

# Scribble Hides Class: Promoting Scribble-Based Weakly-Supervised Semantic Segmentation with Its Class Label

Xinliang Zhang<sup>1,3\*</sup>, Lei Zhu<sup>1-4\*</sup>, Hangzhou He<sup>1-3</sup>, Lujia Jin<sup>1-4</sup>, Yanye Lu<sup>1,3,4†</sup>

<sup>1</sup> Institute of Medical University, Peking University, Beijing, China

<sup>2</sup> Department of Biomedical Engineering, Peking University, Beijing, China

<sup>3</sup> National Biomedical Imaging Center, Peking University, Beijing, China

<sup>4</sup> Institute of Biomedical Engineering, Peking University Shenzhen Graduate School, Beijing, China  
zhangxinliang@tju.edu.cn, {zhulei, zhuang}@stu.pku.edu.cn, {jinlujia, yanye.lu}@pku.edu.cn

## Abstract

Scribble-based weakly-supervised semantic segmentation using sparse scribble supervision is gaining traction as it reduces annotation costs when compared to fully annotated alternatives. Existing methods primarily generate pseudo-labels by diffusing labeled pixels to unlabeled ones with local cues for supervision. However, this diffusion process fails to exploit global semantics and class-specific cues, which are important for semantic segmentation. In this study, we propose a class-driven scribble promotion network, which utilizes both scribble annotations and pseudo-labels informed by image-level classes and global semantics for supervision. Directly adopting pseudo-labels might misguide the segmentation model, thus we design a localization rectification module to correct foreground representations in the feature space. To further combine the advantages of both supervisions, we also introduce a distance entropy loss for uncertainty reduction, which adapts per-pixel confidence weights according to the reliable region determined by the scribble and pseudo-label's boundary. Experiments on the ScribbleSup dataset with different qualities of scribble annotations outperform all the previous methods, demonstrating the superiority and robustness of our method. The code is available at <https://github.com/Zxl19990529/Class-driven-Scribble-Promotion-Network>.

## Introduction

Primarily driven by the availability of extensive pixel-level annotated datasets, the field of semantic segmentation has made remarkable strides in the last decade. However, the challenges of the laborious and time-consuming process of collecting and manually annotating such datasets hinder real-world applications of semantic segmentation. Weakly-supervised semantic segmentation (WSSS) methods utilizing sparse labels have emerged as a prominent trend to overcome this limitation. These methods use annotations at the image, scribble, or bounding box levels as supervision to train the semantic segmentation model. Among them, image-level annotations offer limited spatial supervision, while bounding boxes may lead to overlapping issues when objects are nearby. In comparison, the use of scribble

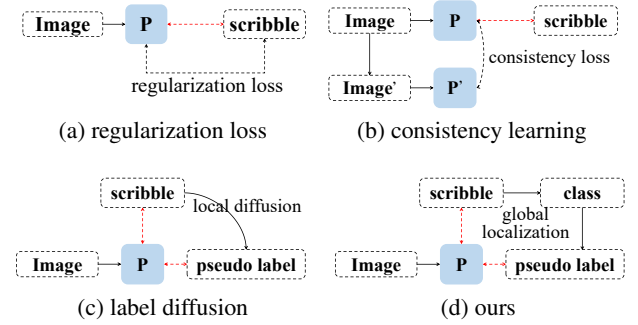


Figure 1: Schematic diagrams of different scribble-based WSSS methods. Existing approaches (a-c) overlooked the class label in scribbles, which provides image-level supervision. “P” represents the model prediction. The red dashed line represents the supervision relationship.

annotations strikes an optimal balance between supervision effectiveness and labor cost (Lin et al. 2016). Consequently, scribble-based WSSS has garnered increasing attention in recent years (Liang et al. 2022; Wu et al. 2023).

The intrinsic challenge in scribble-based WSSS lies in the partial supervision provided by sparse labels. Existing approaches have attempted to address this issue from three perspectives, namely, regularization loss (Tang et al. 2018a,b), consistency learning (Pan et al. 2021; Wang et al. 2022), and label diffusion (Lin et al. 2016; Wu et al. 2023), as illustrated in Figure 1(a-c). Specifically, regularization loss-based methods design specific loss functions to improve the stability of the models. Consistency learning-based approaches aim to capture invariant features to boost fine-grained segmentation performance through consistency loss. However, both methods fail to address the deficiency of pixel-level supervision, leading to limited performance. In contrast, label diffusion-based methods generate pixel-level pseudo-labels by diffusing labeled pixels to unlabeled ones, *i.e.* constructing a graph model on the scribble to generate pseudo-labels for training. However, the diffusion process predominantly relies on local pixel information and fails to exploit the global semantics and class-specific cues of im-

\*These authors contributed equally.

†Corresponding author, yanye.lu@pku.edu.cn

ages, which are important for semantic segmentation. In addition, such pseudo-label generation approaches are heavily dependent on the quantity and quality of the scribbles, where the model performance would be undermined when the scribbles are shrunk or dropped as shown in Figure 7. In fact, sparse scribbles inherently possess class information, which can offer valuable global semantic clues while enriching scribble-based WSSS supervision. However, this advantage has not been extensively explored in existing scribble-based WSSS researches.

In light of this, the present paper is dedicated to promoting the performance of scribble-based WSSS with a globally considered pseudo-label. The image-level class labels could be easily obtained from the scribbles, making it feasible to acquire the globally considered pseudo-label via image-level WSSS methods. Previous image-level WSSS methods have demonstrated that image-level class labels prompt models to focus on discriminative areas within an image, which can be used to compensate for the limitations of local cues provided by scribbles. Drawing inspiration from this, we propose a class-driven scribble promotion (CDSP) network for scribble-based semantic segmentation, which utilizes image-level class labels to generate pseudo-labels.

