

Context Decoupling Augmentation for Weakly Supervised Semantic Segmentation

Yukun Su^{1,2}, Ruizhou Sun^{1,2}, Guosheng Lin^{3†}, and Qingyao Wu^{1,2†}

¹School of Software and Engineering, South China University of Technology

²Key Laboratory of Big Data and Intelligent Robot, Ministry of Education

³School of Computer Science and Engineering, Nanyang Technological University
suyukun666@gmail.com, ruizhousun@foxmail.com, gslin@ntu.edu.sg, qyw@scut.edu.cn

Abstract

Data augmentation is vital for deep learning neural networks. By providing massive training samples, it helps to improve the generalization ability of the model. Weakly supervised semantic segmentation (WSSS) is a challenging problem that has been deeply studied in recent years, conventional data augmentation approaches for WSSS usually employ geometrical transformations, random cropping and color jittering. However, merely increasing the same contextual semantic data does not bring much gain to the networks to distinguish the objects, e.g., the correct image-level classification of “aeroplane” may be not only due to the recognition of the object itself, but also its co-occurrence context like “sky”, which will cause the model to focus less on the object features. To this end, we present a Context Decoupling Augmentation (CDA) method, to change the inherent context in which the objects appear and thus drive the network to remove the dependence between object instances and contextual information. To validate the effectiveness of the proposed method, extensive experiments on PASCAL VOC 2012 and COCO datasets with several alternative network architectures demonstrate that CDA can boost various popular WSSS methods to the new state-of-the-art by a large margin. Code is available at <https://github.com/suyukun666/CDA>

1. Introduction

Semantic segmentation is a foundation in the computer vision field, which aims to predict the pixel-wise classification of the images and it enjoys a wide range of applications. Recently, benefiting from the deep neural net-

[†]Corresponding authors.

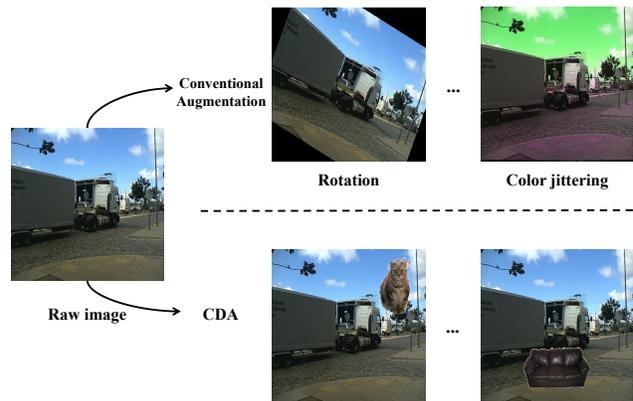


Figure 1. Illustration of the difference between conventional augmentation approaches and our method. Classical data augmentation consists of generating images obtained by basic geometrical transformations or color changes of original training images. Context Decoupling Augmentation (CDA) aims to randomly paste the given object instances into the scenes, so as to decouple the inherent context position of the original objects in the image.

works, modern semantic segmentation models [7, 8, 31, 33] have achieved remarkable progress with massive human-annotated labeled data. However, collecting pixel-level labels is very time-consuming and labor-intensive, which shifts much research attention to weakly supervised semantic segmentation (WSSS). There exist various types of weak supervision for semantic segmentation like using bounding boxes [10, 24], scribbles [30, 40], points [4], and image-level labels [21, 2, 1, 43, 50]. Among them, image-level class labels have been widely used since they demand the least annotation efforts and are already provided in existing large-scale image datasets.

In this paper, we focus on augmentation for WSSS with

image-level labels, which is crucial for deep learning networks. As shown in Figure 1 upper part, given a training image, traditional data augmentation methods utilize some geometrical transformations, such as rotation, scaling, flipping, and even some color conversions to increase the diversity of images to avoid overfitting. However, for weakly supervised semantic segmentation, adjusting the image as a whole and maintain the same contextual semantic relation will not significantly help the networks to mine the object areas. For example, “sofa” always appears in the room in the datasets, therefore, the trained network may not only recognize the objects depending on the instance features but also their co-occurrence context information [29]. Specifically, when object instances often appear at the same time with some accompanying backgrounds, it will cause the networks to yield confounding bias. Namely, the networks can perform classification task well is not due to successfully distinguishing the characteristics of objects, but to being aware of the appearance of certain contextual semantic information, which is harmful to mine the object regions.

Based on this observation, we propose a Context Decoupling Augmentation (CDA) method, designing for disassembling the inherent contextual information of the original image. As shown in Figure 1 bottom half, the “cat” shows in the “sky”, and the “sofa” falls on the “road”. Although some of these scene collocations rarely appear in life, the models can pay more attention to the objects corresponding to the classification labels. Unlike the fully-supervised data augmentation approaches [13], we cannot access the object instance labels to extract the objects under the weakly supervised setting. Therefore, we first adopt off-the-shelf WSSS approaches to obtain the object instances that have been well-segmented. Secondly, we randomly paste the selected foreground instances into the input images to get the new enhanced images and put them into the model for training together with the original ones without augmentation. In this way, we can break the dependency between objects and contextual background, and the models will focus on the internal information of the foreground instances rather than the context information to predict the categories they belong to. Besides, we use an online training technique to conduct data augmentation, which means that the combination of the raw input images and the object instances to be pasted are different each time. This greatly increases the diversity of combinations of various scenes and object instances, and thus enhance the decoupling capability of the networks.

In the proposed context decoupling augmentation framework, we utilize different WSSS networks as our baselines. To verify the effectiveness of our proposed method, extensive experiments show that CDA can improve pseudo-masks more than 2.8% mIoU on average. We achieve new state-of-the-art performance by 66.1% mIoU on the

val set and 66.8% mIoU on the *test* set of PASCAL VOC 2012 [15], and 33.7% mIoU on the *val* set of COCO [32]. The main contributions of our paper can be summarized as follows:

- We present a generally applicable data augmentation approach for weakly supervised semantic segmentation, which, to the best of our knowledge, has not been well explored.
- The proposed context decoupling augmentation (CDA) method does not require additional data and it can remove the correlation between foreground object instances and background context information, which can drive the network focus on object regions rather than the background.
- Experiments on PASCAL VOC 2012 and COCO show the effectiveness of our proposed method and CDA can boost the performance of different WSSS methods to the new state-of-the-art by a large margin.

