# WUDA: Unsupervised Domain Adaptation Based on Weak Source Domain Labels

## Shengjie Liu[1], Chuang Zhu[*1], Wenqi Tang[1]

[1]Beijing University of Posts and Telecommunications, Beijing, China
{shengjie.Liu, czhu, tangwenqi}@bupt.edu.cn

## Abstract

Unsupervised domain adaptation (UDA) for semantic segmentation addresses the cross-domain problem with fine source domain labels. However, the acquisition of semantic labels has always been a difficult step, many scenarios only have weak labels (e.g. bounding boxes). For scenarios where weak supervision and cross-domain problems co-exist, this paper defines a new task: unsupervised domain adaptation based on weak source domain labels (WUDA). To explore solutions for this task, this paper proposes two intuitive frameworks: 1) Perform weakly supervised semantic segmentation in the source domain, and then implement unsupervised domain adaptation; 2) Train an object detection model using source domain data, then detect objects in the target domain and implement weakly supervised semantic segmentation. We observe that the two frameworks behave differently when the datasets change. Therefore, we construct dataset pairs with a wide range of domain shifts and conduct extended experiments to analyze the impact of different domain shifts on the two frameworks. In addition, to measure domain shift, we apply the metric representation shift to urban landscape image segmentation for the first time. The source code and constructed datasets are available at https://github.com/bupt-ai-cz/WUDA.

## Introduction

As one of the most popular computer vision technologies, semantic segmentation has developed very mature and is widely used in many scenarios such as autonomous driving, remote sensing image recognition, and medical image processing. However, like most deep learning models, a serious problem faced by semantic segmentation models in the application is the existence of domain shift. Since the data distribution of the target domain may be very different from the source domain, a model that performs well in the source domain may experience catastrophic performance degradation in the target domain. For cross-domain problems, many studies on domain adaptation have improved the generalization ability of the model. Among these studies, unsupervised domain adaptation (UDA) for semantic segmentation can fine-tune a trained model by using self-training, adversarial training, or image style transfer without any additional annotations of the target domain. Methods of self-training
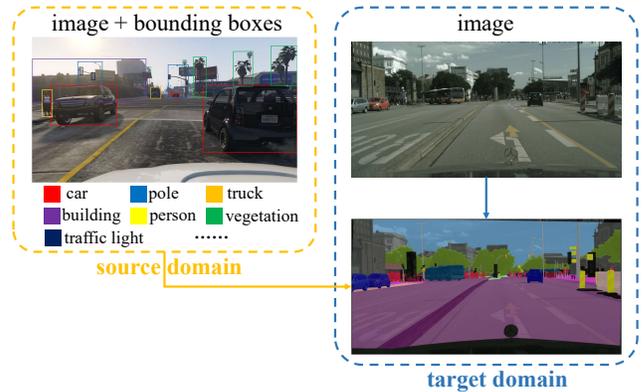
---
[*]Corresponding author.



Figure 1: Schematic diagram of WUDA. The source domain has only bounding boxes and the target domain has no annotations. WUDA achieves semantic segmentation of the target domain images under these conditions.

(Zou et al. 2018, 2019) and adversarial training (Tsai et al. 2018; Luo et al. 2019b) can adapt the model to the distribution of the target domain images. Image style transfer methods (Yang and Soatto 2020; Yang et al. 2020a) can bring the target domain data closer to the distribution of the source domain.

In deep learning tasks, massive amounts of samples are required for training to improve the robustness of the model, therefore, the annotation of large-scale datasets is another difficulty in training deep models. In the semantic segmentation task, in order to obtain pixel-wise mask annotations, it takes a lot of time for a single image (e.g. it takes 1.5 hours to label one image in the Cityscapes (Cordts et al. 2016) dataset), and the annotation of the entire dataset requires huge manpower. Therefore, many computer vision datasets have only weak annotations (e.g. bounding boxes, points etc.)

When the cross-domain problem and the weak labels co-exist (only source domain bounding boxes and target domain images are available), the domain shift and the weak supervision both bring a negative contribution to the pixel-level semantic segmentation. In this case, it is a challenge to achieve accurate semantic segmentation in the target domain. We define this task as Unsupervised Domain Adaptation Based on Weak Source Domain Labels (WUDA). The

schematic diagram of WUDA is shown in Figure 1. For the newly defined task, it is necessary to explore a framework to tackle the problems of transfer learning and weakly supervised segmentation at the same time. The realization of this task can reduce the requirements for source domain labels in future UDA tasks.

In summary, this paper makes the following contributions:

- We define a novel task: unsupervised domain adaptation based on weak source domain labels (WUDA). For this task, we propose two intuitive frameworks: Weakly Supervised Semantic Segmentation + Unsupervised Domain Adaptation (WSSS-UDA) and Target Domain Object Detection + Weakly Supervised Semantic Segmentation (TDOD-WSSS).

- We benchmark typical weakly supervised semantic segmentation, unsupervised domain adaptation, and object detection techniques under our two proposed frameworks, and find that the results of framework WSSS-UDA can reach 83% of the UDA method with fine source domain labels.

- We construct a series of datasets with different domain shifts. To the best of our knowledge, we are the first to use representation shift for domain shift measurement in urban landscape datasets. The constructed dataset will be open for research on WUDA/UDA under multiple domain shifts.

- To further analyze the impact of different degrees of domain shift on our proposed frameworks, we conduct extended experiments using our constructed datasets and find that framework TDOD-WSSS is more sensitive to changes in domain shift than framework WSSS-UDA.

## Related Work

WUDA will involve weakly supervised semantic segmentation, unsupervised domain adaptation, object detection, and the measure of domain shift techniques. In this section, we will review these related previous works.

### Weakly Supervised Semantic Segmentation

In computer vision tasks, pixel-wise mask annotations takes far more time compared to weak annotations (Lin et al. 2014), and the need for time-saving motivates weakly supervised semantic segmentation. Labels for weakly supervised segmentation can be bounding boxes, points, scribbles and image-level tags. Methods (Dai, He, and Sun 2015; Khoreva et al. 2017; Li, Arnab, and Torr 2018; Song et al. 2019; Kulharia et al. 2020) using bounding boxes as supervision usually employ GrabCut (Rother, Kolmogorov, and Blake 2004) or segment proposals techniques to get more accurate semantic labels and can achieve results close (95% or even higher) to fully supervised methods. Point-supervised and scribble-supervised methods (Bearman et al. 2016; Qian et al. 2019; Lin et al. 2016; Vernaza and Chandraker 2017; Tang et al. 2018a,b) take advantage of location and category information in annotations and achieve excellent segmentation results. Tag-supervised methods (Jiang et al. 2019; Wang et al. 2020b; Lee et al. 2021b; Li et al. 2021b) often use class activation mapping (CAM) (Zhou et al. 2016) algorithm to obtain localization maps of the main objects in the images.