The overview of our method is depicted in Figure 1 (d). We begin by extracting image-level class labels from the scribbles and employing them to train a classification model, subsequently generating the globally considered pseudo-label. We then proceed to train a semantic segmentation model with both scribble and pseudo-label for supervision. By doing so, the inclusion of the image-level class label facilitates the acquisition of global semantic information for pseudo-label generation and further benefits the scribble-based WSSS training. Nevertheless, the noisy supervisions in pseudo-labels may affect the model, where we specifically devise a localization rectification module (LoRM) to address this issue, which corrects foreground representations in the latent feature space by referencing other pixels. To further leverage the advantages of both supervisions, we also introduce a distance entropy loss (DEL) for model uncertainty reduction, where the model prediction is assigned with per-pixel confidence based on the reliable region determined by the scribble and the boundary of the pseudo-label. With these integrated components, our method achieves state-of-the-art (SOTA) performance in scribble-based WSSS. Our contributions can be concluded as:

- We present a class-driven scribble promotion network for scribble-based WSSS that utilizes image class information to generate a globally considered pseudo-label. Notably, this is the first approach to exploit image-level class information in the scribble-based WSSS problem.
- A localization rectification module is proposed to correct the foreground representations in the latent feature space that are misled by the noisy pseudo-labels. And a distance entropy loss is proposed to excavate the reliable areas based on proximity to scribbles and pseudo-labels.
- The proposed method outperforms existing state-of-the-art methods. The extensive experiments on the different qualities of scribbles demonstrate the extraordi-

nary robustness of our method.

## Related Works

**Image-level WSSS** The remarkable achievements of early deep learning-based methods in image classification (Simonyan and Zisserman 2014) have spurred numerous works on feature visualization. Zhou et al. (2016) first introduced the class activation map (CAM) technique, which employs global average pooling on deep features to visualize discriminative localization. This technology subsequently catalyzed various efforts to generate semantic pseudo-labels from CAM, facilitating the training of segmentation networks (Kolesnikov and Lampert 2016; Zhang et al. 2021b; Zhu et al. 2023b, 2022). Recently, SEAM (Wang et al. 2020) presented a pixel correlation module that refines current pixel predictions using information on the similar neighbors of the pixel. From another perspective, AFA (Ru et al. 2022) addressed this problem with transformers leveraging multi-head self-attention for effective long-range modeling. Additionally, (Ru et al. 2023) developed patch token contrast and class token contrast modules to capture high-level semantics. The intrinsic capability of image-level supervised semantic segmentation to capture global information makes it a promising approach to promote scribble-supervised semantic segmentation.

**Scribble-based WSSS** Early methods can be traced back to traditional interactive segmentation (Rother, Kolmogorov, and Blake 2004; Grady 2006), which employ graphical models to expand the scribble area and extract foreground regions. These methods typically require multiple continuous interactions to extract foreground masks and generate semantic segmentation results. Recent scribble-based WSSS domain approaches can be categorized into three main groups: regularization loss-based methods (Tang et al. 2018a,b), consistency learning-based methods (Pan et al. 2021; Wang et al. 2022), and label diffusion-based methods (Lin et al. 2016; Vernaza and Chandraker 2017; Xu et al. 2021), as depicted in Figure 1. Regularization loss-based methods aim to enhance network robustness by preventing it from being overconfident. Consistency learning-based methods leverage self-supervised learning strategies to acquire invariant features. While both these two kinds of methods contribute to enhanced network robustness, they still struggle to address the issue of lacking supervision. Specifically, BPG (Wang et al. 2019) utilizes extra boundary data with edge information to improve segmentation performance. Label diffusion-based approaches utilize scribbles to generate pseudo-labels using unsupervised models, such as graph models, and subsequently employ these pseudo-labels to train semantic segmentation models. However, such a diffusion process fails to effectively exploit the global semantic information lurking in the image. More recent works (Liang et al. 2022; Wu et al. 2023) aim to adaptively generate pseudo-labels using a tree filter and a learnable probabilistic model with Gaussian prior, respectively. Despite their advancements, both of these methods still lack image-level supervision, thereby limiting their ability to model global semantic information effectively.

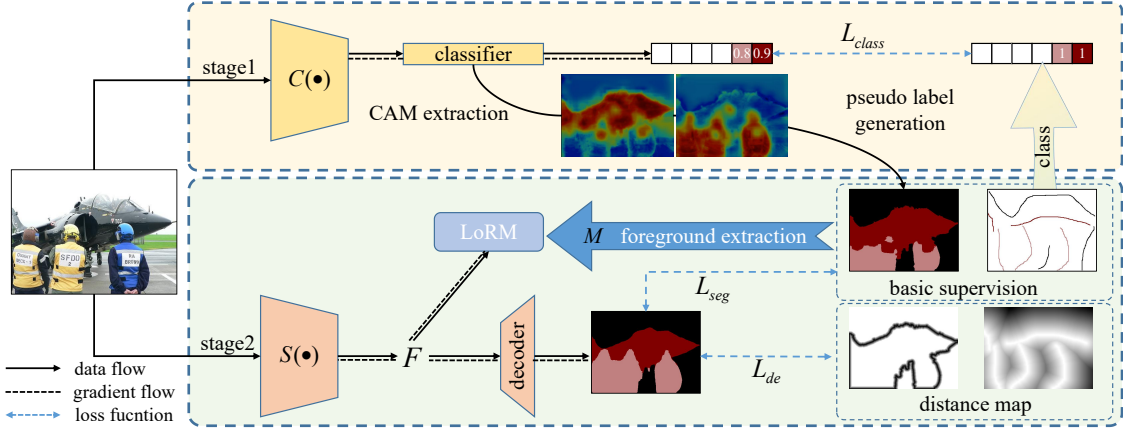


Figure 2: The overview of our method (CDSP). In the first stage, we train a classification model with the image-level class labels extracted from the scribbles to generate the globally considered pseudo-label. Then we train a semantic segmentation model with the globally considered pseudo-label and the scribble label jointly in the second stage. We propose a localization rectification module (LoRM) and a distance entropy loss to assist the training process.

**Other WSSS Methods** Points (Bearman et al. 2016; Chen et al. 2021; Wu et al. 2022, 2023; Liang et al. 2022) and bounding boxes (Dai, He, and Sun 2015; Papandreou et al. 2015; Khoreva et al. 2017; Zhang et al. 2021a) are also common annotations in weakly-supervised semantic segmentation. However, both of them fail to achieve a balance between training supervision and labor costs. The point-level annotation requires less labor, but it provides very limited supervision, hence training a high-accuracy semantic segmentation model is difficult. Bounding boxes suffer from overlapping with each other when encountering multiple objects and provide redundant supervision, which may confuse the model. In comparison, scribbles achieve the best balance between laboring cost and supervision accuracy.