2. Related Work

2.1. WSSS

Image labels as the weak supervision for segmentation have been widely studied in the past few years. Many approaches [44, 2, 1] use CAM [51] to mine the object seed regions by predicting image labels. To solve the problem that only the discriminative regions can be highlighted, researchers designed to expand the object seed regions in various ways. For example, in [47], the target regions are expanded by fusing different discriminative regions generated by convolutional layers with different expansion rates. [44] drives the network to learn the rest parts of the objects by iteratively erasing the target areas. In addition, some previous works [21, 22] use additional data, such as videos and saliency maps, to explore the objects areas.

Although object expansion technologies emerge endlessly, they all use CAM [51] as the cornerstone. The effect of subsequent diffusion depends on the first step of the CAM learning features. As only image-level labels are provided, when objects are closely coupled with contextual backgrounds, such as “boat” and “water”, “aeroplane” and “sky”, “train” and “track”, CAM will mistakenly recognize the background together with foreground objects. As mentioned in [29], the training networks have no incentive to focus attention only on the foreground class as there may be bias towards other contextual factors as a distractor with high correlation. Thus, this is an issue that’s worth thinking about and that needs to be solved.

2.2. Data Augmentation

Data augmentation is a major trick to train deep neural networks, which aims to increase the diversity of the data

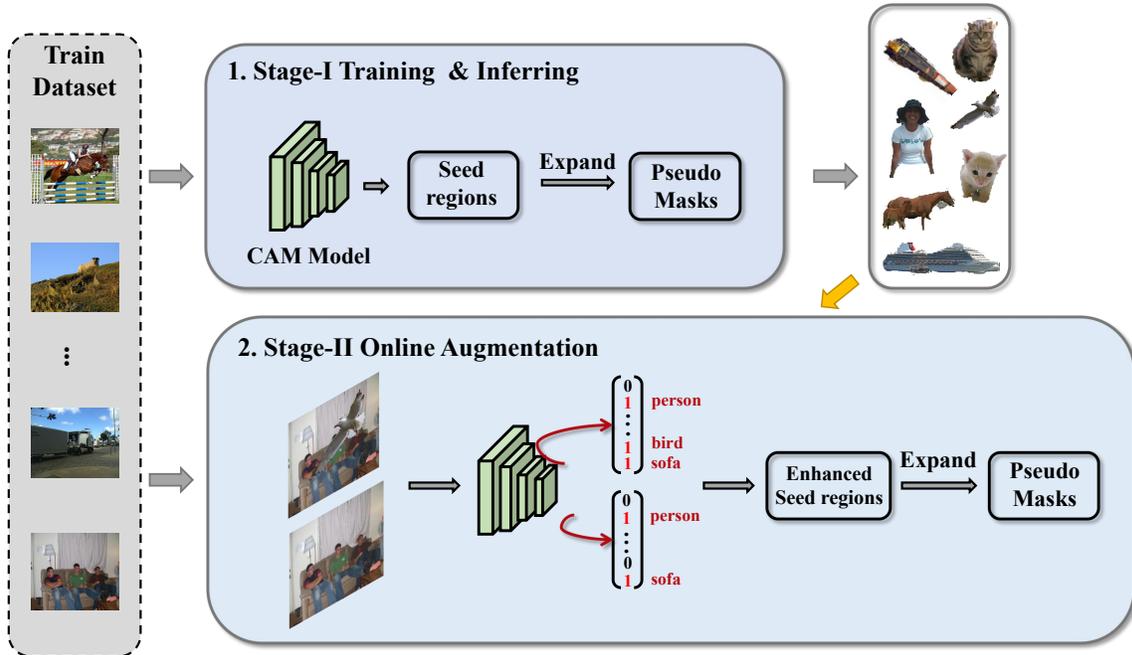


Figure 2. Overview of the proposed augmentation scheme. Stage-I: use the off-the-shelf weakly supervised semantic segmentation methods to obtain some simple object instances with good segmentation. Stage-II: paste the object instances randomly into the raw images to form the new input images, and perform online data augmentation training in a pairwise way with the original input images.

by increasing the training samples and avoid overfitting to a certain extent. Conventional data augmentation approaches perform a series of operations on the basic data, such as rotation, flipping, adding Gaussian noise, etc. Some works have explored synthesizing training data [17, 35] for further generalizability. Generating new training samples by Stylizing ImageNet [18] can lead to better classification performances. Recently, GAN [52] has been employed to transfer the style of the images and to make the content of the images from one domain to another, which can enrich the semantic information of the images to train the deep neural networks. Furthermore, [49] introduced a method to mix two random samples and divide the classification results proportionally to enhance images. [12] conducted augmentation by randomly cutting out some areas in the sample and filled it with 0 pixel value, and keep the result of classification unchanged.

For object detection and segmentation, a popular data augmentation way is “copy-and-paste” [13, 14]. These works pasted real segmented objects into natural images, which is beneficial to increase the object complexity of the internal images and can help to solve the problem of small target detection. However, obtaining these segmented objects requires pixel-wise instance labels. [36] used box-supervision and the off-the-shelf faster-RCNN [37] method to segment and generate masks via cut-and-paste. [3] adopted the unsupervised cut-and-paste learning method to

generate new combined images, but this kind of method is only applicable to the image of single object. It is the first time that we employ copy-and-paste in the WSSS field and it does not require the help of pixel-wise labels and other auxiliary approaches. Thus, for WSSS, such a data augmentation scheme is significant and has not been well explored.

