# Rectifying Pseudo Label Learning via Uncertainty Estimation for Domain Adaptive Semantic Segmentation

**Zhedong Zheng, Yi Yang**

**Abstract** This paper focuses on the unsupervised domain adaptation of transferring the knowledge from the source domain to the target domain in the context of semantic segmentation. Existing approaches usually regard the pseudo label as the ground truth to fully exploit the unlabeled target-domain data. Yet the pseudo labels of the target-domain data are usually predicted by the model trained on the source domain. Thus, the generated labels inevitably contain the incorrect prediction due to the discrepancy between the training domain and the test domain, which could be transferred to the final adapted model and largely compromises the training process.

To overcome the problem, this paper proposes to explicitly estimate the prediction uncertainty during training to rectify the pseudo label learning for unsupervised semantic segmentation adaptation. Given the input image, the model outputs the semantic segmentation prediction as well as the uncertainty of the prediction. Specifically, we model the uncertainty via the prediction variance and involve the uncertainty into the optimization objective. To verify the effectiveness of the proposed method, we evaluate the proposed method on two prevalent synthetic-to-real semantic segmentation benchmarks, *i.e.*, GTA5 → Cityscapes and SYNTHIA → Cityscapes, as well as one cross-city benchmark, *i.e.*, Cityscapes → Oxford RobotCar. We demonstrate through extensive experiments that the proposed approach (1) dynamically sets different confidence thresholds according to the prediction variance, (2) rectifies the learning from noisy pseudo labels, and (3) achieves significant improvements over the conventional pseudo label learning and yields competitive performance on all three benchmarks.

Zhedong Zheng and Yi Yang are with the Australian Artificial Intelligence Institute (AAII), University of Technology Sydney, NSW, Australia. E-mail: zhedong.zheng@student.uts.edu.au, yi.yang@uts.edu.au

## 1 Introduction

Deep neural networks (DNNs) have been widely adopted in the field of semantic segmentation, yielding the state-of-the-art performance [Liang et al., 2017, Wei et al., 2018]. However, recent works show that DNNs are limited in the scalability to the unseen environments, *e.g.*, the testing data collected in rainy days [Hendrycks and Dietterich, 2019, Wu et al., 2019]. One straightforward idea is to annotate more training data of the target environment and then re-train the segmentation model. However, semantic segmentation task usually demands dense annotations and it is unaffordable to manually annotate the pixel-wise label for collected data in new environments. To address the challenge, the researchers, therefore, resort to unsupervised semantic segmentation adaption, which takes one step closer to real-world practice. In unsupervised semantic segmentation adaptation, two datasets collected in different environments are considered: a labeled source-domain dataset where category labels are provided for every pixel, and an unlabeled target-domain dataset where only provides the collected data without annotations. Compared with the annotated data in the target domain, the unlabeled data is usually easy to collect. Semantic segmentation adaptation aims at leveraging the labeled source-domain data as well as the unlabeled target-domain data to adapt the well-trained model to the target environment.

The main challenge of semantic segmentation adaption is the discrepancy of data distribution between the source domain and the target domain. There are two lines of methods for semantic segmentation adaptation. On one hand, sev-

eral existing works focus on the domain alignment by minimizing the distribution discrepancy in different levels, such as pixel level [Wu et al., 2018, Wu et al., 2019, Hoffman et al., 2018], feature level [Huang et al., 2018, Yue et al., 2019, Luo et al., 2019a, Zhang et al., 2019b] and semantic level [Tsai et al., 2018, Tsai et al., 2019, Wang et al., 2019]. Despite great success, this line of work is sub-optimal. Because the alignment objective drives the model to learn the shared knowledge between domains but ignores the domain-specific knowledge. The domain-specific knowledge is one of the keys to the final target, *i.e.*, the model adapted to the target domain. On the other hand, some researchers focus on learning the domain-specific knowledge of the target domain by fully exploiting the unlabeled target-domain data [Zou et al., 2018, Zou et al., 2019, Han et al., 2019]. Specifically, this line of methods usually adopts the two-stage pipeline, which is similar to the traditional semi-supervised framework [Lee, 2013]. The first step is to predict pseudo labels by the knowledge learned from the labeled data, *e.g.*, the model trained on the source domain. The second step is to minimize the cross-entropy loss on the pseudo labels of the unlabeled target-domain data. In the training process, pseudo labels are usually regarded as accurate annotations to optimize the model.

However, one inherent problem exists in the pseudo label based scene adaptation approaches. Pseudo labels usually suffer from the noise caused by the model trained on different data distribution (see Figure 1). The noisy label could compromise the subsequent learning. Although some existing works [Zou et al., 2018, Zou et al., 2019] have proposed to manually set the threshold to neglect the low-confidence pseudo labels, it is still challenging in several aspects: First, the value of the threshold is hard to be determined for different target domain. It depends on the similarity of the source domain and target domain, which is hard to estimate in advance. Second, the value of the threshold is also hard to be determined for different categories. For example, the objectives, such as traffic signs, have rarely appeared in the source domain. The overall confidence score for the rare category is relatively low. The high threshold may ignore the information of rare categories. Third, the threshold is also related to the location of the pixel. For example, the pixel in the center of objectives, such as cars, is relatively easy to predict, while the pixel on the objective edge usually faces ambiguous predictions. It reflects that the threshold should not only consider the confidence score but also the location of the pixel. In summary, every pixel in the segmentation map needs to be treated differently. The fixed threshold is hard to match the demand.

To address the mentioned challenges, we propose one simple and effective method for semantic segmentation adaption via modeling uncertainty, which could provide the pixel-wise threshold for the input image automatically. Without introducing extra parameters or modules, we formulate the uncertainty as the prediction variance. The prediction variance reflects the model uncertainty towards the prediction in a bootstrapping manner. Meanwhile, we explicitly involve the variance into the optimization objective, called variance regularization, which works as an automatic threshold and is compatible with the standard cross-entropy loss. The automatic threshold rectifies the learning from noisy labels and ensures the training in a coherent manner. Therefore, the proposed method could effectively exploit the domain-specific information offered by pseudo labels and takes advantage of the unlabeled target-domain data.

In a nutshell, our contributions are as follows:

– To our knowledge, we are among the first attempts to exploit the uncertainty estimation and enable the automatic threshold to learn from noisy pseudo labels. This is in contrast to most existing domain adaptation methods that directly utilize noisy pseudo labels or manually set the confidence threshold.
– Without introducing extra parameters or modules, we formulate the uncertainty as the prediction variance. Specifically, we introduce a new regularization term, variance regularization, which is compatible with the standard cross-entropy loss. The variance regularization works as the automatic threshold, and rectifies the learning from noisy pseudo labels.
– We verify the proposed method on two synthetic-to-real benchmarks and one cross-city benchmark. The proposed method has achieved significant improvements over the conventional pseudo label learning, yielding competitive performance to existing methods.

## 2 Related work

### 2.1 Semantic Segmentation Adaptation

The main challenge in unsupervised domain adaptation is different data distribution between the source domain and the target domain [Fu et al., 2015, Wang et al., 2018, Li et al., 2020b, Li et al., 2020c, Kang et al., 2020]. To deal with the challenge, some pioneering works [Hoffman et al., 2018, Wu et al., 2018] propose to transfer the visual style of the source-domain data to the target domain. In this way, the model could be trained on the labeled data with the target style. Similarly, some recent works leverage Adversarial Domain Adaptation [Tzeng et al., 2015, Ganin and Lempitsky, 2015, Luo et al., 2020] to transfer the source-domain images or features to multiple domains and intend to learn the domain-invariant features [Wu et al., 2019, Yue et al., 2019]. Furthermore, some works focus on the alignment among the middle activation of neural networks. Luo *et al.* [Luo et al., 2019a, Luo et al., 2019b] utilize the attention mechanism to

| road | sidewalk | building | wall | fence | pole | traffic lgt | traffic sgn | vegetation | terrain |
| sky | person | rider | car | truck | bus | train | motorcycle | bike | unlabeled |

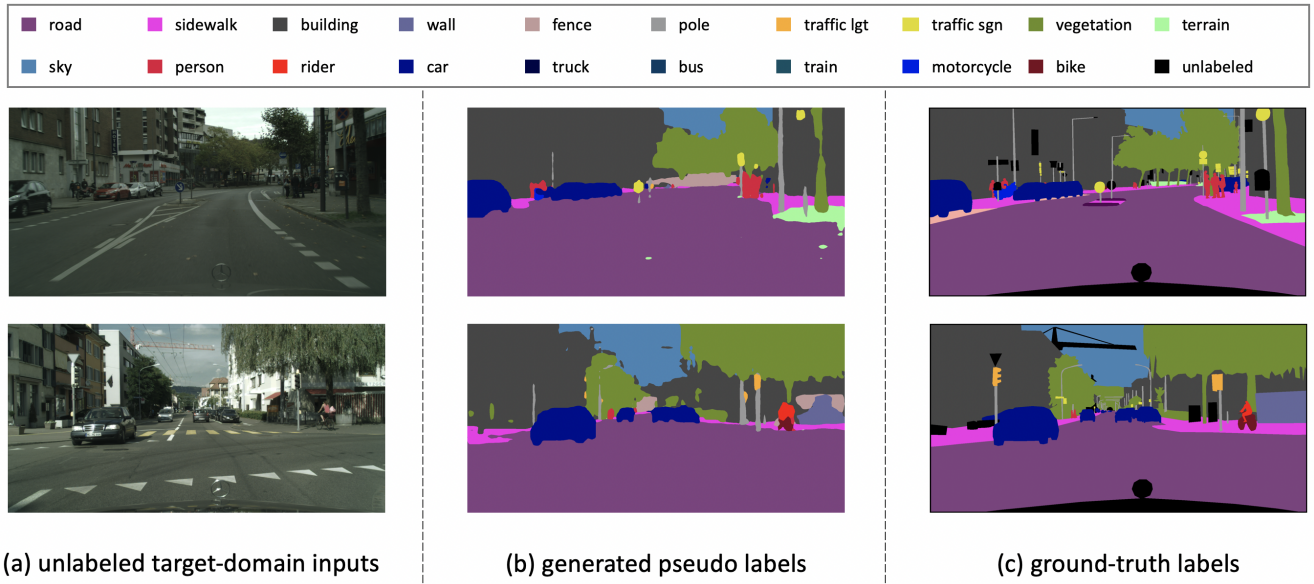(a) unlabeled target-domain inputs   (b) generated pseudo labels   (c) ground-truth labels

