# Salvage of Supervision in Weakly Supervised Object Detection

Lin Sui[1]     Chen-Lin Zhang[1,2]     Jianxin Wu[1*]

[1]State Key Laboratory for Novel Software Technology, Nanjing University, China

[2]4Paradigm Inc., Beijing, China

{suilin0432, zclnjucs, wujx2001}@gmail.com

## Abstract

*Weakly supervised object detection (WSOD) has recently attracted much attention. However, the lack of bounding-box supervision makes its accuracy much lower than fully supervised object detection (FSOD), and currently modern FSOD techniques cannot be applied to WSOD. To bridge the performance and technical gaps between WSOD and FSOD, this paper proposes a new framework, Salvage of Supervision (SoS), with the key idea being to harness every potentially useful supervisory signal in WSOD: the weak image-level labels, the pseudo-labels, and the power of semi-supervised object detection. This paper proposes new approaches to utilize these weak and noisy signals effectively, and shows that each type of supervisory signal brings in notable improvements, outperforms existing WSOD methods (which mainly use only the weak labels) by large margins. The proposed SoS-WSOD method also has the ability to freely use modern FSOD techniques. SoS-WSOD achieves 64.4 $mAP_{50}$ on VOC2007, 61.9 $mAP_{50}$ on VOC2012 and 16.6 $mAP_{50:95}$ on MS-COCO, and also has fast inference speed. Ablations and visualization further verify the effectiveness of SoS.*

## 1. Introduction

Large-scale datasets with precise annotations are critical in developing detection algorithms, but are expensive to obtain. Thus, weakly supervised object detection (WSOD), which only needs image-level labels on training images, is popular these days. WSOD has borrowed ideas from fully supervised object detection (FSOD), such as object proposals [2, 34] and the Fast-RCNN framework [11]. But modern FSOD has discarded external object proposals and has developed better techniques like Faster-RCNN [25] and FPN [21]. Furthermore, current WSOD methods mostly use VGG16 [29] as the backbone and Fast-RCNN [11] as the

detector, which confines both accuracy and speed. That is, due to the lack of detailed box-level annotations, WSOD *cannot* enjoy the progress from FSOD. In fact, it has been shown that modern FSOD techniques such as ResNet backbones and RoIAlign will even *deteriorate* WSOD detectors [28]. The weak image-level label is often *the only supervisory signal* utilized for object detection in WSOD, by resorting to a multi-instance recognition setup [4].

In this paper, we argue that WSOD must *fight hard to harness every potential source of supervisory signal*, and should *find a way to utilize the progress in FSOD*. The proposed Salvage of Supervision (SoS) framework (SoS-WSOD) is illustrated in Fig. 1, which has 3 stages. Stage 1 trains a detector with any WSOD method, and we propose an improved OICR [33] as our stage 1. Stage 2 is pseudo-FSOD, where the difficulty is to generate *good* pseudo box-level annotations in order to boost performance and adopt newer FSOD techniques (e.g., ResNet [14], RoIAlign [13], and FPN [21]), i.e., to *salvage* the supervision. This problem has been largely ignored in WSOD, for which we propose a simple but effective solution. Stage 3 is proposed by us, named as SSOD, in which we split the whole dataset into a "clean" and a "noisy" part, then treat the noisy part as *unlabeled*. That is, we *salvage* additional useful supervisory signals by creating a semi-supervised object detection (SSOD) problem. Hence, we have salvaged supervisory signals out of weak labels by designing novel algorithms to generate high-quality pseudo box-level labels and by creating a semi-supervised learning problem, respectively.

Compared to existing WSOD methods, our SoS-WSOD not only harnesses every potentially useful supervisory signal, but also enables WSOD to fully enjoy both accuracy and speed benefits of modern FSOD methods. Although pseudo FSOD has been tried [7, 10, 33, 43], we will show that SoS-WSOD salvages pseudo-supervision of much higher quality. Hence, our contributions are:

- We propose SoS-WSOD, a new WSOD framework, showing that we must harness *all* potential supervisory signals in WSOD: generate *high-quality* pseudo-annotations for FSOD, and treat the generated pseudo-
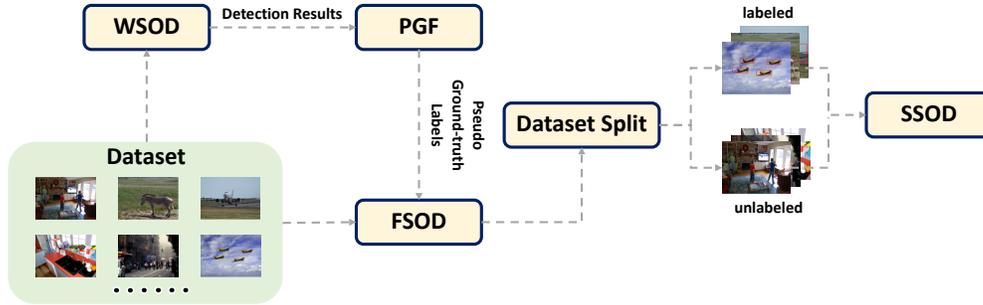
Figure 1. The SoS-WSOD pipeline. Stage 1 trains a weakly supervised detector with only image-level labels. We design PGF to filter its detection results and to generate high-quality pseudo box-level annotations in stage 2, which enables us to train a fully supervised detector. Stage 3 splits the training set into "clean" and unlabeled "noisy" parts, and trains a detector in a semi-supervised manner.

label dataset as a noisy dataset to utilize SSOD.

- We show that although existing WSOD methods lag far behind FSOD in terms of both accuracy and *technique*, it is very beneficial and feasible to fill this gap. Our pseudo-FSOD enjoys the benefits of all modern FSOD techniques in WSOD, and achieves both higher accuracy and faster speed.
- We improve WSOD accuracy by large margins, with 64.4 $mAP_{50}$ on VOC2007, 61.9 $mAP_{50}$ on VOC2012, and 16.6 $mAP_{50:95}$ on MS-COCO. Besides, SoS-WSOD also has fast detection speed.

## 2. Related Work

**Weakly supervised object detection (WSOD).** Weakly supervised object detection (WSOD) seeks to detect the location and type of multiple objects given only image-level labels during training. WSOD methods often utilize object proposals and the multi-instance learning (MIL) framework. WSDDN [4] was the first to integrate MIL into end-to-end WSOD. OICR [33] proposed pseudo groundtruth mining and an online instance refinement branch. PCL [32] clustered proposals to improve pseudo groundtruth mining, and C-MIL [35] improved the MIL loss. Recently, MIST [26] changed the pseudo groundtruth mining rule of OICR, and proposed a Concrete DropBlock module. Zeng *et al.* [28] made the ResNet [14] backbones working in WSOD. CASD [15] proposed self-distillation along with attention to improve WSOD. Some methods [6, 8, 44] proposed to boost WSOD detector performance with the help of fully annotated COCO-60 dataset.

Some methods tried to adopt modern FSOD techniques into WSOD [27, 28]. Some methods have used the output of WSOD methods (pseudo box annotations) to retrain WSOD models with FSOD. W2F [43] proposed a pseudo groundtruth excavation and a pseudo groundtruth adaptation module to mine large and complete objects for retraining. However, they directly retrain WSOD models without considering any noise in the generated pseudo-labeled

dataset which is bound to be very noisy. In contrast, we propose to reconsider the pseudo-labeled dataset with noisy label training perspective and harness the semi-supervised learning paradigm to squeeze better pseudo-labels.

**Semi-supervised object detection (SSOD).** SSOD trains a detector with a small set of images with box-level annotations plus many images without any labels. Compared to WSOD, fewer methods have been proposed for SSOD. SSM [36] stitched high-confidence patches from unlabeled to labeled data. CSD [16] used consistency and background elimination. Recently, STAC [30] used strong data augmentation for unlabeled data. Liu *et al.* [23] used a teacher-student framework, and ISMT [39] used mean teacher. However, these methods need *an exact split* of labeled and unlabeled data, and *noisy-free box-level annotations* for labeled images, but all these are not available in WSOD. We will generate them from the noisy output of the previous stage in SoS-WSOD.

**Learning with noisy labels.** As deep neural networks are annotation-hungry, training DNN with noisy labels also attracts much attention, especially in image classification. Some [31, 40] proposed iterative methods to relabel noisy samples by using network predictions. [1, 24] focused on reweighting. Besides, considering samples with smaller losses as clean ones is also commonly used in many works, such as in [12, 19]. In SoS-WSOD, we adopt the small loss idea to split the noisy output of stage 2 to split the data into "clean" and "noisy" parts.

## 3. Salvage of Supervision

**Notation.** We first define our notation. A training set $\mathcal{D}_w$ consists of training images $I_{tr}$ and image-level annotations $L_{tr}$. Specifically, each image $\boldsymbol{x} \in \mathbb{R}^{h \times w \times 3}$ in $I_{tr}$ has a corresponding label $\boldsymbol{y} = [y_1, y_2, \cdots, y_C] \in [0, 1]^C$, where $C$ is the total number of object categories. We will train a detector $W_{final}$ without using any additional annotations.

**Overview.** Algorithm 1 is the pipeline of the proposed SoS-WSOD. We first train a WSOD detector $W_{wsod}$, which generates pseudo bounding boxes $b_{tr}$. These pseudo su-

**Input**: Training images $I_{tr}$ and image-level class labels $L_{tr}$, test images $I_{te}$

1: Train a WSOD model $W_{wsod}$, and generate pseudo groundtruth bounding boxes $b_{tr}$ for training images
2: Use $I_{tr}$ and $b_{tr}$ to train a fully supervised object detector $W_{full}$
3: Divide $I_{tr}$ into a labeled subset $I_{tr}^*$ with pseudo boxes $b_{tr}^*$ and an unlabeled subset $I_{tr}'$
4: Use $W_{full}$ to initialize, and learn a semi-supervised $W_{final}$ on $L_{tr}$, $I_{tr}^*$ (with $b_{tr}^*$) and $I_{tr}'$
5: **Return:** Use $W_{final}$ to predict the bounding boxes and their class labels for test images

---

pervision signals are used to train an FSOD model $W_{full}$, which can use modern FSOD techniques. Then, we treat the generated pseudo-labeled dataset as a noisy one. With our proposed splitting rule, it is split into an unlabeled subset with $I_{tr}'$ and a "clean" labeled subset (those images with confident pseudo boxes) which consists of $I_{tr}^*$ and $b_{tr}^*$. Finally, we adopt an SSOD method to train the final detector $W_{final}$ on the pseudo labeled dataset.

