# Uncertainty-Aware Consistency Regularization for Cross-Domain Semantic Segmentation

Qianyu Zhou[a,2], Zhengyang Feng[a,2], Qiqi Gu[a], Guangliang Cheng[b], Xuequan Lu[c,**], Jianping Shi[b], Lizhuang Ma[a,**]

[a]*Shanghai Jiao Tong University, 800 Dongchuan Road, Shanghai and 200240, China*
[b]*SenseTime Research, 1900 Hongmei Road, Shanghai and 200233, China*
[c]*Deakin University, 75 Pigdons Rd, Waurn Ponds, VIC 3216, Australia*

arXiv:2004.08878v4 [cs.CV] 19 Aug 2021

## ABSTRACT

Unsupervised domain adaptation (UDA) aims to adapt existing models of the source domain to a new target domain with only unlabeled data. Most existing methods suffer from noticeable negative transfer resulting from either the error-prone discriminator network or the unreasonable teacher model. Besides, the local regional consistency in UDA has been largely neglected, and only extracting the global-level pattern information is not powerful enough for feature alignment due to the abuse use of contexts. To this end, we propose an uncertainty-aware consistency regularization method for cross–domain semantic segmentation. Firstly, we introduce an uncertainty-guided consistency loss with a dynamic weighting scheme by exploiting the latent uncertainty information of the target samples. As such, more meaningful and reliable knowledge from the teacher model can be transferred to the student model. We further reveal the reason why the current consistency regularization is often unstable in minimizing the domain discrepancy. Besides, we design a ClassDrop mask generation algorithm to produce strong class-wise perturbations. Guided by this mask, we propose a ClassOut strategy to realize effective regional consistency in a fine-grained manner. Experiments demonstrate that our method outperforms the state-of-the-art methods on four domain adaptation benchmarks, *i.e.,* GTAV → Cityscapes and SYNTHIA → Cityscapes, Virtual KITTI ⟶ KITTI and Cityscapes ⟶ KITTI.

## 1. Introduction

Semantic segmentation aims to identify the semantic category of each pixel in a given image. Recent studies have shown rapid progress with a variety of CNN-based algorithms trained on a large-scale annotated dataset to tackle this problem [47, 3, 4, 43]. However, due to the time-consuming process of annotating pixel-wise labels [10], building such a large annotated dataset is cost-expensive. Compared with manual annotation, the label of synthetic data is much easier to obtain, and thus it is natural to use synthetic data to supervise the segmentation model instead of real data [58, 59]. However, there always exists a significant performance drop when the learned source models are directly applied to target data, due to the existence of a domain gap between the synthetic images and real images.

To address this issue, various unsupervised domain adaptation (UDA) techniques have been proposed from the domain distribution shift perspective to align the latent feature distributions between the source domain and target domain. Many researchers have exploited additional supervised signals based on the adversarial framework such as depth [40, 5, 68], style [83, 30, 87], category constraint [31, 8], decision boundary [60, 39] and other domain-invariant information [50] to promote the feature alignment. However, due to the fact that it always requires a domain classifier (discriminator) during the training procedure, these adversarial-based approaches often suffer from training instability and the phenomenon of negative transfer [52, 9].

Consistency regularization is one of the non-adversarial methods exploited in cross-domain segmentation to cope with the negative effect caused by adversarial training [9, 73]. This kind of consistency-based methods usually perform the feature-level domain alignment between a student model and a teacher model. The teacher model is an exponential moving average (EMA) of the student model, and then the teacher model could
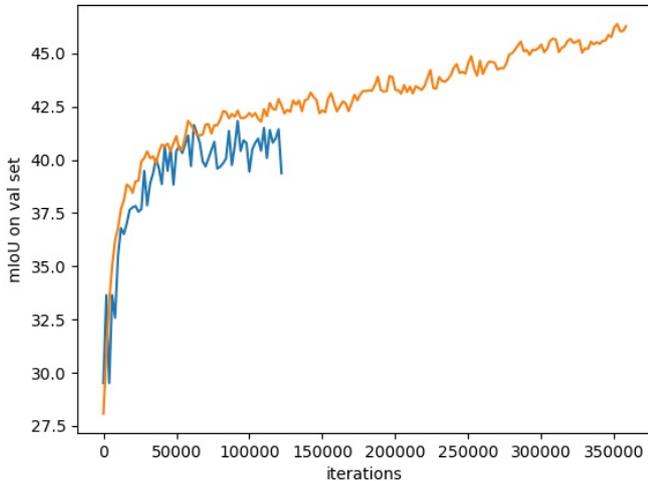
---

Fig. 1: mIoU comparison on the validation set of Cityscapes by adapting from GTA5 dataset to Cityscapes dataset. The blue line corresponds to the conventional Mean Teacher strategy [9]. The orange line corresponds to the consistency-based adaptation combined with our proposed uncertainty guided module.

transfer the learned knowledge to the student. The target predictions of the student and teacher model under different perturbations are penalized by a consistency constraint.

In the previous consistency-based works [9, 57], a common consistency loss, *i.e.,* Mean Square Error, is used to ensure the consistency between the student's prediction and the teacher's prediction. We observe that such a simple consistency constraint is usually weak for domain adaptive semantic segmentation, which is reflected in two respects. Firstly, this kind of alignment did not consider the reliability of the teacher predictions, and not all pixel-wise predictions are highly confident for knowledge transfer. Directly imposing a consistency constraint onto all pixels is inappropriate, which could harm the learning process by generating unreasonable guidance for the student model. Secondly, although the whole training of consistency-based adaptation is more stable than adversarial-based adaptation, it is still insufficient. Due to the fact that the basic Mean-Teacher structure may trigger the "error accumulation", it could take more training iterations to converge and even may lead to early performance degradation during the adaptation process. The performance curve on the target domain images is shown in Fig. 1.

In the existing consistency regularization methods, *e.g.,* [9, 57], the inconsistent penalty is usually adopted on the global level for prediction map, while the region-wise consistency on the local level is ignored, *i.e.,* some contextual object occurrence should be consistent wherever the environments are. Only extracting the global-level pattern information is not powerful enough for the feature-level representation alignment. Without this alignment, the performance will drop significantly in the target domain. We attempt to learn the robust representations to varying environments by exploring the fine-grained regional consistency, to prevent the model from abusing the contexts.

Motivated by the above facts, we propose a novel uncertainty-aware consistency regularization scheme to address the domain shift for cross-domain segmentation. Firstly, we introduce a dynamic weighting scheme with an uncertainty-guided consistency loss to capture the understanding of hidden epistemic uncertainty of target predictions for UDA in semantic segmentation. Secondly, we design a ClassDrop mask generation algorithm to produce strong class-wise perturbations. Guided by this mask, we present an innovative ClassOut strategy to keep the local regional consistency in a fine-grained manner. The whole architecture includes a student model, a teacher model, and our proposed uncertainty module.

In detail, our uncertainty-guided consistency constraints are imposed between the Mean-Teacher system and our proposed uncertainty module, which motivates both the student model and teacher model to alternately promote each other by providing positive feedback, thus leading to the domain gap to be gradually reduced. To cope with the instability of the conventional consistency regularization framework, we introduce a dynamic weighting scheme of the consistency loss, which is to calculate a time-dependent threshold for filtering out the unreasonable predictions along with mining the highly confident pixel-wise predictions of the target sample. In this manner, the adaptation is realized in a more accurate direction, instead of the rough distribution matching. To address the issue of local regional consistency in UDA, we propose a ClassOut strategy to learn more robust region-wise features under varying environments. Our main idea is that the same input image should be invariant under the perturbations by randomly dropping some categories. We design a ClassDrop mask generation algorithm to generate such strong class-wise perturbations. This mask is utilized to filter out the regions of the input target image and the uncertainty mask at the same time to ensure regional consistency on the local level.

Our main contributions are summarized as follows.

- We propose an uncertainty-aware consistency regularization framework for cross-domain semantic segmentation, which is a practical, intuitive and elegant contribution to the field. It is also a simple yet effective method for UDA in semantic segmentation.

- We design an uncertainty-guided consistency loss with a dynamic time-dependent weighting scheme and further reveal the reason why the current consistency regularization is often unstable in minimizing the domain discrepancy. We also show that our method can effectively ease this issue by mining the most reliable and meaningful samples between the source and the target domains.

- We develop a ClassOut strategy for keeping the local regional consistency in UDA. Meanwhile, we propose a ClassDrop mask generation algorithm to produce strong class-wise perturbations for guiding the ClassOut.

- We provide extensive experimental results with two common backbone networks, *i.e.,* VGG16 and ResNet101 and show that our approach achieves outstanding performance on four challenging benchmark datasets including both the synthetic-to-real adaptation and cross-city adaptation, *i.e.,*

GTAV $\longrightarrow$ Cityscapes, SYNTHIA $\longrightarrow$ Cityscapes, Virtual KITTI $\longrightarrow$ KITTI and Cityscapes $\longrightarrow$ KITTI.

## 2. Related Work

### 2.1. Semantic Segmentation

Semantic segmentation is a highly active research field in computer vision. Traditional works of semantic segmentation mainly focused on manually designed image features. With the recent surge of deep learning, a lot of CNN-based methods have been studied and we have witnessed a rapid boost in semantic segmentation performance. Long *et al.* [47] firstly formulated semantic segmentation as a per-pixel classification problem and proposed a fully convolutional network (FCN). With modifications for pixel-wise prediction, many recent approaches have been proposed, such as DeepLab v2 [3], DeepLab v3+ [4], EMANet [43], *etc.* Such models are generally trained on datasets with pixel-wise annotation, *e.g.,* Cityscapes [10], PASCAL [14] and COCO [46]. However, building such large-scale datasets with dense annotations costs expensive human labor. An alternative approach is to train a model on synthetic data generated from virtual 3D environments, for example, GTAV [58], SYNTHIA [59], *etc.* Unfortunately, when directly applying the model trained on the synthetic data to the real-world scenarios, the performance will be noticeably degraded. The main reason lies in the large domain gap or distribution shift between the source domain and target domains.

### 2.2. Domain Adaptation

In conventional machine learning, there holds a basic assumption that the training data and testing data are sampled independently from an identical distribution (*i.i.d*), while this assumption does not always hold in real-world scenarios. Domain Adaptation aims to mitigate the performance drop caused by the distribution mismatch between training and testing data when applying the trained model into the testing data. Unsupervised Domain Adaptation (UDA) refers to the setting when the labeled target data is not available. This question has been well studied in image classification. Please refer to [11] for a comprehensive survey. Conventional methods aim to learn domain-invariant representations through Maximum Mean Discrepancy (MMD) [23, 48, 1, 63], geodesic flow kernel [24], sub-space alignment [16], asymmetric metric learning[35]. Inspired by GAN [26], adversarial learning is successfully applied in UDA to align the feature distributions from different domains. DANN [22] was the pioneering work, it encouraged a generator to enforce the two distributions to be as close as possible, and to fool the domain classifier at the same time. Most of these UDA methods work on simple and small classification datasets (e.g., MNIST [38] and SVHN [54]), and may have limited performance in more challenging tasks, like semantic segmentation.

### 2.3. Domain Adaptation for Semantic Segmentation

Recently many approaches have been proposed to address the domain shift in semantic segmentation. Pioneered by [31], Hoffman *et al.* proposed a domain-adversarial training method by aligning the features between two domains. Following this line, many works have been introduced to address the cross-domain semantic segmentation via the adversarial-based methods, which have achieved great successes in this field. This kind of distribution alignment could be performed at different representation layer, such as pixel-level [30, 61, 7, 25, 77, 78], feature level [8, 50, 31, 82, 6, 2, 55] and output level [66, 67, 5, 52, 76, 71]. Many researchers have exploited additional supervised signal based on the adversarial framework such as depth [40, 5, 68], style [83, 30, 87], category constraint [31, 8], decision boundary [60, 39], and other domain-invariant information to promote the feature alignment. Despite their efforts, these approaches need to maintain an extra discriminator network, thus suffering from training instability and negative transfer [52, 9].

To tackle these issues, another line of non-adversarial methods, *e.g.,* self-training [88, 89, 44, 15] have been recently studied and applied in the field of UDA. However, these methods need to generate pseudo labels and fine-tune the segmentation model iteratively in many stages, they cannot be trained end-to-end. Different from the above self-training approaches, consistency-based methods [9, 57] is a completely different way and a simple online method to learn domain-invariant information in an end-to-end manner.