Fig. 1: Samples of the noisy pseudo labels on Cityscapes [Cordts et al., 2016]. We leverage the widely-used baseline model [Tsai et al., 2018] to generate pseudo labels. Despite the large area of correct prediction, the pseudo labels still suffer from the data distribution biases, and inevitably contains incorrect predictions. (Best viewed in *color*)

refine the feature alignment. Instead of modifying the visual appearance, the alignment between the high-level semantic features also attracts a lot of attention. Tsai *et al.* [Tsai et al., 2018, Tsai et al., 2019] propose to utilize the discriminator to demand the similar semantic outputs between two domains. In summary, this line of methods focuses on the alignment, learning the shared knowledge between the source and target domains. However, the domain-specific information is usually ignored, which is one of the keys to the adaptation in the target environment. Therefore, in this paper, we resort to another line of methods, which is based on pseudo label learning.

## 2.2 Pseudo label learning

Another line of semantic segmentation adaptation approaches utilizes the pseudo label to adapt the model to target domain [Zou et al., 2018, Zou et al., 2019, Zheng and Yang, 2020]. The main idea is close to the conventional semi-supervised learning approach, entropy minimization, which is first proposed to leverage the unlabeled data [Grandvalet and Bengio, 2005]. Entropy minimization encourages the model to give the prediction with a higher confidence score. In practice, Reed *et al.* [Reed et al., 2014] propose bootstrapping via entropy minimization and show the effectiveness on the object detection and emotion recognition. Furthermore, Lee *et al.* [Lee, 2013] exploit the trained model to predict pseudo labels for the unlabeled data, and then fine-tune the model as supervised learning methods to fully leverage the unla-

beled data. Recently, Pan *et al.* [Pan et al., 2019] utilize the pseudo label learning to minimize the distribution of target-domain data with the source-domain prototypes. For unsupervised semantic segmentation, Zou *et al.* [Zou et al., 2019, Zou et al., 2018] introduce the pseudo label strategy to the semantic segmentation adaptation and provide one comprehensive analysis on the regularization terms. In a similar spirit, Zheng *et al.* [Zheng and Yang, 2020] also apply the pseudo label to learn the domain-specific features, yielding competitive results. However, one inherent weakness of the pseudo label learning is that the pseudo label usually contains noisy predictions. Despite the fact that most pseudo labels are correct, wrong labels also exist, which could compromise the subsequent training. If the model is fine-tuned on the noisy label, the error would also be transferred to the adapted model. Different from existing works, we do not treat the pseudo labels equally and intend to rectify the learning from noisy labels. The proposed method explicitly predict the uncertainty of pseudo labels, when fine-tuning the model. The uncertainty could be regarded as an automatic threshold to adjust the learning from noisy labels.

## 2.3 Co-training

Co-training is a semi-supervised learning method, which demands two classifiers to learn complementary information [Blum and Mitchell, 1998]. Some domain adaptation works also explore a similar learning strategy. [Saito et al., 2018, Luo et al., 2019b] explicitly maximizes the discrepancy of

two classifiers by introducing one extra loss, i.e., the $L_{adv}$ in [Saito et al., 2018] and the $L_{weight}$ in [Luo et al., 2019b], to obtain complementary classifiers. [Saito et al., 2018] minimizes the feature discrepancy via adversarial training. Similarly, [Luo et al., 2019b] apply the classifier discrepancy on the discriminator loss to stabilize the training. In contrast, the proposed method enables the classifier discrepancy in nature, since we deploy two classifiers on different intermediate layers. We do not introduce such loss to encourage the classifier discrepancy. Otherwise, every pseudo label will be high-uncertainty. For instance, if the two classifiers output one identical category prediction, we will not punish the network. In contrast, [Saito et al., 2018] will punish the classifiers for enabling adversarial training. Besides, [Saito et al., 2018, Luo et al., 2019b] still use conventional segmentation loss and do not deal with noisy labels, when the proposed method uses the classifier discrepancy to rectify the pseudo label learning on segmentation.

## 2.4 Uncertainty Estimation

To address the noise, existing works have explored the uncertainty estimation from different aspects, such as the input data, the annotation and the model weights. In this work, we focus on the annotation uncertainty. Our target is to learn a model that could predict whether the annotation is correct, and learn from noisy pseudo labels. Among existing works, Bayesian networks are widely used to predict the uncertainty of weights in the network [Nielsen and Jensen, 2009]. In a similar spirit, Kendall *et al.* [Kendall and Gal, 2017] apply the Bayesian theory to the prediction of computer vision tasks, and intend to provide not only the prediction results but also the confidence of the prediction. Further, Yu *et al.* [Yu et al., 2019] explicitly model the uncertainty via an extra auxiliary branch, and involve the random noise into training. The model could explicitly estimate the feature mean as well as the prediction variance. Inspired by the above-mentioned works, we propose to leverage the prediction variance to formulate the uncertainty. There are two fundamental differences between previous works and ours: (1) We do not introduce extra modules or parameters to simulate the noise. Instead, we leverage the prediction discrepancy within the segmentation model. (2) We explicitly involve the uncertainty into the training target and adopt the adaptive method to learn the pixel-wise uncertainty map automatically. The proposed method does not need manually setting the threshold to enforce the pseudo label learning.

## 3 Methodology

In Section 3.1, we first provide the problem definition and denotations. We then revisit the conventional domain adaptation method based on the pseudo label and discuss the limitation of the pseudo label learning (see Section 3.2). To deal with the mentioned limitations, we propose to leverage the uncertainty estimation. In particular, we formulate the uncertainty as the prediction variance and provide one brief definition in Section 3.3, followed by the proposed variance regularization, which is compatible with the standard cross-entropy loss in Section 3.4. Besides, the implementation details are provided in Section 3.5.

### 3.1 Problem Definition

Given the labeled dataset $X_s = \{x_s^i\}_{i=1}^M$ from the source domain and the unlabeled dataset $X_t = \{x_t^j\}_{j=1}^N$ from the target domain, semantic segmentation adaptation intends to learn the projection function $F$, which maps the input image $X$ to the semantic segmentation $Y$. $M$ and $N$ denote the number of the labeled data and the unlabeled data. The source-domain semantic segmentation label $Y_s = \{y_s^i\}_{i=1}^M$ is provided for every labeled data of the source domain $X_s$, while the target-domain label $Y_t = \{y_t^j\}_{j=1}^N$ remains unknown during the training. The aim of unsupervised domain adaptation is to estimate the model parameter $\theta_t$, which could minimize the prediction bias on the target-domain inputs:

$$Bias(p_t) = \mathbb{E}[F(x_t^j|\theta_t) - p_t^j], \tag{1}$$

where $p_t$ is the ground-truth class probability of target data. Ideally, $p_t^j$ is one-hot vector and the maximum value of $p_t^j$ is 1. The ground-truth label $y_t^j = \arg\max p_t^j$. In contrast, $F(x_t^j|\theta_t)$ is the predicted probability distribution of $x_t^j$. When we minimize the prediction bias in Equation 1, the discrepancy between predicted results and the ground-truth probability is minimized.

### 3.2 Pseudo Label Learning Revisit

Pseudo label learning is to leverage the pseudo label to learn from the unlabeled data. The common practice contains two stages. The first stage is to generate the pseudo label for the unlabeled target-domain training data. The pseudo labels could be obtained via the model trained on source-domain data: $\hat{y}_t^j = \arg\max F(x_t^j|\theta_s)$. We note that $\theta_s$ is the model parameters learned from the source-domain training data. Therefore, the pseudo labels $\hat{y}_t$, are not accurate in nature due to different data distribution between $X_s$ and $X_t$. We denote $\hat{p}_t^j$ as the one-hot vector of $\hat{y}_t^j$. If the class index $c$ equals to $\hat{y}_t^j$, $\hat{p}_t^j(c) = 1$ else $\hat{p}_t^j(c) = 0$. The second stage of pseudo learning is to minimize the prediction bias. We could formulate the bias as the similar style of Equation 1:

$$Bias(p_t) = \mathbb{E}[F(x_t^j|\theta_t) - \hat{p}_t^j] + \mathbb{E}[\hat{p}_t^j - p_t^j]. \tag{2}$$
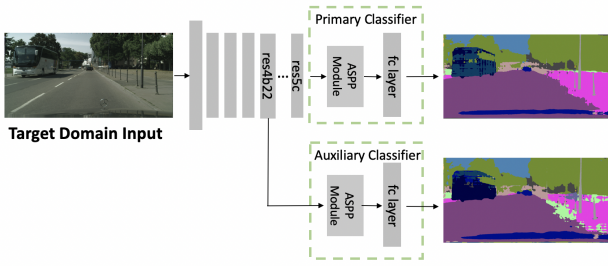
Fig. 2: Illustration of the two-classifier model based on Deeplab-v2 [Chen et al., 2017], which adopts ResNet-101 [He et al., 2016a] as backbone. We follow the previous works [Zhao et al., 2017, Tsai et al., 2018, Tsai et al., 2019, Luo et al., 2019a, Luo et al., 2019b, Zheng and Yang, 2020] to add an auxiliary classifier with the similar structure as the primary classifier. The auxiliary classifier takes the activation of the shallow layer *res4b22* as the input, while the primary classifier leverages that of *res5c*. The ASPP module denotes Atrous Spatial Pyramid Pooling layer [Chen et al., 2017], and the fc layer denotes the fully-connected layer. The original goal of two-classifier model is to evade the problem of gradient vanishing and help the training. In this work, we take one step further to leverage the prediction discrepancy of two classifiers as the uncertainty estimation.

