

International Journal of Computer Vision
ISIM: Iterative Self-Improved Model for Weakly Supervised Segmentation
--Manuscript Draft--

| | |
|--|--|
| Manuscript Number: | VISI-D-21-00725 |
| Full Title: | ISIM: Iterative Self-Improved Model for Weakly Supervised Segmentation |
| Article Type: | Manuscript |
| Keywords: | Weakly Supervised Semantic Segmentation, Self SupervisedLearning, Class Activation Map, Pascal VOC 2012 |
| Corresponding Author: | Cenk Bircanoğlu Bahçeşehir Üniversitesi: Bahcesehir Universitesi İstanbul, TURKEY |
| Corresponding Author Secondary Information: | |
| Corresponding Author's Institution: | Bahçeşehir Üniversitesi: Bahcesehir Universitesi |
| Corresponding Author's Secondary Institution: | |
| First Author: | Cenk Bircanoğlu |
| First Author Secondary Information: | |
| Order of Authors: | Cenk Bircanoğlu Nafiz Arica |
| Order of Authors Secondary Information: | |
| Funding Information: | |
| Abstract: | Weakly Supervised Semantic Segmentation (WSSS) is a challenging task aiming to learn the segmentation labels from class-level labels. In the literature, exploiting the information obtained from Class Activation Maps (CAMs) is widely used for WSSS studies. However, as CAMs are obtained from a classification network, they are interested in the most discriminative parts of the objects, producing non-complete prior information for segmentation tasks. In this study, to obtain more coherent CAMs with segmentation labels, we propose a framework that employs an iterative approach in a modified encoder-decoder-based segmentation model, which simultaneously supports classification and segmentation tasks. As there are no ground-truth segmentation labels given, the same model also generates the pseudo-segmentation labels with the help of dense Conditional Random Fields (dCRF). As a result, the proposed framework becomes an iterative self-improved model. The experiments performed with DeepLabv3 and UNet models show a significant gain on the Pascal VOC12 dataset, and the DeepLabv3 application increases the current state-of-the-art metric by 1%2.5. The implementation associated with the experiments can be found: https://github.com/cenkbircanoglu/isim . |

12
13
14
15
16
17
18
19 [Click here to view linked References](#)



ISIM: Iterative Self-Improved Model for Weakly Supervised Segmentation

24 Cenk Bircanoglu^{1,2*} and Nafiz Arica^{2†}
25

26 ¹*Cognition, Adevinta, 24 rue des Jeûneurs, Paris, 75002, France.
27

28 ²Computer Engineering Department, Bahcesehir University,
29 Besiktas, Istanbul, 34353, Turkey.

30
31 *Corresponding author(s). E-mail(s):
32 cenk.bircanoglu@adevinta.com;
33 cenk.bircanoglu@bahcesehir.edu.tr;

34 Contributing authors: nafiz.arica@eng.bau.edu.tr;

35 [†]These authors contributed equally to this work.

40 Abstract

41 Weakly Supervised Semantic Segmentation (WSSS) is a challenging task
42 aiming to learn the segmentation labels from class-level labels. In the lit-
43 erature, exploiting the information obtained from Class Activation Maps
44 (CAMs) is widely used for WSSS studies. However, as CAMs are obtained
45 from a classification network, they are interested in the most discrim-
46 inative parts of the objects, producing non-complete prior information
47 for segmentation tasks. In this study, to obtain more coherent CAMs
48 with segmentation labels, we propose a framework that employs an iter-
49 ative approach in a modified encoder-decoder-based segmentation model,
50 which simultaneously supports classification and segmentation tasks. As
51 there are no ground-truth segmentation labels given, the same model also
52 generates the pseudo-segmentation labels with the help of dense Con-
53 ditional Random Fields (dCRF). As a result, the proposed framework
54 becomes an iterative self-improved model. The experiments performed
55 with DeepLabv3 and UNet models show a significant gain on the Pas-
56 cal VOC12 dataset, and the DeepLabv3 application increases the current
57 state-of-the-art metric by %2.5. The implementation associated with
58 the experiments can be found: <https://github.com/cenkbircanoglu/isim>.

59
60 **Keywords:** Weakly Supervised Semantic Segmentation, Self Supervised
61 Learning, Class Activation Map, Pascal VOC 2012
62
63

2 Iterative Self-Improved Model for Weakly Supervised Segmentation

1 Introduction

In the last decade, researchers have achieved outstanding results on the semantic segmentation problem, one of the fundamental tasks of Computer Vision. These achievements have been mainly reached within the fully supervised setting where the training dataset contains pixel-level class labels. However, accessing that kind of dataset for a specific problem is a compelling task, and it is well-known that collecting pixel-level class labels is costly and time-consuming. Consequently, this issue diverted the attention of the researchers working on semantic segmentation tasks to apply another approach called Weakly Supervised Semantic Segmentation (WSSS), where weaker labels that are less costly to collect, such as class labels, bounding boxes, scribbles, are employed to obtain pixel-level class labels.

Semantic Segmentation tasks in both fully supervised and weakly supervised settings have several challenges, such as cluttering of the background, occlusions, and intra-class variations. In addition to these challenges, there is another significant challenge in WSSS called the supervision gap. The supervision gap happens when the provided labels contain less information than the desired labels. For example, class labels contain information about the existence of the object classes in the image, but they do not hold any information about its shape, size, color, or any other attributes. In addition, there is no knowledge on if one or multiple instances of a class exist in an image. Therefore, this gap significantly increases the complexity of the segmentation task and has a notable impact on substantial performance differences between WSSS and Fully Supervised Semantic Segmentation (FSSS) solutions.

Most of the advanced studies in WSSS have employed a method called Class Activation Maps (CAMs) to obtain better segmentation labels in various ways. However, there is a significant issue with CAMs; they are interested in the most discriminative parts of the objects in the image. When CAMs are utilized to target segmentation labels, focusing on the most discriminating regions naturally concludes reaching non-complete and incorrect information about the objects and background. Most of the time, some parts of an object may be declared as background, or the other class instances may be missed in CAMs. Therefore, it is more reasonable to expect non-perfect results for the segmentation tasks with CAMs and accept that they require some modifications to reach more coherent results.

The above observations make us focus on boosting the CAMs results with the implicit information pieces of the images to narrow the performance gap between FSSS and WSSS methods. In this study, we propose a general framework having an iterative approach to reveal these implicit information pieces of the images in a multi-task learning setting. The primary task of the framework is to train against the classification and segmentation labels simultaneously within the same network by using only the class labels and images as a source. The training objective against the segmentation label is scraping the implicit information about the objects from the images with the help of segmentation loss.

14 The proposed framework is built on an encoder-decoder-based segmentation
15 model that is modified to support learning classification and segmentation
16 simultaneously. As there are no ground-truth-values for the segmentation task,
17 the same model also generates the segmentation labels. After giving the task
18 of generation of the segmentation labels, fundamentally, the proposed frame-
19 work becomes a self-improved encoder-decoder model supporting multi-task
20 learning. Both encoder and decoder parts of the network are an instance of
21 Fully Convolutional Networks (FCN). The overall network learns the semantic
22 segmentation labels while the CAMs are produced as outcomes of the training
23 of the encoder part and classification branch. The intuition behind this idea
24 is that the encoder should carry the latent features of the objects to make
25 the decoder learn segmentation labels while optimizing itself in class labels
26 training.
27

28 The main objective of the proposed framework can be summarized as pixel-
29 wise propagation of the CAMs by adding a pixel-wise loss to the training
30 phase and using the image content. Additionally, to reinforce this propagation,
31 dense Conditional Random Fields (denseCRF, dCRF) are employed in the
32 production phase of pixel-level labels from CAMs.

33 There are three main contributions of this study,

- 34 • The proposed method is a general framework that can employ any
35 encoder-decoder-based segmentation model aiming to achieve CAMs
36 which is more coherent with segmentation results.
- 37 • To the best of our knowledge, it is the first iterative self-improved
38 segmentation architecture with CAMs and dCRF.
- 39 • It is illustrated with the experiments on PASCAL VOC 2012 that
40 our framework with a version of DeepLab achieves the state-of-the-art
41 performance with only image-level annotations.

42 The paper is organized as follows: after giving a detailed overview of related
43 works on WSSS in section 2, the motivation of the study and proposed architec-
44 ture are explained in the third section. Section 4 describes the implementation
45 details, experimental results. Finally, in Section 5, we conclude the paper by
46 discussing the proposed approach with experimental results.
47

49 2 Related Works 50

51 Segmentation is an essential ingredient of many systems built on image under-
52 standing. It plays a significant role in various applications from different
53 domains such as autonomous driving, analysis of biomedical images, video
54 surveillance, augmented reality, and fashion. Image segmentation algorithms
55 designed to support these applications have led to numerous methods for an
56 extensive history of image segmentation. Even though there are many stud-
57 ies in the literature using traditional Computer Vision and Machine Learning
58 methods, such as thresholding [1], clustering [2], Conditional and Markov
59 Random Fields [3], and others [4–6], our focus in this study is on the latter
60 approaches that apply Deep Learning methods.
61

4 Iterative Self-Improved Model for Weakly Supervised Segmentation

In recent years, FSSS applications have become state-of-the-art on the popular benchmark datasets using Deep Learning methods [7–13]. The main objective of FSSS is assigning the correct class label to each pixel of the image by observing the pixel-level ground truths. To fulfill this objective, the authors [14] introduced the Fully Convolutional Network (FCN) to obtain low-dimensional segmentation masks by converting the fully connected layers to convolutional layers in CNN architecture. Conditional Random Fields (CRFs) were added as a post-processing step by executing it on the results of FCN to enhance their performances [15]. Noh et al. [16] proposed a new model which consists of two parts, encoder, and decoder to obtain relatively high-dimensional segmentation masks. In the network, an instance of FCN takes charge of the encoding part, and the decoder part employs the transposed convolution layers. This study enlarges the dimension of the segmentation result without using interpolation directly on the output of an FCN. One representative example of the encoder-decoder-based methods called UNet [17] employed specific connections between the encoder and decoder blocks and symmetric expanding paths. In addition to these, there are two popular families in FSSS called DeepLab [13, 18, 19] and R-CNN based models [20]. DeepLab introduced the usage of the dilated convolutions in FCN. R-CNN based models introduced Region Proposal Network, which proposes the candidate regions and extracts Region of Interest (RoI) and RoIPool layer, which computes the features from the proposals. There are also other methods, which are popular because of their performance, elegance, or simplicity [21–25].