### 3.1. Stage 1: Improved WSOD

A traditional WSOD detector starts the process. Besides the given image-level annotations $I_{tr}$, most WSOD methods use external object proposals $R$ as extra inputs. Among them, the pipeline of OICR [33] is widely used, which first selects a small number of most confident object proposals $\hat{R}$ as foreground proposals and then refines them by filtering and adding bounding box regression branches.

We propose to improve OICR with two simple changes as our stage 1. First, recent works [20, 26, 32, 41] demonstrate that better proposal mining rules are critical in obtaining higher recall of objects, which are essential for WSOD detectors. For example, MIST [26] proposed to mine more proposals with low overlap between each other. We find that MIST can catch more objects but will also mine a large number of wrong proposals, while OICR is able to mine accurate proposals but ignores many groundtruth instances. Hence, we introduce a mining rule which strikes a balance between recall and precision. In addition, inspired by CASD [15], we find the multi-input technique is also helpful even *without* using inverted attention and CASD's self-attention transfer. More details are in the appendix.

Our proposed WSOD (stage 1) is a strong baseline (cf. Sec. 4). However, we will also show that SoS-WSOD can achieve excellent performance by adopting a weaker WSOD baseline in stage 1, too.

### 3.2. Stage 2: High-quality pseudo boxes for FSOD

If we are able to output pseudo bounding boxes $b_{tr}$ from stage 1's detector $W_{wsod}$ that are *accurate to some extent*, a subsequent FSOD using these boxes can further improve detection. [33] was the first to re-train a WSOD detector by selecting the top-one detection result per class as pseudo groundtruth label, but it will miss a large amount of objects, especially for complicated datasets such as MS-COCO. As



Figure 2. Comparison of W2F [43] (top) and PGF (bottom). W2F tends to generate clustered objects in complicated scenarios.

will be shown in the ablations in Sec. 4, missed objects will be treated as backgrounds and will even deteriorate the detection accuracy. W2F [43] proposed pseudo groundtruth excavation (PGE) and pseudo groundtruth adaption (PGA) to generate pseudo groundtruth from WSOD output. However, W2F only dealt with the VOC datasets, which have a small number of objects per image and the objects are often large in size. Both modules in W2F are designed to mine *large complete* objects, and are not suitable for general detection. Figure 2 shows that W2F tends to cluster multiple objects into one pseudo-box. Instead, we propose a simple but effective algorithm called pseudo groundtruth filtering (PGF) to generate high-quality pseudo-boxes from stage 1's WSOD model, whose pipeline is shown in Algorithm 2.

For each groundtruth class, we only keep the top-scored predictions and those with high confidence ($\geq t_{keep}$, line 6). Then, we remove tiny proposals which are mostly contained inside other proposals in the same category (lines 8-10). After PGF generates pseudo groundtruth $\hat{P}$, in SoS-WSOD, we are able to use $\hat{P}$ to supervise and train an FSOD detector $W_{full}$ using *modern* FSOD methods (e.g., Faster-RCNN [25] + FPN [21]). Please note that the impact of our pseudo-FSOD phase is two-fold. First, the retrained WSOD detector gets accuracy *and speed* gains from these salvaged supervisory signals. Besides, now we are able to use *almost all modern FSOD techniques which are previously not applicable in WSOD*. In other words, a WSOD detector now has the *flexibility to select most backbones and architecture as needed in WSOD*, without resorting to extensive efforts (e.g., as in [28]).

**Algorithm 2** Pseudo Groundtruth Filtering (PGF)

---
**Input**: boxes $P$ with scores $S$ for an input image (output of stage 1) and its active labels $y_1, \ldots, y_m$, keep threshold $t_{keep}$, containment threshold $t_{con}$

**Output**: Pseudo groundtruth boxes $\hat{P}$

1:　　$\hat{P} = \varnothing$
2:　　**for** $i = 1, \ldots, m$ **do**
3:　　　　$S_i = S[i, :]$　　// get scores for the $i$-th active class
4:　　　　$ind_{max}, S_i^{max} = \max(S_i)$　　// get index and score of the top proposal
5:　　　　$P_i^{max} = P[ind_{max}, :]$　　// get bounding box for the top proposal
6:　　　　Remove all proposals whose scores $< t_{keep}$, and the remaining boxes form a set $P_i$
7:　　　　$P_i = P_i \bigcup P_i^{max}$
8:　　　　**for** any two different bounding boxes $u, v$ remaining in $P_i$ **do**
9:　　　　　　**if** $\frac{|u \cap v|}{|v|} \geq t_{con}$ **then** $P_i = P_i \setminus v$
10:　　　　**end for**
11:　　　　$\hat{P} = \hat{P} \bigcup P_i$
12:　　**end for**

---

We intentionally designed PGF to be very simple in order to achieve both generality and simplicity. In practice, it is also possible to tailor the pseudo groundtruth mining algorithm to the characteristics of the dataset (e.g., as in [43].)

### 3.3. Stage 3: Split noisy data for SSOD

FSOD detectors can bring performance gains to WSOD detectors if a high percentage of the generated pseudo groundtruth are correct. However, noisy or wrong pseudo groundtruth (e.g., missing instances, wrong classification results or inaccurate bounding boxes) are inevitable in WSOD. To deal with this issue, we propose to further salvage supervision by treating the generated pseudo-labeled dataset from the perspective of noisy label learning. After splitting the dataset into "clean" labeled part and unlabeled "noisy" part, the semi-supervised learning paradigm can be used.

**Data split.** Many works [12, 17] have shown that noisy annotations will deteriorate the performance. The pseudo groundtruth boxes $\hat{P}$ generated by PGF will inevitably have many noisy ones. As shown in [12, 18, 19], a deep network tends to fit clean data first, then gradually memorize noisy ones. Thus, we use the FSOD detector $W'_{full}$ (the detector before performing learning rate decay in the pseudo-FSOD stage) to divide training images $I_{tr}$ into labeled $I_{tr}^*$ (with relatively clean pseudo groundtruth boxes $b_{tr}^*$) and unlabeled ones $I_{tr}'$ (whose pseudo groundtruth boxes are more noisy). In a classification problem, the split is simple [12]: calculate the loss of each training image, and those with smaller loss values are "clean" ones. But, in object detec-

tion, it is hard to decide whether an image is clean simply based on the sum of all losses of all proposals. We follow the small loss idea but revise it for object detection.

Surely we want to focus on foreground objects, hence we propose the following simple splitting process. In Faster-RCNN, regions of interest (RoIs, denoted as $R$) are divided into foreground and background RoIs according to the IoU between RoIs and pseudo groundtruth boxes. Then, we do *not* calculate losses for background RoIs, and accumulate the RPN losses and RoI losses (both classification and regression branches) of different foreground RoIs. The aggregated loss is the split loss for an input image:

$$\mathcal{L}_{split}(I) = \frac{1}{N_{pos}} \sum_i \mathbb{1}_{fore}^{R_i} \mathcal{L}_{split}(R_i), \qquad (1)$$

$$\mathcal{L}_{split}(R_i) = \mathcal{L}_{RPN_{Cls}}(R_i) + \mathcal{L}_{RPN_{Reg}}(R_i)$$
$$+ \mathcal{L}_{RoI_{Cls}}(R_i) + \mathcal{L}_{RoI_{Reg}}(R_i), \qquad (2)$$

where $N_{pos}$ is the number of foreground RoIs, $\mathbb{1}_{fore}^{R_i}$ is the indicator function for whether a proposal $R_i$ belongs to foreground RoIs or not, $\mathcal{L}_{RPN}$ and $\mathcal{L}_{RoI}$ are RPN and RoI head losses, respectively, and $Cls$ and $Reg$ stand for classification and regression, respectively. We then rank all training images by $\mathcal{L}_{split}$ and choose images with small loss values as "clean" labeled data. For simplicity, we use a hard threshold $K$ to decide sizes of each part. We believe there exist better but more complicated choices such as dynamically using GMM to fit the loss distribution to divide pseudo-labeled dataset. Fig. 3 shows our labeled "clean" part are indeed cleaner than the "noisy" unlabeled part.

The optimal choice of $K$ varies depends on the size and difficulty of datasets, and we provide a rule-of-thumb for it. We find that traditional WSOD performs well for easy images which will be split as "clean". Thus, we set $K$ around the number of images which has a single class label (or few for datasets with many object categories). Ablation studies in Sec. 4 verify the effectiveness of such a rule-of-thumb.

**Semi-supervised detection.** Now we can perform semi-supervised detection. Unbiased Teacher [23] is an effective semi-supervised detector, whose key idea is a teacher-student pair updated by mutual learning. It first trains a detector using only labeled data (i.e., burn-in) and then uses it to initialize both the student and the teacher detectors. In the mutual learning phase, the teacher will dynamically generate pseudo labels for unlabeled data with weak data augmentation. The student will learn from both well-annotated labeled data and strong augmented unlabeled data with the generated pseudo labels. The teacher will receive updates from the student via exponential moving average.

But, clean data is *not* available in WSOD. We use the Unbiased Teacher pipeline with a few changes and improvements. First, the pseudo-FSOD training is actually a burn-in process, and we do not need an additional burn-in stage.
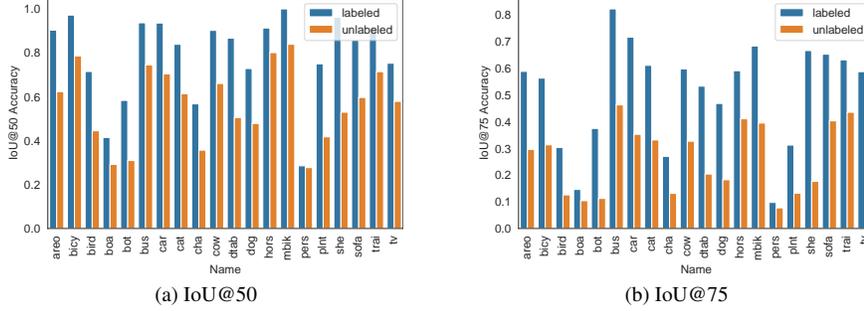
(a) IoU@50

(b) IoU@75

Figure 3. Per-class accuracy of pseudo groundtruth bounding boxes of the labeled and unlabeled subsets on VOC2007 when $K = 2000$.