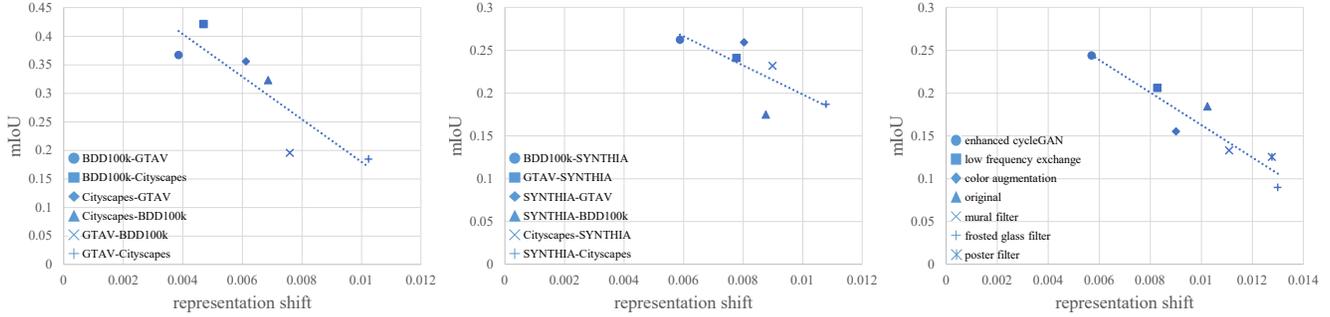### Unsupervised Domain Adaptation for Semantic Segmentation

Unsupervised Domain Adaptation (UDA) is committed to solving the problem of poor model generalization caused by inconsistent data distribution in the source and target domains. Self-training (ST) and adversarial training (AT) are key schemes of UDA: self-training schemes (Zou et al. 2018, 2019; Lian et al. 2019; Li et al. 2020; Lv et al. 2020; Melas-Kyriazi and Manrai 2021; Tranheden et al. 2021; Guo et al. 2021; Araslanov and Roth 2021) typically set a threshold to filter pseudo-labels with high confidence on the target domain, and use the pseudo-labels to supervise target domain training; adversarial training methods (Tsai et al. 2018; Luo et al. 2019b,a; Du et al. 2019; Vu et al. 2019; Tsai et al. 2019; Yang et al. 2020b; Wang et al. 2020a; Li et al. 2021a) usually add a domain discriminator to the model. The adversarial game of the segmenter and the discriminator can make the segmentation results of the source and target domains tend to be consistent. There are also works (Zhang et al. 2019; Pan et al. 2020; Wang et al. 2020c; Yu et al. 2021; Wang, Peng, and Zhang 2021; Mei et al. 2020) that perform both self-training and adversarial training to achieve good segmentation results on the target domain.

### Object Detection

Autonomous driving technology has greatly promoted the development of object detection. There are many pioneering works that can be widely used in various object detection tasks, such as some two-stage methods (Girshick et al. 2014; Girshick 2015; Ren et al. 2015) that first perform object extraction, and then classify the extracted objects. Yolo series of algorithms (Redmon et al. 2016; Redmon and Farhadi 2017, 2018; Bochkovskiy, Wang, and Liao 2020; Jocher et al. 2020) can simultaneously achieve object extraction and classification in one network. The current popular object detection method Yolov5 (Jocher et al. 2020) has been able to achieve 72.7% mean average precision (mAP) on the coco2017 val dataset. Object detection techniques also help to extract bounding boxes in weakly supervised segmentation methods (Lan et al. 2021; Lee et al. 2021a).

### Domain Shift Assessment

Domain shift comes from the difference between the source and target domain data. There are various factors (e.g. image content, view angle, image texture, etc.) that contribute to domain shift. While for Convolutional Neural Networks (CNN), texture is the most critical factor. Many studies (Geirhos et al. 2018; Nam et al. 2019) suggest that the focus of Convolutional Neural Networks (CNN) and human eyes is different when processing images: human eyes are sensitive to the content information of the image (e.g. shapes), while CNN is more sensitive to the style of the image (e.g. texture). Actually, if it involves the calculation of image texture differences, most methods are based on the output features of the middle layer of the neural network. For example,

(a) Linear regression results for datasets with 16 categories. Pearson correlation: -0.874

(b) Linear regression results for datasets with 19 categories. Pearson correlation: -0.741

(c) Linear regression results for augmented Cityscapes, the source domain dataset is GTAV. Pearson correlation: -0.936

Figure 2: Linear regression results of target domain mIoU and representation shift. There is a strong correlation between the two metrics, demonstrating the feasibility of representation shift for semantic segmentation tasks.

in work (Stacke et al. 2020), the metric representation shift is proposed to measure the domain shift of pathological images. In that paper, the mean value of each channel in the feature map is used as the style information of an image. The experimental results show that in the task of pathological image classification, the value of representation shift and the classification accuracy on the test set have a strong correlation, which justifies the use of the intermediate layer features to calculate the domain shift. In the pioneering work on image style transfer (Gatys, Ecker, and Bethge 2015), the style loss is calculated using the Gram matrix of the output features of the two images.

## Construct Datasets with Different Domain Shifts

In section **Experiments**, we will analyze the impact of domain shift changes on the experimental results. Therefore, we first introduce the domain shift metric *representation shift* for semantic segmentation and construct datasets with different representation shifts.

### Representation Shift for Semantic Segmentation

Karin Stacke et al. (Stacke et al. 2020) propose the metric representation shift to measure the domain shift between pathological image datasets. Since the capture of pathological images is affected by many factors (e.g. slide preparation, staining protocol, scanner properties, etc.), domain shift often exists between pathological images scanned at different medical centers. Therefore, for a pathological image classification model trained on the source domain, it may not achieve good results on test images. In work (Stacke et al. 2020), representation shift is highly correlated with the classification accuracy in the target domain: a high representation shift means that the accuracy on the target domain tends to be low. The metric solves this problem: in the absence of ground truth in the target domain, we can estimate whether the predictions made by the trained model are credible by computing the representation shift between the source and target domains.