The first term is the difference between the prediction and the pseudo label, while the second term is the error between the pseudo label and the ground-truth label. When fine-tuning the model in the second stage, we fix the pseudo label. Therefore, the second term is one constant. Existing methods usually optimize the first term as the pretext task. It equals to considering the pseudo labels $\hat{p}_t$ as true labels. Existing methods train the model parameter $\theta_t$ to minimize the bias between the prediction and pseudo labels. In practice, the cross-entropy loss is usually adopted [Zou et al., 2018, Zou et al., 2019, Zheng and Yang, 2020]. The objective could be formulated as:

$$L_{ce} = \mathbb{E}[-\hat{p}_t^j \log F(x_t^j|\theta_t)]. \tag{3}$$

**Discussion.** There are two advantages of pseudo label learning : First, the model is only trained on the target-domain data. The training data distribution is close to the testing data distribution, minoring the input distribution discrepancy. Second, despite the domain discrepancy, most pseudo labels are correct. Theoretically, the fine-tuned model could arrive the competitive performance with the fully-supervised model. However, one inherent problem exists that the pseudo label inevitably contains noise. The wrong annotations are transferred from the source model to the final model. Noisy pseudo label could largely compromise the training.

### 3.3 Uncertainty Estimation

To address the label noise, we model the uncertainty of the pseudo label via the prediction variance. Intuitively, we could formulate the variance of the prediction as:

$$Var(p_t) = \mathbb{E}[(F(x_t^j|\theta_t) - p_t^j)^2]. \tag{4}$$

Since $p_t$ remains unknown, one naive way is to utilize the pseudo label $\hat{p}_t$ to replace the $p_t$. The variance could be approximated as:

$$Var(p_t) \approx \mathbb{E}[(F(x_t^j|\theta_t) - \hat{p}_t^j)^2]. \tag{5}$$

However, in Equation 2, we have pushed $F(x_t^j|\theta_t)$ to $\hat{p}_t$. When optimizing the prediction bias, the variance in Equation 5 will also be minimized. It could not reflect the real prediction variance during training. In this paper, therefore, we adopt another approximation as:

$$Var(p_t) \approx \mathbb{E}[(F(x_t^j|\theta_t) - F_{aux}(x_t^j|\theta_t))^2], \tag{6}$$

where $F_{aux}(x_t|\theta_t)$ denotes the auxiliary classifier output of the segmentation model. As shown in Figure 2, we adopt the widely-used two-classifier model, which contains one primary classifier as well as one auxiliary classifier. We note that the extra auxiliary classifier could be viewed as a free lunch since most segmentation models, including PSPNet [Zhao et al., 2017] and the modified DeepLab-v2 in [Tsai et al., 2018, Tsai et al., 2019, Luo et al., 2019a, Zheng and Yang, 2020], contain the auxiliary classifier to solve the gradient vanish problem [He et al., 2016b] and help the training. In this paper, we further leverage the auxiliary classifier to estimate the variance. In practice, we utilize the KL-divergence of two classifier predictions as the variance:

$$D_{kl} = \mathbb{E}[F(x_t^j|\theta_t) \log(\frac{F(x_t^j|\theta_t)}{F_{aux}(x_t^j|\theta_t)})], \tag{7}$$

If two classifiers provide two different class predictions, the approximated variance will obtain one large value. It reflects the uncertainty of the model on the prediction. Besides, it is worthy to note that the proposed variance in Equation 7 is independent with the pseudo label $\hat{p}_t$.

**Discussion:** *What leads to the discrepancy of the primary classifier and the auxiliary classifier?* First of all, the main reason is different receptive fields. As shown in Figure 2, the auxiliary classifier is located at the relatively shallow layer, when the primary classifier learns from the deeper layer. The input activation is different between two classifiers, leading to the prediction difference. Second, the two classifiers have not been trained on the target-domain data. Therefore, both classifiers may have different biases to the target-domain data. Third, we apply the dropout function [Srivastava et al., 2014] to two classifiers, which also could lead to the different prediction during training. The prediction discrepancy helps us to estimate the uncertainty.

**Algorithm 1** Training Procedure of the Proposed Method

---

**Require:** The target domain dataset $X_t = \{x_t^j\}_{j=1}^N$; The generated
   pseudo label $\hat{Y}_t = \{\hat{y}_t^j\}_{j=1}^N$;
**Require:** The source-domain parameter $\theta_s$; The iteration number $T$.
 1: Initialize $\theta_t = \theta_s$;
 2: **for** $iteration = 1$ to $T$ **do**
 3:     Input $x_t^j$ to $F(\cdot|\theta_t)$, extract the prediction of two classifiers,
        calculate the prediction variance according to Equation 7:

$$D_{kl} = \mathbb{E}[F(x_t^j|\theta_t)\log(\frac{F(x_t^j|\theta_t)}{F_{aux}(x_t^j|\theta_t)})]. \qquad (8)$$

 4:     We fix the prediction variance, and calculate the original cross-
        entropy loss according to Equation 2, where $\hat{p}_t^j$ is the one-hot vec-
        tor of the pseudo label $\hat{y}_t^j$:

$$L_{ce} = \mathbb{E}[-\hat{p}_t^j \log F(x_t^j|\theta_t)]. \qquad (9)$$

 5:     We combine the prediction variance with the conventional ob-
        jective to obtain the rectified objective. Update the $\theta_t$ according to
        Equation 12:

$$L_{rect} = \mathbb{E}[exp\{-D_{kl}\}L_{ce} + D_{kl}] \qquad (10)$$

 6: **end for**
 7: **return** $\theta_t$.

---

## 3.4 Variance Regularization

In this paper, we propose the variance regularization term to
rectify the learning from noisy labels. It leverages the ap-
proximated variance introduced in Section 3.3. The rectified
objective could be formulated as:

$$L_{rect} = \mathbb{E}[\frac{1}{Var(p_t)}Bias(p_t) + Var(p_t)] \qquad (11)$$

It is worthy to note that we do not intend to minimize the
prediction bias under all conditions. If the prediction vari-
ance has received one large value, we will not punish the
prediction bias $Bias(p_t)$. Meanwhile, to prevent that the
model predicts the large variance all the time, as a trade-off,
we introduce the regularization term via adding $Var(p_t)$.
Besides, since $Var(p_t)$ could be zero, it may lead to the
problem of dividing by zero. To stabilize the training, we
adopt the policy in [Kendall and Gal, 2017] that replace
$1/Var$ as $exp(-Var)$. Therefore, the loss term could be
rewritten with the approximated terms as:

$$L_{rect} = \mathbb{E}[exp\{-D_{kl}\}L_{ce} + D_{kl}]. \qquad (12)$$

The training procedure of the proposed method is summa-
rized in Algorithm 1. In practice, we utilize the parameter
$\theta_s$ learned in the source-domain dataset to initialize the $\theta_t$.
In every iteration, we calculate the prediction variance as
well as the cross-entropy loss for the given inputs. We utilize
the $L_{rect}$ to update the $\theta_t$. The training cost of the rectified
objective approximately equals to the conventional pseudo
label learning, since no extra modules are introduced.

**Discussion:** *What are the advantages of the proposed vari-
ance regularization?* First, the proposed variance regular-
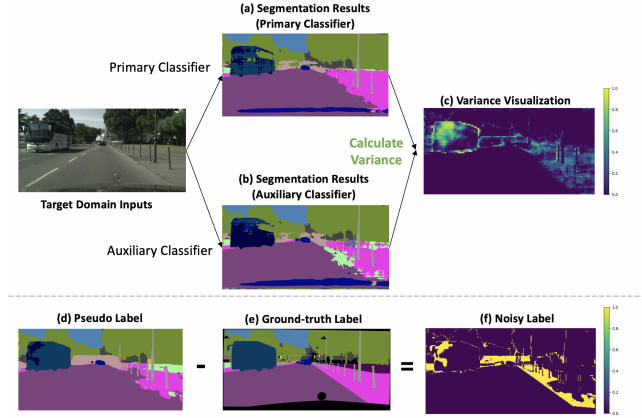ization does not introduce extra parameters or modules to



Fig. 3: Illustration of the prediction variance between two
classifiers, *i.e.*, the primary classifier and the auxiliary clas-
sifier. The areas, where have ambiguous predictions, obtain
large value of the prediction variance. Meanwhile, we could
observe that the high-variance area has considerable over-
laps with the noise in the pseudo label. (Best viewed in
*color*)

model the uncertainty. Different from [Yu et al., 2019], we
do not explicitly introduce the Gaussian noise or extra branches.
Instead, we leverage the prediction variance of the model it-
self. Second, the proposed variance regularization has good
scalability. If the variance equals to zero, the optimization
loss degrades to the objective of the conventional pseudo
learning and the model will focus on minimizing the pre-
diction bias only. In contrast, when the value of variance is
high, the model is prone to neglect the bias and skip am-
biguous pseudo labels; Third, the proposed variance regu-
larization has the same shape of the prediction, and could
works as the pixel-wise threshold of the pseudo label. As
shown in Figure 3, we could observe that the noise usually
exists in the area with high variance. The proposed rectified
loss assigns different thresholds to different areas. For exam-
ple, for the location with coherent predictions, the variance
regularization drives the model trust pseudo labels. For the
area with ambiguous predictions, the variance regularization
drives the model to neglect pseudo labels. Different from
existing works that set the unified threshold for all training
samples, the proposed pseudo label could provide more ac-
curate and adaptive threshold for every pixel.