Then, the student learns by minimizing

$$\mathcal{L} = \mathcal{L}_{sup} + \lambda_u \mathcal{L}_{unsup}, \qquad (3)$$

where the student will learn from both labeled ($\mathcal{L}_{sup}$) and unlabeled ($\mathcal{L}_{unsup}$) data, and $\lambda_u$ is the weight to balance the unsupervised and supervised loss terms. Specifically, $\mathcal{L}_{sup}$ and $\mathcal{L}_{unsup}$ are defined as follows:

$$\mathcal{L}_{sup} = \sum_i \mathcal{L}_{RPN_{Cls}}(\boldsymbol{x}_i^s, \boldsymbol{b}_i^s) + \mathcal{L}_{RPN_{Reg}}(\boldsymbol{x}_i^s, \boldsymbol{b}_i^s)$$
$$+ \mathcal{L}_{RoI_{Cls}}(\boldsymbol{x}_i^s, \boldsymbol{b}_i^s) + \mathcal{L}_{RoI_{Reg}}(\boldsymbol{x}_i^s, \boldsymbol{b}_i^s), \qquad (4)$$

$$\mathcal{L}_{unsup} = \sum_i \mathcal{L}_{RPN_{Cls}}(\boldsymbol{x}_i^u, \boldsymbol{b}_i') + \mathcal{L}_{RoI_{Cls}}(\boldsymbol{x}_i^u, \boldsymbol{b}_i'), \qquad (5)$$

where $\boldsymbol{x}_i^s$ and $\boldsymbol{b}_i^s$ are images and pseudo groundtruth boxes in the "clean" subset $I_{tr}^*$. $\boldsymbol{x}_i^u$ and $\boldsymbol{b}_i'$ are images in the unlabeled subset $I_{tr}'$ and pseudo groundtruth dynamically generated by the teacher. $\mathcal{L}_{sup}$ is for labeled data only. For $\mathcal{L}_{unsup}$, the teacher generates pseudo labels for the student with weak data augmentations, then the student uses strong data augmentations along with pseudo labels to calculate it. We believe the predictions of the teacher are less accurate than annotations for the "clean" data, so $\mathcal{L}_{unsup}$ only contains the classification loss. In other words, all the regression branches are only learned with "clean" labeled data.

Pseudo boxes generated by the teacher for unlabeled images are not always accurate. Hence, different from regular SSOD, we utilize the image-level labels (i.e., another salvage of supervision) by filtering out false positive pseudo labels, which brings additional benefits to WSOD. Finally, the student detector updates its weights according to Eq. 3, and the teacher receives its update from the student through exponential moving average (EMA).

## 4. Experiments

We evaluated our SoS-WSOD on three standard WSOD benchmark datasets: VOC2007 [9], VOC2012 [9] and MS-COCO [22]. VOC2007 has 2501 training, 2510 validation and 4952 test images. VOC2012 contains 5717 training, 5823 validation, and 10991 test images. MS-COCO contains around 110,000 training and 5000 validation images. Following the common WSOD evaluation protocol, we use training and evaluation images to train our model on VOC2007 and VOC2012, and evaluate on the test images. For MS-COCO, we train our model on the training images and evaluate on the validation images. We use $mAP_{50:95}, mAP_{50}$ and $mAP_{75}$ as evaluation metrics for both MS-COCO and VOC2007. For VOC2012, since labels for test images are not released, we report $mAP_{50}$ results returned by the official evaluation server.

### 4.1. Implementation details

We use PyTorch on RTX3090 GPUs, and our code will be released soon. Backbone models are pretrained on ImageNet. It is worth noting that WSOD methods lag behind FSOD in terms of backbone and other techniques. For example, state-of-the-art WSOD methods still use VGG16 as the backbone, while FSOD methods use better architectures. Extra efforts are needed in order to adapt modern backbones to WSOD [28]. Instead, in stage 2 and 3 our SoS-WSOD has the freedom to choose backbones and detection architectures. For simplicity and efficiency, we use FPN with ResNet50 backbone as the FSOD detector in our main experiments *without any extra handling*. We also use VGG16 as backbone and remove FPN for fair comparisons.

Details of our improved OICR for the WSOD training stage (stage 1) are available in the appendix. In PGF (Algorithm 2), we set $t_{keep} = 0.2$ and $t_{con} = 0.85$ for *all* datasets. Although generating pseudo groundtruth labels with TTA (Test Time Augmentation) leads to higher accuracy, the high computational cost (1.5/3/33 hours on VOC2007/2012/MS-COCO) makes it hard to use in large-scale datasets. To keep the same setting in all experiments, we do *not* use TTA in Algorithm 2.

In both stages 2 and 3, we keep *all* hyperparameters except $K$ the same as the *default* hyperparameters in [37] and [23], respectively. We reduce the total training iterations to reduce the training cost. As for $K$ in the data splitting process, we use our rule-of-thumb to set it as 2000 and 30000 for VOC2007 and MS-COCO, respectively. More details can be found in the appendix.

| Method | Backbone | VOC07 $m\text{AP}_{50}$ | VOC12 $m\text{AP}_{50}$ |
|---|---|---|---|
| Pure WSOD | | | |
| PCL [32] | VGG16 | 43.5 | 40.6 |
| W2F [43] | VGG16 | 52.4 | 47.8 |
| Pred Net [3] | VGG16 | 52.9 | 48.4 |
| C-MIDM + FR [10] | VGG16 | 53.6 | 50.3 |
| SLV + FR [7] | VGG16 | 53.9 | - |
| WSOD2 [42] | VGG16 | 53.6 | 47.2 |
| IM-CFB [41] | VGG16 | 54.3 | 49.4 |
| MIST [26] | VGG16 | 54.9 | 52.1 |
| CASD [15] | VGG16 | 56.8 | 53.6 |
| SoS-WSOD (stage 1) | VGG16 | 55.0 | 52.5 |
| SoS-WSOD (stage 1+2+3) | VGG16 | **60.3** | **57.7** |
| SoS-WSOD (stage 1+2+3) | ResNet50 | **64.4** | **61.9** |
| WSOD with transfer (using fully annotated COCO-60) | | | |
| OCUD + FR [44] | ResNet50 | 60.2 | - |
| LBBA [8] | VGG16 | 56.6 | 55.4 |
| CaT [6] | VGG16 | 59.2 | - |

Table 1. Comparison on PASCAL VOC.

| Method | Backbone | MS-COCO $m\text{AP}_{50:95}$ | $m\text{AP}_{50}$ | $m\text{AP}_{75}$ |
|---|---|---|---|---|
| PCL [32] | VGG16 | 8.5 | 19.4 | - |
| C-MIDN [10] | VGG16 | 9.6 | 21.4 | - |
| WSOD2 [42] | VGG16 | 10.8 | 22.7 | - |
| MIST [26] | VGG16 | 12.4 | 25.8 | 10.5 |
| CASD [15] | VGG16 | 12.8 | 26.4 | - |
| SoS-WSOD (stage 1) | VGG16 | 11.9 | 24.2 | 10.7 |
| SoS-WSOD (stage 1+2+3) | VGG16 | **15.5** | **30.5** | **14.3** |
| SoS-WSOD (stage 1+2+3) | ResNet50 | **16.6** | **32.8** | **15.2** |

Table 2. Comparison on MS-COCO.

## 4.2. Comparison with state-of-the-art methods

We compare our method with state-of-the-art WSOD methods, with the results reported in Tables 1 and 2. All results are reported with TTA. Our improved WSOD baseline (stage 1 of SoS-WSOD) reaches $55.0\%$ $m\text{AP}_{50}$, $52.5\%$ $m\text{AP}_{50}$ and $11.9\%$ $m\text{AP}_{50:95}$ on VOC2007, VOC2012 and MS-COCO, respectively, which is already pretty strong.

For a fair comparison, we used VGG16 as backbone and did not use modern FPN architecture in stages 2 and 3. By harnessing all possible supervision signals, SoS-WSOD finally reaches $60.3\%$ and $57.7\%$ $m\text{AP}_{50}$ on VOC2007 and VOC2012, which outperforms previous methods by large margins. On MS-COCO, SoS-WSOD reaches $15.5\%$ $m\text{AP}_{50:95}$, $30.5\%$ $m\text{AP}_{50}$ and $14.3\%$ $m\text{AP}_{75}$, which outperforms previous methods significantly, too.

When further adopting modern techniques in FSOD, with the help of ResNet50 backbone and FPN architecture, SoS-WSOD further reaches $64.4\%$ and $61.9\%$ $m\text{AP}_{50}$ on VOC2007 and VOC2012, respectively. On MS-COCO, accuracy is boosted to $16.6\%$ $m\text{AP}_{50:95}$, $32.8\%$ $m\text{AP}_{50}$ and $15.2\%$ $m\text{AP}_{75}$.

Recently, some methods [6, 8, 44] leverage the well-annotated MS-COCO-60 dataset (removing the 20 cate-

| WSOD | PGF | SSOD | $m\text{AP}_{50:95}$ | $m\text{AP}_{50}$ | $m\text{AP}_{75}$ |
|---|---|---|---|---|---|
| ✓ | | | 26.2 | 54.1 | 22.8 |
| ✓ | ✓ | | 27.3 | 57.6 | 22.5 |
| ✓ | ✓ | ✓ | **31.6** | **62.7** | **28.1** |

Table 3. Ablations of SoS-WSOD stages on VOC2007.

| WSOD | PGF | SSOD | $m\text{AP}_{50}$ |
|---|---|---|---|
| ✓ | | | 51.8 |
| ✓ | ✓ | | 53.9 |
| ✓ | ✓ | ✓ | **59.6** |

Table 4. Ablations of SoS-WSOD stages on VOC2012.

gories in VOC). As shown in Table 1, they have higher accuracy than pure WSOD methods because of the additional cross-domain data. However, SoS-WSOD achieves higher accuracy than them without resorting to these additional data.

## 4.3. Ablation studies and visualization

**Are salvaged supervision signals useful?** Tables 1 and 2 already show that both pseudo boxes (stage 2) and semi-supervised detection (stage 3) notably improve detection accuracy on all 3 datasets. Furthermore, Tables 3 to 5 show results on VOC2007, VOC2012 and MS-COCO, respectively. Our improved WSOD (stage 1 of SoS-WSOD) reaches $54.1\%$, $51.8\%$ and $23.6\%$ $m\text{AP}_{50}$ on VOC2007, VOC2012 and MS-COCO, respectively. After pseudo-FSOD (stage 2), $m\text{AP}_{50}$ is improved by $3.5\%$, $2.1\%$ and $3.9\%$, respectively. Finally, another $5.1\%$, $5.7\%$ and $3.1\%$ higher $m\text{AP}_{50}$ are boosted by stage 3, respectively. For the stricter $m\text{AP}_{50:95}$ metric on MS-COCO, stage 2 and 3 bring $18.1\%$ and $15.5\%$ relative improvements, respectively.