Similarly, we envision that representation shift is also suitable for measuring domain shift in semantic segmentation tasks: for segmentation models trained on source domain data, the mIoU metric tested on the target domain is also highly correlated with representation shift. We carry out relevant experiments to prove this. Below follows a specific description of the representation shift used in our experiments: consider a semantic segmentation model with a feature extraction module and a classification module. Let $F(x) = \{f_1(x), ..., f_c(x)\}$ denote the extracted feature map of the input image $x$, where $f_c(x) \in \{\mathbb{R}^{h \times w}\}$ represents the $c$ th channel of the feature map. The average value of $f_c(x)$ is denoted as $a_c(x)$:

$$a_c(x) = \frac{1}{h}\frac{1}{w}\sum_{i,j}^{h,w} f_c(x)_{i,j}. \tag{1}$$

Let $p_{a_c}^{\mathcal{S}}$ denotes the continuous distribution of $a_c(x)$ values across the source domain dataset $X^{\mathcal{S}} = \{x_1^{\mathcal{S}}, ..., x_{n_s}^{\mathcal{S}}\}$, where $n_s$ is the number of images in source domain. Similarly, the distribution $p_{a_c}^{\mathcal{T}}$ of the target domain dataset $X^{\mathcal{T}} = \{x_1^{\mathcal{T}}, ..., x_{n_t}^{\mathcal{T}}\}$ can be obtained.

Then the representation shift $R$ for semantic segmentation is defined as follows:

$$R\left(X^{\mathcal{S}}, X^{\mathcal{T}}\right) = \frac{1}{c}\sum_{i=1}^{c} W(p_{a_i}^{\mathcal{S}}, p_{a_i}^{\mathcal{T}}), \tag{2}$$

where $W$ denotes the Wasserstein distance between the two distributions.

### Feasibility Analysis

To verify the rationality of applying representation shift in image segmentation tasks, we use four common semantic segmentation datasets for our research: Cityscapes (Cordts et al. 2016), BDD100k (Yu et al. 2020), GTAV (Richter et al. 2016) and SYNTHIA (Ros et al. 2016). Each measurement of representation shift uses one of the datasets as the source domain and the other three datasets as target domains (e.g. GTAV as the source domain, cityscapes, SYNTHIA and BDD100k as the target domain).
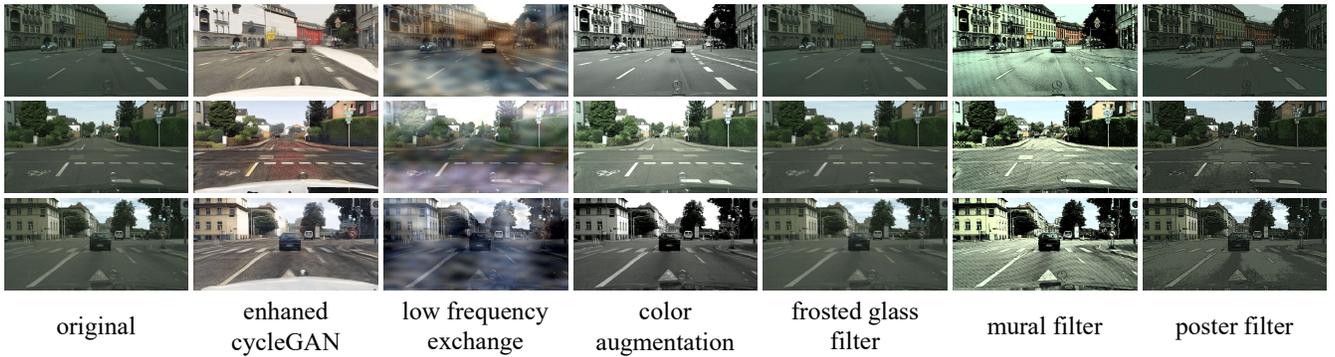
Figure 3: Visualization of different image augmentation methods for adjusting for domain shift. Best viewed in color.



Figure 4: Compared with the original cycleGAN, enhanced cycleGAN can preserve more content information.

Our data collection consists of the following 3 steps: 1) Train a semantic segmentation model using source domain data. 2) Use the trained model to extract the features of the source and target domains, then calculate the representation shift defined in Eq. 2. 3) Use the trained model to make predictions in the target domain and calculate mIoU. During the training process in step 1, we use a Deeplabv3 for semantic segmentation, set the batch size to 2, and perform 10,000 iterations without using any data enhancement methods such as resize, flip, and clipping to ensure that the model fits the original images in the source domain. The collected representation shift-mIoU pairs are used for linear regression analysis.

As can be seen from Figures 2(a) and 2(b), the value of representation shift and mIoU have a high degree of negative correlation. Pearson correlations of representation shift and mIoU reach -0.874 and -0.741 for 16 and 19 categories semantic segmentation, respectively. This conclusion is similar to that of the pathological image classification task in the article (Stacke et al. 2020), which proves that for the semantic segmentation task, using the representation shift to measure the domain shift is also applicable.

## Data Construction

For the popularly used dataset pair GTAV-Cityscapes in UDA methods, we perform different data augmentation processes on the target domain (Cityscapes) to simulate different domain shifts. This helps in section **Experiments** to analyze the impact of different domain shifts on our proposed frameworks. The visualization of image augmentations is shown in Figure 3. For possible future needs to construct dataset pairs with specific domain shifts, we propose a data construction process (Algorithm 1). Assume that data augmentation operations are sufficient.