## 3.5 Implementation

**Network Architecture.** In this work, we utilize the widely-
used Deeplab-v2 [Chen et al., 2017] as the baseline model,
which adopts the ResNet-101 [He et al., 2016a] as the back-
bone model. We follow most existing works [Tsai et al.,
2018, Tsai et al., 2019, Luo et al., 2019a, Luo et al., 2019b,

| Datasets | GTA5 | SYNTHIA | Cityscapes | Oxford RobotCar |
|----------|------|---------|------------|-----------------|
| #Train | 24,966 | 9,400 | 2,975 | 894 |
| #Test | - | - | 500 | 271 |
| #Category | 19 | 16 | 19 | 9 |
| Synthetic | ✓ | ✓ | × | × |

Table 1: List of categories and number of images in four datasets, *i.e.*, GTA5 [Richter et al., 2016], SYNTHIA [Ros et al., 2016], Cityscapes [Cordts et al., 2016] and Oxford RobotCar [Maddern et al., 2017].

Zheng and Yang, 2020] to add one auxiliary classifier. The auxiliary classifier has similar structure with the primary classifier, including one Atrous Spatial Pyramid Pooling (ASPP) module [Chen et al., 2017] and one fully-connected layer. The auxiliary classifier is added after the *res4b22* layer. We also insert the dropout layer [Srivastava et al., 2014] before the fully-connected layer, and the dropout rate is 0.1.

**Pseudo Label.** To verify the effectiveness of the proposed method, we deploy two existing methods, *i.e.*, AdaptSegNet [Tsai et al., 2018] and MRNet [Zheng and Yang, 2020], to generate the pseudo labels of the target-domain dataset.

– AdaptSegNet [Tsai et al., 2018] is one widely-adopted baseline model, which utilize the adversarial training to align the semantic outputs.
– MRNet [Zheng and Yang, 2020] is one recent work, which leverages the memory module to regularize the model training, especially for the target-domain data.

Specifically, MRNet arrives superior performance to Adapt-SegNet in terms of mIoU on three benchmarks. Therefore, if not specific, we adopt the pseudo label generated by the stronger baseline, *i.e.*, MRNet. **It is worth mentioning that we do not use source-domain training data. In practice, we fine-tune the model only on the target-domain training data with pseudo labels.**

**Training Details.** The input image is resized to $1280 \times 640$ with scale jittering from $[0.8, 1.2]$, and then we randomly crop $512 \times 256$ for training. Horizontal flipping is applied with the possibility of $50\%$. We train the model with mini-batch size of 9, and the parameters of batch normalization layers are also fine-tuned. The learning rate is set to $0.0001$. Following [Zhao et al., 2017, Zhang et al., 2019a, Zhang et al., 2020], we deploy the ploy learning rate policy by multiplying the factor $(1 - \frac{iter}{total-iter})^{0.9}$. The total iteration is set as $100k$ iterations and we adopt the early-stop strategy. We stop the training after 50k iterations. When inference, we follow [Zheng and Yang, 2020] to combine the output of both classifier as the final result. $Output = \arg\max(F(x_t^j|\theta_t) + 0.5F_{aux}(x_t^j|\theta_t))$. Our implementation is based on Pytorch [Paszke et al., 2017].

## 4 Experiment

### 4.1 Datasets and Evaluation Metric

**Datasets.** To simplify, we denote the test setting as A → B, where A represents the labeled source domain and B denotes the unlabeled target domain. We evaluate the proposed method on two widely-used synthetic-to-real benchmarks: *i.e.*, GTA5 [Richter et al., 2016]→Cityscapes [Cordts et al., 2016] and SYNTHIA5 [Ros et al., 2016]→Cityscapes [Cordts et al., 2016]. Both source dataset, *i.e.*, GTA5 and SYNTHIA are the synthetic datasets, and the corresponding annotation is easy to obtain. Specifically, the GTA5 dataset is collected from a video game, which contains $24,966$ images for training. The SYNTHIA dataset is rendered from a virtual city and comes with pixel-level segmentation annotations, containing $9,400$ training images. The realistic dataset, Cityscapes, collect street-view scenes from 50 different cities, which contains $2,975$ training images and $500$ images for validation. Besides, we also evaluate the performance on the cross-city benchmark, *i.e.*, Cityscapes [Cordts et al., 2016]→Oxford RobotCar [Maddern et al., 2017]. We utilize the annotation of Cityscapes training images in this setting. The Oxford RobotCar dataset serves as the unlabeled target domain, containing $894$ training images and $271$ validation images. We note that this setting is challenging in different weather conditions. Oxford RobotCar is collected in the rainy days, while the Cityscapes dataset is mostly collected in the sunny days. The differences between datasets are listed in Table 1.

**Evaluation Metric.** We report pre-class IoU and mean IoU over all classes. For SYNTHIA → Cityscapes, due the limited annotated classes in the source dataset, we report the results based on 13 categories as well as 16 categories with three small-scale categories. For Cityscapes → Oxford Robot-Car, we follow the setting in [Tsai et al., 2019] and report 9 pre-class IoU as well as the mIoU accuracy.

### 4.2 Comparisons with state-of-the-art methods

**Synthetic-to-real.** We compare the proposed method with other recent semantic segmentation adaptation methods that have reported the results or can be re-implemented by us on three benchmarks. For a fair comparison, we mainly compare the results based on the same network structure, *i.e.*, DeepLabv2. The competitive methods cover a wide range of approaches and could be roughly categorised according to the usage of pseudo label: AdaptSegNet [Tsai et al., 2018], SIBAN [Luo et al., 2019a], CLAN [Luo et al., 2019b], APODA [Yang et al., 2020] and PatchAlign [Tsai et al., 2018] do not leverage the pseudo labels and focus on aligning the distribution between the source domain and the target domain; CBST [Zou et al., 2018], MRKLD [Zou et al., 2019], and our implemented MRNet+Pseudo are based on the pseudo

| Method | Road | SW | Build | Wall | Fence | Pole | TL | TS | Veg. | Terrain | Sky | PR | Rider | Car | Truck | Bus | Train | Motor | Bike | mIoU |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Source | 75.8 | 16.8 | 77.2 | 12.5 | 21.0 | 25.5 | 30.1 | 20.1 | 81.3 | 24.6 | 70.3 | 53.8 | 26.4 | 49.9 | 17.2 | 25.9 | 6.5 | 25.3 | 36.0 | 36.6 |
| AdaptSegNet [Tsai et al., 2018] | 86.5 | 36.0 | 79.9 | 23.4 | 23.3 | 23.9 | 35.2 | 14.8 | 83.4 | 33.3 | 75.6 | 58.5 | 27.6 | 73.7 | 32.5 | 35.4 | 3.9 | 30.1 | 28.1 | 42.4 |
| SIBAN [Luo et al., 2019a] | 88.5 | 35.4 | 79.5 | 26.3 | 24.3 | 28.5 | 32.5 | 18.3 | 81.2 | 40.0 | 76.5 | 58.1 | 25.8 | 82.6 | 30.3 | 34.4 | 3.4 | 21.6 | 21.5 | 42.6 |
| CLAN [Luo et al., 2019b] | 87.0 | 27.1 | 79.6 | 27.3 | 23.3 | 28.3 | 35.5 | 24.2 | 83.6 | 27.4 | 74.2 | 58.6 | 28.0 | 76.2 | 33.1 | 36.7 | 6.7 | 31.9 | 31.4 | 43.2 |
| APODA [Yang et al., 2020] | 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 |
| PatchAlign [Tsai et al., 2019] | **92.3** | 51.9 | 82.1 | 29.2 | 25.1 | 24.5 | 33.8 | 33.0 | 82.4 | 32.8 | 82.2 | 58.6 | 27.2 | 84.3 | 33.4 | 46.3 | 2.2 | 29.5 | 32.3 | 46.5 |
| AdvEnt [Vu et al., 2019] | 89.4 | 33.1 | 81.0 | 26.6 | **26.8** | 27.2 | 33.5 | 24.7 | 83.9 | 36.7 | 78.8 | 58.7 | 30.5 | 84.8 | **38.5** | 44.5 | 1.7 | 31.6 | 32.4 | 45.5 |
| Source | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 29.2 |
| FCAN [Zhang et al., 2018] | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 46.6 |
| Source | 71.3 | 19.2 | 69.1 | 18.4 | 10.0 | 35.7 | 27.3 | 6.8 | 79.6 | 24.8 | 72.1 | 57.6 | 19.5 | 55.5 | 15.5 | 15.1 | 11.7 | 21.1 | 12.0 | 33.8 |
| CBST [Zou et al., 2018] | 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 |
| MRKLD [Zou et al., 2019] | 91.0 | **55.4** | 80.0 | 33.7 | 21.4 | 37.3 | 32.9 | 24.5 | 85.0 | 34.1 | 80.8 | 57.7 | 24.6 | 84.1 | 27.8 | 30.1 | **26.9** | 26.0 | 42.3 | 47.1 |
| Source | 51.1 | 18.3 | 75.8 | 18.8 | 16.8 | 34.7 | 36.3 | 27.2 | 80.0 | 23.3 | 64.9 | 59.2 | 19.3 | 74.6 | 26.7 | 13.8 | 0.1 | 32.4 | 34.0 | 37.2 |
| MRNet [Zheng and Yang, 2020] | 89.1 | 23.9 | 82.2 | 19.5 | 20.1 | 33.5 | 42.2 | 39.1 | 85.3 | 33.7 | 76.4 | 60.2 | 33.7 | 86.0 | 36.1 | 43.3 | 5.9 | 22.8 | 30.8 | 45.5 |
| MRNet+Pseudo | 90.5 | 35.0 | 84.6 | 34.3 | 24.0 | 36.8 | 44.1 | 42.7 | 84.5 | 33.6 | **82.5** | 63.1 | 34.4 | 85.8 | 32.9 | 38.2 | 2.0 | 27.1 | 41.8 | 48.3 |
| MRNet+Ours | 90.4 | 31.2 | **85.1** | **36.9** | 25.6 | **37.5** | **48.8** | **48.5** | 85.3 | 34.8 | 81.1 | **64.4** | **36.8** | **86.3** | 34.9 | **52.2** | 1.7 | 29.0 | **44.6** | **50.3** |