## Method

In this part, we first retrospect the general problem formulation of label diffusion-based methods and their limitations. Then we introduce CDSP with the pseudo-label generation, basic supervision, LoRM, and DEL sequentially in detail.

### General Problem Formulation

Denoting  $\Omega = \{y_i | i = 1, \dots, n\}$  as the ground truth label set and  $\Omega_s$  as the sparse scribble label, where  $\Omega_s \subset \Omega$  and  $|\Omega_s| \ll |\Omega|$ . The objective function of the scribble-based WSSS can be formulated as:

$$\min c(\mathbf{P}_{\Omega_s}, \Omega_s), \quad (1)$$

where  $c(\cdot, \cdot)$  denotes the criterion function, which is usually cross-entropy.  $\mathbf{P}_{\Omega_s}$  denotes the model predictions corresponding to the sparse scribble label. Such sparse supervision limits the model’s performance and decreases the certainty of the model. Most existing label diffusion-based methods make efforts on devising a graphical diffusion model or learnable probabilistic model with low-level cues

$\phi$  to generate the pseudo-label  $\tilde{\Omega} = \{\tilde{y}_i | i = 1, \dots, n\}$  by diffusing the labeled pixels to unlabeled ones:

$$\tilde{\Omega} = \phi(\Omega_s). \quad (2)$$

Combined with Eq 1, a complete objective function for scribble-based WSSS can be obtained:

$$\min(c(\mathbf{P}_{\Omega_s}, \Omega_s) + c(\mathbf{P}_{\tilde{\Omega}}, \tilde{\Omega})). \quad (3)$$

As shown in Eq. 2, because only scribble-annotated pixels are considered, it is hard for the diffusion methods to capture the global information from the scribbles, making the diffused label  $\tilde{y}$  provide locally considered supervision. Besides, it is evident that the diffused pseudo-label heavily depends on the scribble, where its quality may be undermined by a shrunk or dropped version of the scribble.

### Class-driven Scribble Promotion

To solve the problems mentioned above, we naturally think of utilizing the class label derived from the sparse scribble to provide global cues for image-supervised segmentation when generating the pseudo-label. Denoting  $\tilde{\phi}$  as the classification model with a fully connected layer, the pseudo-label  $\tilde{\Omega}$  can be obtained from the image  $\mathbf{I} \in \mathbb{R}^{3 \times H \times W}$  with multi-class label  $\mathbf{k} \in \mathbb{R}^{1 \times K}$ :

$$\tilde{\Omega} = \tilde{\phi}(\mathbf{I}, \mathbf{k}), \quad (4)$$

where all the pixels are taken into account to generate the pseudo-label. After that, we further introduce the LoRM and DEL to strike the advantages of both supervisions as shown in Figure 2. In general, the overall loss function for supervision can be formulated as:

$$\mathcal{L} = \mathcal{L}_{seg} + \mathcal{L}_{lorm} + \mathcal{L}_{de}. \quad (5)$$

$\mathcal{L}_{seg}$  represents the basic supervision from the scribble and pseudo-label,  $\mathcal{L}_{lorm}$  represents the supervision from the LoRM, and  $\mathcal{L}_{de}$  is the supervision from DEL. The details of each component will be introduced sequentially in the following parts.

## Pseudo-label Generation and Basic Supervision

To obtain the pseudo-label with Eq. 4, we first train a multi-label classification model  $\mathcal{C}(\cdot)$  followed by a  $K$ -class classifier (e.g. Resnet (He et al. 2016) with an FC layer) with image-level classes extracted from the scribbles. After the model converges, the image  $\mathbf{I}$  is fed into the model to generate the class activate map of the  $k^{th}$  class:

$$CAM_k(\mathbf{I}) = ReLU\left(\sum_{i=1}^C \mathbf{W}_{i,k} \mathbf{F}_i\right), \quad (6)$$

where  $\mathbf{F} = \mathcal{C}(\mathbf{I})$ ,  $\mathbf{F} \in \mathbb{R}^{C \times HW}$  is the feature maps of the last layer,  $\mathbf{W}$  is the weight matrix in the classifier. We follow existing image-supervised semantic segmentation methods to threshold the  $CAM$  into binary masks and integrate them into a single channel multi-class mask (Wang et al. 2020; Chen et al. 2022) to generate the pseudo-label  $\tilde{\Omega}$ . It is also possible to adopt one-stage image-supervised WSSS methods (Ru et al. 2022; Zhu et al. 2023a) as  $\hat{\phi}$  to generate the pseudo-label. With both pseudo-label and scribble, the basic supervision can be summarized as:

$$\mathcal{L}_{seg} = \mathcal{L}_{segs} + \mathcal{L}_{segc}. \quad (7)$$

In detail,  $\mathcal{L}_{segs}$  denotes the sparse supervision from the scribble label in the form of a partial cross-entropy:

$$\mathcal{L}_{segs} = \frac{1}{|\Omega_s|} \sum_{\mathbf{y}_i \in \Omega_s} c(\mathbf{y}_i, \mathbf{p}_i), \quad (8)$$

where  $c(\mathbf{y}_i, \mathbf{p}_i) = -\sum_{k=1}^K \mathbf{y}_{i,k} \log(\mathbf{p}_{i,k})$ ,  $K$  is the class number,  $\mathbf{p}_i$  is the prediction from the model,  $\mathbf{y}_i$  is the one-hot label.  $\mathcal{L}_{segc}$  denotes the supervision from the pseudo-label, which can be formulated as a smoothed cross-entropy:

$$\mathcal{L}_{segc} = \frac{1}{|\tilde{\Omega}|} \sum_{\mathbf{y}_i \in \tilde{\Omega}} [(1 - \epsilon)c(\mathbf{y}_i, \mathbf{p}_i) + \epsilon c(\frac{1}{K}, \mathbf{p}_i)], \quad (9)$$

where  $\epsilon = 0.1$  is a regularization item of label smoothing (Müller, Kornblith, and Hinton 2019) to prevent the model from being over confident.