**Enhanced cycleGAN** CycleGAN (Zhu et al. 2017) is a

---

**Algorithm 1: Dataset Construction Algorithm**

**Input**: expected domain shift interval $(a - \delta, a + \delta)$, source domain dataset $X^{\mathcal{S}}$, target domain dataset $X^{\mathcal{T}}$, a list of image augmentation operations $\{O_1, ...O_n\}$, domain shift calculation formula $R$.
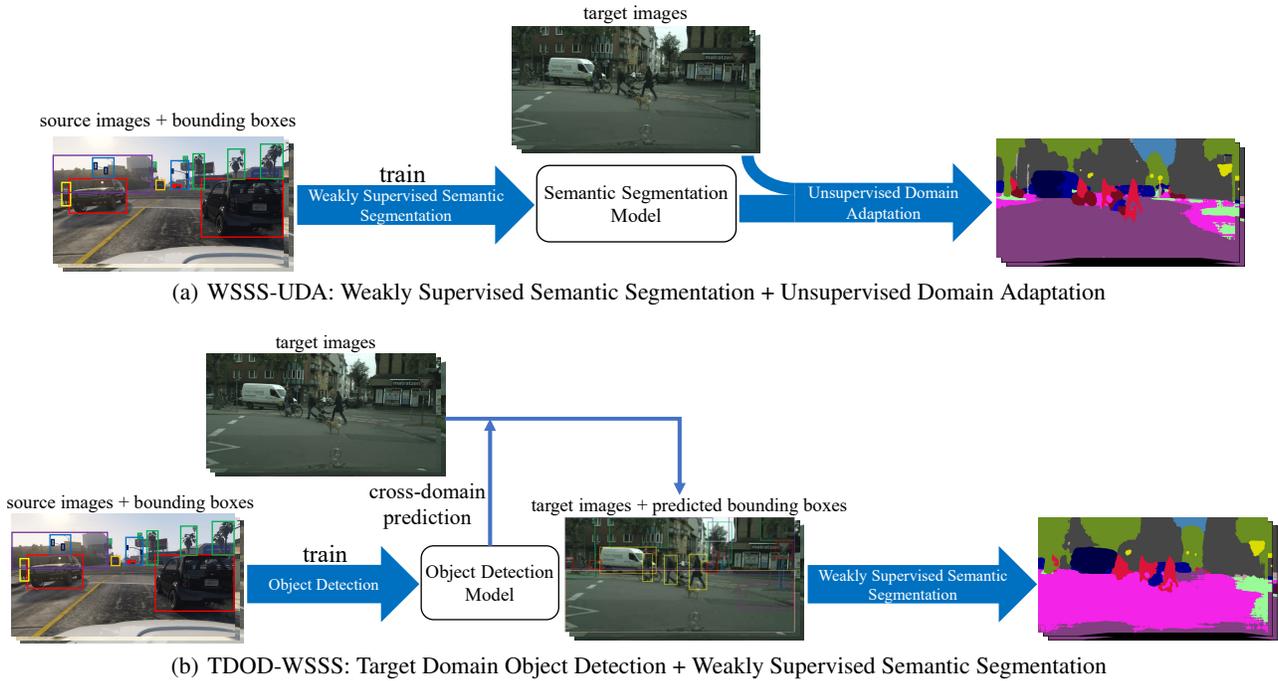
**Output**: expected dataset $\hat{X}^{\mathcal{T}}$.

1: Let $i = 1$.
2: **while** $i <= n$ **do**
3:     $\hat{X}^{\mathcal{T}} = O_i(X^{\mathcal{T}})$.
4:     **if** $R(\hat{X}^{\mathcal{T}}, X^{\mathcal{S}})$ **in** $(a - \delta, a + \delta)$ **then**
5:         **break**
6:     **else**
7:         $i = i + 1$.
8:     **end if**
9: **end while**
10: **return** $\hat{X}^{\mathcal{T}}$

---

common method for unsupervised image-to-image translation, which can realize the conversion of image style between two domains. In order to improve the generalization ability of the model, cycleGAN is used in many studies to reduce the domain shift. Therefore, we implement an enhanced cycleGAN to transfer Cityscapes images into GTAV style. Original cycleGAN only performs adversarial training at the image level, which may result in missing content in the style-transferred Cityscapes images. We add a discriminator on each transfer branch (source to target/target to source) to conduct class-level adversarial training: for generated images, the additional discriminator identifies each part (e.g. vegetation, building, etc.) of the image is real or not. As shown in Figure 4, our enhanced version of cycleGAN can further preserve the content information of the image while changing the style of the image.

**Low frequency exchange** FDA (Yang and Soatto 2020) introduces the method of low frequency exchange in unsupervised domain adaptation for semantic segmentation. Perform fast Fourier transform (FFT) on the source domain images and target domain images separately and transplants the low frequency parts of the source domain to the target domain, so that the target domain has the low-frequency features of the source domain. Theoretically, low frequency ex-

(a) WSSS-UDA: Weakly Supervised Semantic Segmentation + Unsupervised Domain Adaptation



(b) TDOD-WSSS: Target Domain Object Detection + Weakly Supervised Semantic Segmentation

Figure 5: Frameworks for WUDA task. (a) First, perform box-supervised semantic segmentation in the source domain. With the segmentation model pre-trained on the source domain, UDA methods can be performed in the target domain to enhance the generalization ability of the model. (b) Train an object detection model in the source domain. With a pre-trained object detection model in the source domain, we can predict bounding boxes on the target domain and then implement weakly supervised segmentation.

change can bring the distribution of the target domain closer to the source domain.

**Other image augmentations** We also made other augmentations to target domain images: 1) Color augmentation. Randomly adjust the brightness, contrast and color of the image, as shown in Figure 3, column 4; 2) Image filters. Use image processing software to filter images with the following effects: frosted glass, mural and poster, as shown in Figure 3, columns 5-7. These image augmentation methods change the domain shift between the source and target domains to varying degrees.

The constructed datasets cover a wide range of domain shifts as shown in Figure 2(c). The relationship between mIoU and representation shift further justifies the use of representation shift to measuring domain shift in semantic segmentation tasks.

## Frameworks

Considering that WUDA contains both weakly supervised and cross-domain problems, we first propose the framework WSSS-UDA: In order to take advantage of the fine bounding box labels in the source domain, first perform box-supervised segmentation in the source domain, then the problem is transformed into a UDA task. The schematic diagram of the framework WSSS-UDA is shown in Figure 5(a).

However, the cross-domain process is not necessarily

done on the semantic segmentation task. According to the study (Redmon et al. 2016), Yolo has a strong generalization ability: when trained on natural images and tested on the artwork, Yolo outperforms the two-stage detection methods (e.g. R-CNN (Girshick et al. 2014)) by a wide margin. Therefore, we implement the cross-domain process on the object detection task and propose the framework TDOD-WSSS (Figure 5(b)): First, use the bounding box labels of the source domain to train the object detection model, then predict bounding boxes on the target domain. Finally, implement box-supervised segmentation in the target domain.

We benchmark typical methods under our proposed frameworks. For weakly supervised semantic segmentation, we first use GrabCut to obtain pseudo-labels, then perform three epochs of self-training to gradually improve the accuracy of the segmentation results. For unsupervised domain adaptation, we benchmark CBST (Zou et al. 2018) and enhanced IAST (Mei et al. 2020), where CBST is an ST-based method and IAST combines ST and AT. Our enhanced IAST optimizes the self-training process of the original IAST. For the object detection method, we adopt the widely used Yolov5 (Jocher et al. 2020).

