

# Prototypical Pseudo Label Denoising and Target Structure Learning for Domain Adaptive Semantic Segmentation

Pan Zhang<sup>1</sup>\*, Bo Zhang<sup>2</sup>, Ting Zhang<sup>2</sup>, Dong Chen<sup>2</sup>, Yong Wang<sup>1</sup>, Fang Wen<sup>2</sup>

<sup>1</sup>University of Science and Technology of China <sup>2</sup>Microsoft Research Asia

## Abstract

*Self-training is a competitive approach in domain adaptive segmentation, which trains the network with the pseudo labels on the target domain. However inevitably, the pseudo labels are noisy and the target features are dispersed due to the discrepancy between source and target domains. In this paper, we rely on representative prototypes, the feature centroids of classes, to address the two issues for unsupervised domain adaptation. In particular, we take one step further and exploit the feature distances from prototypes that provide richer information than mere prototypes. Specifically, we use it to estimate the likelihood of pseudo labels to facilitate online correction in the course of training. Meanwhile, we align the prototypical assignments based on relative feature distances for two different views of the same target, producing a more compact target feature space. Moreover, we find that distilling the already learned knowledge to a self-supervised pretrained model further boosts the performance. Our method shows tremendous performance advantage over state-of-the-art methods. We will make the code publicly available.*

## 1. Introduction

Despite the remarkable success of deep learning in computer vision, attaining high performance requires vast quantities of data. It is usually expensive to obtain labels for dense prediction tasks, *e.g.*, semantic segmentation. Therefore, people think of leveraging abundant photo-realistic synthetic images with freely generated labels [45, 46]. However, deep neural networks are notoriously sensitive to the domain misalignment that any nuanced unrealism in rendered images will induce poor generalization to real data. To address this issue, domain adaption techniques aim to transfer the knowledge learned from the synthetic images (source domain) to real ones (target domain) with minimal performance loss. In this work, we focus on the challenging

case, *unsupervised domain adaptation* (UDA), where there are no accessible labels in the target domain. Specifically, we solve the UDA problem for semantic segmentation.

Rather than explicitly aligning the distributions of the source and target domains as most predominant solutions [6, 28, 49, 55, 59], self-training [35, 69, 75, 76] has recently emerged as a simple yet competitive approach in the UDA task. This is achieved by iteratively generating a set of pseudo labels based on the most confident predictions on the target data and then relying on these pseudo labels to retrain the network. In this way, the network gradually learns the adaptation in the self-paced curriculum learning. However, the performance still lags far behind the supervised learning or semi-supervised learning using a few labeled samples, making unsupervised domain adaptation impractical in real scenarios.

After dissecting the self-training, we find two key ingredients are lacking in previous works. First, typical practice [75, 76] suggests selecting the pseudo labels according to a strict confidence threshold, while high scores are not necessarily correct, making the network fail to learn reliable knowledge in the target domain. Second, due to the domain gap, the network is prone to produce dispersed features in the target domain. It is likely that for target data, the closer to the source distribution, the higher the confidence score. As a result, data lying far from the source distribution (*i.e.* low scores) will never be considered during the training.

In this paper, we propose to online denoise the pseudo labels and learn a compact target structure to address the above two issues respectively. We resort to prototypes, *i.e.*, the class-wise feature centroids, to accomplish the two tasks. (1) We rectify the pseudo labels by estimating the class-wise likelihoods according to its relative feature distances to all class prototypes. This depends on a practical assumption that the prototype lies closer to the true centroid of the underlying cluster, implying that false pseudo labels are in the minority. It is worth noting that the prototypes are computed on-the-fly, and thus the pseudo labels are progressively corrected throughout the training. (2) We draw inspiration from the Deepcluster [4] to learn the intrinsic structure of the target domain. Instead of directly learning

\*This work is done during the first author’s internship at Microsoft Research Asia.

from the cluster assignment, we propose to align soft prototypical assignments for different views of the same target, which produces a more compact target feature space. We refer to our method *ProDA* as we heavily rely on prototypes for domain adaption.