3. Framework

Our approach mainly consists of two stages : (1) we first collect the easy examples of well-segmented objects by using off-the-shelf WSSS methods; (2) then we train the network in a pairwise manner with online augmentation. In this section, we will describe these two stages in details.

3.1. Object Instances Collecting

We aim to apply data augmentation on one of the WSSS models (*i.e.*, IRNet [1]). To some extent, the WSSS method can successfully predict good masks for some easy objects with class labels. Therefore, as shown in Figure 2, in the first stage, we train the original network and we are able to select qualified object instances through the scene complexity of the image, the scope of the object and the semantic relevance by setting some criteria.

Specifically, for the inferring phase after training the network, we follow two main criteria for collecting object instances: (i) the current image should only have a single

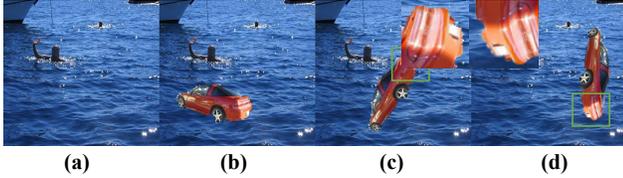


Figure 3. Different kinds of pasting methods used in experiments. (a) Raw input, (b) Random rescale pasting, (c) Random rescale + rotation pasting, (d) Random rescale + rotation + Gaussian smoothing pasting.

class. The intuition behind this is that in the case of only a single class, the image information should be simple and without a complex semantic environment, the segmentation results of the model should be more accurate; (ii) the segmentation result of the current image should meet the condition, $\epsilon_1 < \frac{m}{n} < \epsilon_2$, where ϵ_1 and ϵ_2 are two threshold factors, respectively. m is the number of pixels belonging to the foreground object, n is the number of pixels of the entire image. The reason lies that if the scale value of $\frac{m}{n}$ is too large, it should be that the background is incorrectly identified as the foreground. In contrast, if the scale value is too small, it should be that the model has not been able to recognize enough foreground object pixel information. Different from existing synthesis approaches [13, 14], our method is based on self-provided masks to obtain qualified object instances images.

3.2. Online Augmentation Training

Blending. Before we take a step to train the network in the second stage, we first introduce how to blend the object instances into the natural images. As shown in Figure 3, we show different types of pasting skills in our experiments. It’s worth mentioning that we only paste objects that have not appeared in the original images. The significance of this is that we can increase the diversity of objects of the images, while also reducing the dependence of the same objects in the inherent scene. By randomly rescaling the objects, we can paste them into the images appropriately to prevent them from being too large or too small. The addition of random rotation can change the inherent orientation properties of the objects. Adding Gaussian smoothing can help the added objects boundary blend more naturally.

In some cases, the blending may not be ideal, we elaborate on several possibilities for random pasting. As shown in Figure 4, we have listed several augmented images of random pasting and we call them “perfect”, “good” and “noise” examples. As for the “good” example, the new object “bird” covers part of the “dog” in the original image, however, we argue that this could help to erase the discriminative regions and force the network to discover more object regions like the function in [44]. The noise example shows that the

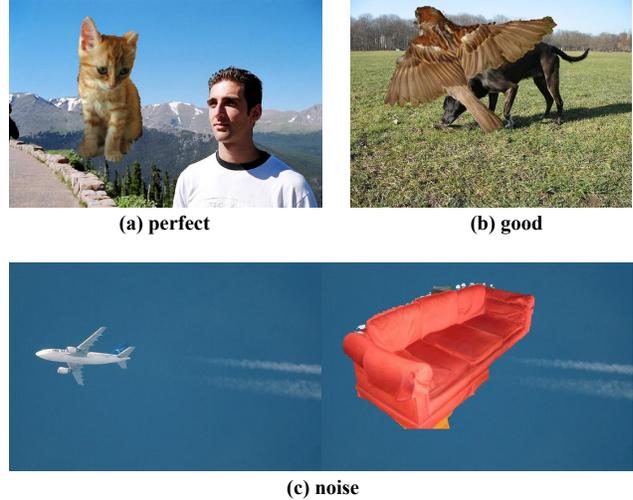


Figure 4. Examples of the input augmented images with varying degrees of occlusion.

“sofa” completely covers the “aeroplane” in the original image, which will cause confusion to network classification. However, we consider that such hard examples do not account for the majority. Most objects occupy in the middle or prominent location of the natural images. The random blending method we employ tends to paste the new objects into the off-center position of the images. Thus, this case does not affect learning. Hence, our framework is robust to the quality of augmentation. According to our experiments, this simple random blending method performs well in boosting the performance.

Online Training. The augmentation scheme is conducted online to enhance the network trained in stage-I to improve the ability to distinguish object features. Formally, in each batch, we sample $N/2$ images from the training dataset and the same number object instances images from the subset which is provided from stage-I. Then we randomly paste the segmented objects into the input images, which creates a $N/2$ batch new images. Thus, a batch of size N is generated online for each augmentation iteration. The construction process of the online augmentation learning is summarized in Algorithm 1.

Note that we train the online augmentation method in a pairwise manner as shown in Figure 2 stage-II left. We consider this can further help the networks to recognize the objects for the reason that some images have new blended objects, while some do not, which can help the classifier find more discriminative features. The motivation behind this is similar to “finding the differences” with the human visual system. When the two images have a different object but with a duplicated background, which can often leave a deep impression. For the same reason, this can make the network classifier learn better features of this kind of object.

Algorithm 1 Stage-II: Online Augmentation.

Input:

The training dataset images \mathcal{I} and the corresponding labels \mathcal{L} ;

The object instances \mathcal{O} and the corresponding labels \mathcal{T} .