### 2.4. Consistency Regularization

Consistency Regularization is applied in the field of semi-supervised learning, which employs unlabeled data to produce consistent predictions under different perturbations [64]. Tarvainen *et al.* [64] firstly encouraged consistency between the predictions of a student network and a teacher network. The teacher's weights are an exponential moving average of those of the student, leading to faster convergence and improved results. French *et al.* [18] then applied the Mean-Teacher framework to the unsupervised domain adaptation for image classification. To address the domain shift for magnetic resonance imaging (MRI), Perone *et al.* [57] applied the self-ensembling method to the medical imaging segmentation task. Considering the UDA task for urban scenes, Choi *et al.* [9] proposed a self-ensembling with the GAN-based data augmentation method for cross-domain segmentation. Our work is mostly related to [9]. Inspire by the work [80] designed for semi-supervised 3D left atrium segmentation, we propose to capture the latent uncertainty understanding of the teacher model, and encourage the student model to learn from that reliable knowledge.

## 3. Methodology

In this section, we present our uncertainty-aware consistency regularization method for unsupervised domain adaptive segmentation. Following the unsupervised domain adaptation protocol [8, 31, 6], the synthetic data is utilized as the source domain $S$, and the real data as target domain $T$. In the source domain, we have access to the synthetic images $x_s \in S$ along with their corresponding ground-truth labels $y_s$. In the target domain, only unlabeled images $x_t \in T$ are available.
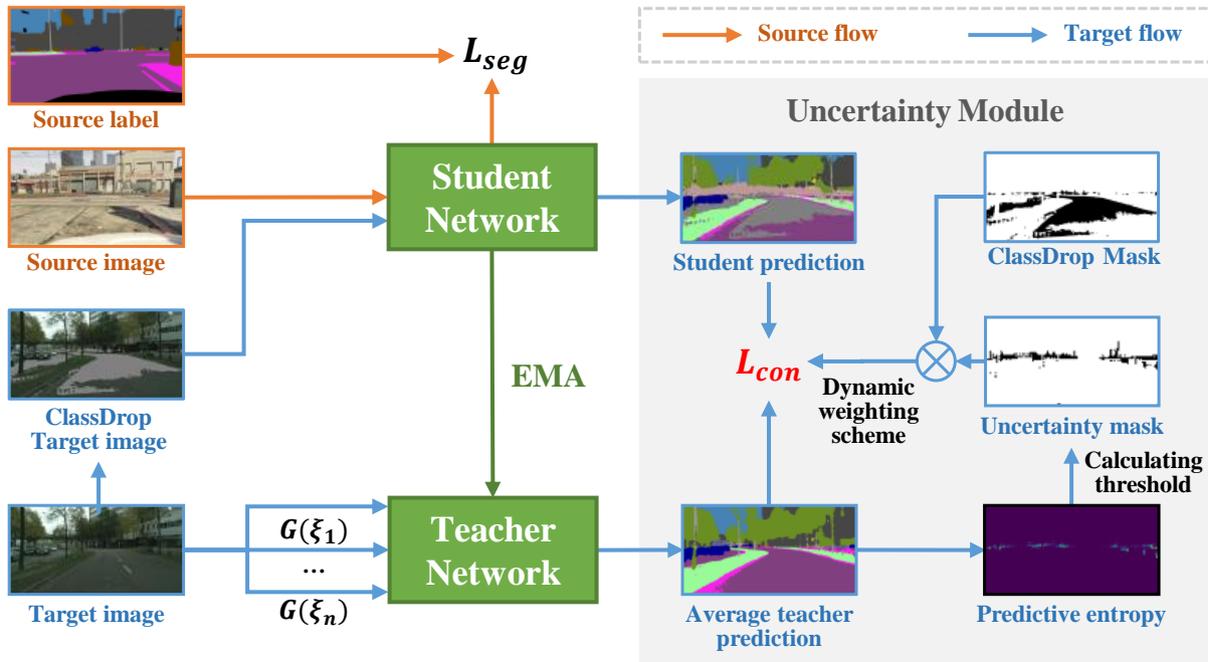
Fig. 2: An overview of the proposed framework. The whole framework includes a student network, a teacher network updated by exponential moving average (EMA), and our uncertainty module. A ClassDrop Mask is generated by the target image and then used to filter the local regions of the target image. The ClassDrop target image is fed into the student network to get student prediction. We employ different augmentations $G(\xi_i)$ for the input target sample, and they are fed into the teacher model. In our uncertainty module, we perform $N$ times stochastic forward passes to get an average teacher prediction. Then, with the estimation of predictive entropy and the proposed dynamic threshold, we could get the uncertainty mask. Thus, the ClassDrop Mask is element-wise multiplied with the uncertainty mask for filtering out the unreasonable predictions. Guided by the proposed dynamic weighting scheme and ClassOut strategy, our uncertainty-guided consistency loss $L_{con}$ could encourage the teacher model to transfer more reliable knowledge to the student.

## 3.1. Overview

The overview of our proposed uncertainty-aware consistency regularization method is illustrated in Fig. 2. The whole framework includes three modules: a student model $f_S$, a teacher model $f_T$, and our uncertainty module. The key idea is to decrease the uncertainty of the error-prone teacher model as training progress thus leading the adaptation process in a more accurate and stable way.

Specifically, a ClassDrop mask is generated by the target image to provide strong class-wise perturbations by randomly dropping some classes that are presented in the target image. This mask will be utilized to filter out the local regions of both the target image and the uncertainty mask (a mask we defined to indicate the uncertain pixels). For the former, a ClassDrop target image is fed into the student network to get student prediction. For the latter, we will explain the data flow in detail. Firstly, we employ data augmentation, *e.g.,* Gaussian Noise, for the input target samples. In our proposed uncertainty module, we perform stochastic forward passes to calculate the mean of target predictions. In this way, we are able to employ our teacher model as a Bayesian network to estimate the latent uncertainty information of the teacher predictions. We formulate the uncertainty as the pixel-wise predictive entropy. Then, we calculate a time-dependent threshold for filtering out those unreasonable predictions along with mining the high confident pixel-wise predictions of the target sample. Thus, the ClassDrop Mask is element-wise multiplied with the uncertainty mask for filtering out the unreasonable predictions.



Fig. 3: Previous work vs. our approach.

With the help of the proposed dynamic weighting scheme, an uncertainty-guided consistency loss is penalized to target predictions under different perturbations, which could lead the student model to gradually learn from the more meaningful and reliable predictions of the teacher model during the training process.

## 3.2. Uncertainty Module

As shown in Fig. 3, the uncertainty module serves as a bridge for connecting the teacher model and the student model. According to the uncertainty estimation method in Bayesian networks [33], we are motivated to capture the understanding of epistemic uncertainty using stochastic forward passes. In step 1, we perform the stochastic forward pass and then extract the uncertainty information of the error-prone teacher model. We

Fig. 4: The main idea of the ClassOut strategy is that the same input image should be invariant under the ClassDrop perturbations. Specifically, we firs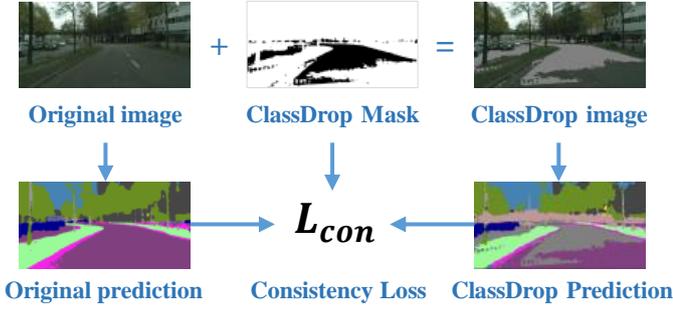tly generate a ClassDrop mask from an original target image. Then, guided by this mask, we calculate a consistency loss between the original prediction from the teacher and the ClassDrop prediction from the student. Therefore, we can keep the local regional consistency on a fine-grained level.

formulate the uncertainty as the pixel-wise predictive entropy. In step 2, we calculate the uncertainty mask given our time-dependent threshold, and a ClassDrop mask given the target image. Guided by these masks, we enforce an uncertain-aware consistency loss and a ClassOut Strategy onto the student predictions and the teacher predictions, thus the student model could learn credible knowledge from the teacher.

The teacher's weights $\Phi'_t$ at training step $t$ are updated by the student's weights $\Phi_t$ with a smoothing coefficient $\alpha \in [0, 1]$, which can be formulated as follows:

$$\Phi'_t = \alpha \cdot \Phi'_{t-1} + (1 - \alpha) \cdot \Phi_t, \tag{1}$$

where $\alpha$ refers to the EMA decay that controls the updating rate.

Specifically, we make $N$ copies of the target image and inject a Gaussian noise for the target predictions following prior works [9]. Then, we perform $N$ stochastic forward passes for the target teacher sample to get the average teacher prediction. Given a set of pixel-wise predicted class scores $\{\boldsymbol{P}_i^{(h,w,c)}(x_t)\}_{i=1}^N$ of the target samples, the average teacher prediction is formulated as:

$$\hat{P}_c = \frac{1}{N} \sum_{i=1}^N \boldsymbol{P}_i^{(h,w,c)}(x_t), \tag{2}$$

where $\hat{P}_c$ denotes the mean of the predictive probability of the $c$-th class after $N$ times stochastic forward passes. Thus, the pixel-wise predictive entropy is as follows:

$$\mu^{(h,w)} = -\sum_{c=1}^C \hat{P}_c \cdot log(\hat{P}_c), \tag{3}$$

where $\zeta$ refers to the predictive entropy in pixel level. All the volumes of each pixel's uncertainty forms a set $Z = \{\zeta\}_{i=1}^N$.

### 3.3. Dynamic Weighting Scheme

With the help of the uncertainty of each pixel, we could calculate a dynamic threshold to filter out the unreliable pixel-wise prediction. On top of that, certain pixels with high confident

probabilities will be left and the student model could gradually learn the reliable target predictions from the teacher model.

In particular, we first calculate the uncertainty threshold $R$ to select the confident pixels according to the uncertainty map we have estimated. Inspired by the ramp-up function of consistency weight in [9], we came up with the Eq. 4, which dynamically increases the threshold while the uncertainty is decreased during the training process. This design is a time-dependent ramp-up function, which changes dynamically over time:

$$R = \alpha + (1 - \alpha) \cdot e^{\beta(1 - t/t_{max})^2} \cdot Z_{sup} \tag{4}$$

where $t$ denotes the current training step and $t_{max}$ is the maximum training step. $Z_{sup}$ means the upper bound of the volumes' self-information, which is denoted by $Z_{sup} = sup\{\mu\}_{i=1}^N$. And $\alpha$ and $\beta$ are two hyper-parameters.

The uncertainty-aware consistency loss $L_{con}$ is imposed between the prediction maps extracted from the student and the predictions from the teacher network.

$$L_{con}(f_S, f_T) = \sum_{h=1}^H \sum_{w=1}^W I(\mu^{(h,w)} < R) \\ \cdot \|f_{\theta(x_{T_1})^{(h,w)}} - f_{\theta'}(x_{T_2})^{(h,w)}\|^2, \tag{5}$$

where $I$ is an indicator function, and $x_{T_1}$ and $x_{T_2}$ are two input target samples with different augmentations. $f_{\theta(x_T)}$ and $f_{\theta'}(x_T)$ are the student and teacher prediction map after the softmax function, respectively. Note that the prediction map $f_{\theta'}(x_{T_2})$ used for consistency regularization is the stochastic one rather than the average one. Our uncertainty mask $M_{uncertainty} = I(\mu^{(h,w)} < R)$ can reweigh not only the Mean Squared Error (MSE) loss but also the Cross-Entropy Loss. For simplicity, we use the Mean Squared Error (MSE) in this paper.

### 3.4. ClassOut Strategy

Previous consistency regularization methods, e.g., [9, 57], usually impose the inconsistent penalty on the global level for prediction map, while the region-wise consistency on the local level is largely ignored, i.e., some contextual object occurrence should be consistent whatever the environments are. Only extracting the global-level pattern information is not powerful enough for the feature-level representation alignment. Due to the lack of local regional consistency, the performance will drop significantly in the target domain. Our goal is to learn the robust representations to varying environments by exploring the fine-grained regional consistency, to prevent the model from abusing the contexts.

Firstly, we propose an innovative ClassDrop mask generation algorithm to provide strong class-wise perturbations, as shown in Algorithm 1. To be specific, we firstly get the pseudo labels $\tilde{Y}_T$ from the target predictions $f_{\theta'}(X_T)$. The set of the classes presented in $\tilde{Y}_T$ are noted as $C$. We get a class ratio $\delta$ sampled from a uniform distribution. Then, we randomly select $\delta|C|$ classes in $C$. A binary mask M is generated by setting the pixels from those classes to 1 in M, whereas all others will have a value 0. This mask is utilized to filter out the local regions in both the target image and the uncertainty mask by an element-wise multiplication.