Supercharged with the above techniques, our ProDA can demonstrate clear superiority over prior works. Moreover, we find that the domain adaptation can also benefit from the task-agnostic pretraining — distilling the knowledge to a self-supervised model [10, 24] further boosts the performance to a record high. Our contributions can be summarized as follows:

- We propose to online correct the soft pseudo labels according to the relative feature distances to the prototypes, whereas the prototypes are also updated on-the-fly. The network thereby learns from denoised pseudo labels throughout the training.
- We propose to rely on the soft prototypical assignment to teach the learning of an augmented view so that a compact target feature space can be obtained.
- We show that distilling the already-learned knowledge to a self-supervised pretrained model further improves the performance significantly.
- The proposed ProDA substantially outperforms state-of-the-art. With the Deeplabv2 [8] network, our method achieves the Cityscapes [14] segmentation mIOU by 57.5 and 55.5 when adapting from the GTA5 [45] and SYNTHIA [46] datasets, improving the adaption gain<sup>1</sup> by 52.6% and 58.5% respectively over the prior leading approach.

## 2. Related Work

**Unsupervised domain adaptation.** As suggested by the theoretical analysis [3], domain alignment methods focus on reducing the distribution mismatch by optimizing some divergence [31, 36] or adopting adversarial training [21, 41] at either the image level [1, 12, 13, 20, 27, 50, 64], the intermediate feature level [6, 28, 49, 59] or the output level [55]. However, aligning global distribution cannot guarantee a small expected error on the target domain [7, 29, 72]. Recent approaches [18, 37, 60] attempt to align distribution in a class-wise manner, aiming to promote fine-grained feature alignment. In fact, it is unnecessary to rigorously align the distribution as long as the features are well-separated.

On the other hand, techniques originated from semi-supervised learning (SSL) offer competitive performance. Entropy minimization and its variants [9, 47, 58] encourage the network to make sharp predictions on the unlabeled target data, and the resulting network is prone to be over-confident on false predictions. To address this, self-training [75] that leverages iteratively generated pseudo la-

els has been proposed. However, the pseudo labels are inevitably noisy. Hence, [76] adds confidence regularization terms to the network, while [73] explicitly estimates a prediction confidence map to reduce the side-effect of unreliable labels. In [35], self-training and image translation are found mutually beneficial. A recent work [69] generates pseudo labels based on categorical centroids and enforces feature alignment in category level. However, these self-training approaches are optimized in an alternative manner — labels are fixed over the course of representation learning, and only get updated after the entire training stage. In contrast, we propose an online pseudo label updating scheme where the false predictions are rectified according to the prototypical context estimated in the target domain.

**Unsupervised representation learning.** A surge of research interest has been recently attracted to unsupervised learning [43]. Early efforts dedicate to designing pretext tasks [17, 19, 30, 70], which are proven beneficial for UDA when utilized as auxiliary tasks on target data [48, 54, 67]. The gap with supervised learning is considerably closed by a few prominent works [10, 24] that build on contrastive learning. A series of recent works [2, 5, 22, 40] find that the network is able to learn rich semantic features as long as they are consistent under different augmented views. Yet these methods assume image-wise discrimination [65], making them unsuitable for learning pixel-level semantics for segmentation tasks. In this work, we find the marriage of consistent learning and cluster-based representation learning fits remarkably well with the UDA problem and learn a compact target feature space inspired from Deepcluster [4]. Differently, we align relative feature distances rather than cluster assignments for different augmented views.

**Learning from noisy labels.** Self-training even with careful thresholding still gives noisy pseudo labels. Therefore, this work is also motivated by emerging techniques [53] of learning from noisy labels. A straightforward way is to design robust losses [63, 71], but these methods fail to handle real-world noisy data. Self-label correction [52, 62, 73, 74] is a more appealing approach. A typical manner [32, 39] under this category is to train two or multiple learners simultaneously and exploit their agreement of predictions to measure the label reliability. Our proposed pseudo label denoising is more close to [23] which online corrects the incorrect labels according to the prototypes determined by some complex heuristics. In contrast, we are able to compute prototypes on-the-fly. On the other hand, knowledge distillation (KD) [26, 33, 34, 66] is proven effective to transfer clean knowledge from the teacher model to the student even when the network learns from itself [10, 33]. In this work, we demonstrate that the knowledge distillation to a self-supervised pretrained model further pushes the performance limit in our task.

<sup>1</sup>The mIoU gain relative to the model without domain adaption.