## Localization Rectification Module

Adopting the pseudo-label directly for supervision can lead to absurd predictions (Wang et al. 2018), particularly evident when foreground objects are nearby, as illustrated in Figure 3(c). Rather than correcting the pseudo-label itself, we are motivated to refine the feature representations of the model so that the model can adopt pseudo-labels with different qualities. To achieve this goal, we propose a novel module namely LoRM. The primary concept behind the LoRM is to leverage the inherent similarity of representations among foreground pixels belonging to the same semantic class. By doing so, mispredicted pixels can be refined through a weighted combination of representations from other pixels. Let  $\mathbf{F} \in \mathbb{R}^{C \times H \times W}$  denotes the feature map generated by the last layer of the segmentation backbone  $\mathcal{S}(\cdot)$ , and  $\mathbf{M} \in \mathbb{R}^{H \times W}$  denotes the pseudo mask, as depicted in Figure 2. The LoRM takes  $\mathbf{F}$  and  $\mathbf{M}$  as inputs and operates accordingly to rectify the representations.

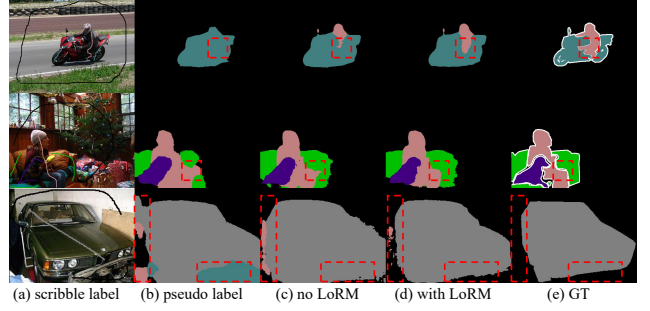


Figure 3: Visualization results employing resnet50 backbone and deeplabV2 segmentor. (a) is the original image with scribble label, (b) is the pseudo-label for training, (c) is the prediction trained with  $\mathcal{L}_{seg}$ , (d) is the prediction trained with  $\mathcal{L}_{seg} + \mathcal{L}_{lorm}$ . (e) is the ground truth label.

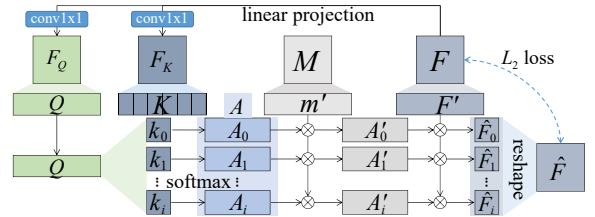


Figure 4: The illustration of LoRM.

As detailed in Figure 4, the feature map  $\mathbf{F}$  is firstly linear projected into  $\mathbf{F}_Q \in \mathbb{R}^{C \times H \times W}$  and  $\mathbf{F}_K \in \mathbb{R}^{C \times H \times W}$  with a single convolution, then flattened along the row axis into  $\mathbf{Q} \in \mathbb{R}^{C \times HW}$  and  $\mathbf{K} \in \mathbb{R}^{C \times HW}$ . Taking  $\mathbf{K}$  as the key set to be refined, and  $\mathbf{Q}$  as the query set for similarity matching, we calculate the weighted similarity matrix  $\mathbf{A}$  by:

$$\mathbf{A} = \text{softmax}\left(\frac{\mathbf{Q}^T \mathbf{K}}{\|\mathbf{Q}^T\|_2^C \|\mathbf{K}\|_2^C}\right), \quad (10)$$

where  $\mathbf{A} \in \mathbb{R}^{HW \times HW}$ ,  $\text{softmax}$  is implemented along the row axis, the L2-norm operation  $\|\cdot\|_2^C$  of  $\mathbf{Q}^T$  and  $\mathbf{K}$  is implemented along the channel dimension. Each row  $\mathbf{A}_i$  in the matrix  $\mathbf{A}$  describes the similarity between the  $i$ -th feature vector in  $\mathbf{K}$  and all the  $HW$  feature vectors in  $\mathbf{Q}$ . With the help of Eq. 10, the  $i$ -th feature vector can be refined by referencing the feature vectors in other locations. It is worth noting that, the background vectors vary largely, and contribute little to the foreground rectification. Therefore, we extract the foreground mask  $\mathbf{M} \in \mathbb{R}^{H \times W}$  from the pseudo-label and flatten it along the row axis, then element-wise multiply it with  $\mathbf{A}$  leveraging the broadcast technique:

$$\mathbf{A}' = \text{flatten}(\mathbf{M}) * \mathbf{A}, \quad (11)$$

so that the background features in each row  $\mathbf{A}_i$  are largely suppressed in its masked one  $\mathbf{A}'_i$ . Then the original feature map  $\mathbf{F}$  is flattened along the row axis, and it is matrix-multiplied with the masked similarity matrix  $\mathbf{A}'$ :

$$\hat{\mathbf{F}} = \delta * \text{flatten}(\mathbf{F}) \mathbf{A}', \quad (12)$$



where  $\delta$  is a learnable parameter initialized with 1 to control the rectification degree,  $\hat{\mathbf{F}} \in \mathbb{R}^{C \times HW}$  is the refined feature which is finally reshaped back to  $\mathbb{R}^{C \times H \times W}$ . The mean square error loss (MSE) is implemented on the original feature  $\mathbf{F}$  and the refined feature  $\hat{\mathbf{F}}$ :

$$\mathcal{L}_{lorm} = MSE(\mathbf{F}, \hat{\mathbf{F}}). \quad (13)$$

The whole process is realized by efficient matrix operations. With the supervision of Eq. 13, the LoRM achieves the goal of rectifying the misled foreground representations by referencing the representations in other foreground locations.

### Distance Entropy Loss

The LoRM effectively addresses the misalignment in the feature space in the foreground area, but the model remains susceptible to being misled by noisy labels near the object boundary during later training steps. This could undermine the efforts of LoRM and reduce the model’s certainty.

To overcome this challenge, it becomes crucial to identify reliable predictions. We propose that discriminative areas, such as the surroundings of the scribble, are more reliable and should be assigned higher confidence. Conversely, indiscriminative areas like the boundary of the pseudo-label, generated by global class supervision, are less reliable and should be assigned lower confidence. Based on this concept, we introduce a distance map strategy, to assign predictions with different confidence levels according to their distance from the scribble and the pseudo-label boundary respectively, known as the distance entropy loss. By doing so, we can better leverage the advantages of both supervisions during model training.