Likewise, the objective in WSSS is the designation of each pixel of the images with class labels. However, WSSS applications aim to fulfill this objective by learning from image-level ground-truths instead of pixel-level ground-truths. This critical difference significantly affects the literature of WSSS as it requires applying different techniques than FSSS.

In recent years, various studies have been published on achieving segmentation labels by employing weaker labels, and these studies can be broadly categorized into four main approaches [26];

- Expectation-Maximization (EM)
- Multiple Instance Learning (MIL)
- Self-Supervised Learning (SSL)
- Object Proposal Class Inference (OPCI)

In the following sub-sections, we introduce the above categories by giving the details of decisive studies in each of them. However, the studies on natural scene images take more of our attention even though there are WSSS applications [27–34] in other domains such as histopathology and satellite images. They are not considered in scope due to the differences in the characteristics of the solutions and datasets.

2.1 Expectation-Maximization

Theoretically, Expectation-Maximization is an iterative approach that contains two main tasks: optimizing a latent distribution across the image and

learning the segmentation masks from that latent distribution. In practice, the researchers follow these steps,

- 16 1. Generate the segmentation labels from images and class labels with a
17 prior assumption
- 18 2. Train a Fully Convolutional Network to learn segmentation labels
- 19 3. Regenerate the segmentation labels from images, class labels, and the
20 features learned by FCN
- 21 4. Iterate over steps 2 and 3

22 One of the early studies of this category, called CCNN [35], converted the original
23 problem into a biconvex optimization problem and solved it by optimizing
24 the convex latent distribution of fixed FCN outputs with several constraints
25 while training the FCN against the fixed latent distribution. One representative
26 study for this category called EM-Adapt [36] sets the expectation by
27 adding a pixel-level bias to the FCN according to the given class labels, and in
28 the maximization phase, it handles the expectation by optimizing the outputs
29 of the FCN to target these pixel-level labels.

31 32 **2.2 Multiple Instance Learning**

33 Multiple Instance Learning (MIL) is one of the approaches of the supervised
34 learning framework that deals with learning from the labels containing incom-
35 plete information. In practice, this method solves the WSSS task by having a
36 model that aims to learn the class labels of the image and then assign each
37 pixel with one of the given class labels for a given image and its class labels.
38 In general, the WSSS application of MIL comes across as the training of a
39 Convolutional Neural Network (CNN) with image-level loss and inferring the
40 image locations according to the class level prediction. Class Activation Map
41 (CAM), which is extensively used in the literature for several reasons and ways
42 by other studies like this study, is an early approach of attention mechanism
43 on CNN. Mainly, it was proposed to understand, explain, and explore CNN
44 architectures better by achieving the activated areas of the image according to
45 the model's predictions. In CAM study [37], there are two constraints, (i) CNN
46 architecture should be an instance of a Fully Convolutional Network (FCN)
47 (ii) There should be a Global Pooling layer before the classifier layer, which is a
48 convolutional layer. Within these constraints, the calculation of the CAMs can
49 be accessed by the dot product between the last two convolutional layers. As
50 the calculation of CAMs requires some modifications in the well-known CNN
51 architectures, the training step is mandatory for this initial study. There are
52 other studies to extend the work of CAMs to improve the results or ease the
53 process. For example, in Grad-CAM [38], the authors used guided backpropa-
54 gation to reach CAMs without changing any part of the original network and
55 training the network. WILDCAT [39] offered a more complex pooling layer to
56 enrich the activated areas by following mainly the idea of the original CAM
57 study [37]. The authors [40] converted the VGG-16 [41] network to an FCN
58 by replacing the Dense layers with a 1x1 Convolutional layer also employed
59 Global Average Pooling as the pooling layer and trained the modified network
60
61
62

6 *Iterative Self-Improved Model for Weakly Supervised Segmentation*

to learn class-level annotations with a special MIL loss offered. In the end, the top predictions at each location were upsampled bilinearly to achieve the segmentation results. The authors [42] utilized a specific method called guided backpropagation to reach the pixel-level top-class predictions. First, with the help of guided-backpropagation, authors created the coarse class activations for multiple convolution layers of the network and aggregated them after unifying into one scale. Additionally, the process ended with the employing of dense-CRF. The authors [43] trained an FCN by targeting the class labels to generate prior foreground/background masks from the intermediate convolution layers; then, the authors used these masks to learn segmentation masks.

2.3 Self-Supervised Learning

Lately, Self-Supervised Learning has become the most popular approach with its diverse and performant examples for WSSS. The applications in this category mostly contain two steps; one to solve pretext tasks to use the outputs of it in the other task, the other one is a more complex downstream task which is most of the time the actual task. Heaps of times, the pipeline starts with learning pseudo semantic segmentation labels from the given class labels with one network and continues with a segmentation model to target the generated pseudo labels. In most studies, the first network aims to produce CAMs, and then the obtained CAMs are used in the training phase of the segmentation network.

In the literature, there are a vast amount of examples in this category as they reach good performances. SEC [44], which can be described as the best matching study to the given pipeline above, trained a VGG-16 on class labels to produce pseudo-pixel-level labels with the outcome of CAM and dense-CRF [45]. After the training of the first model, DeepLab-LargeFOV [15] was employed against the produced labels with a constraint loss as a segmentation model. Another study called MDC [46] followed the same steps with SEC with one modification; they assigned multi-dilated convolutional layers to the pseudo-ground-truths and generated the CAMs from these multi-dilated layers. AE-PSL [47] is another study that follows SEC's steps; with a significant novelty and a better performance, an adversarial erasing mechanism was integrated into the model to erase the activated regions of the previous steps to encourage the network to learn less discriminative parts. FickleNet [48] employed Grad-CAM instead of CAM to achieve a better performer segmentation model and added a center-fixed spatial dropout to the later convolutional layers to produce the pseudo-ground-truth labels to train an FCN against them. However, Grad-CAM also suffered from the same issue as it focuses on discriminative regions of the objects and the inputs of the second model are not precise.

Later, the researchers started to focus on how to take the CAM results as seed points to propagate from them. The idea behind this is that, hypothetically, adjacent pixels of the activated ones have more possibility to be a part of the same object than the others. Therefore, propagating from CAMs

14 to the adjacent regions can be efficient in producing better pseudo-ground-
15 truths. With this in mind, the authors of DSRG [49] proposed a region-growing
16 approach from CAMs to produce pixel-level labels. In the literature, there is
17 one important study called Pixel-level Semantic Affinity (PSA) [50], which
18 dramatically influences other studies, including ours, with its proposed archi-
19 tecture. The pipeline of PSA has three steps to follow and starts with a
20 CNN training to generate class activations as in SEC. In the second step of
21 the pipeline, instead of using the CAMs as ground-truth segmentation labels,
22 the CAMs were used as seed points, and one more model was added to the
23 pipeline, which performs a random walk to propagate from the seeds to achieve
24 the pseudo-ground-truth for training a segmentation model. IRNet [51] can
25 be called a version of PSA, which also targets a more complex task called
26 instance-level semantic segmentation, with additional pixel-wise clustering to
27 the random walk process. It performs the random walk from low-displacement
28 field centroids in the CAM seeds until the class boundaries to produce pseudo-
29 ground truths. Recently there are more advanced studies influenced by PSA
30 were published. One of them, PuzzleCAM [52], improved the overall perfor-
31 mance significantly with obtaining CAMs on image patches. The network to
32 obtain the CAMs transformed to a siamese network without changing its aim.
33 One branch of this siamese network was kept the same, but the other was orga-
34 nized to create tiles from the input image and merge the CAMs of these tiles.
35 In addition to these, a reconstruction loss was proposed to regulate the results
36 of both branches. There are also some studies to improve CAMs using saliency
37 maps as additional information, such as [53, 54]. Even though most of the stud-
38 ies focus on improving the first model, CAMs, in the pipeline offered in PSA,
39 one study [55] focused on the feature propagation frameworks and examined
40 the Graph Convolutional Network (GCN) instead of AffinityNet to propagate
41 from CAMs. To expand object activation regions like the others, DRS [56] used
42 a different approach than others, and it suppressed the attention on discrimi-
43 native regions and spread the attention to adjacent non-discriminative regions
44 by generating dense localization maps.

47 48 2.4 Object Proposal Class Inference

49 Most of the applications in this category first employ a low-level feature extrac-
50 tor, mainly with one of the Traditional Computer Vision methods, and use the
51 features to generate the pseudo-ground-truths. In the applications, researchers
52 assembly methods from MIL and SSL to solve the WSSS problem. The authors
53 [57] aggregate the superpixels and the features extracted from a CNN to gener-
54 ate the pseudo-ground-truths to train an FCN against them. In PRM [58], the
55 low-level object proposals were secured by employing Multi-scale Combinato-
56 rial Grouping (MCG) [59], and an FCN was trained with the loss they proposed
57 called peak stimulation loss. As the next step, they applied peak backpropa-
58 gation to convert Class Response Map to Peak Response Map for each peak.
59 By performing non-maximum suppression on the class labels, the object pro-
60 posal extracted from PRM peaks. SPML [60] employed a Contrastive Learning
61

8 Iterative Self-Improved Model for Weakly Supervised Segmentation

approach with its four types of relationships between pixels and segments in the feature space, capturing low-level image similarity, semantic annotation, co-occurrence, and feature affinity to access the segmentation masks.

3 Methodology

In most WSSS approaches, CAMs are used to obtain pseudo-segmentation labels as it is or as prior information for the subsequent stages. We also utilize CAMs [37] and propose a generic framework, which iteratively improves the semantic segmentation model in an encoder-decoder-based architecture. This new framework contains elements from all the WSSS categories mentioned in section 2. The proposed method is an example of EM due to its iterative approach, and the overall pipeline is an application of SSL. To generate the pseudo-segmentation labels, CAMs, which is an application of MIL approach, are employed. As CAMs are improved with the help of dCRF, it also contains the elements from OPCI.

In this section, after revealing our motivation, the proposed framework is described in detail.

3.1 Motivation

The utilization of CAMs to attain pseudo-segmentation labels is highly popular in the literature, and they are widely used in WSSS applications. These applications achieved relatively great results on the benchmark datasets; however, they have two main issues that can be improved from our perspective.