- 1: **while** not done **do**
 - 2: $(\mathcal{I}_i, \mathcal{L}_i) \leftarrow$ Draw one sample from training dataset;
 - 3: $(\mathcal{O}_j, \mathcal{T}_j) \leftarrow$ Draw one sample from object instances subset;
 - 4: **while** \mathcal{T}_j in \mathcal{L}_i **do**
 - 5: $(\mathcal{O}_j, \mathcal{T}_j) \leftarrow$ Resample;
 - 6: **end while**
 - 7: $\mathcal{I}'_i \leftarrow$ Blend \mathcal{O}_j into \mathcal{I}_i ;
 - 8: $\mathcal{L}'_i \leftarrow$ Append \mathcal{T}_j in \mathcal{L}_i ;
 - 9: Train CAM \leftarrow Loss($\mathbb{C}(\mathcal{I}_i), \mathcal{L}_i$) + Loss($\mathbb{C}(\mathcal{I}'_i), \mathcal{L}'_i$);
 - 10: **end while**
 - 11: Expansion.
-

3.3. Discussion

The proposed CDA framework contributes a new data augmentation learning strategy. Unlike the previous “copy-and-paste” works, we do not use additional pixel-wise labels. Specifically, by using the self-provided initial segmentation masks of the models, we can obtain the object instances for the next phase augmentation training. Furthermore, since our goal is to decouple the high correlation between objects and their contextual background, we don’t need to consider much about visual context [13, 9], which can greatly improve the efficiency of pasting objects into the images. Besides, we adopt online augmentation training skills. Compared with static offline data augmentation, which merely enlarges the scale of the training dataset in linear-level. Namely, once a new dataset is formed, the number of images will remain unchanged. However, our method is able to obtain exponential-level augmentation, because the combination of object instances and natural images can be ever-changing in each round of training.

4. Experiments

To demonstrate the contributions of the proposed method, we conduct several ablation studies to show the effectiveness of CDA and compare different baselines models to the state-of-the-arts. We will give the details of the datasets, evaluation metric, and baseline models in the following.

4.1. Dataset

All the networks in our framework are trained and evaluated on the PASCAL VOC 2012 [15] and COCO [32] segmentation benchmark for a fair comparison to previous

approaches. As for PASCAL VOC, the official dataset separation has 1464 images for training, 1449 for validation and 1456 for testing. Following the common practice, we take additional annotations to build an augmented training set with 10582 images presented in [19]. COCO is a more challenging benchmark with 81 semantic classes (one background class), 80k, and 40k images for training and validation. We use the standard mean Intersection-over-Union (mIoU) as the evaluation metric for all experiments.

4.2. Implementation Details

To validate the applicability of CDA, we deploy it on three popular WSSS models including IRNet [1], AffinityNet [2] and SEAM [43]. The general training architecture components include a multi-label image classification step, a pseudo-mask generation step, and the final segmentation model (DeepLab-v2 [7]). We strictly follow the same settings as reported in the official codes. Specially, for SEAM [43] and AffinityNet [2] baselines, ResNet38 [20] that pre-trained on ImageNet [11] is adopted as backbone with batch size as 8 and 16, respectively. When training the networks, multi-scale and data augmentation techniques like horizontal flip, random cropping, and color jittering are deployed in both architectures. Following the poly policy $lr_{init} = lr_{init}(1 - itr / max_itr)^\rho$ with $\rho = 0.9$ for decay, the models are trained with a fix input size as 448×448 using Adam optimizer [25]. Besides, online hard example mining [39] is employed on the training loss in SEAM. As for IRNet [1], ResNet50 [20] is used as the backbone network (pretrained on ImageNet). The batch size is set to 16 for the image classification model and 32 for the inter-pixel relation model. The input image is cropped into a fix size of 512×512 using zero padding if needed. The model is trained with the same polynomial decay strategy as in AffinityNet [2] using stochastic gradient descent (SGD) for optimization with 8, 000 iterations. The fully-connected CRF [27] is used in three baselines to refine CAM, pseudo-mask, and segmentation mask with the default parameters in the public code. We set the threshold $\epsilon_1 = 0.1$ and $\epsilon_2 = 0.7$ by experience.

4.3. Ablation Studies

To verify the effectiveness of our CDA, we evaluate CAM seed regions, pseudo-masks, and segmentation masks, respectively. In our experiments, the standard mean Intersection over Union (mIoU) is used on the training set for evaluating CAM seed area masks and pseudo-masks, and on the PASCAL VOC 2012 *val* and *test* sets for evaluating segmentation masks. For the sake of simplicity, since the three WSSS models are all based on CAM [51], we use one of the representative models (IRNet [1]) as a baseline to conduct several ablation studies on CAM in mIoU to illustrate the role of each component of our approach.

Random pasting vs. Other sophisticated augmenta-

Method	operation	mIoU (%)
Conventional Augmentation	Rotation	48.5
	Translation	48.4
Mixup [49]	$\alpha = 0.3$	48.7
	$\alpha = 0.5$	48.5
	$\alpha = 0.8$	49.0
CutOut [12]	Random	48.9
CutMix [48]	Random	49.2
Random pasting (ours)	Rescale	49.8

Table 1. Experiments of different augmentation methods. Here α is the intensity of the interpolation between the eigenvector and the target vector.

Baseline	Rescale	Rotation	Gaussian	mIoU (%)
✓				48.3
✓	✓			49.8
✓	✓	✓		50.8
✓	✓		✓	49.6
✓	✓	✓	✓	50.4

Table 2. The ablation study of the effect on different pasting methods. Baseline indicates the original CAM method without pasting new objects for augmentation.

Training manner	mIoU (%)
Pairwise	50.8
None-pairwise	50.1

Table 3. Experiments of augmentation training manner.

tion methods: As for the traditional augmentation methods, we adopt the random rotation and translation to expand the dataset to three times the original size, however, they can not bring significant boost for the performance. We also compare Mixup [49], CutOut [12] and CutMix [48] methods to generate new augmented images. As shown in Table 1, random rescale pasting outperforms the other three methods achieving **49.8%** mIoU. These results demonstrate that random pasting is suitable for our CDA framework. We consider that proper occlusion helps the network to better mine the features of other areas of the objects, and the situation of complete occlusion is relatively rare which will not affect our learning process.