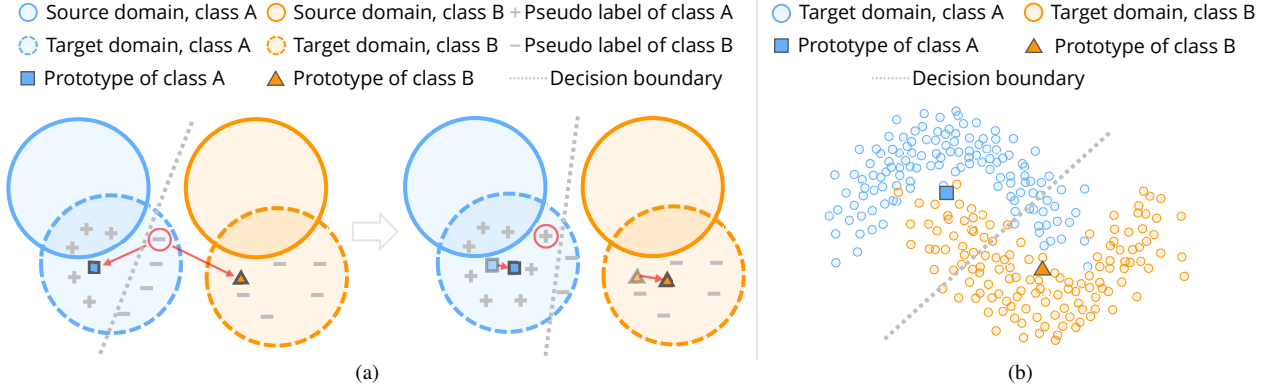


Figure 1: **We illustrate the existing issues of self-training by visualizing the feature space.** (a) The decision boundary (dashed line) crosses the distribution of the target data and induces incorrect pseudo label predictions. This is because the network is unaware of the target distribution when generating pseudo labels. To solve this, we calculate the prototypes of each class on-the-fly and rely on these prototypes to online correct the false pseudo labels. (b) The network may induce dispersed feature distribution in the target domain which is hardly differentiated by a linear classifier.

### 3. Preliminary

Given the source dataset  $\mathcal{X}_s = \{x_s\}_{j=1}^{n_s}$  with segmentation labels  $\mathcal{Y}_s = \{y_s\}_{j=1}^{n_s}$ , we aim to train a segmentation network that learns the knowledge from the source and is expected to achieve low risk on the target dataset  $\mathcal{X}_t = \{x_t\}_{j=1}^{n_t}$  without accessing its ground truth labels  $\mathcal{Y}_t$ , where  $\mathcal{Y}_s$  and  $\mathcal{Y}_t$  share the same  $K$  classes. Generally, the network  $h = g \circ f$  can be regarded as a composite of a feature extractor  $f$  and a classifier  $g$ .

Typically, the networks trained on the source data, *i.e.*, the source model, cannot generalize well to the target data due to the domain gap. To transfer the knowledge, traditional self-training techniques [75, 76] optimize the categorical cross-entropy (CE) with pseudo labels  $\hat{y}_t$ :

$$\ell_{ce}^t = - \sum_{i=1}^{H \times W} \sum_{k=1}^K \hat{y}_t^{(i,k)} \log(p_t^{(i,k)}), \quad (1)$$

where  $p_t = h(x_t)$  and  $p_t^{(i,k)}$  represents the softmax probability of pixel  $x_t^{(i)}$  belonging to the  $k$ th class. Typically, one can use the most probable class predicted by the source network as pseudo labels:

$$\hat{y}_t^{(i,k)} = \begin{cases} 1, & \text{if } k = \arg \max_{k'} p_t^{(i,k')} \\ 0, & \text{otherwise} \end{cases} \quad (2)$$

Here we denote this conversion from the soft predictions to hard labels by  $\hat{y}_t = \xi(p_t)$ . In practice, since the pseudo labels are noisy, only the pixels whose prediction confidence is higher than a given threshold are accounted for the pseudo label retraining. In this way, the network in the target domain is bootstrapped by learning from pseudo labels that only get update till convergence, and then the updated labels are employed for the next training stage.

### 4. Method

#### 4.1. Prototypical pseudo label denoising

We conjecture that updating the pseudo label after one training stage is too late as the network has already overfitted the noisy labels. While simultaneously updating the pseudo labels and the network weights, on the other hand, is prone to give trivial solutions.