**Compatibility with other WSOD methods.** SoS-WSOD is compatible with various WSOD methods in stage 1 and can also improve a weaker WSOD method. We tested the original OICR method in stage 1, and results with TTA are in Table 6. The basic OICR model gets $50.2$ $m\text{AP}_{50}$ on VOC2007. With SoS-WSOD, such a model finally reaches $59.9$ $m\text{AP}_{50}$. These results demonstrate the flexibility and effectiveness of our SoS-WSOD.

**Details about accuracy gains in stage 3.** As stated in [23, 30], a mixture of weak and strong data augmentation and EMA updating are essential for SSOD methods. As shown in Table 7, in addition to gains from strong augmentation and EMA updating, our method always brings significant gains. As for models use EMA and strong data augmentation but do not adopt the splitting rule (i.e., the third stage), we use the same (mix weak and strong) data augmentation used in the student branch of stage 3 and maintain another model by EMA in stage 2. SoS-WSOD works well for VGG16 w/o FPN, even when strong augmentation and EMA updating have minor improvements.

**Effectiveness of PGF.** As shown in Table 1, some

| WSOD | PGF | SSOD | $m\text{AP}_{50:95}$ | $m\text{AP}_{50}$ | $m\text{AP}_{75}$ | $m\text{AP}_S$ | $m\text{AP}_M$ | $m\text{AP}_L$ |
|---|---|---|---|---|---|---|---|---|
| ✓ | | | 11.6 | 23.6 | 10.4 | 2.3 | 11.9 | 20.2 |
| ✓ | ✓ | | 13.7 | 27.5 | 12.2 | 3.8 | 15.1 | 22.0 |
| ✓ | ✓ | ✓ | **15.5** | **30.6** | **14.4** | **5.4** | **16.8** | **24.6** |

Table 5. Ablations of SoS-WSOD stages on MS-COCO.

| Method | $m\text{AP}_{50:95}$ | $m\text{AP}_{50}$ | $m\text{AP}_{75}$ |
|---|---|---|---|
| baseline (OICR) | 24.1 | 50.2 | 19.5 |
| baseline (OICR+) | 27.1 | 55.0 | 24.8 |
| SoS-WSOD (OICR) | 28.5 | 59.9 | 23.8 |
| SoS-WSOD (OICR+) | 32.6 | 64.4 | 29.6 |

Table 6. Ablations of adopting different WSOD model in stage 1. OICR+ is our improved OICR.

| Backbone | FPN | EMA & Aug | Split Rule | $m\text{AP}_{50:95}$ | $m\text{AP}_{50}$ | $m\text{AP}_{75}$ |
|---|---|---|---|---|---|---|
| VGG16 | × | × | × | 26.8 | 56.2 | 22.3 |
| VGG16 | × | ✓ | × | 27.3 | 56.6 | 22.8 |
| VGG16 | × | ✓ | ✓ | **28.8** | **59.0** | **24.3** |
| ResNet50 | ✓ | × | × | 27.3 | 57.6 | 22.5 |
| ResNet50 | ✓ | ✓ | × | 29.6 | 59.6 | 25.7 |
| ResNet50 | ✓ | ✓ | ✓ | **31.6** | **62.7** | **28.1** |

Table 7. Detailed accuracy gains in stage 3. The columns mean whether FPN, strong data augmentation, EMA, and the proposed data splitting rule are used or not.

| Method | Dataset | $m\text{AP}_{50:95}$ | $m\text{AP}_{50}$ | $m\text{AP}_{75}$ |
|---|---|---|---|---|
| Top-One [33] | VOC07 | 25.8 | 53.4 | 22.2 |
| W2F* [43] | VOC07 | 27.2 | 57.0 | 22.4 |
| PGF (Ours) | VOC07 | **27.3** | **57.6** | **22.5** |
| W2F* | MS-COCO | 12.6 | 24.1 | 11.7 |
| PGF (Ours) | MS-COCO | **13.7** | **27.5** | **12.2** |

Table 8. Results of different pseudo groundtruth mining algorithms. * means that we implemented W2F because its code is not available.

WSOD methods tried to add a re-train stage. Most of them follow [33] to re-train a WSOD detector by selecting the top-one detection result per class as pseudo groundtruth labels. However, the performance of the retrained detector model starts to saturate with the increasing performance of the WSOD model, and pseudo-label generation method is one of the most important reasons. Experimental results in Table 8 show the effectiveness of the proposed PGF. Besides, the widely adopted method (Top-One in Table 8) failed to get benefits from modern techniques. When training with RPN and FPN, missed objects will be treated as backgrounds, which will deteriorate localization. Besides, conspicuous objects also have more chance to become the top score proposal which would cause imbalanced anchor allocation and inadequate training in FPN. As for W2F [43], we have slightly better results on VOC07, and are far superior than it on the more complicated MS-COCO, as W2F is designed for large complete objects in VOC.

**Splitting rule vs. random splitting.** To demonstrate the effectiveness of the proposed splitting rule in stage 3, we

| Method | $m\text{AP}_{50:95}$ | $m\text{AP}_{50}$ | $m\text{AP}_{75}$ |
|---|---|---|---|
| Stage 2 w/ EMA & Aug | 29.6 | 59.6 | 25.7 |
| Stage 3 w/ random splitting | 30.2 | 61.1 | 26.3 |
| Stage 3 w/ ours | **31.6** | **62.7** | **28.1** |

Table 9. Comparison of our splitting rule and random splitting

| K | $m\text{AP}_{50:95}$ | $m\text{AP}_{50}$ | $m\text{AP}_{75}$ |
|---|---|---|---|
| 1000 | 31.2 | **63.2** | 26.8 |
| 2000 | **31.6** | 62.7 | **28.1** |
| 3000 | 31.0 | 62.3 | 27.2 |

Table 10. Effects of $K$ in stage 3 on VOC2007.

compare it with the random splitting strategy. As shown in Table 9, adopting random splitting is notably worse than our proposed method. However, random splitting can still surpass simply adopting EMA update and strong data augmentations in stage 2 by a clear margin, which demonstrates the importance of salvaging useful supervisory signals.

**Longer training schedule for WSOD.** Counting all 3 stages in, SoS-WSOD does require more training iterations. Hence, we double the training iteration of the WSOD stage for a further fair comparison. However, we find that $m\text{AP}_{50}$ will drop from $54.1\%$ to $52.7\%$ due to overfitting.

**Size of the labeled subset in SSOD.** In the SSOD stage (stage 3), we split a dataset into labeled and unlabeled subsets. The number of pseudo labeled images, $K$, is a hyperparameter. When we treat a small number of images as "clean" labeled ones, severe class imbalance will deteriorate the performance. However, when splitting most images as labeled, the performance will collapse using fully pseudo annotated labels. As shown in Table 10, $K = 2000$ is a suitable choice for VOC2007. These results also demonstrate the effectiveness of the rule-of-thumb we proposed. We use $K = 2000$ in all our experiments on VOC2007, and double the size to 4000 on VOC2012. Following the proposed rule-of-thumb, for MS-COCO, we use $K = 30000$.

**Hyperparameters in PGF.** Figure 5 shows the effects of hyperparameters $t_{keep}$ and $t_{con}$ introduced in PGF (Algorithm 2). These two hyperparameters are robust and $t_{keep} = 0.2, t_{con} = 0.85$ works best for $m\text{AP}_{50}$ on VOC2007. We tune these two hyperparameters on the smallest VOC2007 dataset and keep them fixed on all other datasets following previous works [15, 32, 35].

**Inference speed.** SoS-WSOD also enjoys speed benefits from modern FSOD methods. We compare the inference speed in Table 11 (on single RTX3090 GPU). Please note

Figure 4. Visualization of SoS-WSOD results on MS-COCO. Top row: groundtruth annotations. 2nd to 4th rows: detection results from stages 1, 2 and 3, respectively. Last column: a failure case.
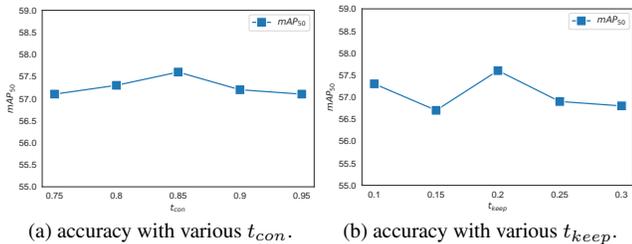


(a) accuracy with various $t_{con}$.    (b) accuracy with various $t_{keep}$.

Figure 5. Effect of hyperparameters in PGF.

Table 11. Inference speed comparison. "Reg" means the bounding box regression branch.

| Method | Pro. Time (s / img) | Inf. Time (s / img) |
|---|---|---|
| OICR (+Reg.) [33] | > 0.2 | 0.101 |
| SoS-WSOD | 0 | 0.031 |

that the time for generating proposals is always far longer than 0.2 seconds per image, e.g., 8.3 s/img for Selective Search [34], while our SoS-WSOD does not need to generate external proposals. Hence, SoS-WSOD not only is significantly faster than baseline WSOD methods, but also eliminates the time to generate external proposals.

Finally, we provide visualization of detection results on MS-COCO in Fig. 4. These results show that SoS-WSOD can mine more correct objects even in complicated environments. Additional visualization results on VOC2007 and MS-COCO are shown in the appendix.

## 5. Conclusions and Remarks

In this paper, we proposed a new three-stage framework called Salvage of Supervision for the weakly supervised ob-ject detection task (SoS-WSOD). SoS-WSOD tackles the WSOD problem from a new perspective, which advocates harnessing all potentially useful supervisory signals (i.e., salvage of supervision) and successfully adopted modern fully supervised detection techniques in WSOD.