Comparison with baseline: We further explore the impact of different pasting methods on data augmentation. Table 2 shows that using random rescale pasting has a 1.5% improvement compared to baseline. After combining rescale and rotation, we can get the best performance to **50.8%** mIoU on PASCAL VOC training set. The results show that applying Gaussian smoothing can not help to

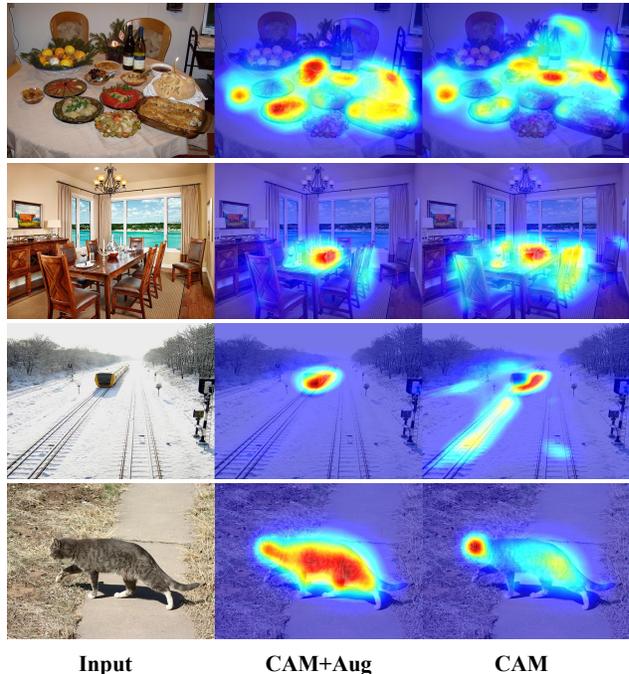


Figure 5. Qualitative visualization of CAMs. Our CDA framework not only suppresses over-activation (1st, 2nd, 3rd row) of the high correlation contextual backgrounds of the objects and expands CAMs to cover the whole object regions (4th row).

prove the performance. Therefore, in subsequent experiments, unless otherwise specified, we will use the random rescale combining with the rotation method.

Figure 5 shows the qualitative comparison between our CAM+Aug by CDA method and the original CAM. As shown in the first and second rows in the figure and the labels of objects are “table”. The original CAM will activate background semantic information that is strongly related to the “table”, such as “chair”. However, by employing the decoupling augmentation training strategy, our method can focus on the target areas. For the image with the label of “train”, CAM even pays attention not to the object itself, but the “track”, which will be detrimental to the subsequent segmentation task. Moreover, CDA can also help the network expand and discover more comprehensive object features but not only the most discriminative regions like the “cat” shown in the last row.

The effect on pairwise training: Compared to merely using the augmented images to train the networks, we use the none-augmented images with the augmented images as pair images to jointly train the models as shown in Figure 2 stage-II. The results shown in Table 3 show that applying pairwise training strategy outperforms the one in single augmented images, which illustrates that this helps the network classifier to learn more discriminative features.

Network	Backbone	CAM	Pseudo-Masks	Seg. Masks (val-set)	Seg. Masks (test-set)
AffinityNet [2] + CDA	ResNet-38	48.0	59.7	61.7	63.7
	ResNet-38	48.9 ^{+0.9}	63.3 ^{+3.6}	64.2 ^{+2.5}	65.8 ^{+2.1}
IRNet* [1] + CDA	ResNet-50	48.3	65.9	63.5	64.8
	ResNet-50	50.8 ^{+2.5}	67.7 ^{+1.8}	65.8 ^{+2.3}	66.4 ^{+1.6}
SEAM [43] + CDA	ResNet-38	55.4	63.4	64.5	65.7
	ResNet-38	58.4 ^{+3.0}	66.4 ^{+3.0}	66.1 ^{+1.6}	66.8 ^{+1.1}

Table 4. Different baselines with our CDA framework performance in mIoU on PASCAL VOC. *denotes our reimplemented results since the original code does not provided pre-trained weights.

Number of pasted objects	Same category objects	mIoU (%)
1	×	50.8
2	×	48.9
3	×	47.8
1	✓	50.2
2	✓	48.6
3	✓	47.4

Table 5. Experiments of different number of pasted objects for augmentation.

The effect on objects numbers: Under the default settings of our experiment, we only paste one new instance that does not exist in the original images. We further explore the effect of pasting multiple objects into the images to conduct augmentation. As shown in Table 5 above the solid line, when the number of object to be pasted increases from one to two, the mIoU performance will decrease. As the number of pasted objects changes to three, it will even worse than the baseline. The results show that over-pasted objects may cover the objects in the original image, making the noise sample dominant. This will confuse the classifier, which will bring negative effects. In addition, as depicted below the solid line in Table 5, when we allow the pasted object to be consistent with the object category in the original image, their general performance is worse than the former. This shows that forcing objects of different categories to be pasted into images can decouple the strong contextual dependence of objects in the original semantic environment.

Analysis of pseudo labels and Segmentation masks: The overall results are shown in Table 4. We can observe that deploying CDA on different weakly supervised semantic segmentation models can improve all their performances. Specifically, SEAM [43] can achieve the best performance in Segmentation Masks on both validation set and testing set. Figure 6 shows that we can obtain more accurate and complete masks covering the object areas.