For the pseudo-label, the pixels around its boundary are indiscriminative, and such an area is probable to provide uncertain supervision. Denoting the coordinates of the  $i^{th}$  point in the image as  $(m, n)$ , and the coordinates of the  $j^{th}$  point on the foreground pseudo-label boundary as  $(m', n')$ , the distance maps of the pseudo-label is designed as:

$$d_c(i) = \min_{\forall j} \left( \frac{\lfloor \sqrt{e^{\lambda_c} [(m - m')^2 + (n - n')^2]} \rfloor_{255}}{255} \right), \quad (14)$$

where  $d_c$  is a probability ranges in  $[0, 1]$  that describes the minimum Euclidean distance between a point and the set of pseudo-label boundary points with the distance value truncated to 255 for normalization and the efficiency of data storage.  $\lambda_c$  is a coefficient to control the scope of the pseudo-label distance map as shown in Figure 5 (f-h). Denoting  $N_c$  as the number of non-zero elements in  $d_c$ , the distance entropy of the pseudo-label is formulated as:

$$\mathcal{L}_{dc} = \frac{1}{N_c} \sum_{i=1}^{N_c} d_c(i) \mathbf{p}_i \log(\mathbf{p}_i). \quad (15)$$

Compared with the pseudo-label, the scribble is certain and correct, the pixels lying around the scribble may largely belong to the same semantic class as the scribble. Moreover, the scribble lying in the foreground’s inner area provides correct supervision, which could suppress the noisy supervision in pseudo-label. But this confidence should decrease

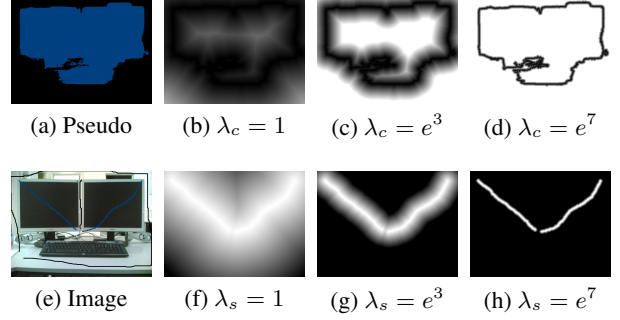


Figure 5: Visualization of distance maps with different coefficients for pseudo label boundary (b-d) and scribble (f-h)

with the increment of the distance. Therefore, denoting the coordinates of the  $i^{th}$  point in the image as  $(m, n)$ , and the  $j^{th}$  foreground scribble point coordinates as  $(m', n')$ , the distance map of the scribble is designed as:

$$d_s(i) = 1 - \min_{\forall j} \left( \frac{\lfloor \sqrt{e^{\lambda_s} [(m - m')^2 + (n - n')^2]} \rfloor_{255}}{255} \right), \quad (16)$$

where  $d_s$  is a probability ranges in  $[0, 1]$  that describes the minimum Euclidean distance between a point and the set of scribble points.  $\lambda_s$  is a coefficient to control the scope of the scribble distance map as shown in Figure 5(b-d). Denoting  $N_s$  as the number of nonzero elements in  $d_s$ , the distance entropy of the scribble is formulated as:

$$\mathcal{L}_{ds} = \frac{1}{N_s} \sum_{i=1}^{N_s} d_s(i) \mathbf{p}_i \log(\mathbf{p}_i), \quad (17)$$

Finally, the overall distance entropy can be formulated as:

$$\mathcal{L}_{de} = \mathcal{L}_{ds} + \mathcal{L}_{dc}. \quad (18)$$

Figure 5 presents visualizations of the distance maps for the scribble and pseudo-label boundaries at different coefficients of  $\lambda_s$  and  $\lambda_c$ . As  $\lambda_s$  increases, the reliable area determined by the scribble becomes more prominent. Conversely, a higher  $\lambda_c$  endows more weights to the pseudo-label in determining the reliable area. Through the distance entropy loss, we effectively excavate the reliable areas and reinforce the prediction certainty of the model by leveraging information from both the scribble and the pseudo-label boundaries.

## Experiments

**Dataset** Our experiments were carried out on the widely used ScribbleSup dataset (Lin et al. 2016), which combines PASCAL VOC2012 and SBD (Hariharan et al. 2011) datasets with scribble annotations. The dataset includes 10,582 training images and 1,449 validation images. To ensure fairness, we used the same scribble generation code as previous works (Lin et al. 2016; Tang et al. 2018b; Pan et al. 2021), maintaining uniform scribble thickness. Additionally, we validated our method on *scribble-shrink* and *scribble-drop* introduced by URSS (Pan et al. 2021) to assess its robustness in diverse scenarios.

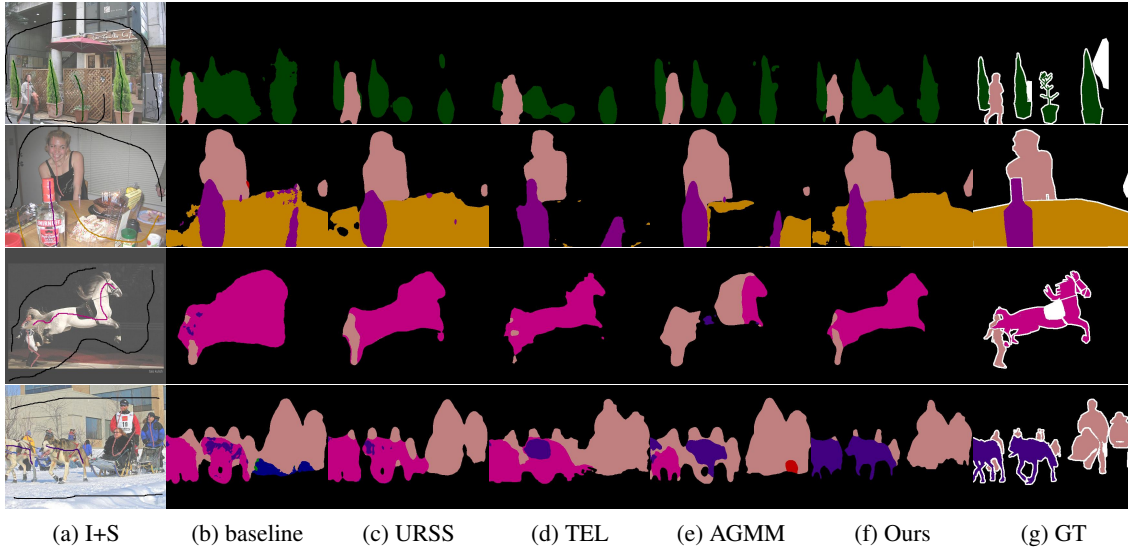


Figure 6: Visualization results comparison. (a) is the image with its scribble annotations. The baseline (b) is deeplabV3+ trained with only scribble annotations. (c) to (e) are recent methods, and (g) is the ground truth label.