The first stage is a WSOD training stage, in which we train a detector with any WSOD method. Pseudo-FSOD, the second stage, improves the WSOD detector by harnessing the pseudo groundtruth generated by PGF and then freely using techniques from modern FSOD. Finally, stage 3 treats the generated pseudo-labeled dataset as a dataset with noisy labels and proposes a novel criterion to split images into labeled and unlabeled subsets, so semi-supervised detection can be used to squeeze useful supervisory signals to further improve the detection performance. Extensive experiments and visualization on VOC2007, VOC2012 and MS-COCO proved both the effectiveness of our SoS-WSOD and extra supervision signals. By successfully utilizing modern FSOD methods, SoS-WSOD can also have faster detection speed than previous WSOD methods.

As for the limitation, SoS-WSOD still suffers from a large performance gap compared to FSOD, especially on COCO. Due to the lack of fully annotated box-level annotations, we need to salvage more supervisory signals in the future. And, SOS-WSOD still suffers problems like part domination, missing instances and clustered instances, which are widely occurred in WSOD. In the future, we will continue to explore to solve the common WSOD problem and develop better rules to split datasets and stronger SSOD methods for the WSOD task.

# References

[1] Eric Arazo, Diego Ortego, Paul Albert, Noel O'Connor, and Kevin McGuinness. Unsupervised label noise modeling and loss correction. In *ICML*, pages 312–321. PMLR, 2019. 2

[2] Pablo Arbeláez, Jordi Pont-Tuset, Jonathan T Barron, Ferran Marques, and Jitendra Malik. Multiscale combinatorial grouping. In *CVPR*, pages 328–335, 2014. 1

[3] Aditya Arun, CV Jawahar, and M Pawan Kumar. Dissimilarity coefficient based weakly supervised object detection. In *CVPR*, pages 9432–9441, 2019. 6, 13

[4] Hakan Bilen and Andrea Vedaldi. Weakly supervised deep detection networks. In *CVPR*, pages 2846–2854, 2016. 1, 2, 11, 13

[5] Zhaowei Cai and Nuno Vasconcelos. Cascade R-CNN: Delving into high quality object detection. In *CVPR*, pages 6154–6162, 2018. 12

[6] Tianyue Cao, Lianyu Du, Xiaoyun Zhang, Siheng Chen, Ya Zhang, and Yan-Feng Wang. CaT: Weakly supervised object detection with category transfer. In *ICCV*, page in press, 2021. 2, 6

[7] Ze Chen, Zhihang Fu, Rongxin Jiang, Yaowu Chen, and Xian-Sheng Hua. SLV: Spatial likelihood voting for weakly supervised object detection. In *CVPR*, pages 12995–13004, 2020. 1, 6, 13

[8] Bowen Dong, Zitong Huang, Yuelin Guo, Qilong Wang, Zhenxing Niu, and Wangmeng Zuo. Boosting weakly supervised object detection via learning bounding box adjusters. In *ICCV*, page in press, 2021. 2, 6, 13

[9] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The PASCAL visual object classes (VOC) challenge. *IJCV*, 88(2):303–338, 2010. 5, 12

[10] Yan Gao, Boxiao Liu, Nan Guo, Xiaochun Ye, Fang Wan, Haihang You, and Dongrui Fan. C-MIDN: Coupled multiple instance detection network with segmentation guidance for weakly supervised object detection. In *ICCV*, pages 9834–9843, 2019. 1, 6, 13

[11] Ross Girshick. Fast R-CNN. In *ICCV*, pages 1440–1448, 2015. 1

[12] Bo Han, Quanming Yao, Xingrui Yu, Gang Niu, Miao Xu, Weihua Hu, Ivor Tsang, and Masashi Sugiyama. Co-teaching: Robust training of deep neural networks with extremely noisy labels. *arXiv preprint arXiv:1804.06872*, 2018. 2, 4

[13] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask R-CNN. In *ICCV*, pages 2961–2969, 2017. 1

[14] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778, 2016. 1, 2

[15] Zeyi Huang, Yang Zou, B. V. K. Vijaya Kumar, and Dong Huang. Comprehensive attention self-distillation for weakly-supervised object detection. In *NeurIPS*, pages 16797–16807, 2020. 2, 3, 6, 7, 11, 13

[16] Jisoo Jeong, Seungeui Lee, Jeesoo Kim, and Nojun Kwak. Consistency-based semi-supervised learning for object detection. In *NeurIPS*, pages 1–9, 2019. 2

[17] Lu Jiang, Zhengyuan Zhou, Thomas Leung, Li-Jia Li, and Li Fei-Fei. MentorNet: Learning data-driven curriculum for very deep neural networks on corrupted labels. In *ICML*, pages 2304–2313, 2018. 4

[18] Hengduo Li, Zuxuan Wu, Chen Zhu, Caiming Xiong, Richard Socher, and Larry S Davis. Learning from noisy anchors for one-stage object detection. In *CVPR*, pages 10588–10597, 2020. 4

[19] Junnan Li, Richard Socher, and Steven CH Hoi. DivideMix: Learning with noisy labels as semi-supervised learning. In *ICLR*, pages 1–13, 2020. 2, 4

[20] Chenhao Lin, Siwen Wang, Dongqi Xu, Yu Lu, and Wayne Zhang. Object instance mining for weakly supervised object detection. In *AAAI*, volume 34, pages 11482–11489, 2020. 3, 11

[21] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature pyramid networks for object detection. In *CVPR*, pages 2117–2125, 2017. 1, 3

[22] 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 *ECCV*, volume 8693 of *LNCS*, pages 740–755, 2014. 5

[23] Yen-Cheng Liu, Chih-Yao Ma, Zijian He, Chia-Wen Kuo, Kan Chen, Peizhao Zhang, Bichen Wu, Zsolt Kira, and Peter Vajda. Unbiased teacher for semi-supervised object detection. In *ICLR*, pages 1–13, 2021. 2, 4, 5, 6, 12

[24] Mengye Ren, Wenyuan Zeng, Bin Yang, and Raquel Urtasun. Learning to reweight examples for robust deep learning. In *ICML*, pages 4334–4343. PMLR, 2018. 2

[25] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster R-CNN: Towards real-time object detection with region proposal networks. In *NeurIPS*, pages 91–99, 2015. 1, 3

[26] Zhongzheng Ren, Zhiding Yu, Xiaodong Yang, Ming-Yu Liu, Yong Jae Lee, Alexander G Schwing, and Jan Kautz. Instance-aware, context-focused, and memory-efficient weakly supervised object detection. In *CVPR*, pages 10598–10607, 2020. 2, 3, 6, 11, 13

[27] Yunhang Shen, Rongrong Ji, Zhiwei Chen, Yongjian Wu, and Feiyue Huang. UWSOD: Toward fully-supervised-level capacity weakly supervised object detection. In *NeurIPS*, volume 33, 2020. 2, 12

[28] Yunhang Shen, Rongrong Ji, Yan Wang, Zhiwei Chen, Feng Zheng, Feiyue Huang, and Yunsheng Wu. Enabling deep residual networks for weakly supervised object detection. In *ECCV*, volume 12353 of *LNCS*, pages 118–136, 2020. 1, 2, 3, 5

[29] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, pages 1–14, 2015. 1

[30] 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. 2, 6

[31] Daiki Tanaka, Daiki Ikami, Toshihiko Yamasaki, and Kiyoharu Aizawa. Joint optimization framework for learning with noisy labels. In *CVPR*, pages 5552–5560, 2018. 2

[32] Peng Tang, Xinggang Wang, Song Bai, Wei Shen, Xiang Bai, Wenyu Liu, and Alan Yuille. PCL: Proposal cluster learning for weakly supervised object detection. *IEEE TPAMI*, 42(1):176–191, 2018. 2, 3, 6, 7, 11, 13

[33] Peng Tang, Xinggang Wang, Xiang Bai, and Wenyu Liu. Multiple instance detection network with online instance classifier refinement. In *CVPR*, pages 2843–2851, 2017. 1, 2, 3, 7, 8, 11, 13

[34] Jasper RR Uijlings, Koen EA Van De Sande, Theo Gevers, and Arnold WM Smeulders. Selective search for object recognition. *IJCV*, 104(2):154–171, 2013. 1, 8

[35] Fang Wan, Chang Liu, Wei Ke, Xiangyang Ji, Jianbin Jiao, and Qixiang Ye. C-MIL: Continuation multiple instance learning for weakly supervised object detection. In *CVPR*, pages 2199–2208, 2019. 2, 7

[36] Keze Wang, Xiaopeng Yan, Dongyu Zhang, Lei Zhang, and Liang Lin. Towards human-machine cooperation: Self-supervised sample mining for object detection. In *CVPR*, pages 1605–1613, 2018. 2

[37] Yuxin Wu, Alexander Kirillov, Francisco Massa, Wan-Yen Lo, and Ross Girshick. Detectron2. `https://github.com/facebookresearch/detectron2`, 2019. 5

[38] Ke Yang, Dongsheng Li, and Yong Dou. Towards precise end-to-end weakly supervised object detection network. In *ICCV*, pages 8372–8381, 2019. 11

[39] Qize Yang, Xihan Wei, Biao Wang, Xian-Sheng Hua, and Lei Zhang. Interactive self-training with mean teachers for semi-supervised object detection. In *CVPR*, pages 5941–5950, 2021. 2

[40] Kun Yi and Jianxin Wu. Probabilistic end-to-end noise correction for learning with noisy labels. In *CVPR*, pages 7017–7025, 2019. 2

[41] Yufei Yin, Jiajun Deng, Wengang Zhou, and Houqiang Li. Instance mining with class feature banks for weakly supervised object detection. In *AAAI*, pages 3190–3198, 2021. 3, 6, 11, 13

[42] Zhaoyang Zeng, Bei Liu, Jianlong Fu, Hongyang Chao, and Lei Zhang. WSOD$^2$: Learning bottom-up and top-down objectness distillation for weakly-supervised object detection. In *ICCV*, pages 8292–8300, 2019. 6, 11, 13

[43] Yongqiang Zhang, Yancheng Bai, Mingli Ding, Yongqiang Li, and Bernard Ghanem. W2F: A weakly-supervised to fully-supervised framework for object detection. In *CVPR*, pages 928–936, 2018. 1, 2, 3, 4, 6, 7, 13

[44] Yuanyi Zhong, Jianfeng Wang, Jian Peng, and Lei Zhang. Boosting weakly supervised object detection with progressive knowledge transfer. In *ECCV*, volume 11216 of *LNCS*, pages 615–631, 2020. 2, 6, 13

## A.1. Introducing the pipeline of OICR

In this part, we will introduce the details of OICR [33], a widely used framework in WSOD. OICR is composed of two parts, a multiple instance detection network (MIDN) and several online instance classifier refinement (OICR) branches. There are different choices to implement the MIDN part. WSDDN [4], the first work to integrate the MIL process into an end-to-end detection model, is the most commonly used one. As for the OICR branch, originally it only contained one classifier and a softmax function. [38] started to introduce the bounding box regressor into OICR branches, which was proved to be effective in many works [15, 26, 41, 42].