Table 2: Quantitative results on GTA5 → Cityscapes. We present pre-class IoU and mIoU. The best accuracy in every column is in **bold**.

| Method | Road | SW | Build | Wall* | Fence* | Pole* | TL | TS | Veg. | Sky | PR | Rider | Car | Bus | Motor | Bike | mIoU* | mIoU |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Source | 55.6 | 23.8 | 74.6 | – | – | – | 6.1 | 12.1 | 74.8 | 79.0 | 55.3 | 19.1 | 39.6 | 23.3 | 13.7 | 25.0 | 38.6 | – |
| SIBAN [Luo et al., 2019a] | 82.5 | 24.0 | 79.4 | – | – | – | 16.5 | 12.7 | 79.2 | 82.8 | 58.3 | 18.0 | 79.3 | 25.3 | 17.6 | 25.9 | 46.3 | – |
| PatchAlign [Tsai et al., 2019] | 82.4 | 38.0 | 78.6 | 8.7 | 0.6 | 26.0 | 3.9 | 11.1 | 75.5 | 84.6 | 53.5 | 21.6 | 71.4 | 32.6 | 19.3 | 31.7 | 46.5 | 40.0 |
| AdaptSegNet [Tsai et al., 2018] | 84.3 | 42.7 | 77.5 | – | – | – | 4.7 | 7.0 | 77.9 | 82.5 | 54.3 | 21.0 | 72.3 | 32.2 | 18.9 | 32.3 | 46.7 | – |
| CLAN [Luo et al., 2019b] | 81.3 | 37.0 | 80.1 | – | – | – | 16.1 | 13.7 | 78.2 | 81.5 | 53.4 | 21.2 | 73.0 | 32.9 | 22.6 | 30.7 | 47.8 | – |
| CCM [Li et al., 2020a] | 79.6 | 36.4 | 80.6 | 13.3 | 0.3 | 25.5 | 22.4 | 14.9 | 81.8 | 77.4 | 56.8 | 25.9 | 80.7 | 45.3 | **29.9** | 52.0 | 52.9 | 45.2 |
| APODA [Yang et al., 2020] | 86.4 | 41.3 | 79.3 | – | – | – | 22.6 | 17.3 | 80.3 | 81.6 | 56.9 | 21.0 | 84.1 | **49.1** | 24.6 | 45.7 | 53.1 | – |
| AdvEnt [Vu et al., 2019] | 85.6 | 42.2 | 79.7 | 8.7 | 0.4 | 25.9 | 5.4 | 8.1 | 80.4 | **84.1** | 57.9 | 23.8 | 73.3 | 36.4 | 14.2 | 33.0 | 48.0 | 41.2 |
| Source | 64.3 | 21.3 | 73.1 | 2.4 | 1.1 | 31.4 | 7.0 | 27.7 | 63.1 | 67.6 | 42.2 | 19.9 | 73.1 | 15.3 | 10.5 | 38.9 | 40.3 | 34.9 |
| CBST [Zou et al., 2018] | 68.0 | 29.9 | 76.3 | 10.8 | 1.4 | 33.9 | 22.8 | 29.5 | 77.6 | 78.3 | 60.6 | 28.3 | 81.6 | 23.5 | 18.8 | 39.8 | 48.9 | 42.6 |
| MRKLD [Zou et al., 2019] | 67.7 | 32.2 | 73.9 | 10.7 | 1.6 | **37.4** | 22.2 | **31.2** | 80.8 | 80.5 | 60.8 | **29.1** | 82.8 | 25.0 | 19.4 | 45.3 | 50.1 | 43.8 |
| Source | 44.0 | 19.3 | 70.9 | 8.7 | 0.8 | 28.2 | 16.1 | 16.7 | 79.8 | 81.4 | 57.8 | 19.2 | 46.9 | 17.2 | 12.0 | 43.8 | 40.4 | 35.2 |
| MRNet [Zheng and Yang, 2020] | 82.0 | 36.5 | 80.4 | 4.2 | 0.4 | 33.7 | 18.0 | 13.4 | 81.1 | 80.8 | 61.3 | 21.7 | 84.4 | 32.4 | 14.8 | 45.7 | 50.2 | 43.2 |
| MRNet+Pseudo | 83.1 | 38.2 | 81.7 | 9.3 | 1.0 | 35.1 | 30.3 | 19.9 | **82.0** | 80.1 | 62.8 | 21.1 | 84.4 | 37.8 | 24.5 | **53.3** | 53.8 | 46.5 |
| MRNet+Ours | **87.6** | 41.9 | **83.1** | **14.7** | **1.7** | 36.2 | **31.3** | 19.9 | 81.6 | 80.6 | **63.0** | 21.8 | **86.2** | 40.7 | 23.6 | 53.1 | **54.9** | **47.9** |

Table 3: Quantitative results on SYNTHIA → Cityscapes. We present pre-class IoU, mIoU and mIoU*. mIoU and mIoU* are averaged over 16 and 13 categories, respectively. The best accuracy in every column is in **bold**.

| Method | road | sidewalk | building | light | sign | sky | person | automobile | two-wheel | mIoU |
|---|---|---|---|---|---|---|---|---|---|---|
| Source | 79.2 | 49.3 | 73.1 | 55.6 | 37.3 | 36.1 | 54.0 | 81.3 | 49.7 | 61.9 |
| AdaptSegNet [Tsai et al., 2018] | 95.1 | 64.0 | 75.7 | 61.3 | 35.5 | 63.9 | 58.1 | 84.6 | 57.0 | 69.5 |
| PatchAlign [Tsai et al., 2019] | 94.4 | 63.5 | 82.0 | 61.3 | 36.0 | 76.4 | 61.0 | 86.5 | 58.6 | 72.0 |
| MRNet [Zheng and Yang, 2020] | **95.9** | 73.5 | 86.2 | 69.3 | 31.9 | 87.3 | 57.9 | 88.8 | **61.5** | 72.5 |
| MRNet+Pseudo | 95.1 | 72.5 | 87.0 | 72.2 | 37.4 | 87.9 | **63.4** | **90.5** | 58.9 | 73.9 |
| MRNet+Ours | **95.9** | **73.7** | **87.4** | **72.8** | **43.1** | **88.6** | 61.7 | 89.6 | 57.0 | **74.4** |

Table 4: Quantitative results on the cross-city benchmark: Cityscapes → Oxford RobotCar. The best accuracy in every column is in **bold**.

label learning to fully exploit the unlabeled target-domain data.

First of all, we consider the widely-used GTA5 → Cityscapes benchmark. Table 2 shows that: (1) The proposed method arrives the state-of-the-art results 50.3% mIoU, which surpasses other methods. Besides, the proposed method also yields the competitive performance in terms of the pre-class IoU. (2) Comparing to our baseline, i.e., MRNet+Pseudo (48.3% mIoU), which adopts the conventional pseudo learning, the proposed method (50.3% mIoU) gains +2.0% mIoU improvement. It verifies the effectiveness of the proposed method in rectifying the learning from the noisy pseudo label. The variance regularization plays an important role in achieving this result; (3) Meanwhile, we could observe that the proposed method outperforms the source-domain model, i.e., MRNet (45.5% mIoU), which provides the pseudo label, 4.8 mIoU. It verifies the effectiveness of the pseudo label learning that push the model to be confident about the prediction. If most pseudo labels are correct, the pseudo label learning could effectively boost the target-domain performance. (4) The proposed method also surpasses the other domain alignment method by a relatively large margin. For example, the modified AdaptSegNet, i.e., PatchAlign [Tsai et al., 2018], leverages the patch-level information, yielding 46.5%, which is inferior to ours. (5) Without using the prior knowledge, the proposed method is also superior to other pseudo label learning works, i.e., CBST [Zou et al., 2018] and MRKLD [Zou et al., 2019]. CBST [Zou et al., 2018] introduces the location knowledge, e.g., sky is always in the upper bound of the image. In this work, we do not apply such prior knowledge, but we note that the prior knowledge is compatible with our method.