---

**Algorithm 1:** ClassDrop Mask Generation Algorithm

   **Input:** teacher model $f_{\theta'}$, target image $X_T$, min class ratio $a$, and max class ratio $b$.
   **Output:** ClassDrop mask $M$

1   $\hat{f}_{\theta'} \leftarrow f_{\theta'}(X_T)$;
2   $\tilde{Y}_T \leftarrow \arg\max_{c'} \ \hat{f}_{\theta'}(i, j, c')$;
3   $C \leftarrow$ Set of the classes present in $\tilde{Y}_T$;
4   $\delta \leftarrow U(a, b)$ ;
5   $c \leftarrow$ Randomly select $\delta|C|$ classes in $C$;
6   **for** *each i, j* **do**
7      $M(i, j) = \begin{cases} 1, & \text{if } \tilde{Y}_T(i, j) \in c \\ 0, & \text{otherwise} \end{cases}$
8   **return** $M$;

---

The main idea of the ClassOut strategy is that the same input image should be invariant under the ClassDrop perturbations. Thus, guided by the ClassDrop mask, we calculate a consistency loss between the original prediction from the teacher and the ClassDrop prediction from the student. Therefore, we can keep the local regional consistency:

$$L_{con} = \|M \odot (f_{\theta'}(M \odot X_T) - f_{\theta'}(X_T))\|^2, \qquad (6)$$

### 3.5. Unified Training

**Consistency Loss:** By unifying the ClassOut strategy and the dynamic weighting scheme into the same framework to realize end-to-end training, we reformulate the consistency loss as follows:

$$L_{con} = \left\| M_{classdrop} \odot M_{uncertainty} \odot (f_\theta(M_{classout} \odot X_T) - f_{\theta'}(X_T)) \right\|^2, \qquad (7)$$

where the final consistency loss is reweighted by the uncertainty mask $M_{uncertainty}$ (defined in section 3.3) and the ClassDrop mask $M_{classout}$ (defined in section 3.4). In other words, the reweighted mask of the consistency loss is the element-wise multiplication between the uncertainty mask $M_{uncertainty}$ and the ClassDrop mask $M_{classout}$. We simplify the previous definition and reformulate it as follows:

$$M_{classdrop}(i, j) = \begin{cases} 1, & \text{if } \tilde{Y}_T(i, j) \in c_{remain} \\ 0, & \text{otherwise} \end{cases} \qquad (8)$$

$$M_{uncertainty}(i, j) = \begin{cases} 1, & \text{if } \mu(i, j) < R \\ 0, & \text{otherwise} \end{cases} \qquad (9)$$

where $c_{remain} = \delta|C|$ is the selected classes from a class set, $\mu$ is the predicted entropy defined in Eq. 3 and $R$ is the dynamic threshold defined in Eq. 4.

**Supervised Loss:** The segmentation loss $L_{seg}$ is the cross-entropy loss for optimizing the images from the source domain, which can be defined as:

$$L_{seg} = -\sum_{h=1}^{H} \sum_{w=1}^{W} \sum_{c=1}^{C} y_s^{(h,w,c)} log(P_s^{(h,w,c)}), \qquad (10)$$

where $y_s$ is the ground truth for source images and $P_s = f_S((\hat{x}_s)^{(h,w,c)}$ is the segmentation output of source-translated input images.

**Total Loss:** The total loss $L_{total}$ is the weighted sum of the segmentation loss $L_{seg}$ and the consistency loss $L_{con}$, and can be written as:

$$L_{total} = L_{seg} + \lambda_{con} L_{con}, \qquad (11)$$

where $L_{con}$ is the combination of Equation 5 and Equation 7. $\lambda_{con}$ is the dynamic weight of the consistency loss. To balance the segmentation loss and the consistency loss, we use the same ramp-up function $\lambda_{con}$ as [65]. It is to increase the dominance of $L_{seg}$ during the early training steps and to increase the dominance of $L_{con}$ during the late training steps.

### 3.6. Discussion

In this subsection, we will discuss the main differences between the existing research and our proposed method.

There exist some works [72, 36] which used uncertainty estimation in domain adaptation; however, those methods [72, 36] always need to maintain a Bayesian Discriminator in adversarial training, thus suffering the drawbacks of negative transfer and remarkable instability of training. Besides, their methods only work well on the simple and small classification dataset, and can hardly work well in structured tasks, *e.g.,* semantic segmentation. Therefore, we do not compare the experimental results with these methods in Section 4. Our uncertainty-aware consistency regularization shows that a non-adversarial approach can achieve the state-of-the-art as well without the need of maintaining an extra discriminator network or carefully tuning the optimization procedure for min-max problems during the domain adaptation procedure.

Different from [28, 84], we focus on investigating the problem of "error accumulation" in consistency regularization, rather than self-training. In contrast to [80] that targets the semi-supervised learning for the 3d left atrium segmentation task, while we target the unsupervised domain adaptation for the image semantic segmentation task. *Our method differs from these approaches in several aspects.* Firstly, we propose a dynamic weighting scheme and a ClassOut strategy for the uncertainty-consistency loss. The uncertainty mask, Classdrop mask are employed in a completely different way from previous works [72, 36, 28, 80]. We further reveal the reason why the current consistency regularization is often unstable in minimizing the distribution discrepancy in Section 1 and Section 4.1. Besides, we also show that our method can effectively ease this issue by mining the most reliable and meaningful samples between the source and the target domains. To sum up, our uncertainty-aware consistency regularization framework is a practical, intuitive and elegant contribution to the field, and it is also a simple yet effective unsupervised domain adaptation method for semantic segmentation. To our best knowledge, there are no such domain adaptive segmentation methods published before.

# 4. Experiments

In this section, we verify the effectiveness of our method with two common backbone networks, *i.e.,* VGG16 and ResNet 101, on both the synthetic-to-real adaptation and cross-city adaptation on four challenging benchmark datasets, *i.e.,* GTAV ⟶ Cityscapes, SYNTHIA ⟶ Cityscapes, Virtual KITTI ⟶ KITTI and Cityscapes ⟶ KITTI.

## 4.1. Datasets

**Cityscapes** [10] is a dataset focused on autonomous driving, which consists of 2,975 images in the training set, and 500 images in the validation set. The images have a fixed spatial resolution of 2048 × 1024 pixels. For the sake of the fairness of experimental results, we follow the same evaluation protocol [66, 69, 52], i.e. we train the model on the unlabeled training set and report the results on the validation set.
**GTAV** [58] is a synthetic dataset including 24,966 photo-realistic images rendered by the gaming engine Grand Theft Auto V (GTAV). The resolution of images is 1914 × 1051 pixels which is similar to Cityscapes that the semantic categories are also compatible between the two datasets. We use all the 19 official training classes in our experiments.
**SYNTHIA** [59] is another synthetic dataset composed of 9,400 annotated synthetic images with the resolution 1280 × 960. Like GTAV, it has semantically compatible annotations with Cityscapes. Following the prior works [8, 82, 6], we use the SYNTHIA-RAND-CITYSCAPES subset [59] as our training set.
**KITTI** [20] is a real-world dataset containing 7,481 images with bounding boxes and another 200 images with pixel-level labels. In the detection task, we split the training set and the validation set manually with a ratio of 9 : 1 following [27]. In the segmentation task, it is used as the target domain only due to the lack of pixel-level annotations.
**Virtual KITTI** [21] is a synthetic dataset which clones the scenes from the KITTI with 21,260 images. Each image is densely annotated at pixel level with category and depth information. It is designed to mimic the conditions of KITTI dataset and has similar scene layouts, camera viewpoints and image resolution to KITTI dataset.

## 4.2. Implementation details

Following common UDA protocols [9, 73, 66, 52], we employ the VGG-16 [62] and ResNet 101 [29] as the backbone of the DeepLab-v2 [3] in our implementations, and the backbone model is pre-trained on ImageNet [12]. For the DeepLab-v2 network, we use Adam as the optimizer. The initial learning rate is $1 \times 10^{-5}$, and the weight decay is $5 \times 10^{-5}$. In our uncertainty module, we perform $N = 8$ times stochastic forward passes to capture the understanding of latent epistemic uncertainty. We set the EMA decay $\alpha$ to 0.999 during the training process. Following [37, 64, 9],the consistency weight is a ramp-up function: $\lambda_{con} = \lambda_0 \times e^{-5(1-t/t_{max})^2}$, where $\lambda_0$ is an initial constant. This time-dependent threshold function is used to increase the certainty at later training steps. We set $\alpha = 0.75$ and $\beta = -5$ in all experiments. Our method is implemented in Pytorch on a single NVIDIA GTX 3090 Ti.

## 4.3. Comparisons with the State-of-the-art Techniques

We compare the results between our method and the state-of-the-art methods on four challenging benchmarks, which includes the synthetic-to-real adaptation, *i.e.,* "GTAV → Cityscapes" and "SYNTHIA → Cityscapes", "Virtual KITTI → KITTI" and cross-city adaptation, *i.e.,* "Cityscapes → KITTI". With VGG16 backbone, our proposed method significantly outperforms the state-of-the-art methods by 5% ∼ 8% on GTAV → Cityscapes, and 2% ∼ 7% on SYNTHIA → Cityscapes. Besides, it is superior to the non-adaptive baseline by 19.5% on GTA5 → Cityscapes and 20% ∼ 24% on SYNTHIA → Cityscapes. With ResNet101 backbone, our proposed method outperforms the state-of-the-art methods by 1% ∼ 3% on GTAV → Cityscapes, and 2% ∼ 6% on SYNTHIA → Cityscapes.

### 4.3.1. Results on GTAV → Cityscapes

As shown in Table 1 and Table 3, we present the adaptation results from GTAV to Cityscapes with VGG16 and ResNet 101, respectively. Source-only denotes the baseline Deeplab-v2 [3] is trained with only source domain data. In the works [82, 66, 6, 52, 76, 55, 70], they mainly focused on distribution alignment via different adversarial mechanisms. But promoting feature alignment only on the high representation level is not enough, *i.e.,* feature level [82, 6] or output level [66, 52, 71]. The best results of mIoU among them are still about 7% worse than our results. To further reduce the domain gap, Hoffman *et. al* [30] introduced an image-to-image translation model to perform a style transfer process on the low appearance level. Such techniques are further integrated into [71, 44, 61, 87, 9] to achieve higher performance, while they are still about 5% ∼ 10% worse than our results. Another line of non-adversarial methods [82, 88, 69] were proposed to address the negative effect of adversarial training. The self-ensembling with GAN-based augmentation [9] has been recently proposed and surpassed most of the previous works. In Table 1, our method could get about 5.3% improvements compared to this work [9]. Extensive experiments in Table 1 and Table 3 show that our approach achieves a new top performance.

### 4.3.2. Results on SYNTHIA → Cityscapes

As shown in Table 2 and Table 9, we list the adaptation results on the task "SYNTHIA → Cityscapes" with VGG16 and ResNet 101, respectively. Due to the fact that the baselines [66, 52] only calculate the results using 13 categories, we also list results for the 13 categories for a fair comparison. Although the domain gap between SYNTHIA and Cityscapes is much larger than that of GTAV to Cityscapes, we could observe in Table 2 that our uncertainty-aware consistency regularization still performs well in terms of both mIoU and per-class IoU. In some semantic categories, such as large objects, *e.g.,* road, building, wall, vegetation, sky, *etc.*, our method could capture the understanding of epidemic uncertainty and remarkably increase the certainty of these categories during the training procedure. In Table 2, the proposed method significantly outperforms the state-of-the-art techniques by 2.5% in mIoU16 and 2% in mIoU13 with VGG16 backbone. It is superior to the non-adaptive baseline by 18.9% in mIoU16 and 24.5% in mIoU13.