In this work, we propose a simple yet effective approach to online update the pseudo labels while avoiding trivial solutions. The key is to fix the soft pseudo labels and progressively weight them by class-wise probabilities, with the update in accordance with the freshly learned knowledge. Formally, we propose to use the weighted pseudo labels for self-training:

$$\hat{y}_t^{(i,k)} = \xi(\omega_t^{(i,k)} p_{t,0}^{(i,k)}), \quad (3)$$

where  $\omega_t^{(i,k)}$  is the weight we propose for modulating the probability and changes as the training proceeds. The  $p_{t,0}^{(i,k)}$  is initialized by the source model and remains fixed throughout the learning process, thus serving as a boilerplate for the subsequent refinement.

We exploit the distances from the prototypes to gradually rectify the pseudo labels. Let  $f(x_t)^{(i)}$  represent the feature of  $x_t$  at index  $i$ . If it is far from the prototype  $\eta^{(k)}$ , the feature centroids of class  $k$ , it is more probable that the learned feature is an outlier, hence we downweight its probability of being classified into  $k$ th category. Concretely, we define the modulation weight in Equation 3 as the softmax over feature distances to prototypes:

$$\omega_t^{(i,k)} = \frac{\exp(-\|f(x_t)^{(i)} - \eta^{(k)}\|/\tau)}{\sum_{k'} \exp(-\|f(x_t)^{(i)} - \eta^{(k')}\|/\tau)}, \quad (4)$$

where  $\tilde{f}$  denotes the momentum encoder [24] of the feature extractor  $f$ , as we desire a reliable feature estimation for  $x_t$ , and  $\tau$  is the softmax temperature empirically set to  $\tau = 1$ . In this equation,  $\omega_t^{(i,k)}$  actually approximates the trust confidence of  $x_t^{(i)}$  belonging to the  $k$ th class. Note that this equation has a close formulation in [44, 51] that shows promise in few-shot learning. Instead of relying on the prototypes for classifying new samples, we attempt to correct the mistakes in the pre-generated pseudo labels according to this prototypical context.

**Prototype computation.** The proposed label updating scheme requires to compute the prototypes on-the-fly. At first, prototypes are initialized according to the predicted pseudo labels  $\hat{y}_t$  for target domain images, which is

$$\eta^{(k)} = \frac{\sum_{x_t \in \mathcal{X}_t} \sum_i f(x_t)^{(i)} * \mathbb{1}(\hat{y}_t^{(i,k)} == 1)}{\sum_{x_t \in \mathcal{X}_t} \sum_i \mathbb{1}(\hat{y}_t^{(i,k)} == 1)}, \quad (5)$$

where  $\mathbb{1}$  is the indicator function. However, such prototype calculation is computational-intensive during training. To address this, we estimate the prototypes as the moving average of the cluster centroids in mini-batches, so that we can track the prototypes that slowly move. Specifically, in each iteration, the prototype is estimated as,

$$\eta^{(k)} \leftarrow \lambda \eta^{(k)} + (1 - \lambda) \eta'^{(k)}, \quad (6)$$

where  $\eta'^{(k)}$  is the mean feature of class  $k$  calculated within the current training batch from the momentum encoder, and  $\lambda$  is the momentum coefficient which we set to 0.9999.

**Pseudo label training loss.** With the online refined pseudo labels, we are able to train the network for target domain segmentation. Instead of using a standard cross-entropy, we adopt a more robust variant, symmetric cross-entropy (SCE) loss [63], to further enhance the noise-tolerance to stabilize the early training phase. Specifically, we enforce

$$\ell_{sce}^t = \alpha \ell_{ce}(p_t, \hat{y}_t) + \beta \ell_{ce}(\hat{y}_t, p_t), \quad (7)$$

where  $\alpha$  and  $\beta$  are balancing coefficients and set to 0.1 and 1 respectively. Following [63], we clamp the one-hot label  $\hat{y}_t$  to  $[1e-4, 1]$ , so as to avoid the numerical issue of  $\log 0$ .

**Why are prototypes useful for label denoising?** First, the prototypes are less sensitive to the outliers that are assumed to be the minority. Second, prototypes treat different classes equally regardless of their occurrence frequency, which is particularly useful to semantic segmentation as class imbalance poses a challenging issue. Experiments show that the proposed label denoising considerably improves the performance for hard classes. More importantly, prototypes help to gradually rectify the incorrect pseudo labels, which we illustrate using a toy example. As shown in Figure 1(a), the classifier  $g$ , may still give a decision boundary crossing the target clusters and yields false pseudo labels. Training

against them cannot further improve the classifier. The prototypes, on the other hand, downweight the false pseudo labels near the decision boundary of  $g$  as they are far from the prototypes. In this way, the network improves, and in turn, makes the prototypes closer to the true cluster centroid.