## Experiments

This section will introduce the datasets and metrics, the implementation details of the experiments, and the presentation and analysis of the results.

| | Road | SW | Build | Wall | Fence | Pole | TL | TS | Veg. | Terrain | Sky | Person | Rider | Car | Truck | Bus | Train | Motor | Bike | mIoU |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| CBST (fine label) | 91.8 | 53.5 | 80.5 | 32.7 | 21.0 | 34.0 | 28.9 | 20.4 | 83.9 | 34.2 | 80.9 | 53.1 | 24.0 | 82.7 | 30.3 | 35.9 | 16.0 | 25.9 | 42.8 | 45.9 |
| E-IAST (fine label) | 94.6 | 65.9 | 87.6 | 41.4 | 25.8 | 36.4 | 52.1 | 54.8 | 82.7 | 23.0 | 89.1 | 68.0 | 23.6 | 88.2 | 42.4 | 52.5 | 26.6 | 50.8 | 62.0 | 56.2 |
| WSSS-UDA (CBST) | 27.3 | 20.3 | 48.7 | 11.9 | 9.2 | 10.1 | 20.1 | 11.3 | 65.1 | 23.6 | 46.0 | 50.1 | 15.0 | 73.1 | 16.5 | 6.2 | 0.0 | 23.2 | 26.1 | 26.5 |
| WSSS-UDA (E-IAST) | 82.5 | 39.1 | 82.4 | 34.4 | 34.2 | 21.0 | 44.6 | 43.4 | 76.7 | 20.4 | 63.4 | 63.8 | 6.5 | 83.2 | 31.7 | 47.5 | 0.0 | 46.5 | 60.0 | 46.4 |
| TDOD-WSSS (Yolov5) | 61.9 | 21.1 | 76.3 | 28.3 | 26.8 | 6.7 | 38.2 | 28.4 | 74.3 | 14.4 | 59.5 | 57.0 | 11.4 | 79.3 | 20.7 | 23.0 | 0.0 | 0.0 | 0.3 | 33.0 |

Table 1: mIoU (%) results of semantic segmentation on Cityscapes with weak GTAV labels (synthesis→real). The remarks behind Framework WSSS-UDA represent which UDA method is used. The remark behind Framework TDOD-WSSS represents which object detection method is used. 'Fine label' means that precise semantic labels are used in the source domain.

| | Road | SW | Build | Wall | Fence | Pole | TL | TS | Veg. | Terrain | Sky | Person | Rider | Car | Truck | Bus | Train | Motor | Bike | mIoU |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| E-IAST (fine label) | 95.1 | 67.2 | 86.7 | 34.0 | 28.4 | 39.9 | 46.9 | 52.0 | 86.6 | 45.0 | 85.1 | 66.4 | 38.6 | 89.4 | 45.2 | 52.0 | 0.0 | 49.5 | 61.7 | 56.3 |
| WSSS-UDA (E-IAST) | 78.5 | 35.0 | 81.8 | 31.4 | 21.0 | 16.0 | 40.7 | 51.8 | 82.3 | 27.6 | 76.1 | 60.7 | 10.8 | 84.7 | 35.8 | 41.1 | 0.0 | 44.4 | 58.7 | 46.9 |
| TDOD-WSSS (Yolov5) | 71.1 | 30.5 | 77.7 | 26.6 | 31.8 | 8.0 | 18.0 | 39.8 | 78.6 | 25.0 | 59.9 | 57.0 | 30.2 | 80.2 | 48.2 | 45.1 | 0.0 | 34.2 | 54.9 | 43.0 |

Table 2: mIoU (%) results of semantic segmentation on Cityscapes with weak BDD100k labels (real→real).

## Datasets and Evaluation

GTAV (Richter et al. 2016) and Cityscapes (Cordts et al. 2016) are the most widely used autonomous driving datasets in UDA methods. For WUDA, we also adopt the GTAV-Cityscapes dataset pair in our experiments. In addition, in order to simulate the real-to-real scene, we carried out further experiments on the dataset pair BDD100k-Cityscapes.

GTAV dataset as the source domain does not have the annotation of the bounding boxes. We use the python package scikit-image to perform class-wise connected domain detection on the semantic labels of GTAV, and then obtain the box of each connected domain.

We use the GrabCut method in the OpenCV package to extract pseudo-labels from bounding boxes. When an occlusion occurs between objects, we assume that the smaller target is in front, as in the method (Khoreva et al. 2017).

For metrics, we use mean intersection over union (mIoU) for evaluation in all experiments.

## Implementation

All our experiments are implemented with Pytorch on a single NVIDIA Tesla V100 with 32 GB memory. In framework WSSS-UDA, to ensure model compatibility, the model used for weakly supervised segmentation is consistent with the corresponding UDA method: CBST uses Deeplabv2 (Chen et al. 2017) with ResNet-38 (He et al. 2016) as the backbone and our enhanced IAST uses Deeplabv2 with ResNet-101 as the backbone. All other settings remain the same as the default setting for CBST and IAST. In framework TDOD-WSSS, the weakly supervised segmentation step uses Deeplabv2 with ResNet-101 as the backbone. The object detection step uses a randomly initialized Yolov5l.

## Main Results

The results of the two frameworks for WUDA are shown in Table 1 and 2. On the synthesis-to-real dataset, the results show that framework WSSS-UDA with E-IAST has a higher potential, its highest mIoU can reach 46.4%, and this value reaches 82.6% of E-IAST with fine source domain labels. However, for WSSS-UDA with CBST, the weak source domain labels bring a relatively large attenuation to the results (45.9% → 26.5% ). Framework TDOD-WSSS only achieves a mIoU of 33.0%. This is because the Yolov5 trained on the source domain has misdetections on the target domain.

On the real-to-real dataset, the mIoU result of framework WSSS-UDA reaches 46.9%, which is similar to that on the synthesis-to-real dataset. For framework TDOD-WSSS, the result is 43%, which is a 10 percentage point increase compared to the result on the GTAV-Cityscapes dataset. As the dataset changes, the results of the two frameworks have very different trends. According to Figure 2(a), the domain shifts of GTAV-Cityscapes and BDD100k-Cityscapes are significantly different. We can reasonably hypothesize that the two frameworks have different sensitivities to changes of domain shifts. Therefore, we conduct extended experiments to verify our hypothesis in the next subsection.

Figure 6 shows the qualitative results of semantic segmentation on Cityscapes. With the support of accurate labels in the source domain, the enhanced IAST can achieve very accurate segmentation, and the segmentation error area is very small. When using bounding boxes as the source domain labels, framework WSSS-UDA also achieves generally satisfactory results, however, it is not as detailed as the fully supervised UDA. When it comes to framework TDOD-WSSS, the segmentation results become coarser and have large areas of mis-segmented regions. This is because the misdetection of Yolov5 brings more noise to the subsequent weakly supervised segmentation.