The first issue appears as a result of the standard process of CAMs generation. The typical procedure of obtaining CAMs contains a classification model optimized with a classification loss that mainly focuses on the most discriminative parts of the objects in the image that causes to ignore the less discriminating parts. Ignoring less discriminative regions or concentrating on some parts of the objects is reasonable for a classification model as its primary interest is the information of the object's existence in the image. However, when it comes to the segmentation models, every single pixel is equally essential. As a result of this contradiction between the nature of classification and segmentation, employing standard CAMs procedure to generate pseudo segmentation labels in WSSS applications can perform well; however, it can not achieve the optimal solution by itself. In light of this information, injecting a pixel-level loss to the training phase may help to achieve more coherent CAMs hypothetically by trying to equalize the importance of the pixels and learn from the image content. This pixel-level loss will help activate the less discriminative regions of the objects in the CAMs.

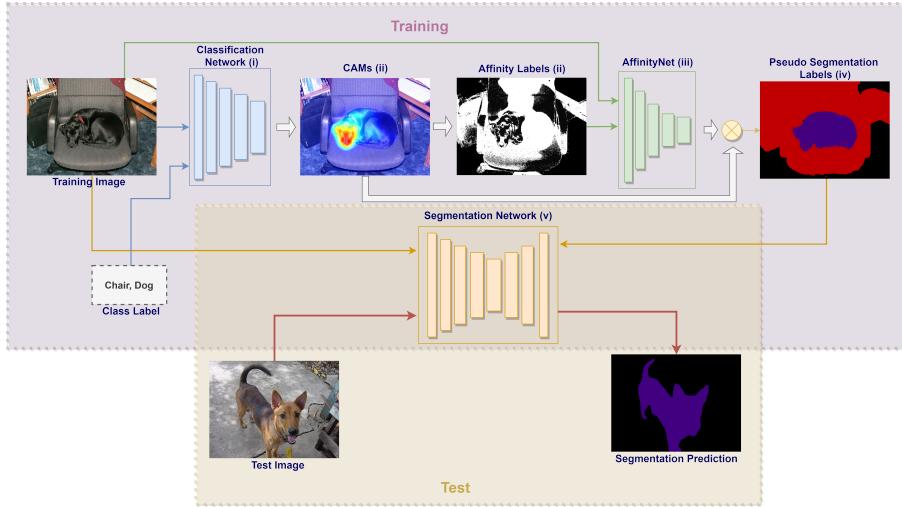
The other downside of using CAMs appears in the transition step of creating the pseudo-segmentation labels from them by applying a threshold. This thresholding process on the obtained CAMs has vital importance as the decision about each pixel belonging to an object or background is done here. This transition process mainly results in having non-complete, non-precise

14 results. The most identified problem is the wrong categorization of a pixel
15 of the foreground object as a background pixel. Therefore, working on these
16 pseudo-segmentation labels as ground truth significantly affects the next steps.
17 Moreover, in previous applications, there is no link between the threshold value
18 and the training steps of the classification network. Therefore, we propose an
19 iterative process with a pixel-level loss by setting this threshold at the begin-
20 ning of the training and implicitly optimizing CAMs and the classification
21 network according to this threshold value.

22 These observations above bring us to the point that the CAM-based WSSS
23 model should aim to learn from the whole body of objects and be boosted with
24 a generated pixel-level loss. We think that it can be achieved by iteratively
25 improving the CAMs with the help of a pixel-level segmentation loss using
26 pseudo-ground-truths at each step of the training phase. In addition, deciding
27 the threshold value at the beginning of the process and optimizing the model,
28 CAMs can be crucial to have more coherent pseudo-segmentation labels. With
29 this motivation, we propose a framework based on an encoder-decoder seg-
30 mentation model that generates its segmentation labels from the classification
31 branch with the help of denseCRF and thresholding process and learns these
32 labels with the decoder part while learning the classification labels simultane-
33 ously and iterating over this process multiple times to improve itself at each
34 iteration.

35 3.2 Proposed Framework

36 We follow the general pipeline of the PSA method [50] previously proposed
37 and used by many other studies such as PuzzleCAM, CDA, and SEAM [52,
38 61, 62]. This section first presents the PSA pipeline and then describes the
39 modifications in the pipeline as a proposed framework in detail.

10 *Iterative Self-Improved Model for Weakly Supervised Segmentation*31 **Fig. 1:** End-to-End Pipeline of PSA
32
33

34 The PSA pipeline contains several steps, visualized in Figure 1 in the training
 35 phase. It is a representative application of a two-stage framework with a
 36 Self-Supervised Learning approach to achieve a semantic segmentation model
 37 using class labels. Generation of the pseudo-ground-truths in the first stage
 38 and targeting them in the second stage are the two stages of Self-Supervised
 39 approaches. In the first stage, after obtaining the CAMs, PSA offers a new
 40 model called AffinityNet, which propagates CAMs by using class-agnostic sim-
 41 ilarities to generate segmentation labels. In the second stage, the segmentation
 42 network is trained by taking these labels as pseudo-ground-truths. Briefly, PSA
 43 contains five steps in the training phase to fulfill this two-stage framework: (i)
 44 training of the classification model to obtain CAMs, (ii) generation of pixel-
 45 level affinity labels from CAMs with applying a threshold and denseCRF,
 46 (iii) training AffinityNet to learn class agnostic pixel-similarities targeting
 47 pixel-level affinity labels, (iv) generation of the segmentation labels from the
 48 combination of CAMs and predictions of AffinityNet, (v) training of the seg-
 49 mentation model using segmentation labels. In the beginning, the classification
 50 model takes the training images with their class labels and produces CAMs.
 51 Then these CAMs are processed with threshold and denseCRF to reach affin-
 52 ity labels. AffinityNet is trained with the images and their affinity labels, and
 53 then segmentation labels are generated by applying a random walk algorithm
 54 on CAMs based on the affinity matrix produced from AffinityNet. Finally,
 55 the segmentation network is trained by taking those segmentation labels as
 56 pseudo-ground-truths. In the test phase, Segmentation Network is used to
 57 obtain segmentation results in the test images.
 58
 59

In PSA-based approaches, adjustments are primarily made in the first step to improve its overall performance by reaching more accurate CAMs and affinity labels. As the problem of the CAMs mainly comes from the loss function and employing a classification network to obtain them, the most straightforward alternative is employing a segmentation network with a segmentation loss. However, as there are no ground truths for the segmentation task, this requires extra supplementary updates.

This study proposes an iterative approach to train a model that concurrently targets both the classification and segmentation tasks while generating pseudo-segmentation labels. The proposed model is an encoder-decoder-based segmentation model, which enhances the pseudo-segmentation labels and CAMs at each iteration.

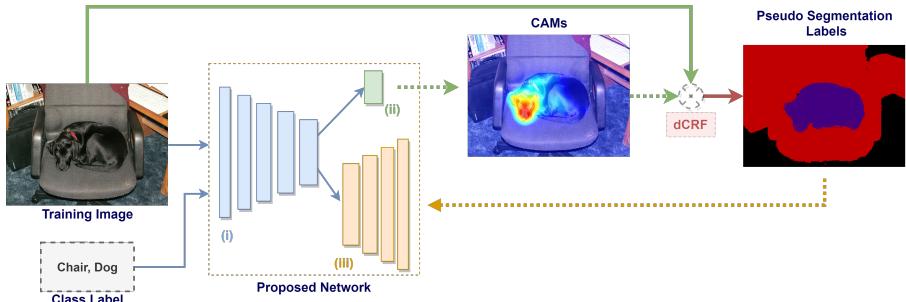


Fig. 2: Proposed Architecture: (i) Encoder network, (ii) Classification Branch, (iii) Decoder network

In the proposed encoder-decoder-based architecture, the tasks for WSSS learning are shared by different parts of the model. The encoder network learns the class-level annotations while discovering the latent features in CAMs to support the segmentation task. The decoder network aims the pixel-level annotations. The training steps of the proposed architecture are given in the pseudo-code Algorithm-1 and visualized in Figure 2. The overall training phase is an iterative approach containing three main steps. Even though it can be divided into three steps to make it easy to explain and formalize, in practice, the training is done in one step without stopping or starting it. Given the training set of images and their class labels, the first step is to train the encoder network, from which we can extract CAM of each training image. This step can be considered as the initialization of the whole model. In the second step, pseudo-segmentation labels are produced with the help of the denseCRF algorithm. In the last step, the whole model is retrained by taking the pseudo-segmentation labels as ground truth. Two different methods are employed for the pseudo-segmentation label generation and encoder-decoder training steps. In the first one, iteration continues until convergence, and in the second one,

10
11 12 *Iterative Self-Improved Model for Weakly Supervised Segmentation*
1314 fixed-size iterations are used. The performances of these two methods are given
15 in the experimental results.
1617 **Algorithm 1** Training steps of the Proposed Architecture
18

```

1: procedure TRAIN(images, class_labels)
2:   x  $\leftarrow$  images
3:   z  $\leftarrow$  class_labels
4:   y  $\leftarrow$  null
5:   Model  $\leftarrow$  SegmentationModelwithPretrainedImageNet
6:   Model = TrainEncoder(Model, x, z)
7:   for i = 1 to N do
8:     y = GeneratePseudoGroundTruths(Model, x, z)
9:     Model = TrainEncoderDecoder(Model, x, y, z)
10:   end for
11:   return Model
12: end procedure

```

32
33 The only architectural change in the encoder-decoder-based segmentation
34 model is that a new classification branch, marked as (ii) in Figure 2,
35 is added after the encoder network with Global Average Pooling (GAP) and
36 (1x1) Convolutional layer to support classification and to produce CAMs. The
37 proposed architecture is generic to practice with any encoder-decoder-based
38 segmentation model.
3940 In the training procedure, we utilize denseCRF to refine CAMs and practice
41 the outputs as pseudo-ground-truths for the segmentation task. The utilization
42 of denseCRF over CAMs and accessing the segmentation labels have
43 two implicit benefits. First, to improve the CAMs, the power of denseCRF
44 is embedded in the encoder network. Second, the information pieces of the
45 images become more extractable for the model. After generating and using
46 segmentation labels in the model, the problem becomes multi-task learning.
47 With these changes, the encoder network is simultaneously optimized to learn
48 the discriminative regions and the other regions related to the object with the
49 help of classification and segmentation loss, respectively. Hypothetically, learning
50 other features should positively affect the CAMs, and to challenge this
51 hypothesis, we extensively operate experiments on the proposed framework.
5253 However, in practice, the proposed framework has two potential issues
54 which can affect the experiments significantly and make the training unstable.
55 The first one is observed when each pixel of the pseudo segmentation masks is
56 marked as the background. The second issue is that classification or segmentation
57 loss suppresses the other one, and as a consequence, both or one of them
58 starts to diverge. The first issue is suspected to be seen mostly when the clas-
59 sification network is not confident enough to distinguish between the classes.
60 In addition that, discriminating between the classes is not enough as there is a
61 threshold applied to CAMs. The model should have a higher confidence score
62

14 in the relevant pixels than the threshold value. Also, it is reported in Section
 15 that the classification models are not 100% accurate on predictions. Therefore,
 16 there is a possibility of inconsistency in the pseudo-ground-truths, which
 17 will affect the overall training process. The decoder part of the model becomes
 18 confused when there are all background pixels as targets for the segmentation
 19 task. As a result, it will affect the encoder part and may make the classifica-
 20 tion loss diverge. The function given in Eq 2 is proposed to ease this issue.
 21 The aim of the proposed function is that when there is a segmentation mask
 22 that involves only background pixels, ignore that mask and do not calculate
 23 any loss for that specific image.