## 4.2. Structure learning by enforcing consistency

Pseudo labels can be denoised when the feature extractor  $f$  generates compact target features. However, due to the domain gap, the generated target distribution is more likely to be dispersed, as shown in Figure 1(b). In this case, the prototypes fail to rectify the labels of the data whose features lie in the far end of the cluster even when the target features from the source model are well-separated. A recent work [42] has identified this issue in semi-supervised learning, but the issue becomes worse in domain adaptation since a few pseudo labeled data cannot cover the entire distribution in the target domain.

To this end, we aim to learn the underlying structure of target domain, and hope to obtain more compact features that are friendly to the pseudo label refinement. Motivated by the recent success of unsupervised learning, we perform simultaneously clustering and representation learning. As opposed to learning against the prototypical assignment, we use the prototypical assignment under weak augmentation to guide the learning for the strong augmented view. Specifically, let  $\mathcal{T}(x_t)$  and  $\mathcal{T}'(x_t)$  respectively denote the weak and strong augmented views for  $x_t$ . We make use of the momentum encoder  $\tilde{f}$  to generate a reliable prototypical assignment for  $\mathcal{T}(x_t)$  which is,

$$z_{\mathcal{T}}^{(i,k)} = \frac{\exp(-\|\tilde{f}(\mathcal{T}(x_t))^{(i)} - \eta^{(k)}\|/\tau)}{\sum_{k'} \exp(-\|\tilde{f}(\mathcal{T}(x_t))^{(i)} - \eta^{(k')}\|/\tau)}, \quad (8)$$

where  $\tau = 1$  by default. Likewise, the soft assignment  $z_{\mathcal{T}'}$  for  $\mathcal{T}'(x_t)$  can be obtained in the same manner except that we use the original trainable feature extractor  $f$ . Since  $z_{\mathcal{T}}$  is more reliable as the feature is given by a momentum encoder  $\tilde{f}$  and the input  $x_t$  suffers from less distortion, we use it to teach  $f$  to produce consistent assignments for  $\mathcal{T}(x_t)$ . Hence, we minimize the Kullback–Leibler (KL) divergence between the prototypical assignments under two views:

$$\ell_{kl}^t = \text{KL}(z_{\mathcal{T}} \| z_{\mathcal{T}'}). \quad (9)$$

Intuitively, this equation enforces the network to give consistent prototypical labeling for the adjacent feature points, which results in more compact feature space in the target domain.

Similar to the previous works that simultaneously learn the clustering and representation, the proposed prototypical consistent learning may suffer from degeneration issue, *i.e.*, one cluster becomes empty. To amend this, we use a regularization term from [76], which encourages the output to

be evenly distributed to different classes,

$$\ell_{reg}^t = - \sum_{i=1}^{H \times W} \sum_{j=1}^K \log p_i^{(i,k)}. \quad (10)$$

Now equipped with the online label correction and the prototypical consistent learning, we train the domain adaptation network with the following total loss:

$$\ell_{total} = \ell_{ce}^s + \ell_{sce}^t + \gamma_1 \ell_{kl}^t + \gamma_2 \ell_{reg}^t. \quad (11)$$

By default, the weighting coefficients  $\gamma_1$  and  $\gamma_2$  are set to 10, 0.1 respectively.

### 4.3. Distillation to self-supervised model

After the training with Equation 11 converges, we further transfer knowledge from the learned target model to a student model with the same architecture but pretrained in a self-supervised manner. To be concrete, we initialize the feature extractor of the student model  $h^\dagger$  with SimCLRv2 [11] pretrained weights, and we apply a knowledge distillation (KD) loss, which lets the student mimic the teacher by minimizing the KL divergence of their predictions on the unlabeled target images. Besides, following the self-training paradigm, we rely on the teacher model  $h$  to generate one-hot pseudo labels  $\xi(p_t)$  so as to teach the student model. To prevent the model from forgetting the source domain, the source images are also utilized. Altogether, we train the student model using the following loss,

$$\ell_{KD} = \ell_{ce}^s(p_s, y_s) + \ell_{ce}^t(p_t^\dagger, \xi(p_t)) + \beta \text{KL}(p_t \| p_t^\dagger), \quad (12)$$

where  $p_t^\dagger = h^\dagger(x_t)$  is the output of the student model, and we set  $\beta = 1$ . In practice, such self-distillation can be applied multiple times once the model converges, which helps the domain adaptation to achieve even higher performance.