4.4. Comparison with State-of-the-arts

Finally, we compare our framework with state-of-the-art methods on the PASCAL VOC 2012 and COCO dataset in-

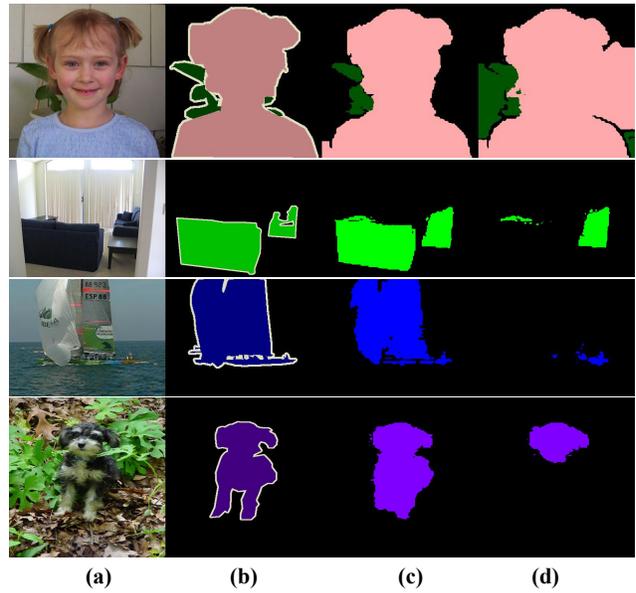


Figure 6. Visualization of pseudo-masks (baseline: IRNet [1]). (a) Input images. (b) Ground-Truth labels. (c) Our CAM+Aug. (d) Original CAM.

cluding both the validation set and the testing set. For a fair comparison, we adopt the same DeepLab [6, 7] architectures as reported in the original papers. On PASCAL VOC 2012, as is shown in Table 6, although different baselines already boosts performance compared to previous methods, when CDA is deployed in the models, SEAM [43] can achieve the best performance and outperform other state-of-the-arts by a large margin. IRNet [1] yield the second best performance and can beat its later published works. On COCO, CDA deployed on IRNet achieves 33.7% mIoU on the val set, which surpasses the previous best model by 1.1% mIoU. Figure 7 presents qualitative results of our CDA approach applying on IRNet baseline and compares them to itself. We can observe that CDA can make more accurate predictions on objects, which shows better demarcations in some coherent areas. Meanwhile, CDA can help to expand and discover more comprehensive object regions.

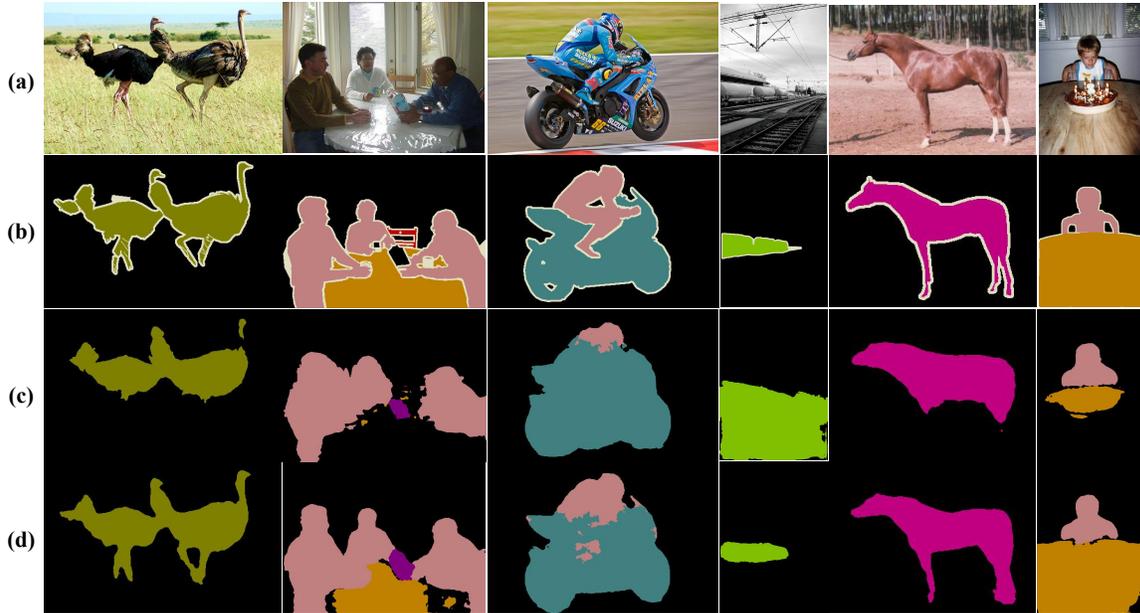


Figure 7. Qualitative results on the PASCAL VOC 2012 val set. (a) Input images. (b) Ground-truth labels. (c) Results obtained by IRNet [1] baseline. (d) Results of our IRNet + CDA. More results can be found in the supplementary material.

Methods	Backbone	Saliency	<i>val</i>	<i>test</i>
CCNN [34] _{ECCV'15}	VGG16	-	35.3	35.6
SEC [26] _{ECCV'16}	VGG16	-	50.7	51.1
STC [45] _{TPAMI'17}	VGG16	✓	49.8	51.2
AdvEra [44] _{CVPR'17}	VGG16	✓	55.0	55.7
DCSP [5] _{BMVC'17}	ResNet101	✓	60.8	61.9
MDC [46] _{CVPR'18}	VGG16	✓	60.4	60.8
MCOF [42] _{CVPR'18}	ResNet101	✓	60.3	61.2
DSRG [23] _{CVPR'18}	ResNet101	✓	61.4	63.2
AffinityNet [2] _{CVPR'18}	ResNet-38	-	61.7	63.7
IRNet [1] _{CVPR'19}	ResNet50	-	63.5	64.8
FickleNet [28] _{CVPR'19}	ResNet101	✓	64.9	65.3
SEAM [43] _{CVPR'20}	ResNet38	-	64.5	65.7
ICD [16] _{CVPR'20}	ResNet101	-	64.1	64.3
IRNet + CDA (ours)	ResNet50	-	65.8	66.4
SEAM + CDA (ours)	ResNet38	-	66.1	66.8

Table 6. Performance comparisons with other state-of-the-art WSSS methods on PASCAL VOC 2012 dataset. The **best** and **second best** performance under each set are marked with corresponding formats.