24 Supposed that y_i is the threshold applied CAM for the instance i and it is
 25 defined as,

$$27 \quad 28 \quad 29 \quad 30 \quad 31 \quad 32 \quad y_i = \begin{pmatrix} a_{1,1} & a_{1,2} & \cdots & a_{1,n} \\ a_{2,1} & a_{2,2} & \cdots & a_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m,1} & a_{m,2} & \cdots & a_{m,n} \end{pmatrix} \quad (1)$$

33 the formulation of the segmentation loss for one instance becomes as below
 34 with our basic modifications,

$$35 \quad 36 \quad \text{ModifiedCrossEntropyLoss}(\hat{y}, y) = \text{CrossEntropyLoss}(\hat{y} \cdot r_y, y) \quad (2)$$

37 where r_y is defined as ,

$$38 \quad 39 \quad 40 \quad 41 \quad 42 \quad r_y = \begin{cases} 0, & \text{if } \sum_{j=0,k=0} a_{jk} = 0 \\ 1, & \text{else } \sum_{j=0,k=0} a_{jk} > 0 \end{cases} \quad (3)$$

43 On the second issue, instead of generating the pseudo-ground-truths in each
 44 epoch, we create them with a different frequency to give some room to the
 45 network to learn class and segmentation labels. On this point, we investigate
 46 and report the outcomes in the following sections.

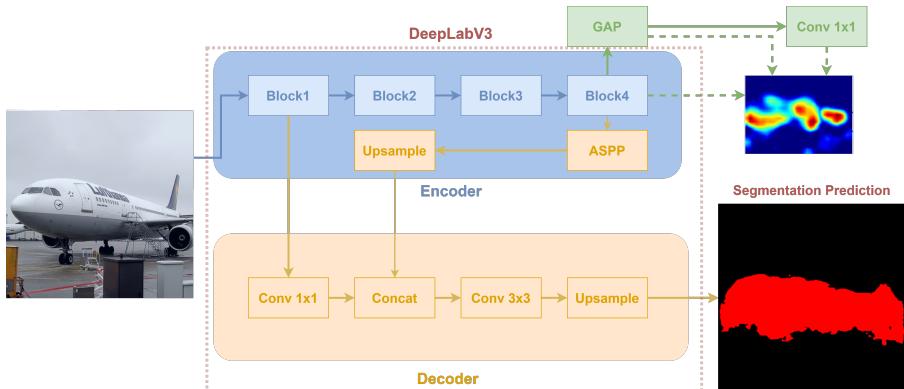
47 To challenge the proposed framework and process, we concentrate on two
 48 common and well-known segmentation models, one from the DeepLab family,
 49 DeepLabv3, and the other is the UNet model. The original architectures of
 50 DeepLabv3 and UNet are modified according to the recommended method.
 51 The modified architectures are called DeepLabCAM and UNetCAM to refer
 52 in our study.

53 3.2.1 DeepLabCAM

56 DeepLab is one of the most popular image segmentation approaches in the
 57 literature. The networks belonging to this family achieved marvelous results on
 58 various tasks in fully supervised settings. In this study, one version of DeepLab
 59 called DeepLabV3 is preferred as a segmentation network to transform for a
 60 strong reason as the primary network. The literature shows that the dimension
 61 of the CAMs has an essential role in the results. Most of the applications in
 62

14 *Iterative Self-Improved Model for Weakly Supervised Segmentation*

the literature adopt a network from the ResNet family to obtain the CAMs. When the original ResNet is applied for that task, the output dimension of the CAMs becomes 16×16 , and with experiments, it is shown that using 32×32 dimensioned CAMs perform better. Moreover, the dimension is doubled by changing the dilation rate of the last block of the ResNet family. In addition, one of the critical features of the DeepLab family is denoted as using the dilated convolutional layers. With these in mind, we add the classification branch after the fourth convolutional block of the encoder part of the DeepLabv3, as shown in Figure 3.

37 **Fig. 3:** DeepLabCAM Architecture

3.2.2 UNetCAM

UNet is also one of the most popular architectures in segmentation studies with unique characteristics and features. Two different forms of the UNet are modified, and they are executed in the experiments to prove once more the hypothesis that adding segmentation loss affects the CAMs positively.

In both UNetCam versions, ResNet50 is employed as the encoder part, and the decoder part is implemented concerning the constraints in the original UNet architecture. In UNet, there are skip connections between the encoder and the decoder blocks, and these skip connections are kept as they are. The difference between the UNetCAM versions comes from the implementation details of the upsampling layers. The first version contains non-learning upsampling layers, and the second one has learnable deconvolution layers. Also, in both networks, the classification branch is injected after the encoder network as proposed.

In the following section, the implementation details and the experimentation settings are given. The experimentation settings are organized to challenge the proposed framework from multiple perspectives such as genericity, stability, performance.

4 Experiments

The proposed framework is assessed under numerous settings from four different perspectives. First of all, to reveal its primary effect on the performance metric, the proposed framework is compared extensively with its original form, CAMs. As the second, two different segmentation models are transformed according to the proposed method, and they are executed to support the genericity of the approach. As a third perspective, we study the effects of the pseudo-segmentation regeneration frequency as the proposed training procedure possesses iterations over the network. The fourth collection of experiments are arranged related to the impact of the size of the images as image resolution has an apparent effect on the segmentation part.

After investigations introduced above are concluded on the proposed method, another experimentation set is practiced to challenge the other conjecture improving CAMs is the key to enhancing WSSS results. In that direction, the pipeline mentioned previously is applied by employing the CAM results of our proposed framework. All the details of the experiments and the results are given in the following parts of the section.

Additionally, another comparison is made between the end-to-end pipeline with the proposed framework and the state-of-the-art methods in the literature.

All the experiments are performed on PASCAL VOC 2012 dataset [63] with 20 class annotations and additional background annotation. The splits of the official dataset have 1464, 1449, and 1456 images for training, validation, and testing, respectively. We obey the conventional experimental protocol to make fair comparisons and take the additional annotations from SBD [64] to build an augmented training set. There are 10582 images in this extended training set, and it is named the trainaug set in this study. In addition to the images, we only use image-level classification labels during training phases.

Furthermore, segmentation results are assessed on the mean intersection over union (mIoU) metric. In addition to metrics, visualizations of the pseudo-segmentation labels are represented to investigate the effects of the proposed method. The dataset does not contain the ground-truths values for the test split, and to evaluate our results, we use the evaluation server maintained by the owners of the dataset. This evaluation server calculates the metrics on the test split, and they are used as it is. Even though there are ground-truth values for the validation dataset, the same method is preferred among the validation split, aiming for consistency in the results.

Additionally, all the experiments are implemented in PyTorch [65], and a computer that contains 4 RTX 2080 GPU, is used to operate the end-to-end pipeline.

4.1 Experimentations on the Proposed Framework

The proposed framework is challenged from several perspectives. Before presenting the experimental results, we introduce all the crucial functions and

10
11 16 *Iterative Self-Improved Model for Weakly Supervised Segmentation*12
13 parameters applied in these experiments by grouping them as preprocessing,
14 generating pseudo-ground-truths, and training.
1516 **Preprocessing:** In the training phase, the images are resized according to
17 these rules as follows,

- 18 • If the longer side of the image is larger than 640, it is scaled down to 640.
-
- 19 • If the longer side of the image is smaller than 320, it is scaled up to 320.
-
- 20

21 Also, random horizontal flip and random cropping are applied to the images
22 with the given resolution value, which is 320×320 if not otherwise declared,
23 and images are normalized between 0 and 1. Also, the exact same cropping
24 and scaling operations are followed for the pseudo-segmentation labels. In the
25 CAMs prediction phase, four different scaled versions and their horizontally
26 flipped images are used separately to generate CAMs, and the outputs are
27 concatenated into one. 1, 0.5, 1.5, and 2 values are used as scales.28 Although there are possible other data augmentation options, we employ
29 the same steps with the previous studies to make the comparison fair.30 **Generation of Pseudo-ground-truths:** Three steps are applied to gen-
31 erate pseudo-ground-truths, respectively (i) Producing CAMs as a result of
32 encoder and classification branch, (ii) Applying a threshold value to decide
33 which pixel involves foreground and background, and (iii) Executing dCRF on
34 the thresholded CAMs. The threshold value is chosen as 0.3, although choos-
35 ing another value for this threshold does not have a significant effect if they
36 are not too big or too small, as the network feeds itself with the outcomes and
37 reaches the balance with the result of the iterations. Also, this behavior can
38 be pointed out as a significant difference from the previous studies as it is a
39 new addition to the CAM process. Additionally, Unary potentials are used in
40 dCRF calculations.
4142 **Training Process:** To support the classification and the segmentation
43 simultaneously, the model employs two different loss functions. The multi-label
44 classification task calculates the loss with the Multi-Label Soft Margin loss
45 function, and modified pixel-wise Cross-Entropy loss, which is described in
46 Section 3.2, is adopted for segmentation tasks. Additionally, Poly optimizer, an
47 extension of Stochastic Gradient Descent (SGD), is employed with momentum
48 value 0.9 with different initial learning rates for the encoder, decoder, and
49 classifier layers as 0.1, 0.01, 1, respectively.50 The training is performed for 50 epochs in the default settings, and the first
51 five epochs are only trained among classification labels. After the fifth epoch,
52 the network is trained with both class and segmentation labels for each epoch.
53 In default settings, segmentation labels are regenerated every ten epochs after
54 the fifth epoch.
5556 4.1.1 CAM vs. DeepLabCAM
5758 To prove our hypothesis and evaluate our framework, an extensive set of exper-
59 imental configurations is designed. First of all, to perceive the effects of the
60

proposed method, several variations of it and the corresponding implementation of CAM formed. In this direction, seven different backbones are employed for CAM and DeepLabCAM methods to compare them.