Table 1: Comparison results (mIoU) from GTAV to Cityscapes (with VGG16 backbone).

| Method | road | sidewalk | building | wall | fence | pole | light | sign | vegetation | terrain | sky | person | rider | car | truck | bus | train | motocycle | bike | **mIoU** |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Source Only | 61.0 | 18.5 | 66.2 | 18.0 | 19.6 | 19.1 | 22.4 | 15.5 | 79.6 | 28.5 | 58.0 | 44.5 | 1.7 | 66.6 | 14.1 | 1.1 | 0.0 | 3.2 | 0.7 | 28.3 |
| SIBAN [50] | 83.4 | 13.0 | 77.8 | 20.4 | 17.5 | 24.6 | 22.8 | 9.6 | 81.3 | 29.6 | 77.3 | 42.7 | 10.9 | 76.0 | 22.8 | 17.9 | 5.7 | 14.2 | 2.0 | 34.2 |
| CyDADA [30] | 85.2 | 37.2 | 76.5 | 21.8 | 15.0 | 23.8 | 22.9 | 21.5 | 80.5 | 31.3 | 60.7 | 50.5 | 9.0 | 76.9 | 17.1 | 28.2 | 4.5 | 9.8 | 0.0 | 35.4 |
| AdaptSegNet [66] | 87.3 | 29.8 | 78.6 | 21.1 | 18.2 | 22.5 | 21.5 | 11.0 | 79.7 | 29.6 | 71.3 | 46.8 | 6.5 | 80.1 | 23.0 | 26.9 | 0.0 | 10.6 | 0.3 | 35.0 |
| ROAD [6] | 85.4 | 31.2 | 78.6 | 27.9 | 22.2 | 21.9 | 23.7 | 11.4 | 80.7 | 29.3 | 68.9 | 48.5 | 14.1 | 78.0 | 19.1 | 23.8 | 9.4 | 8.3 | 0.0 | 35.9 |
| CLAN [52] | 88.0 | 30.6 | 79.2 | 23.4 | 20.5 | 26.1 | 23.0 | 14.8 | 81.6 | 34.5 | 72.0 | 45.8 | 7.9 | 80.5 | 26.6 | 29.9 | 0.0 | 10.7 | 0.0 | 36.6 |
| AdaptPatch [67] | 87.3 | 35.7 | 79.5 | 32.0 | 14.5 | 21.5 | 24.8 | 13.7 | 80.4 | 32.0 | 70.5 | 50.5 | 16.9 | 81.0 | 20.8 | 28.1 | 4.1 | 15.5 | 4.1 | 37.5 |
| APODA [76] | 88.4 | 34.2 | 77.6 | 23.7 | 18.3 | 24.8 | 24.9 | 12.4 | 80.7 | 30.4 | 68.6 | 48.9 | 17.9 | 80.8 | 27.0 | 27.2 | 6.2 | 19.1 | 10.2 | 38.0 |
| CrCDA [32] | 86.8 | 37.5 | 80.4 | 30.7 | 18.1 | 26.8 | 25.3 | 15.1 | 81.5 | 30.9 | 72.1 | 52.8 | 19.0 | 82.1 | 25.4 | 29.2 | 10.1 | 15.8 | 3.7 | 39.1 |
| SWD [39] | 91.0 | 35.7 | 78.0 | 21.6 | 21.7 | 31.8 | 30.2 | 25.2 | 80.2 | 23.9 | 74.1 | 53.1 | 15.8 | 79.3 | 22.1 | 26.5 | 1.5 | 17.2 | 30.4 | 39.9 |
| DCAN [42] | 82.3 | 26.7 | 77.4 | 23.7 | 20.5 | 20.4 | 30.3 | 15.9 | 80.9 | 25.4 | 69.5 | 52.6 | 11.1 | 79.6 | 24.9 | 21.2 | 1.3 | 17.0 | 6.7 | 36.2 |
| CrDoCo [7] | 89.1 | 33.2 | 80.1 | 26.9 | 25.0 | 18.3 | 23.4 | 12.8 | 77.0 | 29.1 | 72.4 | 55.1 | 20.2 | 79.9 | 22.3 | 19.5 | 1.0 | 20.1 | 18.7 | 38.1 |
| CDA [82] | 72.9 | 30.0 | 74.9 | 12.1 | 13.2 | 15.3 | 16.8 | 14.1 | 79.3 | 14.5 | 75.5 | 35.7 | 10.0 | 62.1 | 20.6 | 19.0 | 0.0 | 19.3 | 12.0 | 31.4 |
| CBST [42] | 66.7 | 26.8 | 73.7 | 14.8 | 9.5 | 28.3 | 25.9 | 10.1 | 75.5 | 15.7 | 51.6 | 47.2 | 6.2 | 71.9 | 3.7 | 2.2 | 5.4 | 18.9 | **32.4** | 30.9 |
| ADVENT [69] | 86.8 | 28.5 | 78.1 | 27.6 | 24.2 | 20.7 | 19.3 | 8.9 | 78.8 | 29.3 | 69.0 | 47.9 | 5.9 | 79.8 | 25.9 | 34.1 | 0.0 | 11.3 | 0.3 | 35.6 |
| PyCDA [45] | 86.7 | 24.8 | 80.9 | 21.4 | **27.3** | 30.2 | 26.6 | 21.1 | 86.6 | 28.9 | 58.8 | 53.2 | 17.9 | 80.4 | 18.8 | 22.4 | 4.1 | 9.7 | 6.2 | 37.2 |
| LSD-seg [61] | 88.0 | 30.5 | 78.6 | 25.2 | 23.5 | 16.7 | 23.5 | 11.6 | 78.7 | 27.2 | 71.9 | 51.3 | 19.5 | 80.4 | 19.8 | 18.3 | 0.9 | 20.8 | 18.4 | 37.1 |
| SSF-DAN [13] | 88.7 | 32.1 | 79.5 | 29.9 | 22.0 | 23.8 | 21.7 | 10.7 | 80.8 | 29.8 | 72.5 | 49.5 | 16.1 | 82.1 | 23.2 | 18.1 | 3.5 | 24.4 | 8.1 | 37.7 |
| Conservative Loss [87] | 85.6 | 38.3 | 78.6 | 27.2 | 18.4 | 25.3 | 25.0 | 17.1 | 81.5 | 31.3 | 70.6 | 50.5 | 22.3 | 81.3 | 25.5 | 21.0 | 0.1 | 18.9 | 4.3 | 38.1 |
| PIT [27] | 86.2 | 35.0 | 82.1 | 31.1 | 22.1 | 23.2 | 29.4 | 28.5 | 79.3 | 31.8 | 81.9 | 52.1 | 23.2 | 80.4 | 29.5 | 26.9 | 30.7 | 20.5 | 1.2 | 41.8 |
| BDL [44] | 89.2 | 40.9 | 81.2 | 29.1 | 19.2 | 14.2 | 29.0 | 19.6 | 83.7 | 35.9 | 80.7 | 54.7 | 23.3 | 82.7 | 25.8 | 28.0 | 2.3 | **25.7** | 19.9 | 41.3 |
| SIM [71] | 88.1 | 35.8 | 83.1 | 25.8 | 23.9 | 29.2 | 28.8 | 28.6 | 83.0 | 36.7 | 82.3 | 53.7 | 22.8 | 82.3 | 26.4 | 38.6 | 0.0 | 19.6 | 17.1 | 42.4 |
| TGCF-DA + SE [9] | 90.2 | 51.5 | 81.1 | 15.0 | 10.7 | **37.5** | **35.2** | 28.9 | 84.1 | 32.7 | 75.9 | **62.7** | 19.9 | 82.6 | 22.9 | 28.3 | 0.0 | 23.0 | 25.4 | 42.5 |
| **Ours** | **95.1** | **66.5** | **84.7** | **35.1** | 19.8 | 31.2 | 35.0 | **32.1** | **86.2** | **43.4** | **82.5** | 61.0 | **25.1** | **87.1** | **35.3** | **46.1** | 0.0 | 24.6 | 17.5 | **47.8** |

Table 2: Comparison results (mIoU) from SYNTHIA to Cityscapes (with VGG16 backbone).

| Method | road | sidewalk | building | wall | fence | pole | light | sign | vegetation | sky | person | rider | car | bus | motocycle | bike | **mIoU** | **mIoU*** |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Source Only | 6.8 | 15.4 | 56.8 | 0.8 | 0.1 | 14.6 | 4.7 | 6.8 | 72.5 | 78.6 | 41.0 | 7.8 | 46.9 | 4.7 | 1.8 | 2.1 | 22.6 | 24.1 |
| Cross-city [8] | 62.7 | 25.6 | 78.3 | - | - | - | 1.2 | 5.4 | 81.3 | 81.0 | 37.4 | 6.4 | 63.5 | 16.1 | 1.2 | 4.6 | - | 35.7 |
| SIBAN [50] | 70.1 | 25.7 | 80.9 | - | - | - | 3.8 | 7.2 | 72.3 | 80.5 | 43.3 | 5.0 | 73.3 | 16.0 | 1.7 | 3.6 | - | 37.2 |
| ROAD [6] | 77.7 | 30.0 | 77.5 | 9.6 | 0.3 | 25.8 | 10.3 | 15.6 | 77.6 | 79.8 | 44.5 | 16.6 | 67.8 | 14.5 | 7.0 | 23.8 | 36.2 | - |
| AdaptSegNet [66] | 78.9 | 29.2 | 75.5 | - | - | - | 0.1 | 4.8 | 72.6 | 76.7 | 43.4 | 8.8 | 71.1 | 16.0 | 3.6 | 8.4 | - | 37.6 |
| CLAN [52] | 80.4 | 30.7 | 74.7 | - | - | - | 1.4 | 8.0 | 77.1 | 79.0 | 46.5 | 8.9 | 73.8 | 18.2 | 2.2 | 9.9 | - | 39.3 |
| AdaptPatch [67] | 72.6 | 29.5 | 77.2 | 3.5 | 0.4 | 21.0 | 1.4 | 7.9 | 73.3 | 79.0 | 45.7 | 14.5 | 69.4 | 19.6 | 7.4 | 16.5 | 33.7 | 39.6 |
| SPIGAN [40] | 71.1 | 29.8 | 71.4 | 3.7 | 0.3 | 33.2 | 6.4 | 15.6 | 81.2 | 78.9 | 52.7 | 13.1 | 75.9 | **25.5** | 10.0 | 20.5 | 36.8 | - |
| CrCDA [32] | 74.5 | 30.5 | 78.6 | 6.6 | 0.7 | 21.2 | 2.3 | 8.4 | 77.4 | 79.1 | 45.9 | 16.5 | 73.1 | 24.1 | 9.6 | 14.2 | 35.2 | 41.1 |
| APODA [76] | 82.9 | 31.4 | 72.1 | - | - | - | 10.4 | 9.7 | 75.0 | 76.3 | 48.5 | 15.5 | 70.3 | 11.3 | 1.2 | 29.4 | - | 41.1 |
| SWD [39] | 83.3 | 35.4 | 82.1 | - | - | - | 12.2 | 12.6 | 83.8 | 76.5 | 47.4 | 12.0 | 71.5 | 17.9 | 1.6 | 29.7 | - | 43.5 |
| CrDoCo [7] | 62.2 | 21.2 | 72.8 | 4.2 | 0.8 | 30.1 | 4.1 | 10.7 | 76.3 | 73.6 | 45.6 | 14.9 | 69.2 | 14.1 | 12.2 | 23.0 | 33.4 | - |
| DCAN [42] | 79.9 | 30.4 | 70.8 | 1.6 | 0.6 | 22.3 | 6.7 | 23.0 | 76.9 | 73.9 | 41.9 | 16.7 | 61.7 | 11.5 | **10.3** | 38.6 | 35.4 | - |
| CDA [82] | 65.2 | 26.1 | 74.9 | 0.1 | 0.5 | 10.7 | 3.7 | 3.0 | 76.1 | 70.6 | 47.1 | 8.2 | 43.2 | 20.7 | 0.7 | 13.1 | 29.0 | 34.8 |
| CBST [88] | 69.6 | 28.7 | 69.5 | **12.1** | 0.1 | 25.4 | 11.9 | 13.6 | 82.0 | **81.9** | 49.1 | 14.5 | 66.0 | 6.6 | 3.7 | 32.4 | 35.4 | 36.1 |
| ADVENT [69] | 67.9 | 29.4 | 71.9 | 6.3 | 0.3 | 19.9 | 0.6 | 2.6 | 74.9 | 74.9 | 35.4 | 9.6 | 67.8 | 21.4 | 4.1 | 15.5 | 31.4 | 36.6 |
| PyCDA [45] | 80.6 | 26.6 | 74.5 | 2.0 | 0.1 | 18.1 | **13.7** | 14.2 | 80.8 | 71.0 | 48.0 | 19.0 | 72.3 | 22.5 | 12.1 | 18.1 | 35.9 | 42.6 |
| Conservative Loss [87] | 80.0 | 31.4 | 72.9 | 0.4 | 0.0 | 22.4 | 8.1 | 16.7 | 74.8 | 72.2 | 50.9 | 12.7 | 53.9 | 15.6 | 1.7 | 33.5 | 34.2 | 40.3 |
| LSD-seg [61] | 80.1 | 29.1 | 77.5 | 2.8 | 0.4 | 26.8 | 11.1 | 18.0 | 78.1 | 76.7 | 48.2 | 15.2 | 70.5 | 17.4 | 8.7 | 16.7 | 36.1 | - |
| GIO-Ada [5] | 78.3 | 29.2 | 76.9 | 11.4 | 0.3 | 26.5 | 10.8 | 17.2 | 81.7 | **81.9** | 45.8 | 15.4 | 68.0 | 15.9 | 7.5 | 30.4 | 37.3 | 43.0 |
| SSF-DAN [13] | 87.1 | 36.5 | 79.7 | - | - | - | 13.5 | 7.8 | 81.2 | 76.7 | 50.1 | 12.7 | 78.0 | 35.0 | 4.6 | 1.6 | - | 43.4 |
| PIT [27] | 81.7 | 26.9 | 78.4 | 6.3 | 0.2 | 19.8 | 13.4 | 17.4 | 76.7 | 74.1 | 47.5 | 22.4 | 76.0 | 21.7 | 19.6 | 27.7 | 38.1 | 44.9 |
| BDL [44] | 72.0 | 30.3 | 74.5 | 0.1 | 0.3 | 24.6 | 10.2 | **25.2** | 80.5 | 80.0 | **54.7** | **23.2** | 72.7 | 24.0 | 7.5 | **44.9** | 39.0 | - |
| TGCF-DA + SE [9] | 90.1 | 48.6 | 80.7 | 2.2 | 0.2 | 27.2 | 3.2 | 14.3 | **82.1** | 78.4 | 54.4 | 16.4 | 82.5 | 12.3 | 1.7 | 21.8 | 38.5 | 46.6 |
| **Ours** | **93.1** | **53.2** | **81.1** | 2.6 | **0.6** | 29.1 | 7.8 | 15.7 | 81.7 | 81.6 | 53.6 | 20.1 | **82.7** | 22.9 | 7.7 | 31.3 | **41.5** | **48.6** |