Specifically, we denote $\boldsymbol{I} \in \mathbb{R}^{h \times w \times 3}$ as an RGB image, $\boldsymbol{y} = [y_1, y_2, \ldots, y_C] \in [0,1]^C$ as its corresponding groundtruth class labels, and $\boldsymbol{R} \in \mathbb{R}^{4 \times N}$ as the pre-computed object proposals. $C$ is the total number of object categories and $N$ is the number of proposals. With the help of a pre-trained backbone model, we can extract the feature map for $\boldsymbol{I}$, and proposal feature vectors are extracted by an RoI pooling layer and two FC layers. Following WSDDN, proposal feature vectors are branched into two streams to produce classification logits $\boldsymbol{x}^c \in \mathbb{R}^{C \times N}$ and detection logits $\boldsymbol{x}^d \in \mathbb{R}^{C \times N}$. Then $\boldsymbol{x}^c$ and $\boldsymbol{x}^d$ will be normalized by passing through two softmax layers along the category direction and the proposal direction, respectively, as shown in Equation 6. $[\sigma(\boldsymbol{x}^c)]_{ij}$ represents the probability of proposal $j$ belonging to class $i$ and $[\sigma(\boldsymbol{x}^d)]_{ij}$ represents the likelihood of proposal $j$ to contain an informative part of class $i$ among all proposals in image $\boldsymbol{I}$.

$$[\sigma(\boldsymbol{x}^c)]_{ij} = \frac{\exp^{x^c_{ij}}}{\sum_{k=1}^{C} \exp^{x^c_{kj}}} , \ [\sigma(\boldsymbol{x}^d)]_{ij} = \frac{\exp^{x^d_{ij}}}{\sum_{k=1}^{N} \exp^{x^d_{ik}}} . \tag{6}$$

The final proposal scores of a multiple instance detection network are computed by element-wise product: $\boldsymbol{x}^R = \sigma(\boldsymbol{x}^c) \odot \sigma(\boldsymbol{x}^d)$. During the training process, image score of the $c^{th}$ category $\phi_c$ can be obtained by summing over all proposals: $\phi_c = \sum_{r=1}^{N} \boldsymbol{x}^R_{c,r}$. Then the MIL classification loss is calculated by Equation 7.

$$\mathcal{L}_{mil} = -\sum_{c=1}^{C} [y_c \log \phi_c + (1 - y_c \log(1 - \phi_c))] . \tag{7}$$

As to the online instance classifier refinement (OICR) branches, they are added on top of MIDN, i.e., WSDDN here. Proposal feature vectors are fed into another $K$ refinement stages and to generate classification logits $x^k \in \mathbb{R}^{(C+1) \times N}, k \in \{1, 2, \ldots, K\}$. The $k^{th}$ branch is supervised by pseudo labels $\boldsymbol{y}^k \in [0,1]^{(C+1) \times N}$, which are generated by top-score proposals of each category from the previous branch. One proposal will be encouraged to be classified as the $c$-th class only if it has high overlap with any top-score proposal of the previous OICR branch. The loss for the classifier of the $k^{th}$ branch is defined as Equation 8, where $w_r^k$ is the loss weight of proposal $r$:

$$\mathcal{L}_r^k = -\frac{1}{N} \sum_{r=1}^{N} \sum_{c=1}^{C+1} w_r^k y_{c,r}^k \log x_{c,r}^k . \tag{8}$$

The loss for bounding box regressor of the $k^{th}$ OICR branch is defined as Equation 9, $N_{pos}$ is the number of positive proposals in the $k^{th}$ branch, $\lambda_{reg}$ is a scalar weight of the regression loss, $t_r^k, \hat{t}_r^k$ are the predicted and pseudo groundtruth offsets of the $r^{th}$ positive proposal in the $k^{th}$ branch, respectively:

$$\mathcal{L}_{reg}^k = \frac{1}{N_{pos}} \sum_{r=1}^{N_{pos}} \lambda_{reg} \mathcal{L}_{smooth-L1}(t_r^k, \hat{t}_r^k) . \tag{9}$$

## A.2. Details of our improved OICR

In this part, we provide details of our improved OICR [33], which is used in stage 1. As we claimed in Sec. 3.1, we proposed an improved OICR as the baseline in our main experiments.

**Mining Rules.** Recent works [20, 26, 32, 41] demonstrate that better mining rules are critical in obtaining higher recall of objects. OICR mines proposals that have high overlap with top-scoring proposals. MIST [26] mines more proposals with low overlap between each other but mines many wrong proposals, too. We notice that recall and precision are both essential for mining proposals. Hence, we introduce a mining rule (Algorithm A.3) to strike a balance between the two factors. In Line 6, the rule to only retain the top $p$ percent of proposals is learned from MIST, but we remove low score proposals to keep the precision.

**Multi-Input.** A very recent paper CASD [15] showed that the self-attention transfer between different versions of an input image is beneficial for boosting performance in WSOD. We find that adopting the multi-input technique alone is also helpful for performance and stability of the training process even without using inverted attention, CASD's self-attention transfer and other tricks. We randomly select inputs with two different scales and their flipped versions, feed them into the model to obtain RoI scores for different inputs, and average the scores of each proposal to get the final RoI scores.

## A.3. Implementation Details

In this section, we provide additional implementation details for completeness.

In the WSOD training stage, we set the maximum iteration numbers to 50k, 60k and 200k for VOC2007, VOC2012 and MS-COCO, respectively. Batch size is set

**Algorithm A.3** Mining Rules in SoS-WSOD

---

**Input**: An input image $I$, class labels $y_1, \ldots, y_m$ that are active in $I$, a set of proposals $R$ with size $n$, maximum percent $p$, score threshold $s_t$

**Output**: Pseudo groundtruth seed boxes $\hat{R}$ for $I$

1:     $\hat{R} = \varnothing$
2:     Feed $I$ and $R$ into the model to obtain RoI scores $S$ for each proposal in $R$
3:     **for** $i = 1, \ldots, m$ **do**
4:         $S_i = S[i, :]$     // get scores for the $i$-th active class
5:         $R_i = \text{SORTED}_{S_i}(R)$   // sort the proposals according to the scores in $S_i$
6:         Pick top $n \times p$ proposals, but *remove those whose scores are low* $(< s_t)$. Denote them as $R'_i$
7:         $R'_i = \text{NMS}(R'_i, 0.01)$   // remove those proposals having overlap with higher scored ones
8:         $\hat{R} = \hat{R} \bigcup R'_i$
9:     **end for**

---

to be 4 for the basic OICR model. as we input 4 images with 4 different input transformations, the actual batch size is 16 when we use the improved OICR. When training the improved OICR model, $p = 0.1, s_t = 0.05$ are set for all datasets. When training the FSOD model with pseudo ground-truth, maximum iteration numbers are 12k, 18k, 50k for VOC2007, VOC2012 and MS-COCO, respectively. Learning rate and batch size are 0.01 and 8 for VOC2007 and VOC2012. For MS-COCO, we double the batch size to 16 and adjust the learning rate to 0.02 based on batch size. The learning rate is decayed with a factor of 10 at (8k, 10.5k), (12k, 16k) and (30k, 40k) for VOC2007, VOC2012 and MS-COCO, respectively. When mining potential useful supervisory signals by the semi-supervised learning paradigm, maximum iteration numbers are 15k, 30k, 50k for VOC2007, VOC2012 and MS-COCO, respectively. Batch sizes for the unlabeled subset and "clean" labeled subset are both 8 on VOC2007 and VOC2012, and doubled to 16 on MS-COCO. Learning rate is set to 0.01 on all datasets. We do not modify any other hyperparameters of object detectors.

As for the data argumentation, following [27], we use random flip and multi-scale training in which scales range from 480 to 1216 with stride 32 in stage 1. In stage 2 and 3, we apply the same data augmentations as [23]. For weak augmentation, only scale transform and random flip are used. Color jittering, grayscale, Gaussian blur, and cutout patches are randomly applied for strong augmentation additionally.

## A.4. Ability to adopt modern backbones

In order to show that SoS-WSOD can readily enjoy the benefits from modern fully supervised object detection

techniques, we conducted experiments using ResNet101 and ResNeXt101, which are widely used in fully supervised object detection, as the backbone of SoS-WSOD in stages 2 and 3. In Table A.1, we show the results on VOC2007. These results demonstrate that our SoS-WSOD can successfully adopt different modern backbones. Note that TTA was *not* used for results in Table A.1.

| Backbone | PGF | SSOD | $mAP_{50:95}$ | $mAP_{50}$ | $mAP_{75}$ |
|---|---|---|---|---|---|
| ResNet50 | ✓ | | 27.3 | 57.6 | 22.5 |
| ResNet50 | ✓ | ✓ | 31.6 | 62.7 | 28.1 |
| ResNet101 | ✓ | | 28.7 | 58.2 | 24.2 |
| ResNet101 | ✓ | ✓ | 32.4 | 63.2 | 29.3 |
| ResNeXt101 | ✓ | | 29.1 | 59.1 | 25.5 |
| ResNeXt101 | ✓ | ✓ | 33.0 | 64.7 | 30.1 |

Table A.1. Results for SoS-WSOD when using ResNet101 and ResNeXt101 as the backbone on VOC2007.

## A.5. Ability to adopt different detector architectures

In order to show that SoS-WSOD can also enjoy benefits from different detector architectures, we conducted experiments using Cascade R-CNN [5] with ResNet50 as the backbone on the VOC2007 dataset. Experiment results in Table A.2 show that SoS-WSOD can successfully adopt different modern detector architectures such as Cascade R-CNN. The experimental results also illustrate that using Cascade R-CNN as the detector, SoS-WSOD can obtain performance gains and more high-quality detection results.

| Detector | PGF | SSOD | $mAP_{50:95}$ | $mAP_{50}$ | $mAP_{75}$ |
|---|---|---|---|---|---|
| Faster R-CNN | ✓ | | 27.3 | 57.6 | 22.5 |
| Faster R-CNN | ✓ | ✓ | 31.6 | 62.7 | 28.1 |
| Cascade R-CNN | ✓ | | 29.9 | 56.7 | 27.6 |
| Cascade R-CNN | ✓ | ✓ | 32.5 | 61.3 | 30.8 |

Table A.2. Results for SoS-WSOD when using Cascade R-CNN as the detector on VOC2007.