In the first experiment set, the classifier backbones of the DeepLabCAM and CAM are adjusted as the versions of the ResNet and ResNeSt network family. Dilated convolution is used in the last blocks of all networks to double the dimension of the CAMs.

Before presenting the segmentation results, in Table 1, classification scores are shared. Even though previous WSSS studies ignored these metrics, they are essential metrics as the accuracy of the classification model affects the CAM performance directly. If the model does not correctly predict the class, there will be no information about it to learn in the following steps of the pipeline. As shown from Table 1, models containing ResNeSt based backbones perform better in the classification task and, most probably, obtain better CAMs.

Table 1: Classification Accuracy

| | ResNet | | | ResNeSt | | | |
|-------|--------|-------|-------|---------|-------|-------|--------|
| | 50 | 101 | 152 | 50 | 101 | 200 | 269 |
| train | 95.86 | 96.35 | 96.58 | 97.46 | 97.96 | 98.49 | 0.9861 |
| val | 93.52 | 93.33 | 93.95 | 95.44 | 96.38 | 96.2 | 96.31 |

Table 2: mIoU on Pseudo Segmentation Labels

| | ResNet | | | ResNeSt | | | |
|-------------------------|--------|-------|-------|---------|--------------|-------|-------|
| | 50 | 101 | 152 | 50 | 101 | 200 | 269 |
| CAM ¹ | 48.93 | 50.33 | 50.83 | 53.49 | 57.78 | 57.33 | 54.75 |
| DeepLabCAM ¹ | 54.89 | 57.11 | 58.35 | 56.91 | 63.09 | 62.48 | 61.11 |

This table contains the mIoU metric for the trainaug split of PASCAL VOC 2012 dataset

The results of the CAM and DeepLabCAM versions are reported in Table 2. Moreover, significant performance improvement is apparent for each backbone from that table as the gap is more than 3% even in the closest form, rising higher than 7%.

The proposed method is proved to perform better than CAMs, with results presented in Table 2, at least for the modified versions of the DeepLabv3. The next set of experiments are done with another segmentation model called UNet.

18 *Iterative Self-Improved Model for Weakly Supervised Segmentation*

4.1.2 CAM vs. UNetCAM

In the second set of configurations, another well-known segmentation model with different characteristics and internals than DeepLab is hired called UNet. There are two different versions of the UNet composed in this experiment set, and they are named UNetCAM Learnable and UNetCAM Non-Learnable, as described in Section 3.2.2. In the architecture of UNetCAM Learnable, the upsampling part of the decoder is constituted of transposed convolution layers, and for the other one, upsampling is done with bilinear interpolation. ResNet50 is selected to be the encoder part of the UNetCAM.

Table 3: CAM vs UNetCAM

| | CAM | UNetCAM Learnable | UNetCAM Non-Learnable |
|-------|-------|-------------------|-----------------------|
| train | 48.93 | 54.30 | 53.75 |
| val | 47.59 | 53.01 | 52.41 |

All the experiments are done by using trainaug split

Table 3 presents the results of both UNetCAM networks and CAM, and it shows that UNetCAM enhances the results by at least 3%, which supports the experimentation results on DeepLabCAM.

Table 4: CAM vs UNetCAM

| | CAM | UNetCAM Learnable | UNetCAM Non-Learnable |
|-------|-------|-------------------|-----------------------|
| train | 48.14 | 52.75 | 51.90 |
| val | 46.60 | 51.32 | 50.56 |

All the experiments are done by using train split

To go a little further, we change the data source with the current settings and use the train split instead of trainaug. Table 4 proves that the UNetCAM is performing better even with fewer training images.

One another exciting examination is arranged to interpret the effects of the usage of dCRF. For this purpose, in Table 5, the metrics are calculated after dCRF is applied to the CAMs. We are investigating if our proposed method embeds the dCRF to the network and the further application of the dCRF becomes somehow useless. The embedding idea matches with the results of Table 5, but it also shows that dCRF still has room to enhance the UNetCAM results.

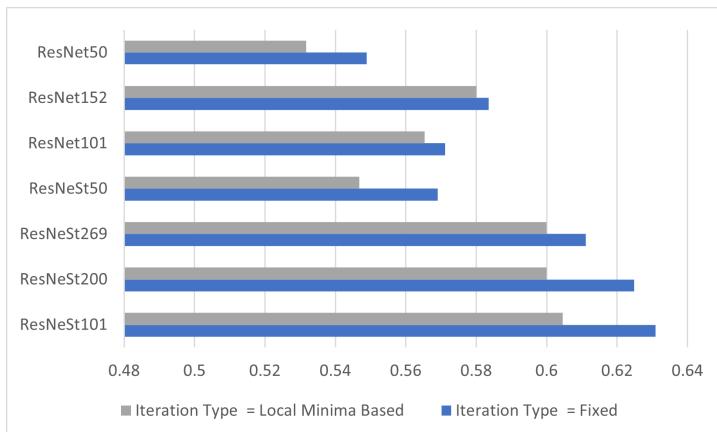
15 **Table 5:** CAM + dCRF vs UNetCAM + dCRF

| | CAM | UNetCAM Learnable | UNetCAM Non-Learnable |
|-------|-------|-------------------|-----------------------|
| train | 53.18 | 59.90 | 58.74 |
| val | 50.72 | 57.11 | 57.01 |

19 All the experiments are done by using trainaug split
20

21 4.1.3 Investigation on the Iteration Procedure

22 The conception of an iterative self-improved training procedure has its
23 challenges, such as stabilization and optimization. In our proposal, the classifi-
24 cation performance directly impacts the segmentation network as it uses the
25 classification results. Furthermore, the segmentation network implicitly affects
26 the features of the classifier network. If things go wrong for one task, they will
27 go the same in the other, and the worsening will continue growing. The prob-
28 lem with optimization is that the overall network can be stuck in early steps
29 without improving its performance.
30

47 **Fig. 4:** mIoU According to Iteration Procedure
48
49
50

51 Intuitively, these two dilemmas may happen in two different ways, (i) when
52 the classifier and the segmentation part could not find the room to optimize
53 itself, (ii) when one of the tasks dominates the other in the training phase. To
54 examine these potential problems, we set experimentations by changing the
55 frequency of pseudo-ground-truth creation. The experiments focus on two main
56 configurations; the first is waiting until the losses converge, and the second
57 is fixing the generation frequency manually at the beginning of the training
58 procedure.
59

60 Figure 4 represents the experimental results on the frequency of the pseudo-
61 ground-truth generation, and from these results, it is clear that fixing the
62 frequency performs better. Intuitively, this makes sense as this is an instance
63
64

20 *Iterative Self-Improved Model for Weakly Supervised Segmentation*

11
12 of the expectation-maximization problem, and giving room to each task and
13 training both in relaxed conditions performs better than strict rules.
14
15

16 **4.1.4 Investigation on Image Resolutions**
17
18

19 There are several differences between the applications of the segmentation
20 and classification problems. One of them is that their input resolution pref-
21 erences are different. Most of the classification architecture in their original
22 forms accept the images which resolutions vary between 224 to 299. However,
23 it increases to 512 or higher for the popular segmentation architectures with
24 viable reasons. As the proposed network supports segmentation and classifi-
25 cation simultaneously, studying the different input sizes makes sense to reveal
26 the framework's stability and robustness.
27

28 In this part of the study, the experimentation reveals the effects of the
29 input resolutions on the proposed method, and four input sizes are selected to
30 run tests as 224x224, 320x320, 448x448, and 512x512.
31
32

Table 6: Experiments with Different Input Sizes

| | 256×256 | 320×320 | 448×448 | 512×512 |
|------------|---------|---------|---------|---------|
| CAM | 57.52 | 57.78 | 55.34 | 54.94 |
| DeepLabCAM | 62.86 | 63.09 | 57.94 | 56.26 |

33
34 This table contains the mIoU metric between predicted segmentation labels and original
35 segmentation labels
36
37

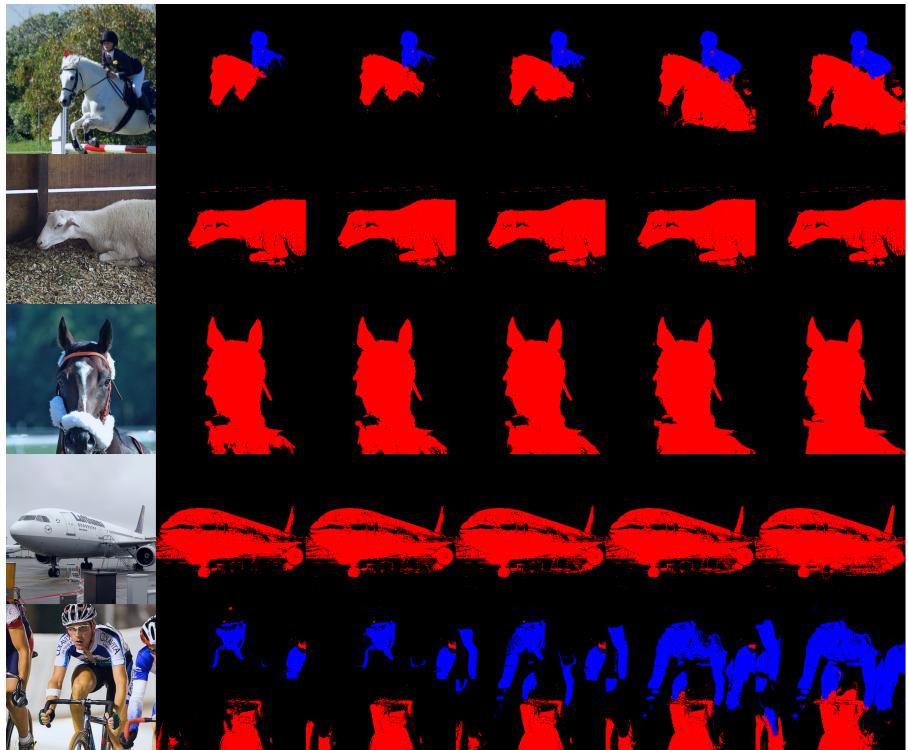
38 According to the reported Table 6, the proposed architecture performs
39 better than the original CAM for each input size. In addition to that, when
40 the input size is 320x320, it reaches its best.
41
42

43 These four experiments prove that the proposed network performs better
44 than CAMs when employed as segmentation labels for the dataset. It engages
45 the presented network results as affinity labels and executes the pipeline to
46 achieve the final segmentation model.
47
48

49 **4.1.5 Qualitative Results**
50
51

52 Qualitative investigations can be done on the visualizations in Figure 5 to
53 understand the effects of the segmentation loss and iterations more clearly. The
54 first column of the figure shows the original image, and the other columns from
55 left to right contain the pseudo-segmentation labels generated after epochs 5,
56 15, 25, 35, and 45, respectively. For the images in the first and the last rows,
57 it is evident that the latter pseudo masks are more accurate than the previous
58 ones. However, in the fourth row, it first begins to perform better, and in the
59 last column, some pixels that are the part of the background start to be marked
60 as a plane, which is evidence of becoming poorer. Even though it is hard to
61 say the effects of the framework are in the positive or negative direction, it is
62
63
64
65

obvious from the images, one thing is clear, the proposed model has effects on the qualitative results.