We also employed the two-stage method faster-RCNN for object detection, however, the misdetection of this model in the target domain is catastrophic. Therefore, we did not implement further experiments with two-stage object detection

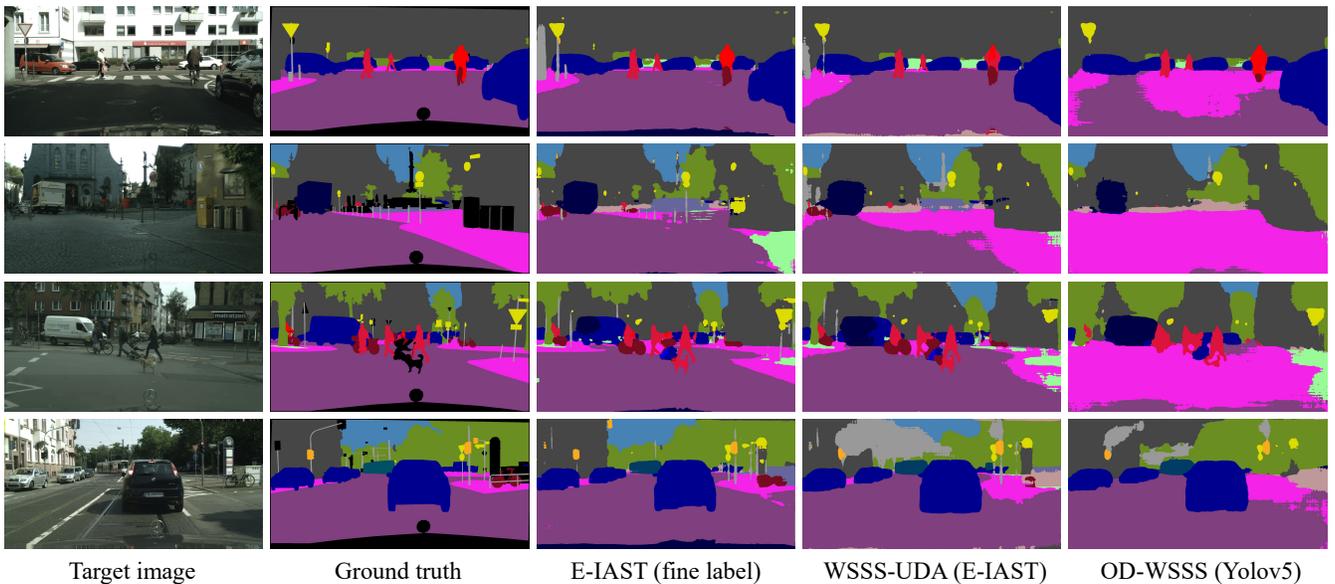| Target image | Ground truth | E-IAST (fine label) | WSSS-UDA (E-IAST) | OD-WSSS (Yolov5) |

Figure 6: Visualization of Semantic Segmentation on Cityscapes with different frameworks. Framework WSSS-UDA can achieve results close to UDA methods with fine source labels. However, the results of framework TDOD-WSSS are rough in detail and have many wrong segmentations. Best viewed in color.

methods.

## Analysis under Multiple Domain Shifts

| target domain style | WSSS-UDA (E-IAST) | WSSS-UDA (CBST) | TDOD-WSSS (Yolov5) |
|---|---|---|---|
| original | 46.4 | 26.5 | 33.0 |
| enhanced cycleGAN | 42.0 | 26.7 | 28.7 |
| low frequency exchange | 41.6 | 24.2 | 33.4 |
| color augmentation | 43.0 | 24.5 | 25.7 |
| mural filter | 37.6 | 23.5 | 6.0 |
| frosted glass filter | 35.2 | 13.4 | 3.7 |
| poster filter | 39.6 | 20.4 | 1.9 |

Table 3: mIoU (%) results of semantic segmentation on datasets with multiple domain shifts (the source domain dataset is GTAV).

Figure 7 shows the variation of mIoU results of different frameworks with domain shift. When the domain shift is moderate or small, the mIoU of WSSS-UDA is more stable than that of TDOD-WSSS, and WSSS-UDA with the E-IAST method can stabilize at a high level. When it comes to large domain shift, the mIoU of both frameworks drops: WSSS-UDA drops slowly and TDOD-WSSS drops significantly. The mIoU values of different frameworks under different domain shifts are shown in Table 3.

In general, the domain shift has a larger impact on TDOD-WSSS, while WSSS-UDA is less sensitive to changes in domain shift. WSSS-UDA has more potential in the applications of complex target domains.
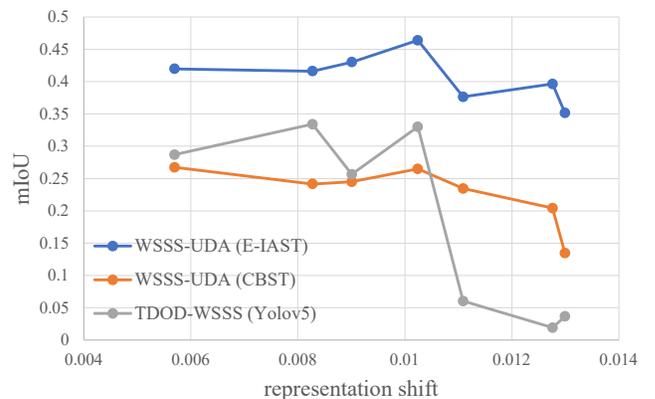


Figure 7: The variation of mIoU results of different frameworks with domain shift.

pervision. Meanwhile, we also propose two intuitive frameworks for this task. The results show that by using a suitable UDA method in the framework WSSS-UDA, mIoU on the target domain can reach 83% of UDA methods with fine source labels. In addition, we apply representation shift to semantic segmentation of urban landscapes for the first time and analyze the impact of different domain shifts on the two proposed frameworks. Experiments prove that framework WSSS-UDA is more tolerant of domain shift.

However, there is still a long way to go for current methods to achieve precise segmentation on WUDA tasks. We hope more excellent solutions will be proposed to tackle WUDA in the future.

## Conclusion

This paper defines a novel unsupervised domain adaptation task that requires only source domain bounding boxes for su-

## References

Araslanov, N.; and Roth, S. 2021. Self-supervised augmentation consistency for adapting semantic segmentation. In

*Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 15384–15394.