## A.6. Result on VOC2012

The results on VOC2012 we reported in Sec. 4 of the main paper were directly returned from the evaluation server of the PASCAL VOC Challenge [9]. The detailed results of SoS-WSOD (using all stages) can be obtained by visiting these two anonymous result links.[1][2]

---

[1] http://host.robots.ox.ac.uk:8080/anonymous/Q4JFTS.html
[2] http://host.robots.ox.ac.uk:8080/anonymous/PDK0Q9.html

| Method | Backbone | aero | bicy | bird | boa | bot | bus | car | cat | cha | cow | dtab | dog | hors | mbik | pers | plnt | she | sofa | trai | tv | $m\text{AP}_{50}$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | | Pure WSOD | | | | | | | | | | | | | |
| WSDDN [4] | VGG16 | 39.3 | 43.0 | 28.8 | 20.4 | 8.0 | 45.5 | 47.9 | 22.1 | 8.4 | 33.5 | 23.6 | 29.2 | 38.5 | 47.9 | 20.3 | 20.0 | 35.8 | 30.8 | 41.9 | 20.1 | 30.2 |
| OICR [33] | VGG16 | 58.0 | 62.4 | 31.1 | 19.4 | 13.0 | 65.1 | 62.2 | 28.4 | 24.8 | 44.7 | 30.6 | 25.3 | 37.8 | 65.5 | 15.7 | 24.1 | 41.7 | 46.9 | 64.3 | 62.6 | 41.2 |
| PCL [32] | VGG16 | 54.4 | 69.0 | 39.3 | 19.2 | 15.7 | 62.9 | 64.4 | 30.0 | 25.1 | 52.5 | 44.4 | 19.6 | 39.3 | 67.7 | 17.8 | 22.9 | 46.6 | 57.5 | 58.6 | 63.0 | 43.5 |
| W2F [43] | VGG16 | 63.5 | 70.1 | 50.5 | 31.9 | 14.4 | 72.0 | 67.8 | 73.7 | 23.3 | 53.4 | 49.4 | 65.9 | 57.2 | 67.2 | 27.6 | 23.8 | 51.8 | 58.7 | 64.0 | 62.3 | 52.4 |
| C-MIDN [10] | VGG16 | 53.3 | 71.5 | 49.8 | 26.1 | 20.3 | 70.3 | 69.9 | 68.3 | 28.7 | 65.3 | 45.1 | 64.6 | 58.0 | 71.2 | 20.0 | 27.5 | 54.9 | 54.9 | 69.4 | 63.5 | 52.6 |
| C-MIDN + FR [10] | VGG16 | 54.1 | 74.5 | 56.9 | 26.4 | 22.2 | 68.7 | 68.9 | 74.8 | 25.2 | 64.8 | 46.4 | 70.3 | 66.3 | 67.5 | 21.6 | 24.4 | 53.0 | 59.7 | 68.7 | 58.9 | 53.6 |
| Pred Net [3] | VGG16 | 66.7 | 69.5 | 52.8 | 31.4 | 24.7 | 74.5 | 74.1 | 67.3 | 14.6 | 53.0 | 46.1 | 52.9 | 69.9 | 70.8 | 18.5 | 28.4 | 54.6 | 60.7 | 67.1 | 60.4 | 52.9 |
| SLV [7] | VGG16 | 65.6 | 71.4 | 49.0 | 37.1 | 24.6 | 69.6 | 70.3 | 70.6 | 30.8 | 63.1 | 36.0 | 61.4 | 65.3 | 68.4 | 12.4 | 29.9 | 52.4 | 60.0 | 67.6 | 64.5 | 53.5 |
| SLV + FR [7] | VGG16 | 62.1 | 72.1 | 54.1 | 34.5 | 25.6 | 66.7 | 67.4 | 77.2 | 24.2 | 61.6 | 47.5 | 71.6 | 72.0 | 67.2 | 12.1 | 24.6 | 51.7 | 61.1 | 65.3 | 60.1 | 53.9 |
| WSOD2 [42] | VGG16 | 65.1 | 64.8 | 57.2 | 39.2 | 24.3 | 69.8 | 66.2 | 61.0 | 29.8 | 64.6 | 42.5 | 60.1 | 71.2 | 70.7 | 21.9 | 28.1 | 58.6 | 59.7 | 52.2 | 64.8 | 53.6 |
| IM-CFB [41] | VGG16 | 64.1 | 74.6 | 44.7 | 29.4 | 26.9 | 73.3 | 72.0 | 71.2 | 28.1 | 66.7 | 48.1 | 63.8 | 55.5 | 68.3 | 17.8 | 27.7 | 54.4 | 62.7 | 70.5 | 66.6 | 54.3 |
| MIST [26] | VGG16 | 68.8 | 77.7 | 57.0 | 27.7 | 28.9 | 69.1 | 74.5 | 67.0 | 32.1 | 73.2 | 48.1 | 45.2 | 54.4 | 73.7 | 35.0 | 29.3 | 64.1 | 53.8 | 65.3 | 65.2 | 54.9 |
| CASD [15] | VGG16 | 70.5 | 70.1 | 57.0 | 45.8 | 29.5 | 74.5 | 72.8 | 71.4 | 25.3 | 67.6 | 49.3 | 64.7 | 65.8 | 72.7 | 23.7 | 25.9 | 56.3 | 60.8 | 65.4 | 66.5 | 56.8 |
| SoS-WSOD (ours) | VGG16 | 67.4 | 83.1 | 56.2 | 20.2 | 44.6 | 80.9 | 82.0 | 78.7 | 30.3 | 76.0 | 49.5 | 56.6 | 74.9 | 76.1 | 30.1 | 29.7 | 64.1 | 56.6 | 76.7 | 72.6 | 60.3 |
| SoS-WSOD (ours) | ResNet50 | 77.9 | 81.2 | 58.9 | 26.7 | 54.3 | 82.5 | 84.0 | 83.5 | 36.3 | 76.5 | 57.5 | 58.4 | 78.5 | 78.6 | 33.8 | 37.4 | 64.0 | 63.4 | 81.5 | 74.0 | 64.4 |
| | | | | | | | | | WSOD with transfer | | | | | | | | | | | | | |
| OCUD + FR [44] | ResNet50 | 65.5 | 57.7 | 65.1 | 41.3 | 43.0 | 73.6 | 75.7 | 80.4 | 33.4 | 72.2 | 33.8 | 81.3 | 79.6 | 63.0 | 59.4 | 10.9 | 65.1 | 64.2 | 72.7 | 67.2 | 60.2 |
| LBBA [8] | VGG16 | 70.3 | 72.3 | 48.7 | 38.7 | 30.4 | 74.3 | 76.6 | 69.1 | 33.4 | 68.2 | 50.5 | 67.0 | 49.0 | 73.6 | 24.5 | 27.4 | 63.1 | 58.9 | 66.0 | 69.2 | 56.6 |

Table A.3. Per-class detection results on the VOC2007 test set.

| Method | Backbone | aero | bicy | bird | boa | bot | bus | car | cat | cha | cow | dtab | dog | hors | mbik | pers | plnt | she | sofa | trai | tv | $\text{CorLoc}_{50}$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | | Pure WSOD | | | | | | | | | | | | | |
| WSDDN [4] | VGG16 | 65.1 | 58.8 | 58.5 | 33.1 | 39.8 | 68.3 | 60.2 | 59.6 | 34.8 | 64.5 | 30.5 | 43.0 | 56.8 | 82.4 | 25.5 | 41.6 | 61.5 | 55.9 | 65.9 | 63.7 | 53.5 |
| OICR [33] | VGG16 | 81.7 | 80.4 | 48.7 | 49.5 | 32.8 | 81.7 | 85.4 | 40.1 | 40.6 | 79.5 | 35.7 | 33.7 | 60.5 | 88.8 | 21.8 | 57.9 | 76.3 | 59.9 | 75.3 | 81.4 | 60.6 |
| W2F [43] | VGG16 | 85.4 | 87.5 | 62.5 | 54.3 | 35.5 | 85.3 | 86.6 | 82.3 | 39.7 | 82.9 | 49.4 | 76.5 | 74.8 | 90.0 | 46.8 | 53.9 | 84.5 | 68.3 | 79.1 | 79.9 | 70.3 |
| Pred Net [3] | VGG16 | 88.6 | 86.3 | 71.8 | 53.4 | 51.2 | 87.6 | 89.0 | 65.3 | 33.2 | 86.6 | 58.8 | 65.9 | 87.7 | 93.3 | 30.9 | 58.9 | 83.4 | 67.8 | 78.7 | 80.2 | 70.9 |
| SLV [7] | VGG16 | 84.6 | 84.3 | 73.3 | 58.5 | 49.2 | 80.2 | 87.0 | 79.4 | 46.8 | 83.6 | 41.8 | 79.3 | 88.8 | 90.4 | 19.5 | 59.7 | 79.4 | 67.7 | 82.9 | 83.2 | 71.0 |
| SLV + FR [7] | VGG16 | 85.8 | 85.9 | 73.3 | 56.9 | 52.7 | 79.7 | 87.1 | 84.0 | 49.3 | 82.9 | 46.8 | 81.2 | 89.8 | 92.4 | 21.2 | 59.3 | 80.4 | 70.4 | 82.1 | 78.8 | 72.0 |
| WSOD2 [42] | VGG16 | 87.1 | 80.0 | 74.8 | 60.1 | 36.6 | 79.2 | 83.8 | 70.6 | 43.5 | 88.4 | 46.0 | 74.7 | 87.4 | 90.8 | 44.2 | 52.4 | 81.4 | 61.8 | 67.7 | 79.9 | 69.5 |
| MIST [26] | VGG16 | 87.5 | 82.4 | 76.0 | 58.0 | 44.7 | 82.2 | 87.5 | 71.2 | 49.1 | 81.5 | 51.7 | 53.3 | 71.4 | 92.8 | 38.2 | 52.8 | 79.4 | 61.0 | 78.3 | 76.0 | 68.8 |
| SoS-WSOD (ours) | VGG16 | 82.4 | 91.8 | 66.4 | 47.5 | 63.5 | 88.7 | 94.8 | 85.8 | 44.7 | 93.6 | 63.5 | 70.6 | 91.6 | 93.5 | 37.8 | 62.0 | 90.6 | 71.6 | 86.6 | 83.2 | 75.5 |
| SoS-WSOD (ours) | ResNet50 | 89.5 | 93.0 | 71.8 | 49.2 | 72.5 | 88.7 | 93.8 | 88.4 | 54.4 | 94.3 | 70.5 | 70.6 | 93.0 | 95.1 | 39.7 | 70.2 | 89.6 | 74.7 | 88.1 | 86.3 | 78.7 |
| | | | | | | | | | WSOD with transfer | | | | | | | | | | | | | |
| OCUD + FR [44] | ResNet50 | 85.8 | 67.5 | 87.1 | 68.6 | 68.3 | 85.8 | 90.4 | 88.7 | 43.5 | 95.2 | 31.6 | 90.9 | 94.2 | 88.8 | 72.4 | 23.8 | 88.7 | 66.1 | 89.7 | 76.7 | 75.2 |
| LBBA [8] | VGG16 | 89.2 | 82.0 | 74.2 | 53.2 | 51.2 | 84.8 | 87.5 | 83.7 | 46.2 | 87.0 | 48.3 | 84.7 | 79.9 | 92.4 | 40.3 | 47.6 | 88.7 | 65.6 | 81.0 | 81.7 | 72.5 |