43 **Fig. 5:** Visualization of Pseudo Segmentation Labels
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

4.1.6 AffinityNet Training

AffinityNet is the network proposed in PSA to learn class-agnostic pixel-wise affinities. The model is trained on the images with a target called the affinity labels, generated from the denseCRF algorithm on CAMs of the DeepLab-CAM network. There are two threshold values to apply on CAMs to generate affinity labels more confidently, one for the background and the other for the foreground, with the values 0.1 and 0.3, respectively. The foreground threshold has the same value as the threshold used in DeepLabCAM to optimize according to it. We followed the same steps to train and evaluate the model as performed in the original study. In addition to the overall procedure, we only adjust the backbone instead of ResNet-38 in our experiments, and we use deeper versions and versions of ResNeSt architecture. As the optimizer, the Poly optimizer is chosen. Moreover, AffinityNet trained for three epochs with batch size 16 on 512×512 resolution images.

22 *Iterative Self-Improved Model for Weakly Supervised Segmentation*

In the end-to-end pipeline, the backbones are held in common generally. It means that we operate the same backbone in DeepLabCAM, AffinityNet, and Segmentation model unless otherwise declared.

19 **Table 7:** mIoU on Pseudo Segmentation Labels

| | ResNet | | | ResNeSt | | | |
|--------------------------------|--------|-------|-------|---------|--------------|-------|--------------|
| | 50 | 101 | 152 | 50 | 101 | 200 | 269 |
| CAM ¹ | 48.93 | 50.33 | 50.83 | 53.49 | 57.78 | 57.33 | 54.75 |
| DeepLabCAM ¹ | 54.89 | 57.11 | 58.35 | 56.91 | 63.09 | 62.48 | 61.11 |
| DeepLabCAM+RW ¹ | 74.25 | 76.40 | 78.54 | 75.65 | 76.78 | 77.47 | 79.10 |
| DeepLabCAM+RW+CRF ¹ | 74.83 | 77.09 | 79.41 | 75.97 | 77.29 | 77.93 | 79.64 |

This table contains the mIoU metric between CAMs and original segmentation labels for the trainaug split of Pascal VOC2012

Table 7 presents the mIoU values of several backbones with different settings. The first two rows contain the results of CAM and DeepLabCAM, respectively. The third row shows the improvements done on DeepLabCAM results by applying Random Walk with AffinityNet predictions, and in addition to these, the fourth row also has one more step as employing denseCRF. The results in table 7 prove that AffinityNet performs better when it has more accurate affinity labels and also deeper networks achieve better performances.

The comparison between our proposed method with AffinityNet and previous studies with AffinityNet is handled in Section 4.1.8.

42 **4.1.7 Segmentation model training**

According to the previous studies, three different configurations are organized on the version of the DeepLab family to make fair comparisons. In the first set of experiments, DeepLabv2 architecture with the ResNet-101 backbone is exercised. The second set of experiments is done by upgrading the DeepLab version to 3 and keeping the same backbone. The last experiments are performed with again DeepLabv3 but with various backbones. For all of them, the input image resolution is 512×512, and the training is stopped after 50 epochs.

53 **Table 8:** mIoU on Segmentation Labels with DeepLabv3

| | ResNet | | | ResNeSt | | | |
|------|--------|-------|-------|---------|-------|--------------|-------|
| | 50 | 101 | 152 | 50 | 101 | 200 | 269 |
| Val | 67.87 | 70.51 | 72.28 | 69.77 | 73.52 | 74.90 | 74.76 |
| Test | 67.89 | 71.45 | 72.41 | 69.88 | 73.25 | 74.98 | 74.90 |

This table contains the mIoU metric between predicted segmentation labels and original segmentation labels

14 The results of the segmentation model with different backbones are
 15 reported in Table 8. In this Table, the backbone of all three models, DeepLab-
 16 CAM, AffinityNet, and DeepLabv3, is kept precisely the same in the end-to-end
 17 pipeline.

20 **Table 9:** mIoU on Segmentation Labels with DeepLabv3 (ResNet101)

| | ResNet | | | ResNeSt | | | |
|------|--------|-------|-------|---------|-------|--------------|-------|
| | 50 | 101 | 152 | 50 | 101 | 200 | 269 |
| Val | 69.09 | 70.51 | 71.54 | 70.39 | 73.55 | 73.70 | 73.37 |
| Test | 69.46 | 71.45 | 71.62 | 70.43 | 73.39 | 74.50 | 73.50 |

27 This table contains the mIoU metric between predicted segmentation labels and original
 28 segmentation labels

30 In Table 9, the backbone of the last segmentation model is kept the same
 31 as ResNet-101 as it is widely used in the previous experiments.

35 **Table 10:** mIoU on Segmentation Labels with DeepLabv2 (ResNet101)

| | ResNet-50 | ResNet-101 | ResNeSt-200 |
|------|-----------|------------|--------------|
| Val | 68.64 | 70.60 | 72.60 |
| Test | 68.83 | 70.38 | 72.94 |

40 This table contains the mIoU metric between predicted segmentation labels and original
 41 segmentation labels

43 In Table 10, to compare our results with the previous studies, we train the
 44 DeepLabv2 segmentation model with the pseudo-segmentation labels obtained
 45 from ResNet-50, ResNet-101 backbones.

48 4.1.8 Comparison with SOTA

49 In Table 11, we compare our best results with the previous studies, which are
 50 state-of-the-art at their time.

52 Table 11 can be interpreted like this, in overall performance, the best model
 53 is our proposed model with ResNeSt-269 backbone without using any addi-
 54 tional data or saliency map with DeepLabv3 with a margin over %2.5. When
 55 the backbone is kept the same between the studies as ResNet-101, the best per-
 56 formance is observed in the EPS study, which uses extra saliency map data in
 57 the training phase. When there is a restriction to not use any additional data
 58 in the training phase, our model becomes the second-best one with %70.38
 59 mIoU with a small margin on SPML which also does not use any additional
 60 information. In addition to these, EPS and DRS studies use saliency map data
 61 as the input in the other results in the table.

Table 11: Comparison of the results with literature

| Method | Backbone | CAM ¹ | Pseudo Masks ¹ | Seg Masks ² | Seg Masks ³ |
|---------------------|-------------|------------------|---------------------------|------------------------|------------------------|
| PSA | ResNet-38 | 48 | 59,7 | 61,7 | 63,7 |
| IRNet | ResNet-50 | 48,3 | 65,9 | 63,5 | 64,8 |
| FickleNet | ResNet-101 | - | - | 64.9 | 65.3 |
| ICD | ResNet-101 | - | - | 64.1 | 64.3 |
| SEAM | ResNet-38 | 55.4 | 63.4 | 64.5 | 65.7 |
| CDA | ResNet-38 | 58.4 | 66.4 | 66.1 | 66.8 |
| PuzzleCAM | ResNeSt-101 | 61.85 | 72.46 | 66.9 | 67.7 |
| WSGCN | ResNet-101 | - | - | 68.7 | 69.3 |
| Ours ⁴ | ResNet-101 | 57,11 | 77.09 | 70.61 | 70.38 |
| DRS v1 ⁴ | ResNet-101 | - | - | 70.4 | 70.7 |
| EPS ⁴ | ResNet-101 | 69.4 | 71.6 | 70.9 | 70.8 |
| DRS v2 ⁴ | ResNet-101 | - | - | 71.2 | 71.4 |
| Ours ⁵ | ResNet-101 | 57,11 | 77.09 | 70.51 | 71.45 |
| SPML ⁴ | ResNet-101 | - | - | 69.5 | 71.6 |
| EPS ⁷ | ResNet-101 | 69.4 | 71.6 | 71.0 | 71.8 |
| PuzzleCAM | ResNeSt-269 | 62.45 | 74.67 | 71.9 | 72.2 |
| Ours ⁵ | ResNeSt-200 | 62.48 | 77.93 | 72.60 | 72.94 |
| Ours ⁶ | ResNeSt-200 | 62.48 | 77.93 | 74.90 | 74.98 |

This table contains the mIoU metric between predicted segmentation labels and original segmentation labels

¹The mIoU metric calculated on train split

²The mIoU metric calculated on val split

³The mIoU metric calculated on test split

⁴DeepLabv2 with backbone ResNet-101 used as Segmentation Model

⁵DeepLabv3 with backbone ResNet-101 used as Segmentation Model

⁶DeepLabv3 with backbone ResNeSt-200 used as Segmentation Model

⁷DeepLabv1 with backbone ResNet-101 used as Segmentation Model

5 Conclusion

In this study, we propose a novel framework for WSSS. It employs an iterative self-improved approach in an encoder-decoder-based segmentation model to obtain coherent CAMs with segmentation labels. The proposed framework contains various elements from the available approaches in the literature. Therefore, it tries to combine the strengths of those methods.