**Implementation Details** With the pseudo-labels generated by BMP (Zhu et al. 2023a), we employed representative segmentation frameworks deeplabV2 (Chen et al. 2017) and deeplabV3+ (Chen et al. 2018) for method validation and generating competitive results, respectively. We conducted a total of 50 epochs with a base learning rate of  $1e^{-3}$  and batch size set to 16 for training. To ensure stable training, we adopted a learning rate warmup strategy, linearly increasing the learning rate to  $1e^{-3}$  over the first 10 epochs, followed by a cosine decay to zero over the next 40 epochs. Validation results were reported using the last checkpoint. The stochastic gradient descent (SGD) optimizer was utilized with a momentum of 0.9 and weight decay of  $5e^{-4}$ . Data augmentation followed the same strategy used in URSS. All experiments were reported with the mIoU metric (%) and conducted on one NVIDIA RTX 4090 24G GPU with an Intel Xeon Gold 6330 CPU.

**Comparison on ScribbleSup** We deploy resnet101 (He et al. 2016) as the backbone and deeplabV3+ as the segmentor with hyper-parameters of ( $\lambda_s = e^2, \lambda_c = e^7$ ) to generate the best result. The comparison details are recorded in Table 1. It is worth noting that, previous works of ScribbleSup, RAWKS (Vernaza and Chandraker 2017), and NCL (Tang et al. 2018a) adopted CRF for postprocessing, which was fairly time-consuming. For recent works of TEL (Liang et al. 2022) and AGMM (Wu et al. 2023), they were designed for general sparsely supervised segmentation, covering point level, scribble level, and box level annotations. To ensure the fairness, we reimplemented them using standard scribbles commonly used in previous works like ScribbleSup, NCL, and URSS. As shown in Table 1, our method outperforms all the previous methods, exceeding the TEL by 0.6% and AGMM by 1.6%. The test results reported in the last column of Table 1 are acquired from PASCAL VOC2012

website (Everingham and Winn 2012). The visualization comparison of our method using deeplabV3+ with previous SOTA methods is shown in Figure 6, where recent methods fail to capture correct global semantics.

**Shrink and Drop** As scribble-based annotations are flexible, it is common that the user annotates the scribbles with different length and sometimes drop some of the objects. Therefore, evaluating the model’s robustness with different shrink or drop ratios is also essential. Some shrunk or dropped samples are presented in Figure 7. Notably, as depicted in the figure, an increase in the drop or shrink ratio leads to a decrease in the model’s performance. Specifically, when the scribbles are shrunk to points (*shrink ratio* = 1), AGMM and TEL experience an approximately 10% performance degradation. In contrast, our method exhibits only a marginal drop within 1%, showcasing its robustness.

**Ablation on Components** We employ resnet50 backbone with deeplabV2 as the segmentor and use the ScribbleSup (Lin et al. 2016) dataset for training and validation. The optimal hyper-parameter combination of the distance entropy loss with all components is found by grid-search, where  $\lambda_s = 1, \lambda_c = 6$ , then we validate the effectiveness of each module by eliminating them one by one. The results are recorded in Table 2. It can be observed from the first three lines that, employing either scribble or pseudo-label as the basic supervision generates an unsatisfactory result (only around 67%), while using both of them produces a much better result (72.13%). This demonstrates that the scribble and pseudo-label provide complementary supervision and they compensate each other. Additionally, only adding  $\mathcal{L}_{dc}$  on the basic supervision degrades the model to almost the same performance as merely using  $\mathcal{L}_{segc}$ . This issue is attributed to the overfitting of the noisy labels in pseudo-labels of the

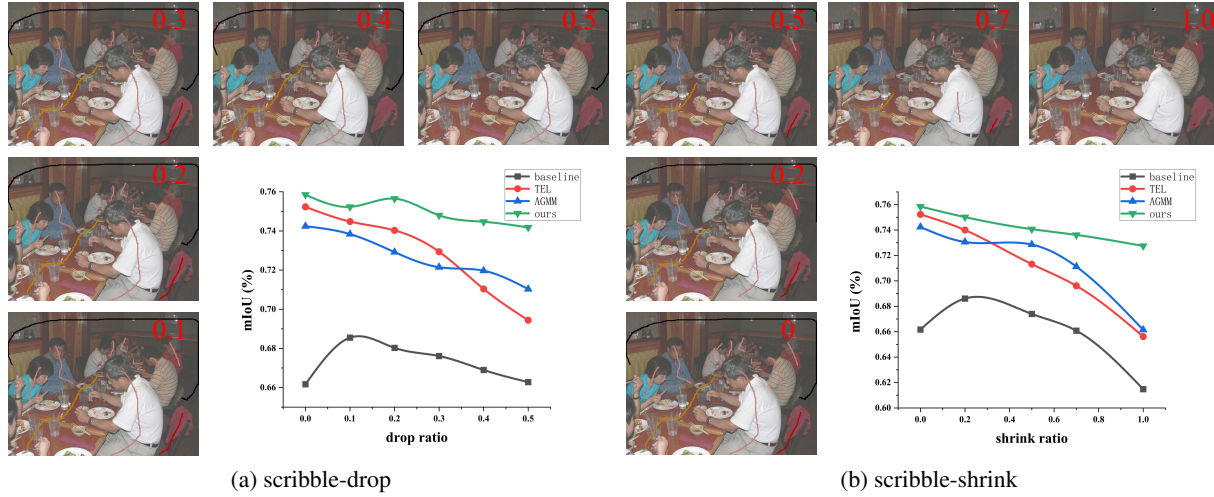


Figure 7: The experiments on scribble-drop and scribble-shrink dataset with different drop or shrink ratios.