Table 3: Comparison results (mIoU) from GTAV to Cityscapes (with ResNet 101 backbone).

| Method | Venue | road | sidewalk | building | wall | fence | pole | light | sign | vegetation | terrain | sky | person | rider | car | truck | bus | train | motocycle | bike | mIoU |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Source Only | - | 63.3 | 15.7 | 59.4 | 8.6 | 15.2 | 18.3 | 26.9 | 15.0 | 80.5 | 15.3 | 73.0 | 51.0 | 17.7 | 59.7 | 28.2 | 33.1 | 3.5 | 23.2 | 16.7 | 32.9 |
| BDL [44] | CVPR'19 | 91.0 | 44.7 | 84.2 | 34.6 | 27.6 | 30.2 | 36.0 | 36.0 | 85.0 | 43.6 | 83.0 | 58.6 | 31.6 | 83.3 | 35.3 | 49.7 | 3.3 | 28.8 | 35.6 | 48.5 |
| APODA [76] | AAAI'20 | 85.6 | 32.8 | 79.0 | 29.5 | 25.5 | 26.8 | 34.6 | 19.9 | 83.7 | 40.6 | 77.9 | 59.2 | 28.3 | 84.6 | 34.6 | 49.2 | 8.0 | 32.6 | 39.6 | 45.9 |
| STAR [49] | CVPR'20 | 88.4 | 27.9 | 80.8 | 27.3 | 25.6 | 26.9 | 31.6 | 20.8 | 83.5 | 34.1 | 76.6 | 60.5 | 27.2 | 84.2 | 32.9 | 38.2 | 1.0 | 30.2 | 31.2 | 43.6 |
| IntraDA [55] | CVPR'20 | 90.6 | 37.1 | 82.6 | 30.1 | 19.1 | 29.5 | 32.4 | 20.6 | 85.7 | 40.5 | 79.7 | 58.7 | 31.1 | 86.3 | 31.5 | 48.3 | 0.0 | 30.2 | 35.8 | 46.3 |
| SIM [71] | CVPR'20 | 90.6 | 44.7 | 84.8 | 34.3 | 28.7 | 31.6 | 35.0 | 37.6 | 84.7 | 43.3 | 85.3 | 57.0 | 31.5 | 83.8 | 42.6 | 48.5 | 1.9 | 30.4 | 39.0 | 49.2 |
| LSE [53] | ECCV'20 | 90.2 | 40.0 | 83.5 | 31.9 | 26.4 | 32.6 | 38.7 | 37.5 | 81.0 | 34.2 | 84.6 | 61.6 | 33.4 | 82.5 | 32.8 | 45.9 | 6.7 | 29.1 | 30.6 | 47.5 |
| WLabel [56] | ECCV'20 | 91.6 | 47.4 | 84.0 | 30.4 | 28.3 | 31.4 | 37.4 | 35.4 | 83.9 | 38.3 | 83.9 | 61.2 | 28.2 | 83.7 | 28.8 | 41.3 | 8.8 | 24.7 | 46.4 | 48.2 |
| CrCDA [32] | ECCV'20 | 92.4 | 55.3 | 82.3 | 31.2 | 29.1 | 32.5 | 33.2 | 35.6 | 83.5 | 34.8 | 84.2 | 58.9 | 32.2 | 84.7 | 40.6 | 46.1 | 2.1 | 31.1 | 32.7 | 48.6 |
| FADA [70] | ECCV'20 | 92.5 | 47.5 | 85.1 | 37.6 | 32.8 | 33.4 | 33.8 | 18.4 | 85.3 | 37.7 | 83.5 | 63.2 | 39.7 | 87.5 | 32.9 | 47.8 | 1.6 | 34.9 | 39.5 | 49.2 |
| LDR [74] | ECCV'20 | 90.8 | 41.4 | 84.7 | 35.1 | 27.5 | 31.2 | 38.0 | 32.8 | 85.6 | 42.1 | 84.9 | 59.6 | 34.4 | 85.0 | 42.8 | 52.7 | 3.4 | 30.9 | 38.1 | 49.5 |
| CCM [41] | ECCV'20 | 93.5 | 57.6 | 84.6 | 39.3 | 24.1 | 25.2 | 35.0 | 17.3 | 85.0 | 40.6 | 86.5 | 58.7 | 28.7 | 85.8 | 49.0 | 56.4 | 5.4 | 31.9 | 43.2 | 49.9 |
| CD-SAM [75] | WACV'21 | 91.3 | 46.0 | 84.5 | 34.4 | 29.7 | 32.6 | 35.8 | 36.4 | 84.5 | 43.2 | 83.0 | 60.0 | 32.2 | 83.2 | 35.0 | 46.7 | 0.0 | 33.7 | 42.2 | 49.2 |
| ASA [85] | TIP'21 | 89.2 | 27.8 | 81.3 | 25.3 | 22.7 | 28.7 | 36.5 | 19.6 | 83.8 | 31.4 | 77.1 | 59.2 | 29.8 | 84.3 | 33.2 | 45.6 | 16.9 | 34.5 | 30.8 | 45.1 |
| CLAN [51] | TPAMI'21 | 88.7 | 35.5 | 80.3 | 27.5 | 25.0 | 29.3 | 36.4 | 28.1 | 84.5 | 37.0 | 76.6 | 58.4 | 29.7 | 81.2 | 38.8 | 40.9 | 5.6 | 32.9 | 28.8 | 45.5 |
| DAST [79] | AAAI'21 | 92.2 | 49.0 | 84.3 | 36.5 | 28.9 | 33.9 | 38.8 | 28.4 | 84.9 | 41.6 | 83.2 | 60.0 | 28.7 | 87.2 | 45.0 | 45.3 | 7.4 | 33.8 | 32.8 | 49.6 |
| Ours | - | 91.3 | 48.6 | **85.5** | 35.8 | 31.4 | 36.7 | 37.5 | 36.8 | **86.3** | 40.3 | 85.7 | **64.3** | 31.1 | **87.7** | 36.7 | 44.9 | **15.9** | **38.9** | **55.4** | **51.9** |

Table 4: Segmentation results of Virtual KITTI → KITTI.

| Method | mIoU |
|---|---|
| GIO-Ada (CVPR'19) [5] | 53.50 |
| Self-Ensembling (ICCV'19) [9] | 55.45 |
| CutMix (BMVC'20) [17] | 55.58 |
| CowMix (Arxiv'20) [19] | 56.07 |
| DACS (WACV'21) [65] | 55.51 |
| PIT + CutMix (ICCV'21) [27] | 56.72 |
| PIT + CowMix (ICCV'21) [27] | 57.24 |
| PIT + DACS (ICCV'21) [27] | 56.57 |
| Ours | **60.16** |

Table 5: Segmentation results of Cityscapes → KITTI.

| Method | mIoU |
|---|---|
| Self-Ensembling (ICCV'19) [9] | 59.54 |
| CutMix (BMVC'20) [17] | 58.78 |
| CowMix (Arxiv'20) [19] | 59.15 |
| DACS (WACV'21) [65] | 59.19 |
| PIT + CutMix (ICCV'21) [27] | 60.09 |
| PIT + CowMix (ICCV'21) [27] | 60.37 |
| PIT + DACS (ICCV'21) [27] | 60.82 |
| Ours | **61.62** |

Table 6: Ablation of each component on SYNTHIA → Cityscapes.

| baseline | $M_{uncertainty}$ | $M_{classout}$ | mIoU | Gain |
|---|---|---|---|---|
| √ | | | 51.5 | - |
| √ | √ | | 53.5 | 2.0 |
| √ | √ | √ | 55.9 | 4.4 |

In Table 9, the proposed method outperforms the state-of-the-art approaches by 2% ~ 5% with ResNet101 backbone.

### 4.3.3. Virtual KITTI → KITTI and Cityscapes → KITTI

In addition to the two commonly-used benchmarks, we also conduct experiments on another synthetic-to-real adaptation, *i.e.,* Virtual KITTI → KITTI, and cross-city adaptation, *i.e.,* Cityscapes → KITTI, to validate the effectiveness of our method. Table 4 shows the results of adapting the model from Virtual KITTI to KITTI. We reproduce Self-Ensembling [9], CutMix [17], CowMix [19], DACS [65] in the same setting. The results of GIO-Ada [5] and PIT [27] are reported in the original papers. We can see that our method significantly improves the mIoU by 3.4% ~ 6.6% compared with the existing UDA methods. In Table 5, we adapt from Cityscapes to KITTI, where the source domains and target domain have different distributions in cross-city road scenes and street views. Our proposed method can outperform the state-of-the-art methods by around 1%. Table 5 demonstrate our effectiveness in cross-city adaptation.

### 4.4. Ablation Study

**Ablation of each component:** In Table 6, we investigate the effects of different design components in SYNTHIA → Cityscapes with ResNet101 backbone. The uncertainty mask $M_{uncertainty}$ and the ClassDrop mask $M_{classout}$ reveals the contribution of the proposed dynamic weighting scheme and the ClassOut strategy are complementary. The consistency regularization baseline is 51.5%. By adding the $M_{uncertainty}$ and $M_{classout}$ sequentially, we boost the mIoU with an additional +2.0% and +2.4%, achieving 53.5% and 55.9%, respectively. These improvements show the effects of individual components of our proposed approach.

**Comparison to the related work [80] :** In Table. 7, we show the experimental comparison on two benchmark datasets with ResNet 101 backbones to demonstrate its effectiveness. Note
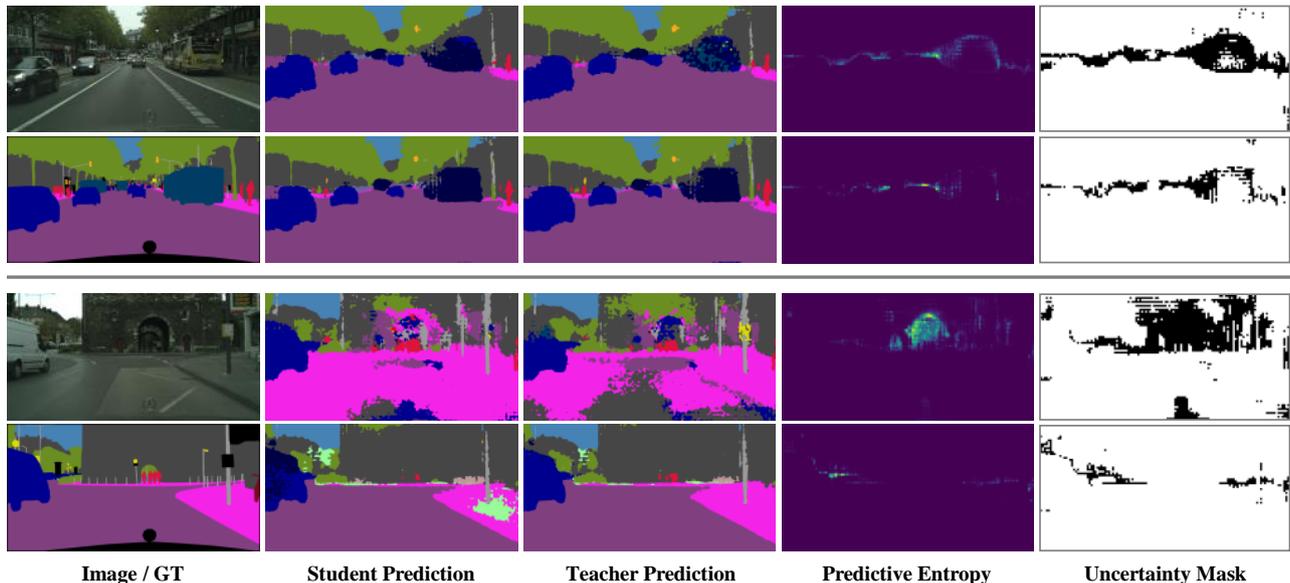
| Image / GT | Student Prediction | Teacher Prediction | Predictive Entropy | Uncertainty Mask |

Fig. 5: Visualization results of GTA5 → Cityscapes (first and second rows) and SYNTHIA → Cityscapes (third and fourth rows). Segmentation results at 10K training steps (first and third rows) and 56K training steps (second and fourth rows). The fourth and fifth columns illustrate the predictive entropy and our uncertainty mask.