The extensive experiments show that the proposed framework is challenged and proven to improve the CAMs with the quantitative and qualitative results. The iterative approach is optimized with the help of a simple modification in the well-known loss function, Cross-Entropy, to achieve the goal. In addition to these, a specific implementation of the proposed framework achieves the state-of-the-art on Pascal VOC 2012 dataset.

As a possible future work, in theory, there is room to improve the loss function to stabilize more the training phase of the proposed framework. And also, combining AffinityNet with the proposed network will have a significant effect on the results as AffinityNet learns from pixel-level information from the

12
13 images, and adding more pixel-level loss to the proposed network may improve
14 the results.
15

16 **Supplementary information.**

17 **References**

- 18
21 [1] Otsu, N.: A threshold selection method from gray-level histograms. IEEE
22 Transactions on Systems, Man, and Cybernetics **9**(1), 62–66 (1979). <https://doi.org/10.1109/TSMC.1979.4310076>
- 24 [2] Dhanachandra, N., Manglem, K., Chanu, Y.J.: Image segmentation using
25 k -means clustering algorithm and subtractive clustering algorithm. Pro-
26 cedia Computer Science **54**, 764–771 (2015). <https://doi.org/10.1016/j.procs.2015.06.090>. Eleventh International Conference on Communication
27 Networks, ICCN 2015, August 21-23, 2015, Bangalore, India Eleventh
28 International Conference on Data Mining and Warehousing, ICDMW
29 2015, August 21-23, 2015, Bangalore, India Eleventh International Con-
30 ference on Image and Signal Processing, ICISP 2015, August 21-23, 2015,
31 Bangalore, India
33
34
- 35 [3] Plath, N., Toussaint, M., Nakajima, S.: Multi-class image segmentation
36 using conditional random fields and global classification. In: Proceedings
37 of the 26th Annual International Conference on Machine Learning. ICML
38 '09, pp. 817–824. Association for Computing Machinery, New York, NY,
39 USA (2009). <https://doi.org/10.1145/1553374.1553479>. <https://doi.org/10.1145/1553374.1553479>
- 42 [4] Neubert, P., Protzel, P.: Compact watershed and preemptive slic: On
43 improving trade-offs of superpixel segmentation algorithms. 2014 22nd
44 International Conference on Pattern Recognition, 996–1001 (2014)
- 46 [5] Kass, M., Witkin, A., Terzopoulos, D.: Snakes: Active contour models.
47 International Journal of Computer Vision **1**, 321–331 (2004)
- 49 [6] Vedaldi, A., Soatto, S.: Quick shift and kernel methods for mode seeking.
50 In: Forsyth, D., Torr, P., Zisserman, A. (eds.) Computer Vision – ECCV
51 2008, pp. 705–718. Springer, Berlin, Heidelberg (2008)
- 53 [7] Zheng, S., Lu, J., Zhao, H., Zhu, X., Luo, Z., Wang, Y., Fu, Y., Feng,
54 J., Xiang, T., Torr, P.H., et al.: Rethinking semantic segmentation from
55 a sequence-to-sequence perspective with transformers. arXiv preprint
56 arXiv:2012.15840 (2020)
- 58 [8] Wu, T., Tang, S., Zhang, R., Cao, J., Zhang, Y.: Cgnet: A light-weight
59 context guided network for semantic segmentation. IEEE Transactions on
60 Image Processing **30**, 1169–1179 (2020)
- 62
63
64
65

Table 12: Class Based Results for the Best Model

| | aeroplane | bicycle | bird | boat | bottle | bus | car | cat | chair | cow |
|--------------|-----------|---------|-------|-------|--------|-------|-------|-------|-------|-------|
| background | 93.07 | 89.51 | 37.46 | 84.24 | 63.12 | 72.56 | 91.70 | 87.14 | 92.16 | 37.49 |
| dining table | 66.95 | 88.32 | 85.73 | 84.89 | 83.48 | 60.10 | 85.77 | 57.01 | 74.88 | 53.77 |

9
10
11 *Iterative Self-Improved Model for Weakly Supervised Segmentation* 27
12
13

- 14 [9] Kirillov, A., Wu, Y., He, K., Girshick, R.: Pointrend: Image segmentation
15 as rendering. In: Proceedings of the IEEE/CVF Conference on Computer
16 Vision and Pattern Recognition, pp. 9799–9808 (2020)
- 17 [10] Yin, M., Yao, Z., Cao, Y., Li, X., Zhang, Z., Lin, S., Hu, H.: Disentangled
18 Non-Local Neural Networks (2020)
- 19 [11] He, J., Deng, Z., Qiao, Y.: Dynamic multi-scale filters for semantic seg-
20 mentation. In: Proceedings of the IEEE/CVF International Conference
21 on Computer Vision (ICCV) (2019)
- 22 [12] He, J., Deng, Z., Zhou, L., Wang, Y., Qiao, Y.: Adaptive pyramid context
23 network for semantic segmentation. In: Proceedings of the IEEE/CVF
24 Conference on Computer Vision and Pattern Recognition (CVPR) (2019)
- 25 [13] Chen, L.-C., Zhu, Y., Papandreou, G., Schroff, F., Adam, H.: Encoder-
26 decoder with atrous separable convolution for semantic image segmenta-
27 tion. In: ECCV (2018)
- 28 [14] Long, J., Shelhamer, E., Darrell, T.: Fully convolutional networks for
29 semantic segmentation. CoRR **abs/1411.4038** (2014) <https://arxiv.org/abs/1411.4038>
- 30 [15] Chen, L.-C., Papandreou, G., Kokkinos, I., Murphy, K., Yuille, A.L.:
31 Semantic Image Segmentation with Deep Convolutional Nets and Fully
32 Connected CRFs (2016)
- 33 [16] Noh, H., Hong, S., Han, B.: Learning deconvolution network for seman-
34 tic segmentation. CoRR **abs/1505.04366** (2015) <https://arxiv.org/abs/1505.04366>
- 35 [17] Ronneberger, O., Fischer, P., Brox, T.: U-net: Convolutional networks for
36 biomedical image segmentation. CoRR **abs/1505.04597** (2015) <https://arxiv.org/abs/1505.04597>
- 37 [18] Chen, L., Papandreou, G., Kokkinos, I., Murphy, K., Yuille, A.L.:
38 Deeplab: Semantic image segmentation with deep convolutional nets,
39 atrous convolution, and fully connected crfs. CoRR **abs/1606.00915**
40 (2016) <https://arxiv.org/abs/1606.00915>
- 41 [19] Chen, L., Papandreou, G., Schroff, F., Adam, H.: Rethinking atrous
42 convolution for semantic image segmentation. CoRR **abs/1706.05587**
43 (2017) <https://arxiv.org/abs/1706.05587>
- 44 [20] He, K., Gkioxari, G., Dollár, P., Girshick, R.B.: Mask R-CNN. CoRR
45 **abs/1703.06870** (2017) <https://arxiv.org/abs/1703.06870>
- 46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

- 9
10
11 28 *Iterative Self-Improved Model for Weakly Supervised Segmentation*
12
13
14 [21] Lin, T., Dollár, P., Girshick, R.B., He, K., Hariharan, B., Belongie, S.J.:
15 Feature pyramid networks for object detection. CoRR **abs/1612.03144**
16 (2016) <https://arxiv.org/abs/1612.03144>
- 17
18 [22] Zhao, H., Shi, J., Qi, X., Wang, X., Jia, J.: Pyramid scene parsing network.
19 CoRR **abs/1612.01105** (2016) <https://arxiv.org/abs/1612.01105>
- 20
21 [23] Ghiasi, G., Fowlkes, C.C.: Laplacian reconstruction and refinement for
22 semantic segmentation. CoRR **abs/1605.02264** (2016) [https://arxiv.](https://arxiv.org/abs/1605.02264)
23 [org/abs/1605.02264](https://arxiv.org/abs/1605.02264)
- 24
25 [24] Visin, F., Kastner, K., Courville, A.C., Bengio, Y., Matteucci, M., Cho,
26 K.: Reseg: A recurrent neural network for object segmentation. CoRR
27 **abs/1511.07053** (2015) <https://arxiv.org/abs/1511.07053>
- 28
29 [25] Chen, L., Yang, Y., Wang, J., Xu, W., Yuille, A.L.: Attention to
30 scale: Scale-aware semantic image segmentation. CoRR **abs/1511.03339**
31 (2015) <https://arxiv.org/abs/1511.03339>
- 32
33 [26] Chan, L., Hosseini, M.S., Plataniotis, K.N.: A comprehensive analysis
34 of weakly-supervised semantic segmentation in different image domains.
35 CoRR **abs/1912.11186** (2019) <https://arxiv.org/abs/1912.11186>
- 36
37 [27] Aresta, G., Araújo, T., Kwok, S., Chennamsetty, S.S., P., M.S.K., Vargh-
38 ese, A., Marami, B., Prastawa, M., Chan, M., Donovan, M.J., Fernandez,
39 G., Zeineh, J., Kohl, M., Walz, C., Ludwig, F., Braunewell, S., Baust, M.,
40 Vu, Q.D., To, M.N.N., Kim, E., Kwak, J.T., Galal, S., Sanchez-Freire, V.,
41 Brancati, N., Frucci, M., Riccio, D., Wang, Y., Sun, L., Ma, K., Fang,
42 J., Koné, I., Boulmane, L., Campilho, A., Eloy, C., Polónia, A., Aguiar,
43 P.: BACH: grand challenge on breast cancer histology images. CoRR
44 **abs/1808.04277** (2018) <https://arxiv.org/abs/1808.04277>
- 45
46 [28] Chan, L., Hosseini, M., Rowsell, C., Plataniotis, K., Damaskinos, S.:
47 Histosegnet: Semantic segmentation of histological tissue type in whole
48 slide images. In: 2019 IEEE/CVF International Conference on Computer
49 Vision (ICCV), pp. 10661–10670 (2019). <https://doi.org/10.1109/ICCV.2019.01076>
- 50
51 [29] Hosseini, M.S., Chan, L., Tse, G., Tang, M., Deng, J., Norouzi, S., Rowsell,
52 C., Plataniotis, K.N., Damaskinos, S.: Atlas of digital pathology: A gen-
53 eralized hierarchical histological tissue type-annotated database for deep
54 learning. In: 2019 IEEE/CVF Conference on Computer Vision and Pat-
55 tern Recognition (CVPR), pp. 11739–11748 (2019). <https://doi.org/10.1109/CVPR.2019.01202>
- 56
57 [30] Lin, H., Chen, H., Dou, Q., Wang, L., Qin, J., Heng, P.: Scannet: A
58 fast and dense scanning framework for metastatic breast cancer detection
- 59
60
61
62
63
64
65