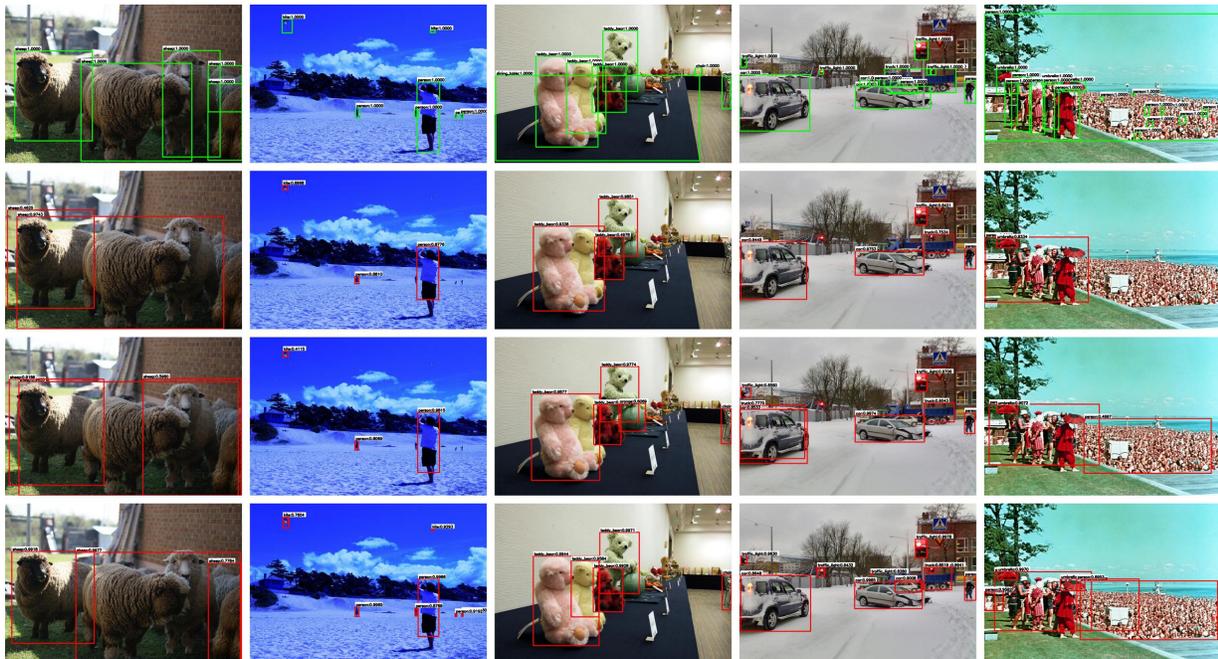Table A.4. Correct localization (CorLoc) results on the VOC2007 trainval set.



Figure A.6. Visualization of SoS-WSOD results on MS-COCO (more examples in addition to Fig. 2 in the main paper). Top row: groundtruth annotations. 2nd to 4th rows: detection results from stages 1, 2 and 3, respectively. Last column: a failure case.
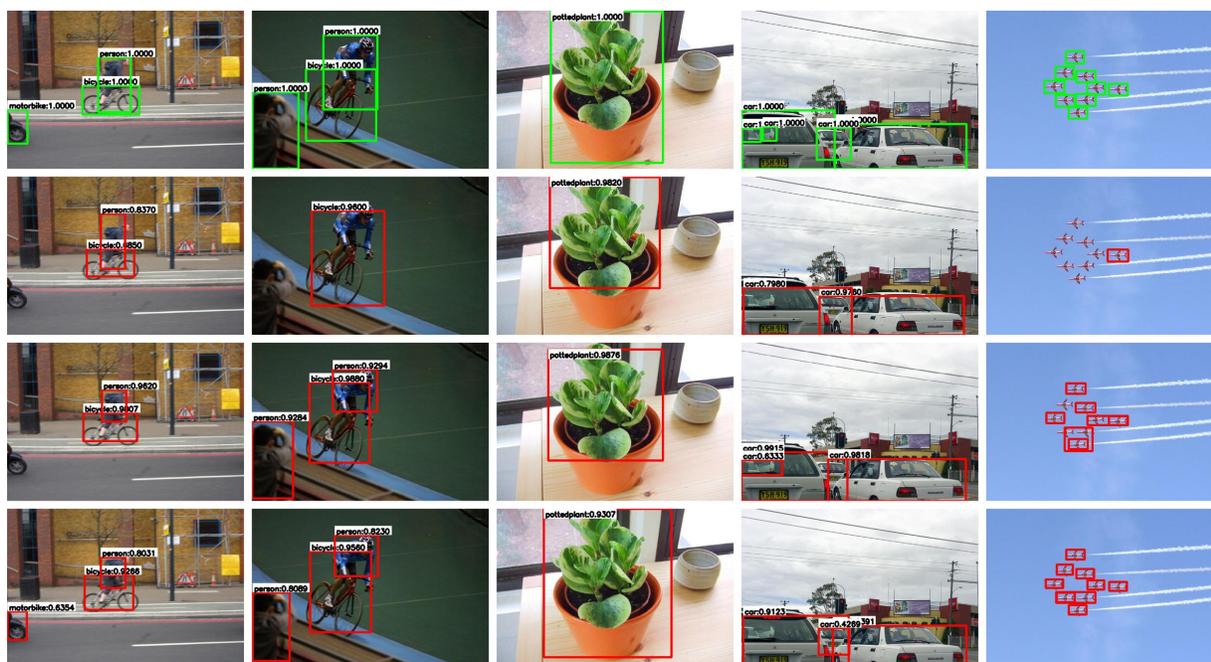
Figure A.7. Visualization of SoS-WSOD results on VOC2007. Top row: groundtruth annotations. 2nd to 4th rows: detection results from stages 1, 2 and 3, respectively.
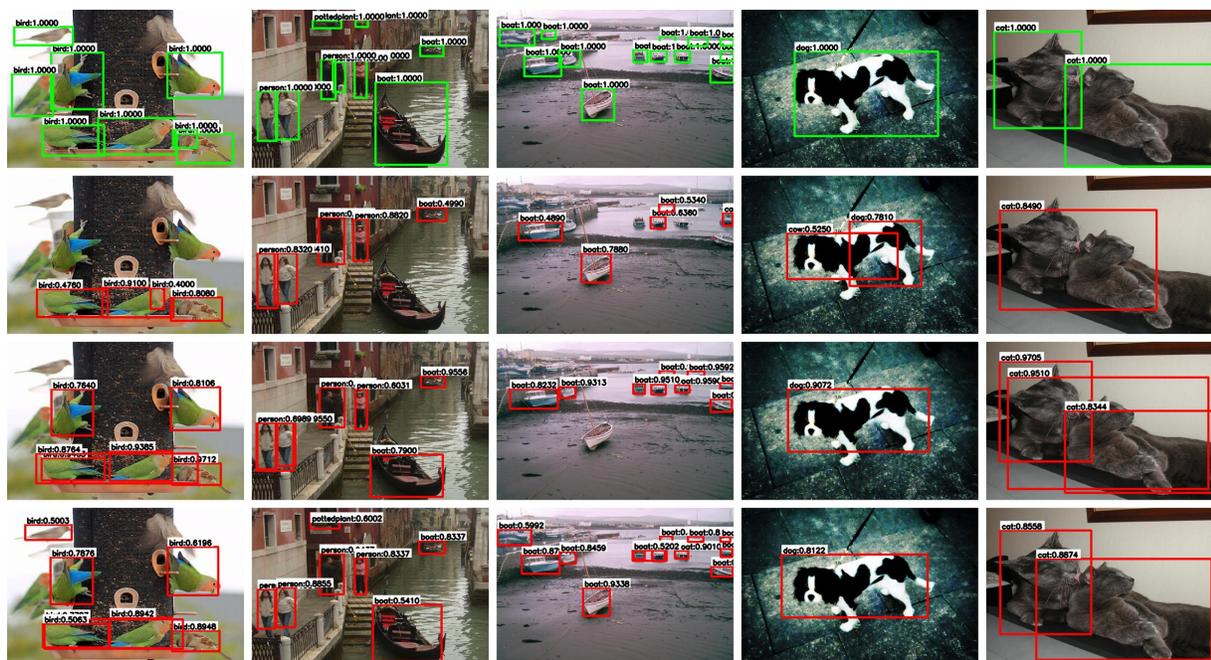


Figure A.8. Visualization of SoS-WSOD results on VOC2007 (more examples in addition to Fig. A.7). Top row: groundtruth annotations. 2nd to 4th rows: detection results from stages 1, 2 and 3, respectively.

## A.7. Per-class detection results

In Table A.3, we report and compare the per-class detection $mAP_{50}$ results on VOC2007. Besides, we also report and compare correct localization (CorLoc) results on VOC2007 trainval set in Table A.4.

## A.8. More visualization results

In Sec. 4 of the main paper, we only show some visualization results on MS-COCO due to the limited space. Here, more visualization results are shown in Fig. A.6 to A.8.