Table 7: Comparisons with the related work [80]. on GTAV ⟶ Cityscapes with ResNet 101 backbone.

| method | mIoU (GTAV) | mIoU$_{13}$ (SYN) |
|---|---|---|
| Mean Teacher [9] | 43.1 | 45.9 |
| + Yu et al. [80] | 44.6 | 47.6 |
| + Ours | **51.9** | **55.9** |

Table 8: Ablation study of each module's improvement from GTA5 to Cityscapes with VGG16 backbone. $L_{seg}$: Segmentation loss, $L_{mse}$: Mean Square Error used in [9], $L_{con}$: Our Uncertainty-Guided Consistency Loss, $IT$: Image-to-Image translation for Style transfer.

| Method | Component | mIoU | Gain |
|---|---|---|---|
| Source Only | $L_{seg}$ | 28.3 | - |
| Choi *et al.* [9] | $L_{seg}+L_{mse}$ | 32.6 | +4.3 |
| Ours (w/o $M_{classout}$) | $L_{seg}+L_{con}$ | 35.6 | +7.3 |
| Choi *et al.* [9] | $L_{seg}+IT_1$ [9] | 35.4 | + 4.1 |
| Ours (w/o $M_{classout}$) | $L_{seg}+IT_2$ [44] | 35.1 | + 3.8 |
| Choi *et al.* [9] | $L_{seg}+L_{mse}+IT_1$ | 42.5 | +14.2 |
| Ours (w/o $M_{classout}$) | $L_{seg}+L_{con}+IT_2$ | 47.8 | +19.5 |

that all the experimental results of Table. 7 are conducted on the same Mean-Teacher baseline with ResNet 101 backbones. We replace the proposed method with the approach [80], and we find that the improvements of [80] are limited over the Mean Teacher baseline, only achieving 44.6 and 47.6 in GTAV → Cityscapes and SYNTHIA → Cityscapes, respectively. Our proposed method outperforms the related work [80] by 7.3 % and 8.3 % on GTAV → Cityscapes and SYNTHIA → Cityscapes, achieving 51.9% and 55.9%, respectively.

**Comparison to the related work [9] :** In Table 8, we compare our method with the non-adaptive baseline and Self-Ensembling (SE) [9] with VGG16 backbone. $L_{seg}$ denotes the supervised segmentation loss, $L_{mse}$ refers to the common Mean Square Error used in [9], and $L_{con}$ is our uncertainty-guided consistency Loss with the dynamic weighting scheme. As we can see, the Source Only baseline achieves 28.3% from GTAV dataset to Cityscapes dataset. We see that in the third row, Choi *et al.* achieves a performance of 32.6% in the original consistency loss ($L_{seg} + L_{mse}$). Our uncertainty-guided consistency loss achieves about 3.0% improvement over directly using the Mean Square Error ($L_{seg} + L_{con}$), reaching 35.6% in mIoU.

As mentioned in Section 2, pixel-level adaptation is also

a key factor in minimizing the discrepancy of data distribution. Therefore, it is helpful to utilize a transferred source domain image dataset whose appearance is more similar to that of the target-domain image dataset. Following common practice [71, 44], we adopt the transferred GTA5 images of [44] which utilizes a CycleGAN[86] structure to adapt the style of GTAV images to the style of Cityscapes images. In the fifth row and sixth row of Table 3, we could find that our Image-to-Image Translation achieves a similar performance compared to [9]. On top of that, as we can see in the last row, our final adaptive performance is superior to the state-of-the-art by 5.3%, resulting in a 19.5% increase in mIoU over the non-adaptive baseline.

### 4.5. Analysis

In this section, we provide visualization results and provide some analysis of our proposed framework.

Fig. 6 shows the comparison results of the per-class IoU gain and comparisons of mIoU between the SE baseline [9] and our proposed method. In many large categories, *i.e.,* road, building,

Table 9: Comparison results (mIoU) from SYNTHIA to Cityscapes (with ResNet 101 backbone).

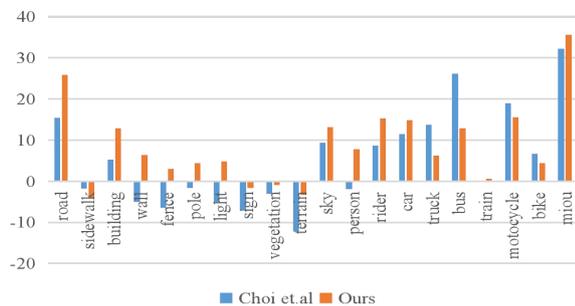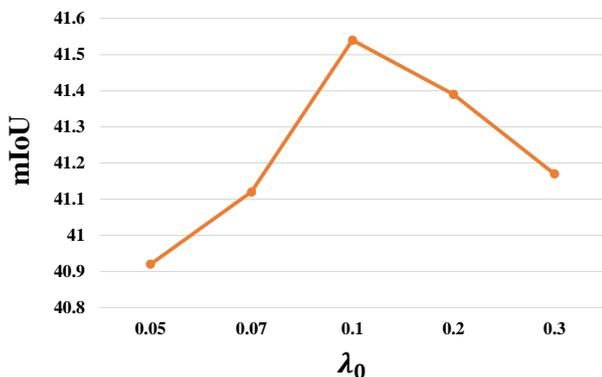| Method | Venue | road | sidewalk | building | light | sign | vegetation | sky | person | rider | car | bus | motocycle | bike | mIoU$_{13}$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Source Only | - | 36.3 | 14.6 | 68.8 | 5.6 | 9.1 | 69.0 | 79.4 | 52.5 | 11.3 | 49.8 | 9.5 | 11.0 | 20.7 | 29.5 |
| BDL [44] | CVPR'19 | 86.0 | 46.7 | 80.3 | 14.1 | 11.6 | 79.2 | 81.3 | 54.1 | 27.9 | 73.7 | 42.2 | 25.7 | 45.3 | 51.4 |
| DADA [68] | ICCV'19 | 89.2 | 44.8 | 81.4 | 8.6 | 11.1 | 81.8 | 84.0 | 54.7 | 19.3 | 79.7 | 40.7 | 14.0 | 38.8 | 49.8 |
| STAR [49] | CVPR'20 | 82.6 | 36.2 | 81.1 | 12.2 | 8.7 | 78.4 | 82.2 | 59.0 | 22.5 | 76.3 | 33.6 | 11.9 | 40.8 | 48.1 |
| IntraDA [55] | CVPR'20 | 84.3 | 37.7 | 79.5 | 9.2 | 8.4 | 80.0 | 84.1 | 57.2 | 23.0 | 78.0 | 38.1 | 20.3 | 36.5 | 48.9 |
| LTIR [34] | CVPR'20 | 92.6 | 53.2 | 79.2 | 1.6 | 7.5 | 78.6 | 84.4 | 52.6 | 20.0 | 82.1 | 34.8 | 14.6 | 39.4 | 49.3 |
| SIM [71] | CVPR'20 | 83.0 | 44.0 | 80.3 | 17.1 | 15.8 | 80.5 | 81.8 | 59.9 | 33.1 | 70.2 | 37.3 | 28.5 | 45.8 | 52.1 |
| LSE [53] | ECCV'20 | 82.9 | 43.1 | 78.1 | 9.1 | 14.4 | 77.0 | 83.5 | 58.1 | 25.9 | 71.9 | 38.0 | 29.4 | 31.2 | 49.4 |
| CrCDA [32] | ECCV'20 | 86.2 | 44.9 | 79.5 | 9.4 | 11.8 | 78.6 | 86.5 | 57.2 | 26.1 | 76.8 | 39.9 | 21.5 | 32.1 | 50.0 |
| WLabel [56] | ECCV'20 | 92.0 | 53.5 | 80.9 | 3.8 | 6.0 | 81.6 | 84.4 | 60.8 | 24.4 | 80.5 | 39.0 | 26.0 | 41.7 | 51.9 |
| CD-SAM [75] | WACV'21 | 82.5 | 42.2 | 81.3 | 18.3 | 15.9 | 80.6 | 83.5 | 61.4 | 33.2 | 72.9 | 39.3 | 26.6 | 43.9 | 52.4 |
| CLAN [51] | TPAMI'21 | 82.7 | 37.2 | 81.5 | 17.1 | 13.1 | 81.2 | 83.3 | 55.5 | 22.1 | 76.6 | 30.1 | 23.5 | 30.7 | 48.8 |
| ASA [85] | TIP'21 | 91.2 | 48.5 | 80.4 | 5.5 | 5.2 | 79.5 | 83.6 | 56.4 | 21.9 | 80.3 | 36.2 | 20.0 | 32.9 | 49.3 |
| DAST [79] | AAAI'21 | 87.1 | 44.5 | 82.3 | 13.9 | 13.1 | 81.6 | 86.0 | 60.3 | 25.1 | 83.1 | 40.1 | 24.4 | 40.5 | 52.5 |
| Ours | - | 85.5 | 42.5 | **83.0** | **20.9** | **25.5** | **82.5** | **88.0** | **63.2** | **31.8** | **86.5** | **41.2** | 25.9 | **50.7** | **55.9** |



Fig. 6: Comparisons of Per-Class IoU Gain between Choi et.al [9] and ours w/o IT with VGG16 backbone in GTAV → Cityscapes.



Fig. 7: Parameter analysis about $\lambda_0$ (SYNTHIA → Cityscapes).

sky, that have long boundaries, we have achieved a per-class IoU performance improvement. In other static categories, such as sidewalk, sign, vegetation and terrain, our method achieves a lower performance degradation than [9]. Besides, as we can see in Table 3 and Table 9, our method shows good performance in some moving objects, *e.g.,* motorcycle, bicycle, etc. Our method obtains an overall better performance than [9].

#### 4.5.1. Visualization

The effectiveness of the uncertainty-aware consistency regularization is shown in Fig. 5. We visualize the student prediction, teacher prediction, the entropy map of the teacher model, and our uncertainty mask. In Fig. 5, as we can see in the fourth column, the predictive entropy captures the latent epidemic uncertainty, especially for some specific large objects, such as car and truck. In the fifth column, the white pixels of the uncertainty mask are the ones with higher confidence. The first and third rows show that our uncertainty mask effectively filters out the unreasonable pixels and guides the teacher to be a good proxy for training the student network in the early stage of the training process. In addition, the second and fourth rows show that our uncertainty module pays attention to the semantic boundary of objects in the later training stage.