10
11 *Iterative Self-Improved Model for Weakly Supervised Segmentation* 2912
13 from whole-slide images. CoRR **abs/1707.09597** (2017) <https://arxiv.org/abs/1707.09597>14
15
16 [31] Nivaggioli, A., Randrianarivo, H.: Weakly supervised semantic segmen-
17 tation of satellite images. CoRR **abs/1904.03983** (2019) <https://arxiv.org/abs/1904.03983>21
22 [32] Robinson, C., Hou, L., Malkin, K., Soobitsky, R., Czawltyko, J.,
23 Dilkina, B., Jovic, N.: Large scale high-resolution land cover map-
24 ping with multi-resolution data. In: Proceedings of the IEEE
25 Conference on Computer Vision and Pattern Recognition (CVPR)
26 (2019). http://openaccess.thecvf.com/content_CVPR_2019/html/Robinson_Large_Scale_High-Resolution_Land_Cover_Mapping_With_Multi-Resolution_Data_CVPR_2019_paper.html29
30 [33] Seferbekov, S.S., Iglovikov, V.I., Buslaev, A.V., Shvets, A.A.: Fea-
31 ture pyramid network for multi-class land segmentation. CoRR
32 **abs/1806.03510** (2018) <https://arxiv.org/abs/1806.03510>33
34 [34] Wang, X., Chen, H., Gan, C., Lin, H., Dou, Q., Huang, Q., Cai, M.,
35 Heng, P.: Weakly supervised learning for whole slide lung cancer image
36 classification. (2018)37
38 [35] Pathak, D., Krähenbühl, P., Darrell, T.: Constrained convolutional neural
39 networks for weakly supervised segmentation. CoRR **abs/1506.03648**
40 (2015) <https://arxiv.org/abs/1506.03648>41
42 [36] Papandreou, G., Chen, L., Murphy, K., Yuille, A.L.: Weakly- and semi-
43 supervised learning of a DCNN for semantic image segmentation. CoRR
44 **abs/1502.02734** (2015) <https://arxiv.org/abs/1502.02734>45
46 [37] Zhou, B., Khosla, A., Lapedriza, À., Oliva, A., Torralba, A.: Learning deep
47 features for discriminative localization. CoRR **abs/1512.04150** (2015)
48 <https://arxiv.org/abs/1512.04150>49
50 [38] Selvaraju, R.R., Das, A., Vedantam, R., Cogswell, M., Parikh, D., Batra,
51 D.: Grad-cam: Why did you say that? visual explanations from deep
52 networks via gradient-based localization. CoRR **abs/1610.02391** (2016)
53 <https://arxiv.org/abs/1610.02391>54
55 [39] Durand, T., Mordan, T., Thome, N., Cord, M.: Wildcat: Weakly super-
56 vised learning of deep convnets for image classification, pointwise local-
57 ization and segmentation. In: 2017 IEEE Conference on Computer Vision
58 and Pattern Recognition (CVPR), pp. 5957–5966 (2017). <https://doi.org/10.1109/CVPR.2017.631>60
61 [40] Pathak, D., Shelhamer, E., Long, J., Darrell, T.: Fully Convolutional
62

9
10
11 30 *Iterative Self-Improved Model for Weakly Supervised Segmentation*
12
13
14 Multi-Class Multiple Instance Learning (2015)
15
16 [41] Simonyan, K., Zisserman, A.: Very Deep Convolutional Networks for
17 Large-Scale Image Recognition (2015)
18
19 [42] Shimoda, W., Yanai, K.: Distinct class-specific saliency maps for weakly
20 supervised semantic segmentation. In: Leibe, B., Matas, J., Sebe, N.,
21 Welling, M. (eds.) Computer Vision – ECCV 2016, pp. 218–234. Springer,
22 Cham (2016)
23
24 [43] Saleh, F., Akbarian, M.S.A., Salzmann, M., Petersson, L., Gould,
25 S., Alvarez, J.M.: Built-in foreground/background prior for weakly-
26 supervised semantic segmentation. CoRR **abs/1609.00446** (2016) <https://arxiv.org/abs/1609.00446>
27
28
29 [44] Kolesnikov, A., Lampert, C.H.: Seed, expand and constrain: Three prin-
30 ciples for weakly-supervised image segmentation. CoRR **abs/1603.06098**
31 (2016) <https://arxiv.org/abs/1603.06098>
32
33 [45] Krähenbühl, P., Koltun, V.: Efficient inference in fully connected crfs with
34 gaussian edge potentials. CoRR **abs/1210.5644** (2012) <https://arxiv.org/abs/1210.5644>
35
36
37 [46] Wei, Y., Xiao, H., Shi, H., Jie, Z., Feng, J., Huang, T.S.: Revisiting dilated
38 convolution: A simple approach for weakly- and semi- supervised seman-
39 tic segmentation. CoRR **abs/1805.04574** (2018) <https://arxiv.org/abs/1805.04574>
40
41
42 [47] Wei, Y., Feng, J., Liang, X., Cheng, M., Zhao, Y., Yan, S.: Object region
43 mining with adversarial erasing: A simple classification to semantic seg-
44 mentation approach. CoRR **abs/1703.08448** (2017) <https://arxiv.org/abs/1703.08448>
45
46
47 [48] Lee, J., Kim, E., Lee, S., Lee, J., Yoon, S.: Ficklenet: Weakly and
48 semi-supervised semantic image segmentation\\using stochastic inference.
49 CoRR **abs/1902.10421** (2019) <https://arxiv.org/abs/1902.10421>
50
51
52 [49] Huang, Z., Wang, X., Wang, J., Liu, W., Wang, J.: Weakly-supervised
53 semantic segmentation network with deep seeded region growing. In:
54 Proceedings of the IEEE Conference on Computer Vision and Pattern
55 Recognition, pp. 7014–7023 (2018)
56
57 [50] Ahn, J., Kwak, S.: Learning pixel-level semantic affinity with image-
58 level supervision for weakly supervised semantic segmentation. CoRR
59 **abs/1803.10464** (2018) <https://arxiv.org/abs/1803.10464>
60
61 [51] Ahn, J., Cho, S., Kwak, S.: Weakly supervised learning of instance
62
63
64
65

9
10
11 *Iterative Self-Improved Model for Weakly Supervised Segmentation* 31
12
13
14
15
16

segmentation with inter-pixel relations. CoRR **abs/1904.05044** (2019)
<https://arxiv.org/abs/1904.05044>

17 [52] Jo, S., Yu, I.-J.: Puzzle-CAM: Improved localization via matching partial
18 and full features (2021)

20 [53] Lee, S., Lee, M., Lee, J., Shim, H.: Railroad is not a train: Saliency as
21 pseudo-pixel supervision for weakly supervised semantic segmentation.
22 CoRR **abs/2105.08965** (2021) <https://arxiv.org/abs/2105.08965>

24 [54] Yao, Q., Gong, X.: Saliency guided self-attention network for weakly-
25 supervised semantic segmentation. CoRR **abs/1910.05475** (2019) <https://arxiv.org/abs/1910.05475>

28 [55] Pan, S., Lu, C., Lee, S., Peng, W.: Weakly-supervised image semantic seg-
29 mentation using graph convolutional networks. CoRR **abs/2103.16762**
30 (2021) <https://arxiv.org/abs/2103.16762>

32 [56] Kim, B., Han, S., Kim, J.: Discriminative region suppression for weakly-
33 supervised semantic segmentation. CoRR **abs/2103.07246** (2021) <https://arxiv.org/abs/2103.07246>

36 [57] Kwak, S., Hong, S., Han, B.: Weakly supervised semantic segmentation
37 using superpixel pooling network. In: AAAI (2017)

39 [58] Zhou, Y., Zhu, Y., Ye, Q., Qiu, Q., Jiao, J.: Weakly supervised instance
40 segmentation using class peak response. CoRR **abs/1804.00880** (2018)
41 <https://arxiv.org/abs/1804.00880>

43 [59] Pont-Tuset, J., Arbelaez, P., Barron, J.T., Marqués, F., Malik, J.: Multi-
44 scale combinatorial grouping for image segmentation and object proposal
45 generation. CoRR **abs/1503.00848** (2015) <https://arxiv.org/abs/1503.00848>

48 [60] Ke, T., Hwang, J., Yu, S.X.: Universal weakly supervised segmentation
49 by pixel-to-segment contrastive learning. CoRR **abs/2105.00957** (2021)
50 <https://arxiv.org/abs/2105.00957>

52 [61] Su, Y., Sun, R., Lin, G., Wu, Q.: Context decoupling augmentation
53 for weakly supervised semantic segmentation. CoRR **abs/2103.01795**
54 (2021) <https://arxiv.org/abs/2103.01795>

56 [62] Wang, Y., Zhang, J., Kan, M., Shan, S., Chen, X.: Self-supervised equiv-
57 ariant attention mechanism for weakly supervised semantic segmentation.
58 CoRR **abs/2004.04581** (2020) <https://arxiv.org/abs/2004.04581>

60 [63] Everingham, M., Van Gool, L., Williams, C.K.I., Winn, J., Zisserman, A.:

10
11
12 32 *Iterative Self-Improved Model for Weakly Supervised Segmentation*
13
14 The pascal visual object classes (voc) challenge. International Journal of
15 Computer Vision **88**(2), 303–338 (2010)
16
17 [64] Hariharan, B., Arbelaez, P., Bourdev, L., Maji, S., Malik, J.: Semantic
18 contours from inverse detectors. In: International Conference on
19 Computer Vision (ICCV) (2011)
20
21 [65] Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G.,
22 Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., Desmaison, A., Kopf,
23 A., Yang, E., DeVito, Z., Raison, M., Tejani, A., Chilamkurthy, S.,
24 Steiner, B., Fang, L., Bai, J., Chintala, S.: Pytorch: An imperative style,
25 high-performance deep learning library. In: Wallach, H., Larochelle,
26 H., Beygelzimer, A., d' Alché-Buc, F., Fox, E., Garnett, R. (eds.)
27 Advances in Neural Information Processing Systems 32, pp. 8024–8035.
28 Curran Associates, Inc., ??? (2019). <http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf>
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65