We observe a similar result on SYNTHIA → Cityscapes (see Table 3). Following the setting in [Zou et al., 2018, Zou et al., 2019], we include the mIoU results of 13 categories as well as 16 categories, which also calculate IoU of other three small-scale objectives, i.e., Wall, Fence and Pole. The

proposed method has achieved $47.9$ mIoU of 16 categories and $54.9$ mIoU* of 13 categories. Comparing to the baseline, MRNet+Pseudo, we yield $+1.4\%$ mIoU and $+1.1\%$ mIoU* improvement. Meanwhile, the proposed method also outperforms the second best method, *i.e.*, APODA [Yang et al., 2020], $1.8\%$ mIoU*.

**Cross-city.** We further evaluate the proposed method on the cross-city benchmark, *i.e.*, Cityscapes → Oxford RobotCar. Both of the source-domain and target-domain datasets are collected in the real-world scenario. We follow the settings in [Tsai et al., 2019] to report IoU of the shared 9 categories between the two datasets. As shown in Table 4, the proposed method arrives $74.4\%$ mIoU. Comparing to the baseline, *i.e.*, MRNet+Pseudo ($73.9\%$), the improvement ($+0.5\%$) on the cross-city benchmark is relatively limited. Therefore, the baseline, MRNet+Pseudo, also could obtain competitive results by directly utilizing all pseudo labels. Besides, it is worthy to note that the proposed method has arrived the 6 of 9 best pre-class IoU accuracy, and achieved $+5.7\%$ on the class of traffic sign, which is a small-scale objective.

**Visualization.** As shown in Figure 4, we provide the qualitative results of the semantic segmentation adaptation on all three benchmarks. Comparing to the source model, the pseudo label learning could significantly improve the performance. Besides, in contrast with the baseline method with conventional pseudo label learning, we observe that the proposed variance regularization has better scalability to small-scale objectives, such as traffic signs and poles. It is because that the noisy pseudo label usually contains the error of predicting the rare category to the common category, *i.e.*, large-scale objectives. The proposed method rectifies the learning from such mistakes, yielding more reasonable segmentation prediction.

## 4.3 Further Evaluations

**Variance Regularization vs. Handcrafted Threshold.** The proposed variance regularization is free from setting the threshold. To verify the effectiveness of the variance regularization, we also compare the conventional pseudo label learning with different thresholds. As shown in Table 5, the proposed regularization arrives the superior performance to the hand-crafted threshold. It is due to that the variance regularization could be viewed as a dynamic threshold, providing different thresholds for different pixels in the same image. For the coherent predictions, the model is prone to learning the pseudo label and maximizing the impact of such labels. For the incoherent results, the model is prone to neglecting the pseudo label automatically and minimizing the negative effect of noisy labels. The best handcrafted threshold is to neglect the label with the prediction score $\leq 0.90$, yielding $48.4\%$ mIoU. In contrast, the proposed method achieves $50.3\%$ mIoU with $+1.9\%$ increment.

| Methods | Threshold | mIoU |
|---|---|---|
| MRNet [Zheng and Yang, 2020] | - | 45.5 |
| Pseudo Learning | > 0.99 | 45.5 |
| Pseudo Learning | > 0.95 | 47.2 |
| Pseudo Learning | > 0.90 | 48.4 |
| Pseudo Learning | > 0.80 | 48.1 |
| Pseudo Learning | > 0.70 | 48.2 |
| Pseudo Learning | > 0.00 | 48.3 |
| Ours | - | 50.3 |

Table 5: Variance Regularization vs. Handcrafted Threshold. The proposed method is free from hand-crafted threshold. '$> k$' denotes that we only utilize the label confidence $> k$ to train the model. We report the mIoU accuracy on GTA5 → Cityscapes.

| Methods | Pseudo Label | mIoU |
|---|---|---|
| AdaptSegNet [Tsai et al., 2018] | - | 42.4 |
| Pseudo Learning | AdaptSegNet | 46.8 |
| Ours | AdaptSegNet | 47.4 |
| MRNet [Zheng and Yang, 2020] | - | 45.5 |
| Pseudo Learning | MRNet | 48.3 |
| Ours | MRNet | 50.3 |

Table 6: Ablation study of the impact of different pseudo labels. The model name in the 'Pseudo Label' column denotes that we deploy the pseudo label generated by the corresponding model.

| Dropout Rate | mIoU |
|---|---|
| Pseudo Learning | 48.3 |
| droprate = 0 | 49.6 |
| droprate = 0.1 | 50.3 |
| droprate = 0.3 | 50.1 |
| droprate = 0.5 | 50.1 |
| droprate = 0.7 | 50.0 |

Table 7: Ablation study of dropout rate on GTA5 → Cityscapes.

**Could the proposed method work on the pseudo label generated by other models (*e.g.*, with more noise)?** To verify the scalability of the proposed method, we adopt the AdaptSegNet [Tsai et al., 2018] to generate pseudo labels. AdaptSegNet is inferior to MRNet in terms of the mIoU on GTA5 → Cityscapes. As shown in Table 6, the proposed method still could learn from the label generated by AdaptSegNet, improving the performance from $42.4\%$ to $47.4\%$. Meanwhile, the proposed method is also superior to the baseline method with the conventional pseudo learning ($46.8\%$ mIoU).

**Training Convergence.** As shown in Figure 6, the conventional pseudo label learning (orange line) is prone to overfit all pseudo labels, including the noisy label. Therefore, the training loss converges to zero. In contrast, the proposed
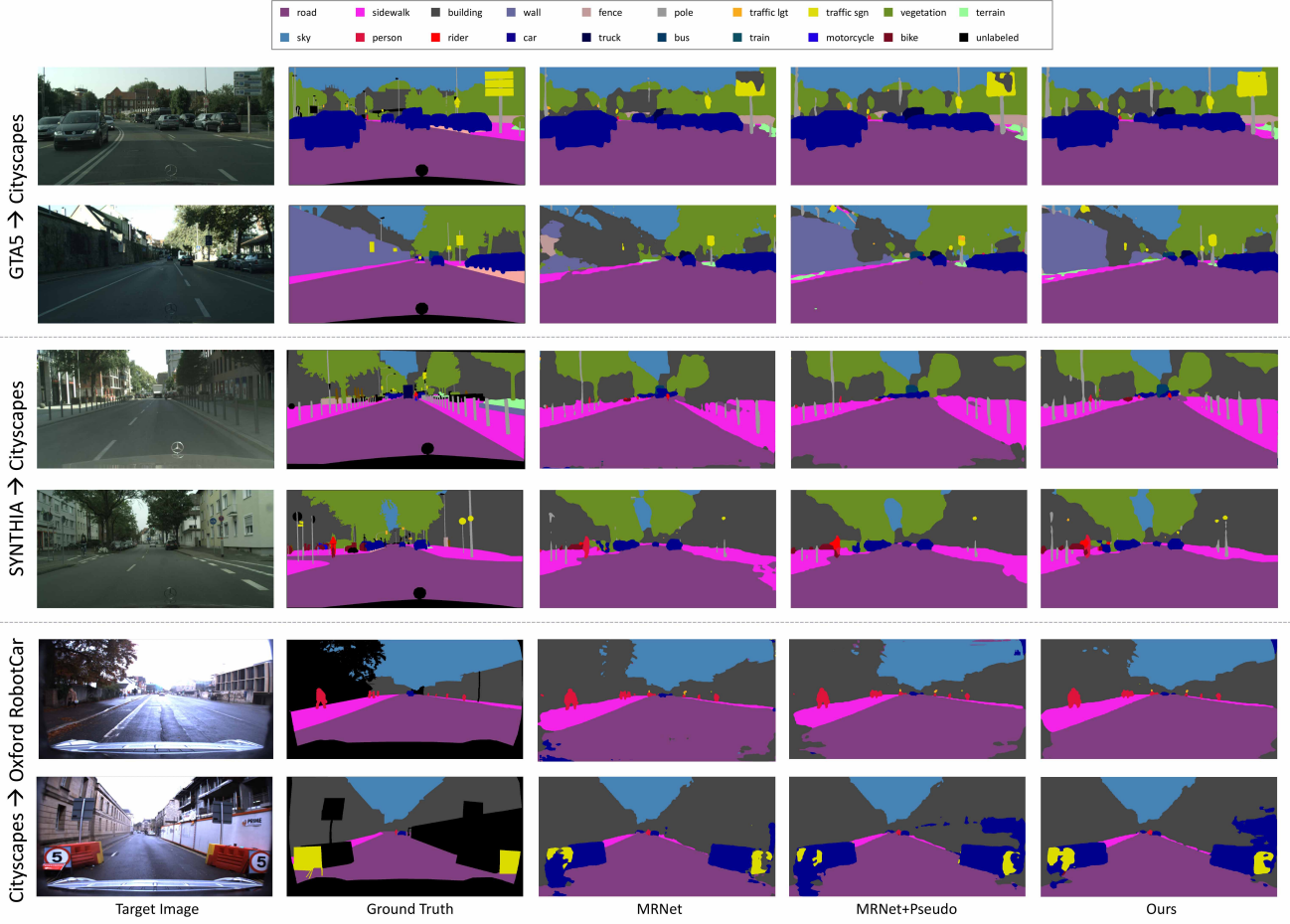
Fig. 4: Qualitative results of semantic segmentation adaptation on GTA5 → Cityscapes, SYNTHIA → Cityscapes and Cityscapes → Oxford RobotCar. We show the original target image, the ground-truth segmentation, the output of the source model, *i.e.*, MRNet, and the baseline, *i.e.*, MRNet+Pseudo. Our results are in the right column. (Best viewed in *color*).