These qualitative results are consistent with our motivations and reveal the reason why current consistency regularization methods are often unstable in minimizing the distribution discrepancy, which lies in two aspects. Firstly, directly imposing a simple MSE constraint as consistency loss onto all pixels could harm the learning process by generating unreasonable guidance from the teacher to the student model. Secondly, due to the "error accumulation" in the teacher model, it could take more training iterations to converge and even may lead to early performance degradation during the adaptation process. In the second row of Figure 4, the entropy is low and not obvious to
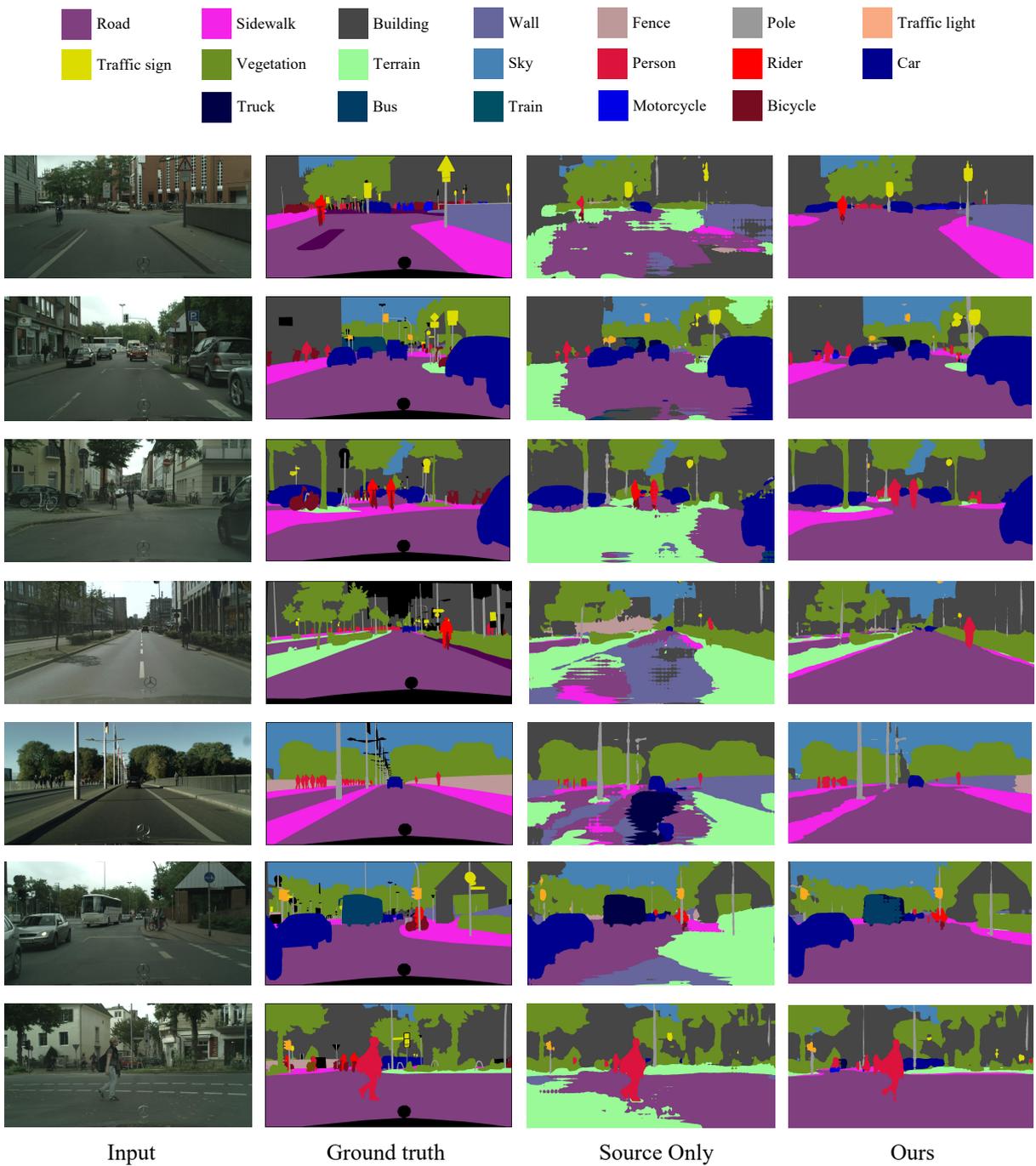
Fig. 8: Semantic segmentation qualitative results from GTA5 to Cityscapes. From left to right: target image, ground truth, source-only prediction, and predictions using our method.

Table 10: Parameter analysis about $\alpha$.

| $\alpha$ | 0.675 | 0.70 | 0.725 | 0.75 | 0.775 | 0.80 |
|---|---|---|---|---|---|---|
| mean IoU | 54.0 | 53.9 | **55.1** | 54.5 | 54.7 | 53.3 |

Table 11: Parameter analysis about $\beta$.

| $\beta$ | -5.3 | -5.2 | -5.1 | -5.0 | -4.9 | -4.8 | -4.7 |
|---|---|---|---|---|---|---|---|
| mean IoU | 54.2 | 54.7 | 53.7 | 54.5 | 54.7 | **55.9** | 54.5 |

Table 12: Parameter analysis about $N$.

| $N$ | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|
| mean IoU | 53.6 | 54.0 | **54.5** | 53.4 | 54.1 |

human naked eyes. However, our uncertainty mask can still filter out these uncertain areas. The main reason is that our dynamic weighting scheme can enable dynamic adjustion during the training. When the entropy is low, the dynamic threshold will be updated accordingly, and thus, we can filter out theses unobvious areas.

In Fig. 8, we illustrate some qualitative results of our models tested on the validation sets of Cityscapes [10] dataset. Following prior works [73, 30, 81], we show the target image, ground truth, source only prediction, and our prediction from left to right. Without domain adaptation, the model trained only on source supervision produces noisy segmentation predictions. With the help of our uncertainty-aware consistency-regularization, our method manages to produce correct predictions at a high level of confidence.

*4.5.2. Parameter Analysis*

In this section, we investigate the sensitivity of hyperparameters $\lambda_0$, $\alpha$, $\beta$, $N$ and augmentations.

**Effect of $\lambda_0$:** $\lambda_0$ means the initial state of consistency weight $\lambda_{con}$, which balances the domain adaptation process among different loss functions, and it is crucial in the training process. In Fig. 7, the best performance from SYNTHIA to Cityscapes (w VGG16 backbone) occurs when the initial value of $\lambda_0$ is 0.1. The results of mean mIoU over the 16 common classes are reported.

**Effect of $\alpha$:** In this experiment, we set $\beta = -5$, $N = 8$ to check the sensitivity of $\alpha$ in SYNTHIA to Cityscapes (w ResNet101 backbone). $\alpha$ is the initial state of the dynamic threshold in Eq 4. In Table 10, we find that when $\alpha$ is set to 0.725, we can obtain the highest performance.

**Effect of $\beta$:** In this experiment, we adapt from SYNTHIA to Cityscapes (w ResNet101 backbone) to discuss the selection of the parameter $\beta$, which controls the exponential speed of the dynamic threshold in Equation 4. We set the other parameters $\alpha = 0.75$, $N = 8$. In Table 11, we observe that the highest mIoU on target domain is achieved when the value of $\beta$ is around $-4.8$, which means that, under such condition, the exponential speed benefits the dynamic threshold of the domain adaptation the most.

**Effect of $N$:** In this part, we analyze the selection of $N$, which is the copy numbers of target image in the stochastic forward pass. We set the other parameters as: $\alpha = 0.75$, $\beta = -5$. In Table 12, we find that when $N$ is set to 8, it achieves the best performance in SYNTHIA to Cityscapes (w ResNet101 backbone). Therefore, $N$ is set to 8 in all experiments.

**Effects of Augmentations:** We investigate the sensitivity of augmentation, *e.g,* Gaussian noise, color jittering, random crop. GN, CJ, RC are the abbreviation of Gaussian noise, color jittering, random crop, respectively. The ablations of each augmentation are shown in Table 13 when adapting from SYNTHIA

to Cityscapes (w ResNet101 backbone). From the Table 13, we can observe that augmentations are complementary and our consistency regularization needs to be conducted under the condition that the student image and the target image are imposed with different augmentations.

*4.6. Limitations*

1) We develop a unified uncertainty-aware consistency regularization in this work. Though our method has achieved very good results, it can hardly treat the stuff regions and the instances of things in a different manner to reduce the uncertainty. 2) The scale-invariant feature across different frames are neglected in this work, which can be utilized as prior knowledge for effective domain adaptation for video semantic segmentation. As future work, these interesting points will be investigated.

## 5. Conclusion

In this paper, we proposed an uncertainty-aware consistency regularization technique to address the domain shift for cross-domain segmentation. Our uncertainty module is capable of estimating the latent uncertainty map for the purpose of a better knowledge transfer. Specifically, We first introduced an uncertainty-guided consistency loss with a dynamic weighting scheme for filtering out the unreasonable pixels and mining the high confident predictions of target samples. Secondly, we present a ClassDrop mask generation algorithm to generate class-wise perturbations. Guided by this mask, we present a ClassOut strategy to keep the local regional consistency in varying environments. Experimental results verify that our method is superior to existing state-of-the-art approaches on four challenging benchmark datasets.

## 6. Acknowledgments

Table 13: Ablation study of each augmentation.

| No Aug | RC | GN | CJ | mIoU |
|--------|----|----|----|------|
| √ | | | | 52.2 |
| √ | √ | | | 51.9 |
| √ | √ | √ | | 53.7 |
| √ | √ | √ | √ | 55.9 |

# References

[1] Cariucci, F.M., Porzi, L., Caputo, B., Ricci, E., Bulo, S.R., 2017. Autodial: Automatic domain alignment layers, in: ICCV.

[2] Chang, W.L., Wang, H.P., Peng, W.H., Chiu, W.C., 2019. All about structure: Adapting structural information across domains for boosting semantic segmentation, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 1900–1909.

[3] Chen, L., Papandreou, G., Kokkinos, I., Murphy, K., Yuille, A.L., 2018a. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. IEEE Trans. Pattern Anal. Mach. Intell. 40, 834–848.

[4] Chen, L.C., Zhu, Y., Papandreou, G., Schroff, F., Adam, H., 2018b. Encoder-decoder with atrous separable convolution for semantic image segmentation, in: Proceedings of the European conference on computer vision (ECCV), pp. 801–818.

[5] Chen, Y., Li, W., Chen, X., Gool, L.V., 2019a. Learning semantic segmentation from synthetic data: A geometrically guided input-output adaptation approach, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 1841–1850.

[6] Chen, Y., Li, W., Van Gool, L., 2018c. Road: Reality oriented adaptation for semantic segmentation of urban scenes, in: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 7892–7901.

[7] Chen, Y.C., Lin, Y.Y., Yang, M.H., Huang, J.B., 2019b. Crdoco: Pixel-level domain transfer with cross-domain consistency, in: CVPR.

[8] Chen, Y.H., Chen, W.Y., Chen, Y.T., Tsai, B.C., Frank Wang, Y.C., Sun, M., 2017. No more discrimination: Cross city adaptation of road scene segmenters, in: Proceedings of the IEEE International Conference on Computer Vision, pp. 1992–2001.

[9] Choi, J., Kim, T., Kim, C., 2019. Self-ensembling with gan-based data augmentation for domain adaptation in semantic segmentation, in: Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 6830–6840.

[10] Cordts, M., Omran, M., Ramos, S., Rehfeld, T., Enzweiler, M., Benenson, R., Franke, U., Roth, S., Schiele, B., 2016. The cityscapes dataset for semantic urban scene understanding, in: Proc. CVPR, pp. 3213–3223.

[11] Csurka, G., 2017. Domain adaptation for visual applications: A comprehensive survey. arXiv preprint arXiv:1702.05374 .

[12] Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L., 2009. Imagenet: A large-scale hierarchical image database, in: CVPR.

[13] Du, L., Tan, J., Yang, H., Feng, J., Xue, X., Zheng, Q., Ye, X., Zhang, X., 2019. Ssf-dan: Separated semantic feature based domain adaptation network for semantic segmentation, in: ICCV.

[14] Everingham, M., Van Gool, L., Williams, C.K., Winn, J., Zisserman, A., 2010. The pascal visual object classes (voc) challenge. International journal of computer vision 88, 303–338.

[15] Feng, Z., Zhou, Q., Cheng, G., Tan, X., Shi, J., Ma, L., 2020. Semi-supervised semantic segmentation via dynamic self-training and class-balanced curriculum. arXiv preprint arXiv:2004.08514v1 .

[16] Fernando, B., Habrard, A., Sebban, M., Tuytelaars, T., 2013. Unsupervised visual domain adaptation using subspace alignment, in: ICCV.

[17] French, G., Laine, S., Aila, T., Laine, S., Mackiewicz, M., Finlayson, G., 2020a. Semi-supervised semantic segmentation needs strong, varied perturbations, in: British Machine Vision Conference.

[18] French, G., Mackiewicz, M., Fisher, M., 2018. Self-ensembling for visual domain adaptation, in: Proceedings of the International Conference on Learning Representations.

[19] French, G., Oliver, A., Salimans, T., 2020b. Milking cowmask for semi-supervised image classification. arXiv preprint arXiv:2003.12022 .