Method	Sup	Segmentor	val	test
AFA (Zhang et al. 2021a)	$\mathcal{I}$	SegFormer	66.0	-
AMN (Lee et al. 2022)	$\mathcal{I}$	r101+v2	70.7	-
BECO (Rong et al. 2023)	$\mathcal{I}$	MiT+v3p	73.7	-
TOKO (Ru et al. 2023)	$\mathcal{I}$	ViT+v2	72.3	-
BoxSup (Dai et al. 2015)	$\mathcal{B}$	vgg16+v1	62.0	-
WSSL (Papandreou et al. 2015)	$\mathcal{B}$	vgg16+v1	67.6	-
SDI (Khoreva et al. 2017)	$\mathcal{B}$	vgg16+v1	65.7	-
BBAM (Lee et al. 2021)	$\mathcal{B}$	r101+v2	63.7	-
ScribbleSup (Lin et al. 2016)	$\mathcal{S}$	vgg16+v1	63.1	-
RAWKS (Vernaza et al. 2017)	$\mathcal{S}$	r101+v1	61.4	-
NCL (Tang et al. 2018a)	$\mathcal{S}$	r101+v1	72.8	-
KCL (Tang et al. 2018b)	$\mathcal{S}$	r101+v2	72.9	-
BPG (Wang et al. 2019)	$\mathcal{S}$	r101+v2	73.2	-
PSI (Xu et al. 2021)	$\mathcal{S}$	r101+v3p	74.9	-
URSS (Pan et al. 2021)	$\mathcal{S}$	r101+v2	74.6	73.3
CCL (Wang et al. 2022)	$\mathcal{S}$	r101+v2	74.4	-
TEL (Liang et al. 2022)	$\mathcal{S}$	r101+v3p	75.2	75.6
AGMM (Wu et al. 2023)	$\mathcal{S}$	r101+v3p	74.2	75.7
Ours	$\mathcal{S}$	r50+v2	73.9	74.2
Ours	$\mathcal{S}$	r101+v2	75.3	75.3
Ours	$\mathcal{S}$	r101+v3p	<b>75.9</b>	<b>76.0</b>
baseline (scribble only)	$\mathcal{S}$	r101+v3p	66.2	69.7

Table 1: Comparison with the state-of-the-arts methods.

model and can be addressed by our LoRM, which improves the model performance from 67.33% to 73.64%. Compared with the baseline, all the components obtain a better performance, and using them all achieves the best performance.

**Ablation on Pseudo-labels** We also conducted experiments with different pseudo-labels to assess their influence, utilizing deeplabV3+ as the segmentor. The results in Table 3 indicate that, as the pseudo-label base accuracy improves, our method exhibits increasing performance.

basic supervision		$\mathcal{L}_{de}$		$\mathcal{L}_{lorm}$	mIoU
$\mathcal{L}_{segs}$	$\mathcal{L}_{segc}$	$\mathcal{L}_{ds}$	$\mathcal{L}_{dc}$		
✓					66.17
	✓				67.23
✓	✓				72.13
✓	✓		✓		67.33
✓	✓	✓			73.38
✓	✓	✓	✓		73.58
✓	✓			✓	73.26
✓	✓	✓		✓	73.51
✓	✓		✓	✓	73.64
✓	✓	✓	✓	✓	<b>73.91</b>

Table 2: The effectiveness of each component.

This demonstrates that our approach directly benefits from image-level WSSS methods, making it a promising avenue for further development.

Method	Base acc	res50	res101
SEAM (Wang et al. 2020)	64.5	69.8	71.8
AFA (Ru et al. 2022)	66.0	71.5	73.3
BMP (Zhu et al. 2023a)	68.1	73.9	75.9

Table 3: Performance adopting different pseudo-labels.

## Conclusion

We propose a class-driven scribble promotion network for the scribble-based WSSS problem. To address the issue of model overfitting to noisy labels, we introduce a localization rectification module. Additionally, a distance entropy loss is incorporated to enhance the robustness of the network. Experimental results show that our method outperforms existing approaches, achieving state-of-the-art performance.

## Acknowledgements

This work was supported in part by the Natural Science Foundation of China (82371112, 62394311, 62394310), in part by Beijing Natural Science Foundation (Z210008), and in part by Shenzhen Science and Technology Program, China (KQTD20180412181221912).