| Methods | Right-prediction Certainty | Wrong-prediction Certainty | Uncertainty Gap |
|---|---|---|---|
| MC-dropout 0.5 | 0.9945 | 0.9733 | 0.0212 |
| MC-dropout 0.7 | 0.9870 | 0.9396 | 0.0474 |
| MC-dropout 0.9 | 0.9486 | 0.8118 | 0.1368 |
| Ours | 0.9767 | 0.8410 | 0.1357 |
| Ours + dropout 0.5 | 0.9673 | 0.8065 | **0.1608** |

Table 8: Comparison with Monte Carlo Dropout.

| $\alpha$ | $\beta$ | mIoU |
|---|---|---|
| 1.0 | 0.0 | 49.3 |
| 0.0 | 1.0 | 47.8 |
| 1.0 | 1.0 | 50.1 |
| 1.0 | 0.5 | 50.3 |

Table 10: Sensitivity of inference weighting.

| Distance Functions | mIoU |
|---|---|
| $\mathbb{E}[(F(x_t^j\|\theta_t) - F_{aux}(x_t^j\|\theta_t))^2]$ | 49.6 |
| $\mathbb{E}[F_{aux}(x_t^j\|\theta_t) \log(\frac{F_{aux}(x_t^j\|\theta_t)}{F(x_t^j\|\theta_t)})]$ | 49.4 |
| $\mathbb{E}[F(x_t^j\|\theta_t) \log(\frac{F(x_t^j\|\theta_t)}{F_{aux}(x_t^j\|\theta_t)})]$ | 50.3 |

Table 9: Ablation study of distance functions on GTA5 → Cityscapes.

method (blue line) also converges, but does not force the loss to be zero. It is because that we provide the variance regularization term, which could punish the wrong prediction for the uncertain pseudo labels with flexibility.

**Effect of Dropout.** The proposed method is not very sensitive to the dropout rate. As shown in Table 7, we could observe two points: 1) The dropout function is not the main reason for variance of the predictions. Without dropout function ($p = 0$), the proposed method still could achieve $49.6\%$ mIoU, which is better than the conventional pseudo label

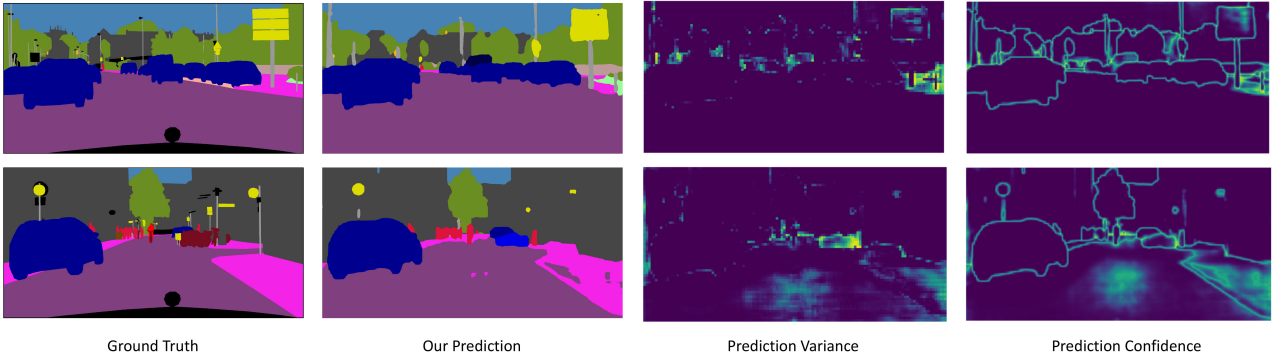| Ground Truth | Our Prediction | Prediction Variance | Prediction Confidence |

Fig. 5: Qualitative results of the discrepancy between the prediction variance and the prediction confidence. We could observe that the prediction variance used in this work has more overlaps with the ambiguous predictions, while the prediction confidence usually focuses on the edge of the two different classes. (Best viewed in *color*).
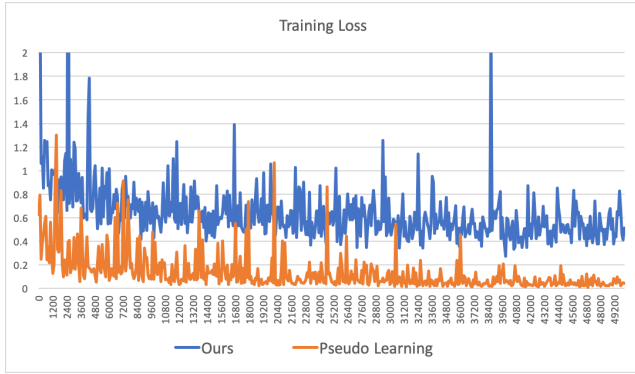


Fig. 6: The training loss of the proposed method and the pseudo label learning. The pseudo label learning is prone to over-fit all pseudo label, and the training loss converges to zero. In contrast, the proposed method would converge to one non-zero constant while training.

learning. 2) With a propose dropout rate, the proposed method could generally achieve better results around $50\%$ mIoU.

**Uncertainty of High-confidence Predictions.** We analyze the variance of high-confidence predictions on Cityscape. Specifically, we calculate the average uncertainty of right-assigned and wrong-assigned prediction with a confidence score$> 0.95$. Here we use the metric $exp\{-D_{kl}\}$ in Equation 12 to report the variance value. The high value means low uncertainty. The average variance of right-assigned high-confidence labels is $0.9901$, when the average variance of wrong-assigned high-confidence labels is $0.9332$. We could see one significant variance gap between the right-assigned labels and wrong-assigned labels, even if they all achieve a high confidence score. The result verifies that the variance value could reflect the difference between wrong-assigned labels and right-assigned labels.

**Comparison with Monte Carlo Dropout.** Monte Carlo Dropout (MC-Dropout) [Gal and Ghahramani, 2016] activates the dropout function when inference to obtain various predictions. Here we compare the ability of representing the uncertainty of the proposed method and MC-Dropout. For a fair comparison, we just replace the prediction of the aux classifier with the main classifier $F_{drop}$ with MC dropout rate of $\{0.5, 0.7, 0.9\}$.

$$D_{mc} = \mathbb{E}[F(x_t^j|\theta_t) \log(\frac{F(x_t^j|\theta_t)}{F_{drop}(x_t^j|\theta_t)})]. \quad (13)$$

Since the prediction score could not reflect the ground-truth uncertainty, we introduce one new metric called uncertainty gap as indicator. Uncertainty gap is the variance difference of right predictions and wrong predictions. Generally, we hope that the right prediction obtains low uncertainty value, while the wrong prediction obtains high uncertainty value. In practice, we use the $exp(-D)$ to keep the value in [0,1]. As shown in Table 8, the proposed method obtains $0.1357$ variance gap, which is competitive to MC-dropout with $0.9$ drop rate. The proposed method is also complementary to MC-dropout. The proposed method with MC-dropout could further boost the uncertainty gap. Meanwhile, it is worth noting that the proposed method directly leverages the variance of both main and auxiliary classifiers without multiple inferences, which can largely save the test time.

**Effect of Distance Functions.** In fact, KL-divergence is an alternative option for variance calculation. We could swap the main and aux classifiers to calculate the distance or use mean-square error (MSE). Here we add one experiment to compare common distance functions (see Table 9). First, we could observe that the model is not very sensitive to the distance metric, since the performances are close. Second, the KL-divergence used in Method is slightly better than swapping the predictions and MSE distance.

**Effect of Inference Weighting.** Inference weighting is one practical trick to combine the predictions of both main and auxiliary classifiers. Generally, the main classifier could achieve better performance, so we give the prediction of the main classifier a larger weight of $\alpha = 1$ and assign $\beta = 0.5$ to the prediction of auxiliary classifier. $Output = \arg\max(\alpha F(x_t^j|\theta_t) + \beta F_{aux}(x_t^j|\theta_t))$. This trick could slightly improve the final performance. Here we provide the ablation study on the sensitivity of inference weighting in Table 10. If we only deploy the main classifier ($\alpha = 1, \beta = 0$), the model could achieve 49.3% mIoU accuracy. When we combine the prediction of two classifiers, the performance could be improved about 1.0% mIoU.

**Uncertainty Visualization.** As a by-product, we also could estimate the prediction uncertainty when inference. We provide the visualization results to show the difference between the uncertainty estimation and the confidence score. As shown in Figure 5, we observe that the model is prone to provide the low confidence score of the boundary pixels, which does not provide the effective cue to the ambiguous prediction. Instead, the proposed prediction variance reflects the label uncertainty, and the highlight area in prediction variance map has lots of overlaps with the wrong prediction.

## 5 Conclusion

We identify the challenge of pseudo label learning in adaptive semantic segmentation and present a simple and effective method to estimate the prediction uncertainty during training. We also involve the uncertainty into the optimization objective as the variance regularization to rectify the training. The regularization helps the model learn from the noisy label, without introducing extra parameters or modules. As a result, we achieve the competitive performance on three benchmarks, including two synthetic-to-real benchmarks and one cross-city benchmark. In the future, we will continue to investigate the usage of uncertainty and the applications to other related tasks, *e.g.*, medical imaging.

## References

Blum and Mitchell, 1998. Blum, A. and Mitchell, T. (1998). Combining labeled and unlabeled data with co-training. In *COLT*.