[20] Gaidon, A., Wang, Q., Cabon, Y., Vig, E., 2016a. Virtual worlds as proxy for multi-object tracking analysis, in: CVPR.

[21] Gaidon, A., Wang, Q., Cabon, Y., Vig, E., 2016b. Virtual worlds as proxy for multi-object tracking analysis, in: CVPR.

[22] Ganin, Y., Lempitsky, V., 2015. Unsupervised domain adaptation by backpropagation, in: ICLR.

[23] Geng, B., Tao, D., Xu, C., 2011. Daml: Domain adaptation metric learning. TIP 20, 2980–2989.

[24] Gong, B., Shi, Y., Sha, F., Grauman, K., 2012. Geodesic flow kernel for unsupervised domain adaptation, in: CVPR.

[25] Gong, R., Li, W., Chen, Y., Gool, L.V., 2019. Dlow: Domain flow for adaptation and generalization, in: CVPR.

[26] Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y., 2014. Generative adversarial nets. Advances in neural information processing systems 27.

[27] Gu, Q., Zhou, Q., Xu, M., Feng, Z., Cheng, G., Lu, X., Shi, J., Ma, L., 2021. Pit: Position-invariant transform for cross-fov domain adaptation, in: Proceedings of the IEEE/CVF International Conference on Computer Vision.

[28] Han, L., Zou, Y., Gao, R., Wang, L., Metaxas, D., 2019. Unsupervised domain adaptation via calibrating uncertainties, in: CVPR Workshops.

[29] He, K., Zhang, X., Ren, S., Sun, J., 2016. Deep residual learning for image recognition, in: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 770–778.

[30] Hoffman, J., Tzeng, E., Park, T., Zhu, J.Y., Isola, P., Saenko, K., Efros, A., Darrell, T., 2018. CyCADA: Cycle-consistent adversarial domain adaptation, in: International conference on machine learning, pp. 1989–1998.

[31] Hoffman, J., Wang, D., Yu, F., Darrell, T., 2016. Fcns in the wild: Pixel-level adversarial and constraint-based adaptation. CoRR abs/1612.02649.

[32] Huang, J., Lu, S., Guan, D., Zhang, X., 2020. Contextual-relation consistent domain adaptation for semantic segmentation, in: European conference on computer vision, pp. 705–722.

[33] Kendall, A., Gal, Y., 2017. What uncertainties do we need in bayesian deep learning for computer vision?, in: NeurIPS.

[34] Kim, M., Byun, H., 2020. Learning texture invariant representation for domain adaptation of semantic segmentation, in: Proc. CVPR, pp. 12975–12984.

[35] Kulis, B., Saenko, K., Darrell, T., 2011. What you saw is not what you get: Domain adaptation using asymmetric kernel transforms, in: CVPR.

[36] Kurmi, V.K., Kumar, S., Namboodiri, V.P., 2019. Attending to discriminative certainty for domain adaptation, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 491–500.

[37] Laine, S., Aila, T., 2016. Temporal ensembling for semi-supervised learning. arXiv preprint arXiv:1610.02242 .

[38] LeCun, Y., Bottou, L., Bengio, Y., Haffner, P., 1998. Gradient-based learning applied to document recognition. IEEE Proc. 86, 2278–2324.

[39] Lee, C.Y., Batra, T., Baig, M.H., Ulbricht, D., 2019a. Sliced wasserstein discrepancy for unsupervised domain adaptation, in: CVPR.

[40] Lee, K.H., Ros, G., Li, J., Gaidon, A., 2019b. Spigan: Privileged adversarial learning from simulation .

[41] Li, G., Kang, G., Liu, W., Wei, Y., Yang, Y., 2020a. Content-consistent matching for domain adaptive semantic segmentation, in: European conference on computer vision, Springer. pp. 440–456.

[42] Li, S., Liu, H.C., Lin, Q., Xie, B., Ding, Z., Huang, G., Tang, J., 2020b. Domain conditioned adaptation network, in: Proc. AAAI, pp. 11386–11393.

[43] Li, X., Zhong, Z., Wu, J., Yang, Y., Lin, Z., Liu, H., 2019a. Expectation-maximization attention networks for semantic segmentation, in: ICCV.

[44] Li, Y., Yuan, L., Vasconcelos, N., 2019b. Bidirectional learning for domain adaptation of semantic segmentation, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 6936–6945.

[45] Lian, Q., Duan, L., Lv, F., 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, pp. 6757–6766.

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

[47] Long, J., Shelhamer, E., Darrell, T., 2015a. Fully convolutional networks for semantic segmentation, in: CVPR.

[48] Long, M., Cao, Y., Wang, J., Jordan, M., 2015b. Learning transferable features with deep adaptation networks, in: ICLR.

[49] Lu, Z., Yang, Y., Zhu, X., Liu, C., Song, Y.Z., Xiang, T., 2020. Stochastic classifiers for unsupervised domain adaptation, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 9111–9120.

[50] Luo, Y., Liu, P., Guan, T., Yu, J., Yang, Y., 2019a. Significance-aware information bottleneck for domain adaptive semantic segmentation, in: Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 6778–6787.

[51] Luo, Y., Liu, P., Zheng, L., Guan, T., Yu, J., Yang, Y., 2021. Category-level adversarial adaptation for semantic segmentation using purified features. IEEE Transactions on Pattern Analysis and Machine Intelligence , 1–1doi:10.1109/TPAMI.2021.3064379.

[52] Luo, Y., Zheng, L., Guan, T., Yu, J., 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, pp. 2507–2516.

[53] Naseer Subhani, M., Ali, M., 2020. Learning from scale-invariant examples for domain adaptation in semantic segmentation, in: European conference on computer vision, Springer. pp. 290–306.

[54] Netzer, Y., Wang, T., Coates, A., Bissacco, A., Wu, B., Ng, A.Y., 2011. Reading digits in natural images with unsupervised feature learning, in: NeurIPS workshop.

[55] Pan, F., Shin, I., Rameau, F., Lee, S., Kweon, I.S., 2020. Unsupervised intra-domain adaptation for semantic segmentation through self-supervision, in: Unsupervised Intra-domain Adaptation for Semantic Segmentation through Self-Supervision, pp. 3764–3773.

[56] Paul, S., Tsai, Y., Schulter, S., Roy-Chowdhury, A.K., Chandraker, M., 2020. Domain adaptive semantic segmentation using weak labels, in: European conference on computer vision, Springer. pp. 571–587.

[57] Perone, C.S., Ballester, P., Barros, R.C., Cohen-Adad, J., 2019. Unsupervised domain adaptation for medical imaging segmentation with self-ensembling. NeuroImage 194, 1–11.

[58] Richter, S.R., Vineet, V., Roth, S., Koltun, V., 2016. Playing for data: Ground truth from computer games, in: ECCV.

[59] Ros, G., Sellart, L., Materzynska, J., Vazquez, D., Lopez, A.M., 2016. The synthia dataset: A large collection of synthetic images for semantic segmentation of urban scenes, in: CVPR.

[60] Saito, K., Watanabe, K., Ushiku, Y., Harada, T., 2018. Maximum classifier discrepancy for unsupervised domain adaptation, in: CVPR.

[61] Sankaranarayanan, S., Balaji, Y., Jain, A., Nam Lim, S., Chellappa, R., 2018. Learning from synthetic data: Addressing domain shift for semantic segmentation, in: CVPR.

[62] Simonyan, K., Andrew, Z., 2014. Very deep convolutional networks for large-scale image recognition. arXiv preprint arXiv:1409.1556 .

[63] Sun, B., Feng, J., Saenko, K., 2016. Return of frustratingly easy domain adaptation, in: AAAI.

[64] Tarvainen, A., Valpola, H., 2017. Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results, in: Advances in Neural Information Processing Systems 30, pp. 1195–1204.

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

[66] Tsai, Y.H., Hung, W.C., Schulter, S., Sohn, K., Yang, M.H., Chandraker, M., 2018. Learning to adapt structured output space for semantic segmentation, in: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 7472–7481.

[67] Tsai, Y.H., Sohn, K., Schulter, S., Chandraker, M., 2019. Domain adaptation for structured output via discriminative patch representations, in: Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 1456–1465.

[68] Vu, T., Jain, H., Bucher, M., Cord, M., Pérez, P., 2019a. DADA: depth-aware domain adaptation in semantic segmentation, in: Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 7363–7372.

[69] Vu, T.H., Jain, H., Bucher, M., Cord, M., Pérez, P., 2019b. Advent: Adversarial entropy minimization for domain adaptation in semantic segmentation, in: CVPR.

[70] Wang, H., Shen, T., Zhang, W., Duan, L., Mei, T., 2020a. Classes matter: A fine-grained adversarial approach to cross-domain semantic segmentation, Springer. pp. 642–659.

[71] Wang, Z., Yu, M., Wei, Y., Feris, R., Xiong, J., Hwu, W.m., Huang, T.S., Shi, H., 2020b. 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, pp. 12635–12644.

[72] Wen, J., Zheng, N., Yuan, J., Gong, Z., Chen, C., 2019. Bayesian uncertainty matching for unsupervised domain adaptation. arXiv preprint arXiv:1906.09693 .

[73] Xu, Y., Du, B., Zhang, L., Zhang, Q., Wang, G., Zhang, L., 2019. Self-ensembling attention networks: Addressing domain shift for semantic segmentation, in: Proceedings of the AAAI Conference on Artificial Intelligence, pp. 5581–5588.

[74] Yang, J., An, W., Wang, S., Zhu, X., Yan, C., Huang, J., 2020a. Label-driven reconstruction for domain adaptation in semantic segmentation, in: European conference on computer vision, Springer. pp. 480–498.

[75] Yang, J., An, W., Yan, C., Zhao, P., Huang, J., 2021. Context-aware domain adaptation in semantic segmentation, in: Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision, pp. 514–524.

[76] Yang, J., Xu, R., Li, R., Qi, X., Shen, X., Li, G., Lin, L., 2020b. An adversarial perturbation oriented domain adaptation approach for semantic segmentation, in: Proceedings of the AAAI Conference on Artificial Intelligence, pp. 12613–12620.

[77] Yang, Y., Lao, D., Sundaramoorthi, G., Soatto, S., 2020c. Phase consistent ecological domain adaptation, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 9011–9020.

[78] Yang, Y., Soatto, S., 2020. Fda: Fourier domain adaptation for semantic segmentation, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 4085–4095.

[79] Yu, F., Zhang, M., Dong, H., Hu, S., Dong, B., 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, pp. 10754–10762.

[80] Yu, L., Wang, S., Li, X., Fu, C.W., Heng, P.A., 2019. Uncertainty-aware self-ensembling model for semi-supervised 3d left atrium segmentation, in: MICCAI.

[81] Zhang, Q., Zhang, J., Liu, W., Tao, D., 2019. Category anchor-guided unsupervised domain adaptation for semantic segmentation, in: Advances in Neural Information Processing Systems, pp. 433–443.

[82] Zhang, Y., David, P., Gong, B., 2017. Curriculum domain adaptation for semantic segmentation of urban scenes, in: Proceedings of the IEEE international conference on computer vision, pp. 2020–2030.

[83] Zhao, S., Li, B., Yue, X., Gu, Y., Xu, P., Tan, Hu, R., Chai, H., Keutzer, K., 2019. Multi-source domain adaptation for semantic segmentation, in: NeurIPS.

[84] Zheng, Z., Yang, Y., 2020. Rectifying pseudo label learning via uncertainty estimation for domain adaptive semantic segmentation. International Journal of Computer Vision (IJCV) doi:10.1007/s11263-020-01395-y.

[85] Zhou, W., Wang, Y., Chu, J., Yang, J., Bai, X., Xu, Y., 2020. Affinity space adaptation for semantic segmentation across domains. IEEE Transactions on Image Processing 30, 2549–2561.

[86] Zhu, J.Y., Park, T., Isola, P., Efros, A.A., 2017. Unpaired image-to-image translation using cycle-consistent adversarial networks, in: Proceedings of the IEEE international conference on computer vision, pp. 2223–2232.

[87] Zhu, X., Zhou, H., Yang, C., Shi, J., Lin, D., 2018. Penalizing top performers: Conservative loss for semantic segmentation adaptation, in: Proceedings of the European Conference on Computer Vision (ECCV), pp. 568–583.

[88] Zou, Y., Yu, Z., Kumar, B., Wang, J., 2018. Unsupervised domain adaptation for semantic segmentation via class-balanced self-training, in: Proceedings of the European conference on computer vision (ECCV), pp. 289–305.

[89] Zou, Y., Yu, Z., Liu, X., Kumar, B., Wang, J., 2019. Confidence regularized self-training, in: Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 5982–5991.