5. Conclusion

In this paper, we propose a Context Decoupling Augmentation (CDA) method for WSSS and to narrow the gap with fully supervision. Specifically, through a two-stage training, the object instances provided by the network itself are copied and pasted into the input images to conduct augmentation. To further improve the ability of network for

Methods	Backbone	<i>val</i>
BFBP [38] _{ECCV'16}	VGG16	20.4
SEC [26] _{ECCV'16}	VGG16	22.4
IRNet [1] _{CVPR'19}	ResNet50	32.6
SEAM [43] _{CVPR'20}	ResNet38	31.9
IAL [41] _{IJCV'20}	VGG16	27.7
IRNet + CDA (ours)	ResNet50	33.7
SEAM + CDA (ours)	ResNet38	33.2

Table 7. Performance comparisons with other state-of-the-art WSSS methods on COCO val in terms of mIoU.

learning object features, we adopt pairwise training manner to help the classifier to distinguish more discriminative features. Experimental results show that CDA can help boost various WSSS methods to the new state-of-the-arts.

Acknowledgement

This work was supported by National Natural Science Foundation of China (NSFC) 61876208, Key-Area Research and Development Program of Guangdong Province 2018B010108002, Central Universities of China under Grant D2192860, and the National Research Foundation, Singapore under its AI Singapore Programme (AISG Award No: AISG-RP-2018-003), and the MOE Tier-1 research grants: RG28/18 (S), RG22/19 (S) and RG95/20.

References

- [1] Jiwoon Ahn, Sunghyun Cho, and Suha Kwak. Weakly supervised learning of instance segmentation with inter-pixel relations. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2209–2218, 2019. [1](#), [2](#), [3](#), [5](#), [7](#), [8](#)
- [2] Jiwoon Ahn and Suha Kwak. Learning pixel-level semantic affinity with image-level supervision for weakly supervised semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4981–4990, 2018. [1](#), [2](#), [5](#), [7](#), [8](#)
- [3] Relja Arandjelović and Andrew Zisserman. Object discovery with a copy-pasting gan. *arXiv preprint arXiv:1905.11369*, 2019. [3](#)
- [4] Amy Bearman, Olga Russakovsky, Vittorio Ferrari, and Li Fei-Fei. What’s the point: Semantic segmentation with point supervision. In *European conference on computer vision*, pages 549–565. Springer, 2016. [1](#)
- [5] Arslan Chaudhry, Puneet K Dokania, and Philip HS Torr. Discovering class-specific pixels for weakly-supervised semantic segmentation. *arXiv preprint arXiv:1707.05821*, 2017. [8](#)
- [6] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille. Semantic image segmentation with deep convolutional nets and fully connected crfs. *arXiv preprint arXiv:1412.7062*, 2014. [7](#)
- [7] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille. 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, 2017. [1](#), [5](#), [7](#)
- [8] Liang-Chieh Chen, Yukun Zhu, George Papandreou, Florian Schroff, and Hartwig Adam. Encoder-decoder with atrous separable convolution for semantic image segmentation. In *Proceedings of the European conference on computer vision (ECCV)*, pages 801–818, 2018. [1](#)
- [9] Wenqing Chu and Deng Cai. Deep feature based contextual model for object detection. *Neurocomputing*, 275:1035–1042, 2018. [5](#)
- [10] Jifeng Dai, Kaiming He, and Jian Sun. Boxesup: Exploiting bounding boxes to supervise convolutional networks for semantic segmentation. In *Proceedings of the IEEE international conference on computer vision*, pages 1635–1643, 2015. [1](#)
- [11] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009. [5](#)
- [12] Terrance DeVries and Graham W Taylor. Improved regularization of convolutional neural networks with cutout. *arXiv preprint arXiv:1708.04552*, 2017. [3](#), [6](#)
- [13] Nikita Dvornik, Julien Mairal, and Cordelia Schmid. Modeling visual context is key to augmenting object detection datasets. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 364–380, 2018. [2](#), [3](#), [4](#), [5](#)
- [14] Debidatta Dwibedi, Ishan Misra, and Martial Hebert. Cut, paste and learn: Surprisingly easy synthesis for instance detection. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1301–1310, 2017. [3](#), [4](#)
- [15] Mark Everingham, SM Ali Eslami, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes challenge: A retrospective. *International journal of computer vision*, 111(1):98–136, 2015. [2](#), [5](#)
- [16] Junsong Fan, Zhaoxiang Zhang, Chunfeng Song, and Tieniu Tan. Learning integral objects with intra-class discriminator for weakly-supervised semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4283–4292, 2020. [8](#)
- [17] Maayan Frid-Adar, Eyal Klang, Michal Amitai, Jacob Goldberger, and Hayit Greenspan. Synthetic data augmentation using gan for improved liver lesion classification. In *2018 IEEE 15th international symposium on biomedical imaging (ISBI 2018)*, pages 289–293. IEEE, 2018. [3](#)
- [18] Robert Geirhos, Patricia Rubisch, Claudio Michaelis, Matthias Bethge, Felix A Wichmann, and Wieland Brendel. Imagenet-trained cnns are biased towards texture; increasing shape bias improves accuracy and robustness. *arXiv preprint arXiv:1811.12231*, 2018. [3](#)
- [19] Bharath Hariharan, Pablo Arbeláez, Lubomir Bourdev, Subhransu Maji, and Jitendra Malik. Semantic contours from inverse detectors. In *2011 International Conference on Computer Vision*, pages 991–998. IEEE, 2011. [5](#)
- [20] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. [5](#)
- [21] Seunghoon Hong, Donghun Yeo, Suha Kwak, Honglak Lee, and Bohyung Han. Weakly supervised semantic segmentation using web-crawled videos. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7322–7330, 2017. [1](#), [2](#)
- [22] Qibin Hou, PengTao Jiang, Yunchao Wei, and Ming-Ming Cheng. Self-erasing network for integral object attention. In *Advances in Neural Information Processing Systems*, pages 549–559, 2018. [2](#)
- [23] Zilong Huang, Xinggang Wang, Jiasi Wang, Wenyu Liu, and Jingdong Wang. Weakly-supervised semantic segmentation network with deep seeded region growing. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7014–7023, 2018. [8](#)
- [24] Anna Khoreva, Rodrigo Benenson, Jan Hosang, Matthias Hein, and Bernt Schiele. Simple does it: Weakly supervised instance and semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 876–885, 2017. [1](#)
- [25] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. [5](#)
- [26] Alexander Kolesnikov and Christoph H Lampert. Seed, expand and constrain: Three principles for weakly-supervised image segmentation. In *European conference on computer vision*, pages 695–711. Springer, 2016. [8](#)