Chen et al., 2017. 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 et al., 2016. 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 *CVPR*.

Fu et al., 2015. Fu, Y., Hospedales, T. M., Xiang, T., and Gong, S. (2015). Transductive multi-view zero-shot learning. *IEEE transactions on pattern analysis and machine intelligence*, 37(11):2332–2345.

Gal and Ghahramani, 2016. Gal, Y. and Ghahramani, Z. (2016). Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *ICML*.

Ganin and Lempitsky, 2015. Ganin, Y. and Lempitsky, V. (2015). Unsupervised domain adaptation by backpropagation. In *ICML*.

Grandvalet and Bengio, 2005. Grandvalet, Y. and Bengio, Y. (2005). Semi-supervised learning by entropy minimization. In *NeurIPS*.

Han et al., 2019. Han, L., Zou, Y., Gao, R., Wang, L., and Metaxas, D. (2019). Unsupervised domain adaptation via calibrating uncertainties. In *CVPR Workshops*.

He et al., 2016a. He, K., Zhang, X., Ren, S., and Sun, J. (2016a). Deep residual learning for image recognition. In *CVPR*.

He et al., 2016b. He, K., Zhang, X., Ren, S., and Sun, J. (2016b). Identity mappings in deep residual networks. In *ECCV*.

Hendrycks and Dietterich, 2019. Hendrycks, D. and Dietterich, T. (2019). Benchmarking neural network robustness to common corruptions and perturbations. In *ICLR*.

Hoffman et al., 2018. Hoffman, J., Tzeng, E., Park, T., Zhu, J.-Y., Isola, P., Saenko, K., Efros, A. A., and Darrell, T. (2018). Cycada: Cycle-consistent adversarial domain adaptation. In *ICML*.

Huang et al., 2018. Huang, H., Huang, Q., and Krahenbuhl, P. (2018). Domain transfer through deep activation matching. In *ECCV*.

Kang et al., 2020. Kang, G., Wei, Y., Yang, Y., and Hauptmann, A. (2020). Pixel-level cycle association: A new perspective for domain adaptive semantic segmentation. In *NeurIPS*.

Kendall and Gal, 2017. Kendall, A. and Gal, Y. (2017). What uncertainties do we need in bayesian deep learning for computer vision? In *NeurIPS*.

Lee, 2013. Lee, D.-H. (2013). Pseudo-label: The simple and efficient semi-supervised learning method for deep neural networks. In *Workshop on challenges in representation learning, ICML*.

Li et al., 2020a. Li, G., Kang, G., Liu, W., Wei, Y., and Yang, Y. (2020a). Content-consistent matching for domain adaptive semantic segmentation. In *ECCV*.

Li et al., 2020b. Li, P., Wei, Y., and Yang, Y. (2020b). Consistent structural relation learning for zero-shot segmentationg. In *NeurIPS*.

Li et al., 2020c. Li, P., Wei, Y., and Yang, Y. (2020c). Meta parsing networks: Towards generalized few-shot scene parsing with adaptive metric learning. In *ACM Multimedia*.

Liang et al., 2017. Liang, X., Lin, L., Wei, Y., Shen, X., Yang, J., and Yan, S. (2017). Proposal-free network for instance-level object segmentation. *IEEE transactions on pattern analysis and machine intelligence*, 40(12):2978–2991.

Luo et al., 2019a. Luo, Y., Liu, P., Guan, T., Yu, J., and Yang, Y. (2019a). Significance-aware information bottleneck for domain adaptive semantic segmentation. In *ICCV*.

Luo et al., 2020. Luo, Y., Liu, P., Guan, T., Yu, J., and Yang, Y. (2020). Adversarial style mining for one-shot unsupervised domain adaptation. In *NeurIPS*.

Luo et al., 2019b. 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 *CVPR*.

Maddern et al., 2017. Maddern, W., Pascoe, G., Linegar, C., and Newman, P. (2017). 1 Year, 1000km: The Oxford RobotCar Dataset. *The International Journal of Robotics Research (IJRR)*, 36(1):3–15.

Nielsen and Jensen, 2009. Nielsen, T. D. and Jensen, F. V. (2009). *Bayesian networks and decision graphs*. Springer Science & Business Media.

Pan et al., 2019. Pan, Y., Yao, T., Li, Y., Wang, Y., Ngo, C.-W., and Mei, T. (2019). Transferrable prototypical networks for unsupervised domain adaptation. In *CVPR*.

Paszke et al., 2017. Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., Lin, Z., Desmaison, A., Antiga, L., and Lerer, A. (2017). Automatic differentiation in pytorch.

Reed et al., 2014. Reed, S., Lee, H., Anguelov, D., Szegedy, C., Erhan, D., and Rabinovich, A. (2014). Training deep neural networks on noisy labels with bootstrapping. *arXiv:1412.6596*.

Richter et al., 2016. Richter, S. R., Vineet, V., Roth, S., and Koltun, V. (2016). Playing for data: Ground truth from computer games. In *ECCV*.

Ros et al., 2016. 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 *CVPR*.

Saito et al., 2018. Saito, K., Watanabe, K., Ushiku, Y., and Harada, T. (2018). Maximum classifier discrepancy for unsupervised domain adaptation. In *CVPR*.

Srivastava et al., 2014. Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., and Salakhutdinov, R. (2014). Dropout: a simple way to prevent neural networks from overfitting. *The journal of machine learning research*, 15(1):1929–1958.

Tsai et al., 2018. 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 *CVPR*.

Tsai et al., 2019. Tsai, Y.-H., Sohn, K., Schulter, S., and Chandraker, M. (2019). Domain adaptation for structured output via discriminative representations. In *ICCV*.

Tzeng et al., 2015. Tzeng, E., Hoffman, J., Darrell, T., and Saenko, K. (2015). Simultaneous deep transfer across domains and tasks. In *ICCV*.

Vu et al., 2019. Vu, T.-H., Jain, H., Bucher, M., Cord, M., and Pérez, P. (2019). Advent: Adversarial entropy minimization for domain adaptation in semantic segmentation. In *CVPR*.

Wang et al., 2018. Wang, J., Zhu, X., Gong, S., and Li, W. (2018). Transferable joint attribute-identity deep learning for unsupervised person re-identification. In *CVPR*.

Wang et al., 2019. Wang, S., Zhang, L., Zuo, W., and Zhang, B. (2019). Class-specific reconstruction transfer learning for visual recognition across domains. *IEEE Transactions on Image Processing*, 29:2424–2438.

Wei et al., 2018. Wei, Y., Xiao, H., Shi, H., Jie, Z., Feng, J., and Huang, T. S. (2018). Revisiting dilated convolution: A simple approach for weakly-and semi-supervised semantic segmentation. In *CVPR*.

Wu et al., 2018. Wu, Z., Han, X., Lin, Y.-L., Gokhan Uzunbas, M., Goldstein, T., Nam Lim, S., and Davis, L. S. (2018). Dcan: Dual channel-wise alignment networks for unsupervised scene adaptation. In *ECCV*.

Wu et al., 2019. Wu, Z., Wang, X., Gonzalez, J. E., Goldstein, T., and Davis, L. S. (2019). Ace: adapting to changing environments for semantic segmentation. In *ICCV*.

Yang et al., 2020. Yang, J., Xu, R., Li, R., Qi, X., Shen, X., Li, G., and Lin, L. (2020). An adversarial perturbation oriented domain adaptation approach for semantic segmentation. In *AAAI*.

Yu et al., 2019. Yu, T., Li, D., Yang, Y., Hospedales, T. M., and Xiang, T. (2019). Robust person re-identification by modelling feature uncertainty. In *ICCV*.

Yue et al., 2019. Yue, X., Zhang, Y., Zhao, S., Sangiovanni-Vincentelli, A., Keutzer, K., and Gong, B. (2019). Domain randomization and pyramid consistency: Simulation-to-real generalization without accessing target domain data. In *ICCV*.

Zhang et al., 2019a. Zhang, L., Li, X., Arnab, A., Yang, K., Tong, Y., and Torr, P. H. (2019a). Dual graph convolutional network for semantic segmentation. In *BMVC*.

Zhang et al., 2019b. Zhang, L., Wang, S., Huang, G.-B., Zuo, W., Yang, J., and Zhang, D. (2019b). Manifold criterion guided transfer learning via intermediate domain generation. *IEEE Transactions on Neural Networks and Learning Systems*, 30(12):3759–3773.

Zhang et al., 2020. Zhang, L., Xu, D., Arnab, A., and Torr, P. H. (2020). Dynamic graph message passing network. In *CVPR*.

Zhang et al., 2018. Zhang, Y., Qiu, Z., Yao, T., Liu, D., and Mei, T. (2018). Fully convolutional adaptation networks for semantic segmentation. In *CVPR*.

Zhao et al., 2017. Zhao, H., Shi, J., Qi, X., Wang, X., and Jia, J. (2017). Pyramid scene parsing network. In *CVPR*.

Zheng and Yang, 2020. Zheng, Z. and Yang, Y. (2020). Unsupervised scene adaptation with memory regularization in vivo. In *IJCAI*.

Zou et al., 2019. Zou, Y., Yu, Z., Liu, X., Kumar, B., and Wang, J. (2019). Confidence regularized self-training. In *ICCV*.

Zou et al., 2018. Zou, Y., Yu, Z., Vijaya Kumar, B., and Wang, J. (2018). Unsupervised domain adaptation for semantic segmentation via class-balanced self-training. In *ECCV*.