Bearman, A.; Russakovsky, O.; Ferrari, V.; and Fei-Fei, L. 2016. What's the point: Semantic segmentation with point supervision. In *European conference on computer vision*, 549–565. Springer.

Bochkovskiy, A.; Wang, C.-Y.; and Liao, H.-Y. M. 2020. Yolov4: Optimal speed and accuracy of object detection. *arXiv preprint arXiv:2004.10934*.

Chen, L.-C.; Papandreou, G.; Kokkinos, I.; Murphy, K.; and Yuille, A. L. 2017. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *IEEE transactions on pattern analysis and machine intelligence*, 40(4): 834–848.

Cordts, M.; Omran, M.; Ramos, S.; Rehfeld, T.; Enzweiler, M.; Benenson, R.; Franke, U.; Roth, S.; and Schiele, B. 2016. The cityscapes dataset for semantic urban scene understanding. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 3213–3223.

Dai, J.; He, K.; and Sun, J. 2015. Boxsup: Exploiting bounding boxes to supervise convolutional networks for semantic segmentation. In *Proceedings of the IEEE international conference on computer vision*, 1635–1643.

Du, L.; Tan, J.; Yang, H.; Feng, J.; Xue, X.; Zheng, Q.; Ye, X.; and Zhang, X. 2019. Ssf-dan: Separated semantic feature based domain adaptation network for semantic segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 982–991.

Gatys, L.; Ecker, A. S.; and Bethge, M. 2015. Texture synthesis using convolutional neural networks. *Advances in neural information processing systems*, 28.

Geirhos, R.; Rubisch, P.; Michaelis, C.; Bethge, M.; Wichmann, F. A.; and Brendel, W. 2018. ImageNet-trained CNNs are biased towards texture; increasing shape bias improves accuracy and robustness. *arXiv preprint arXiv:1811.12231*.

Girshick, R. 2015. Fast R-CNN. In *International Conference on Computer Vision (ICCV)*.

Girshick, R.; Donahue, J.; Darrell, T.; and Malik, J. 2014. Rich feature hierarchies for accurate object detection and semantic segmentation. In *Computer Vision and Pattern Recognition*.

Guo, X.; Yang, C.; Li, B.; and Yuan, Y. 2021. Metacorrection: Domain-aware meta loss correction for unsupervised domain adaptation in semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 3927–3936.

He, K.; Zhang, X.; Ren, S.; and Sun, J. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 770–778.

Jiang, P.-T.; Hou, Q.; Cao, Y.; Cheng, M.-M.; Wei, Y.; and Xiong, H.-K. 2019. Integral object mining via online attention accumulation. In *Proceedings of the IEEE/CVF international conference on computer vision*, 2070–2079.

Jocher, G.; Stoken, A.; Borovec, J.; NanoCode012; ChristopherSTAN; Changyu, L.; Laughing; Hogan, A.; lorenzomammana; tkianai; yxNONG; AlexWang1900; Diaconu, L.; Marc; wanghaoyang0106; ml5ah; Doug; Hatovix; Poznanski, J.; Yu, L.; changyu98; Rai, P.; Ferriday, R.; Sullivan, T.; Xinyu, W.; YuriRibeiro; Claramunt, E. R.; hopesala; pritul dave; and yzchen. 2020. ultralytics/yolov5: v3.0.

Khoreva, A.; Benenson, R.; Hosang, J.; Hein, M.; and Schiele, B. 2017. Simple does it: Weakly supervised instance and semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 876–885.

Kulharia, V.; Chandra, S.; Agrawal, A.; Torr, P.; and Tyagi, A. 2020. Box2seg: Attention weighted loss and discriminative feature learning for weakly supervised segmentation. In *European Conference on Computer Vision*, 290–308. Springer.

Lan, S.; Yu, Z.; Choy, C.; Radhakrishnan, S.; Liu, G.; Zhu, Y.; Davis, L. S.; and Anandkumar, A. 2021. Discobox: Weakly supervised instance segmentation and semantic correspondence from box supervision. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 3406–3416.

Lee, J.; Yi, J.; Shin, C.; and Yoon, S. 2021a. Bbam: Bounding box attribution map for weakly supervised semantic and instance segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2643–2652.

Lee, S.; Lee, M.; Lee, J.; and Shim, H. 2021b. Railroad is not a train: Saliency as pseudo-pixel supervision for weakly supervised semantic segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 5495–5505.

Li, G.; Kang, G.; Liu, W.; Wei, Y.; and Yang, Y. 2020. Content-consistent matching for domain adaptive semantic segmentation. In *European conference on computer vision*, 440–456. Springer.

Li, Q.; Arnab, A.; and Torr, P. H. 2018. Weakly-and semi-supervised panoptic segmentation. In *Proceedings of the European conference on computer vision (ECCV)*, 102–118.

Li, S.; Lv, F.; Xie, B.; Liu, C. H.; Liang, J.; and Qin, C. 2021a. Bi-Classifier Determinacy Maximization for Unsupervised Domain Adaptation. In *AAAI*, volume 2, 5.

Li, Y.; Kuang, Z.; Liu, L.; Chen, Y.; and Zhang, W. 2021b. Pseudo-mask matters in weakly-supervised semantic segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 6964–6973.

Lian, Q.; Lv, F.; Duan, L.; and Gong, B. 2019. Constructing self-motivated pyramid curriculums for cross-domain semantic segmentation: A non-adversarial approach. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 6758–6767.

Lin, D.; Dai, J.; Jia, J.; He, K.; and Sun, J. 2016. Scribblesup: Scribble-supervised convolutional networks for semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 3159–3167.

Lin, T.-Y.; Maire, M.; Belongie, S.; Hays, J.; Perona, P.; Ramanan, D.; Dollár, P.; and Zitnick, C. L. 2014. Microsoft coco: Common objects in context. In *European conference on computer vision*, 740–755. Springer.

Luo, Y.; Liu, P.; Guan, T.; Yu, J.; and Yang, Y. 2019a. Significance-aware information bottleneck for domain adaptive semantic segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 6778–6787.

Luo, Y.; Zheng, L.; Guan, T.; Yu, J.; and Yang, Y. 2019b. Taking a closer look at domain shift: Category-level adversaries for semantics consistent domain adaptation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2507–2516.

Lv, F.; Liang, T.; Chen, X.; and Lin, G. 2020. Cross-domain semantic segmentation via domain-invariant interactive relation transfer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 4334–4343.

Mei, K.; Zhu, C.; Zou, J.; and Zhang, S. 2020. Instance adaptive self-training for unsupervised domain adaptation. In *European conference on computer vision*, 415–430. Springer.