## 5. Experiments

### 5.1. Implementation

**Training.** We use the DeepLabv2 [8] for segmentation with the backbone ResNet-101 [25]. Following [69, 73], we utilize [56] that applies adversarial training at the segmentation output as a warm-up. We apply the SGD solver with the initial learning rate as 1e-4 which is decayed by 0.9 every training epoch, and the training lasts 80 epochs. During the structure learning, the augmentation is composed of random crop, RandAugment [15] and Cutout [16]. For knowledge distillation, we utilize the pretrained SimCLRv2 model with the ResNet-101 backbone as well. An extra batch normalization (BN) layer is introduced after its feature extraction layer so as to accommodate the activation magnitude for our task, with the learning rate set to 6e-4 and 6e-3 respectively before and after this BN layer. During the

distillation stage, we use hard pseudo labels with the selection threshold 0.95. Readers can refer to the supplementary material for more training details and algorithm flow. We conduct all the experiments on 4 Tesla V100 GPUs with PyTorch implementation.

**Dataset.** For evaluation, we adapt the segmentation from game scenes, GTA5 [45] and SYNTHIA [46] datasets, to real scene, the Cityscapes [14] dataset. GTA5 contains 24,966 training images with the resolution of 1914×1052 and we use its 19 categories shared with Cityscapes. SYNTHIA dataset contains 9,400 1280×760 images and we use its 16 common categories with Cityscapes. We also report the results on 13 common categories on this dataset following the protocol of some methods. The Cityscapes dataset contains 2,975 training images and 500 images for validation with the resolution of 2048×1024. Since its testing set does not provide ground truth labeling, we conduct evaluations on its validation set.

### 5.2. Comparisons with state-of-the-art methods

We comprehensively compare our proposed method with the recent leading approaches. These methods could be divided into two categories: 1) domain alignment methods that align the distribution through adversarial training, which include AdaptSeg [55], CyCADA [27], CLAN [37], APODA [68], PatchAlign [57], ADVENT [58], BDL [35] and FADA [61]; 2) self-training approaches, including CBST [75], MRKLD [76], Seg-Uncertainty [73], CAG-UDA [69].

Table 1 shows the comparisons of GTA5 → Cityscapes adaptation. our ProDA arrives at the state-of-the-art mIoU score 57.5, outperforming existing methods by a large margin. Among all the 19 categories, we achieve the best scores on 15 categories. ProDA shows evident advantage in hard classes, *e.g.*, fence, terrace, motor, that cannot be well handled in previous works. Indeed, the performance improvement of ProDA mostly comes from these challenging cases, the small or rare objects, as we regard different categories equally thanks to the prototypes. Note that this is achieved without any heuristic class-balance strategies as [75]. Compared with the non-adapted baseline (*i.e.*, the model purely trained on the source), ProDA offers the mIoU gain by 20.9, outperforming the second-best method by 52.6%.

In Table 2, we show the adaptation results on SYNTHIA → Cityscapes, where ProDA also shows tremendous improvement. The proposed ProDA achieves the mIoU score by 55.5 and 62.0 over the 16 and 13 categories respectively. To be specific, we arrive at the best on 11 out of 16 categories, mostly the hard classes. Relative to the non-adaptive model, the segmentation model after our adaption sees the gain by 20.6, surpassing that of the prior leading approach by 58.4%. While the adaptation from SYNTHIA is more challenging than that from GTA5, our ProDA per-