## References

- Bearman, A.; Russakovsky, O.; Ferrari, V.; and Fei-Fei, L. 2016. What's the point: Semantic segmentation with point supervision. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part VII 14*, 549–565. Springer.
- Chen, H.; Wang, J.; Chen, H. C.; Zhen, X.; Zheng, F.; Ji, R.; and Shao, L. 2021. Seminar learning for click-level weakly supervised semantic segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 6920–6929.
- Chen, L.-C.; Papandreou, G.; Kokkinos, I.; Murphy, K.; and Yuille, A. L. 2017. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *IEEE transactions on pattern analysis and machine intelligence*, 40(4): 834–848.
- Chen, L.-C.; Zhu, Y.; Papandreou, G.; Schroff, F.; and Adam, H. 2018. Encoder-decoder with atrous separable convolution for semantic image segmentation. In *ECCV*, 801–818.
- Chen, Z.; Wang, T.; Wu, X.; Hua, X.-S.; Zhang, H.; and Sun, Q. 2022. Class re-activation maps for weakly-supervised semantic segmentation. In *CVPR*, 969–978.
- Dai, J.; He, K.; and Sun, J. 2015. Boxsup: Exploiting bounding boxes to supervise convolutional networks for semantic segmentation. In *CVPR*, 1635–1643.
- Everingham, M.; and Winn, J. 2012. The PASCAL visual object classes challenge 2012 (VOC2012) development kit. *Pattern Anal. Stat. Model. Comput. Learn., Tech. Rep.*, 2007(1-45): 5.
- Grady, L. 2006. Random walks for image segmentation. *IEEE transactions on pattern analysis and machine intelligence*, 28(11): 1768–1783.
- Hariharan, B.; Arbeláez, P.; Bourdev, L.; Maji, S.; and Malik, J. 2011. Semantic contours from inverse detectors. In *2011 international conference on computer vision*, 991–998. IEEE.
- He, K.; Zhang, X.; Ren, S.; and Sun, J. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 770–778.
- Khoreva, A.; Benenson, R.; Hosang, J.; Hein, M.; and Schiele, B. 2017. Simple does it: Weakly supervised instance and semantic segmentation. In *CVPR*, 876–885.
- Kolesnikov, A.; and Lampert, C. H. 2016. Seed, expand and constrain: Three principles for weakly-supervised image segmentation. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part IV 14*, 695–711. Springer.
- Lee, J.; Yi, J.; Shin, C.; and Yoon, S. 2021. Bbam: Bounding box attribution map for weakly supervised semantic and instance segmentation. In *CVPR*, 2643–2652.
- Lee, M.; Kim, D.; Shim, D.; and Hyun Jung, J. 2022. Threshold matters in wsss: Manipulating the activation for the robust and accurate segmentation model against thresholds. In *CVPR*, 4330–4339.
- Liang, Z.; Wang, T.; Zhang, X.; Sun, J.; and Shen, J. 2022. Tree energy loss: Towards sparsely annotated semantic segmentation. In *CVPR*, 16907–16916.
- Lin, D.; Dai, J.; Jia, J.; He, K.; and Sun, J. 2016. Scribble-sup: Scribble-supervised convolutional networks for semantic segmentation. In *CVPR*, 3159–3167.
- Müller, R.; Kornblith, S.; and Hinton, G. E. 2019. When does label smoothing help? *Advances in neural information processing systems*, 32.
- Pan, Z.; Jiang, P.; Wang, Y.; Tu, C.; and Cohn, A. G. 2021. Scribble-supervised semantic segmentation by uncertainty reduction on neural representation and self-supervision on neural eigenspace. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 7416–7425.
- Papandreou, G.; Chen, L.-C.; Murphy, K. P.; and Yuille, A. L. 2015. Weakly-and semi-supervised learning of a deep convolutional network for semantic image segmentation. In *ICCV*, 1742–1750.
- Rong, S.; Tu, B.; Wang, Z.; and Li, J. 2023. Boundary-Enhanced Co-Training for Weakly Supervised Semantic Segmentation. In *CVPR*, 19574–19584.
- Rother, C.; Kolmogorov, V.; and Blake, A. 2004. "GrabCut" interactive foreground extraction using iterated graph cuts. *ACM transactions on graphics (TOG)*, 23(3): 309–314.
- Ru, L.; Zhan, Y.; Yu, B.; and Du, B. 2022. Learning affinity from attention: end-to-end weakly-supervised semantic segmentation with transformers. In *CVPR*, 16846–16855.
- Ru, L.; Zheng, H.; Zhan, Y.; and Du, B. 2023. Token contrast for weakly-supervised semantic segmentation. In *CVPR*, 3093–3102.
- Simonyan, K.; and Zisserman, A. 2014. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*.
- Tang, M.; Djelouah, A.; Perazzi, F.; Boykov, Y.; and Schroers, C. 2018a. Normalized cut loss for weakly-supervised cnn segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 1818–1827.
- Tang, M.; Perazzi, F.; Djelouah, A.; Ben Ayed, I.; Schroers, C.; and Boykov, Y. 2018b. On regularized losses for weakly-supervised cnn segmentation. In *ECCV*, 507–522.
- Vernaza, P.; and Chandraker, M. 2017. Learning random-walk label propagation for weakly-supervised semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 7158–7166.
- Wang, B.; Qi, G.; Tang, S.; Zhang, T.; Wei, Y.; Li, L.; and Zhang, Y. 2019. Boundary perception guidance: A scribble-supervised semantic segmentation approach. In *IJCAI International joint conference on artificial intelligence*.



- Wang, B.; Qiao, Y.; Lin, D.; Yang, S. D.; and Li, W. 2022. Cycle-consistent learning for weakly supervised semantic segmentation. In *Proceedings of the 3rd International Workshop on Human-Centric Multimedia Analysis*, 7–13.
- Wang, X.; You, S.; Li, X.; and Ma, H. 2018. Weakly-supervised semantic segmentation by iteratively mining common object features. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 1354–1362.
- Wang, Y.; Zhang, J.; Kan, M.; Shan, S.; and Chen, X. 2020. Self-supervised equivariant attention mechanism for weakly supervised semantic segmentation. In *CVPR*, 12275–12284.
- Wu, L.; Fang, L.; Yue, J.; Zhang, B.; Ghamisi, P.; and He, M. 2022. Deep Bilateral Filtering Network for Point-Supervised Semantic Segmentation in Remote Sensing Images. *IEEE Transactions on Image Processing*, 31: 7419–7434.
- Wu, L.; Zhong, Z.; Fang, L.; He, X.; Liu, Q.; Ma, J.; and Chen, H. 2023. Sparsely Annotated Semantic Segmentation With Adaptive Gaussian Mixtures. In *CVPR*, 15454–15464.
- Xu, J.; Zhou, C.; Cui, Z.; Xu, C.; Huang, Y.; Shen, P.; Li, S.; and Yang, J. 2021. Scribble-supervised semantic segmentation inference. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 15354–15363.
- Zhang, B.; Xiao, J.; Jiao, J.; Wei, Y.; and Zhao, Y. 2021a. Affinity attention graph neural network for weakly supervised semantic segmentation. *TPAMI*, 44(11): 8082–8096.
- Zhang, F.; Gu, C.; Zhang, C.; and Dai, Y. 2021b. Complementary patch for weakly supervised semantic segmentation. In *Proceedings of the IEEE/CVF international conference on computer vision*, 7242–7251.
- Zhou, B.; Khosla, A.; Lapedriza, A.; Oliva, A.; and Torralba, A. 2016. Learning deep features for discriminative localization. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2921–2929.
- Zhu, L.; He, H.; Zhang, X.; Chen, Q.; Zeng, S.; Ren, Q.; and Lu, Y. 2023a. Branches Mutual Promotion for End-to-End Weakly Supervised Semantic Segmentation. arXiv:2308.04949.
- Zhu, L.; She, Q.; Chen, Q.; Meng, X.; Geng, M.; Jin, L.; Zhang, Y.; Ren, Q.; and Lu, Y. 2023b. Background-aware classification activation map for weakly supervised object localization. *IEEE Transactions on Pattern Analysis and Machine Intelligence*.
- Zhu, L.; She, Q.; Chen, Q.; You, Y.; Wang, B.; and Lu, Y. 2022. Weakly supervised object localization as domain adaption. In *CVPR*, 14637–14646.