- [27] Philipp Krähenbühl and Vladlen Koltun. Efficient inference in fully connected crfs with gaussian edge potentials. In *Advances in neural information processing systems*, pages 109–117, 2011. 5
- [28] Jungbeom Lee, Eunji Kim, Sungmin Lee, Jangho Lee, and Sungroh Yoon. Ficklenet: Weakly and semi-supervised semantic image segmentation using stochastic inference. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5267–5276, 2019. 8
- [29] Kunpeng Li, Ziyang Wu, Kuan-Chuan Peng, Jan Ernst, and Yun Fu. Tell me where to look: Guided attention inference network. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 9215–9223, 2018. 2
- [30] Di Lin, Jifeng Dai, Jiaya Jia, Kaiming He, and Jian Sun. Scribblesup: Scribble-supervised convolutional networks for semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3159–3167, 2016. 1
- [31] Guosheng Lin, Chunhua Shen, Anton Van Den Hengel, and Ian Reid. Efficient piecewise training of deep structured models for semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3194–3203, 2016. 1
- [32] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014. 2, 5
- [33] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3431–3440, 2015. 1
- [34] Deepak Pathak, Philipp Krahenbuhl, and Trevor Darrell. Constrained convolutional neural networks for weakly supervised segmentation. In *Proceedings of the IEEE international conference on computer vision*, pages 1796–1804, 2015. 8
- [35] Xingchao Peng, Baochen Sun, Karim Ali, and Kate Saenko. Learning deep object detectors from 3d models. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1278–1286, 2015. 3
- [36] Tal Remez, Jonathan Huang, and Matthew Brown. Learning to segment via cut-and-paste. In *Proceedings of the European conference on computer vision (ECCV)*, pages 37–52, 2018. 3
- [37] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. *arXiv preprint arXiv:1506.01497*, 2015. 3
- [38] Fatemehsadat Saleh, Mohammad Sadegh Aliakbarian, Mathieu Salzmann, Lars Petersson, Stephen Gould, and Jose M Alvarez. Built-in foreground/background prior for weakly-supervised semantic segmentation. In *European conference on computer vision*, pages 413–432. Springer, 2016. 8
- [39] Abhinav Shrivastava, Abhinav Gupta, and Ross Girshick. Training region-based object detectors with online hard example mining. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 761–769, 2016. 5
- [40] Paul Vernaza and Manmohan Chandraker. Learning random-walk label propagation for weakly-supervised semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7158–7166, 2017. 1
- [41] Xiang Wang, Sifei Liu, Huimin Ma, and Ming-Hsuan Yang. Weakly-supervised semantic segmentation by iterative affinity learning. *International Journal of Computer Vision*, 128(6):1736–1749, 2020. 8
- [42] Xiang Wang, Shaodi You, Xi Li, and Huimin Ma. Weakly-supervised semantic segmentation by iteratively mining common object features. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1354–1362, 2018. 8
- [43] Yude Wang, Jie Zhang, Meina Kan, Shiguang Shan, and Xilin Chen. Self-supervised equivariant attention mechanism for weakly supervised semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12275–12284, 2020. 1, 5, 7, 8
- [44] Yunchao Wei, Jiashi Feng, Xiaodan Liang, Ming-Ming Cheng, Yao Zhao, and Shuicheng Yan. Object region mining with adversarial erasing: A simple classification to semantic segmentation approach. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1568–1576, 2017. 2, 4, 8
- [45] Yunchao Wei, Xiaodan Liang, Yunpeng Chen, Xiaohui Shen, Ming-Ming Cheng, Jiashi Feng, Yao Zhao, and Shuicheng Yan. Stc: A simple to complex framework for weakly-supervised semantic segmentation. *IEEE transactions on pattern analysis and machine intelligence*, 39(11):2314–2320, 2016. 8
- [46] Yunchao Wei, Huaxin Xiao, Honghui Shi, Zequn Jie, Jiashi Feng, and Thomas S Huang. Revisiting dilated convolution: A simple approach for weakly-and semi-supervised semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7268–7277, 2018. 8
- [47] Fisher Yu and Vladlen Koltun. Multi-scale context aggregation by dilated convolutions. *arXiv preprint arXiv:1511.07122*, 2015. 2
- [48] Sangdoon Yun, Dongyoon Han, Seong Joon Oh, Sanghyuk Chun, Junsuk Choe, and Youngjoon Yoo. Cutmix: Regularization strategy to train strong classifiers with localizable features. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6023–6032, 2019. 6
- [49] Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. *arXiv preprint arXiv:1710.09412*, 2017. 3, 6
- [50] Tianyi Zhang, Guosheng Lin, Weide Liu, Jianfei Cai, and Alex Kot. Splitting vs. merging: Mining object regions with discrepancy and intersection loss for weakly supervised semantic segmentation. In *European Conference on Computer Vision*, 2020. 1
- [51] Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, and Antonio Torralba. Learning deep features for discrimina-

tive localization. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2921–2929, 2016. [2](#), [5](#)

- [52] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *Proceedings of the IEEE international conference on computer vision*, pages 2223–2232, 2017. [3](#)