Melas-Kyriazi, L.; and Manrai, A. K. 2021. Pixmatch: Unsupervised domain adaptation via pixelwise consistency training. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 12435–12445.

Nam, H.; Lee, H.; Park, J.; Yoon, W.; and Yoo, D. 2019. Reducing domain gap via style-agnostic networks. *arXiv preprint arXiv:1910.11645*, 2(7): 8.

Pan, F.; Shin, I.; Rameau, F.; Lee, S.; and Kweon, I. S. 2020. Unsupervised intra-domain adaptation for semantic segmentation through self-supervision. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 3764–3773.

Qian, R.; Wei, Y.; Shi, H.; Li, J.; Liu, J.; and Huang, T. 2019. Weakly supervised scene parsing with point-based distance metric learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, 8843–8850.

Redmon, J.; Divvala, S.; Girshick, R.; and Farhadi, A. 2016. You only look once: Unified, real-time object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 779–788.

Redmon, J.; and Farhadi, A. 2017. YOLO9000: better, faster, stronger. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 7263–7271.

Redmon, J.; and Farhadi, A. 2018. Yolov3: An incremental improvement. *arXiv preprint arXiv:1804.02767*.

Ren, S.; He, K.; Girshick, R.; and Sun, J. 2015. Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks. In *Advances in Neural Information Processing Systems (NIPS)*.

Richter, S. R.; Vineet, V.; Roth, S.; and Koltun, V. 2016. Playing for data: Ground truth from computer games. In *European conference on computer vision*, 102–118. Springer.

Ros, G.; Sellart, L.; Materzynska, J.; Vazquez, D.; and Lopez, A. M. 2016. The synthia dataset: A large collection of synthetic images for semantic segmentation of urban scenes. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 3234–3243.

Rother, C.; Kolmogorov, V.; and Blake, A. 2004. Interactive foreground extraction using iterated graph cuts. *ACM Transactions on Graphics*, (23): 3.

Song, C.; Huang, Y.; Ouyang, W.; and Wang, L. 2019. Box-driven class-wise region masking and filling rate guided loss for weakly supervised semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 3136–3145.

Stacke, K.; Eilertsen, G.; Unger, J.; and Lundström, C. 2020. Measuring domain shift for deep learning in histopathology. *IEEE journal of biomedical and health informatics*, 25(2): 325–336.

Tang, M.; Djelouah, A.; Perazzi, F.; Boykov, Y.; and Schroers, C. 2018a. Normalized cut loss for weakly-supervised cnn segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 1818–1827.

Tang, M.; Perazzi, F.; Djelouah, A.; Ben Ayed, I.; Schroers, C.; and Boykov, Y. 2018b. On regularized losses for weakly-supervised cnn segmentation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 507–522.

Tranheden, W.; Olsson, V.; Pinto, J.; and Svensson, L. 2021. Dacs: Domain adaptation via cross-domain mixed sampling. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 1379–1389.

Tsai, Y.-H.; Hung, W.-C.; Schulter, S.; Sohn, K.; Yang, M.-H.; and Chandraker, M. 2018. Learning to adapt structured output space for semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 7472–7481.

Tsai, Y.-H.; Sohn, K.; Schulter, S.; and Chandraker, M. 2019. Domain adaptation for structured output via discriminative patch representations. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 1456–1465.

Vernaza, P.; and Chandraker, M. 2017. Learning random-walk label propagation for weakly-supervised semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 7158–7166.

Vu, T.-H.; Jain, H.; Bucher, M.; Cord, M.; and Pérez, P. 2019. Advent: Adversarial entropy minimization for domain adaptation in semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2517–2526.

Wang, H.; Shen, T.; Zhang, W.; Duan, L.-Y.; and Mei, T. 2020a. Classes matter: A fine-grained adversarial approach to cross-domain semantic segmentation. In *European conference on computer vision*, 642–659. Springer.

Wang, Y.; Peng, J.; and Zhang, Z. 2021. Uncertainty-aware pseudo label refinery for domain adaptive semantic segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 9092–9101.

Wang, Y.; Zhang, J.; Kan, M.; Shan, S.; and Chen, X. 2020b. Self-supervised equivariant attention mechanism for

weakly supervised semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 12275–12284.

Wang, Z.; Yu, M.; Wei, Y.; Feris, R.; Xiong, J.; Hwu, W.-m.; Huang, T. S.; and Shi, H. 2020c. Differential treatment for stuff and things: A simple unsupervised domain adaptation method for semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 12635–12644.

Yang, J.; An, W.; Wang, S.; Zhu, X.; Yan, C.; and Huang, J. 2020a. Label-driven reconstruction for domain adaptation in semantic segmentation. In *European conference on computer vision*, 480–498. Springer.

Yang, J.; Xu, R.; Li, R.; Qi, X.; Shen, X.; Li, G.; and Lin, L. 2020b. An adversarial perturbation oriented domain adaptation approach for semantic segmentation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, 12613–12620.

Yang, Y.; and Soatto, S. 2020. Fda: Fourier domain adaptation for semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 4085–4095.

Yu, F.; Chen, H.; Wang, X.; Xian, W.; Chen, Y.; Liu, F.; Madhavan, V.; and Darrell, T. 2020. Bdd100k: A diverse driving dataset for heterogeneous multitask learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2636–2645.

Yu, F.; Zhang, M.; Dong, H.; Hu, S.; Dong, B.; and Zhang, L. 2021. Dast: Unsupervised domain adaptation in semantic segmentation based on discriminator attention and self-training. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, 10754–10762.

Zhang, Q.; Zhang, J.; Liu, W.; and Tao, D. 2019. Category anchor-guided unsupervised domain adaptation for semantic segmentation. *Advances in Neural Information Processing Systems*, 32.

Zhou, B.; Khosla, A.; Lapedriza, A.; Oliva, A.; and Torralba, A. 2016. Learning deep features for discriminative localization. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2921–2929.

Zhu, J.-Y.; Park, T.; Isola, P.; and Efros, A. A. 2017. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *Proceedings of the IEEE international conference on computer vision*, 2223–2232.

Zou, Y.; Yu, Z.; Kumar, B. V.; and Wang, J. 2018. Unsupervised Domain Adaptation for Semantic Segmentation via Class-Balanced Self-Training. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 289–305.

Zou, Y.; Yu, Z.; Liu, X.; Kumar, B. V.; and Wang, J. 2019. Confidence Regularized Self-Training. In *The IEEE International Conference on Computer Vision (ICCV)*.