	road	sideway	building	wall	fence	pole	light	sign	vege.	terrace	sky	person	rider	car	truck	bus	train	motor	bike	mIoU	gain
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	+0.0
AdaptSeg [55]	86.5	25.9	79.8	22.1	20.0	23.6	33.1	21.8	81.8	25.9	75.9	57.3	26.2	76.3	29.8	32.1	7.2	29.5	32.5	41.4	+4.8
CyCADA [27]	86.7	35.6	80.1	19.8	17.5	38.0	39.9	41.5	82.7	27.9	73.6	64.9	19.0	65.0	12.0	28.6	4.5	31.1	42.0	42.7	+6.1
CLAN [37]	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	+6.6
APODA [68]	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	+9.3
PatchAlign [57]	<b>92.3</b>	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	+9.9
ADVENT [58]	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	+8.9
BDL [35]	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	+11.9
FADA [61]	91.0	50.6	<b>86.0</b>	43.4	29.8	36.8	43.4	25.0	86.8	38.3	<b>87.4</b>	64.0	38.0	85.2	31.6	46.1	6.5	25.4	37.1	50.1	+13.5
CBST [75]	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	+9.3
MRKLD [76]	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	+10.5
CAG-UDA [69]	90.4	51.6	83.8	34.2	27.8	38.4	25.3	48.4	85.4	38.2	78.1	58.6	34.6	84.7	21.9	42.7	<b>41.1</b>	29.3	37.2	50.2	+13.6
Seg-Uncertainty [73]	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	+13.7
<i>ProDA</i>	87.8	<b>56.0</b>	79.7	<b>46.3</b>	<b>44.8</b>	<b>45.6</b>	<b>53.5</b>	<b>53.5</b>	<b>88.6</b>	<b>45.2</b>	82.1	<b>70.7</b>	<b>39.2</b>	<b>88.8</b>	<b>45.5</b>	<b>59.4</b>	1.0	<b>48.9</b>	<b>56.4</b>	<b>57.5</b>	<b>+20.9</b>

Table 1: Comparison results of GTA5→Cityscapes adaptation in terms of mIoU. The best score for each column is highlighted.

	road	sideway	building	wall*	fence*	pole*	light	sign	vege.	sky	person	rider	car	bus	motor	bike	mIoU	gain	mIoU*	gain*
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	34.9	+0.0	40.3	+0.0
AdaptSeg [55]	79.2	37.2	78.8	-	-	-	9.9	10.5	78.2	80.5	53.5	19.6	67.0	29.5	21.6	31.3	-	-	45.9	+5.6
PatchAlign [57]	82.4	38.0	78.6	8.7	0.6	26.0	3.9	11.1	75.5	<b>84.6</b>	53.5	21.6	71.4	32.6	19.3	31.7	40.0	+5.1	46.5	+6.2
CLAN [37]	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	+7.5
APODA [68]	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	+12.8
ADVENT [58]	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	41.2	+6.3	48.0	+7.7
BDL [35]	86.0	<b>46.7</b>	80.3	-	-	-	14.1	11.6	79.2	81.3	54.1	27.9	73.7	42.2	25.7	45.3	-	-	51.4	+11.1
FADA [61]	84.5	40.1	83.1	4.8	0.0	34.3	20.1	27.2	84.8	84.0	53.5	22.6	85.4	43.7	26.8	27.8	45.2	+10.3	52.5	+12.2
CBST [75]	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	42.6	+7.7	48.9	+8.6
MRKLD [76]	67.7	32.2	73.9	10.7	1.6	37.4	22.2	31.2	80.8	80.5	60.8	<b>29.1</b>	82.8	25.0	19.4	45.3	43.8	+8.9	50.1	+9.8
CAG-UDA [69]	84.7	40.8	81.7	7.8	0.0	35.1	13.3	22.7	84.5	77.6	64.2	27.8	80.9	19.7	22.7	48.3	44.5	+9.6	51.5	+11.2
Seg-Uncertainty [73]	87.6	41.9	83.1	14.7	<b>1.7</b>	36.2	31.3	19.9	81.6	80.6	63.0	21.8	86.2	40.7	23.6	<b>53.1</b>	47.9	+13.0	54.9	+14.6
<i>ProDA</i>	<b>87.8</b>	45.7	<b>84.6</b>	<b>37.1</b>	0.6	<b>44.0</b>	<b>54.6</b>	<b>37.0</b>	<b>88.1</b>	84.4	<b>74.2</b>	24.3	<b>88.2</b>	<b>51.1</b>	<b>40.5</b>	45.6	<b>55.5</b>	<b>+20.6</b>	<b>62.0</b>	<b>+21.7</b>

Table 2: Comparison results of SYNTHIA→Cityscapes adaptation in terms of mIoU. The best score for each column is highlighted. mIoU and mIoU\* denote the averaged scores across 16 and 13 categories respectively.

forms equally well on both datasets.

### 5.3. Discussion

**The effectiveness of pseudo label denoising.** In Table 11, the non-adapted source model only gives 36.6 mIoU on the target domain. After the model warm-up, we get a 5.0 mIoU increase. Initialized with the warm-up, the baseline model, *i.e.*, the vanilla self-training, which is trained using the offline pseudo labels, improves the mIoU to 45.2. Adding symmetric cross-entropy (SCE) brings +0.4 mIoU gain. By contrast, we are able to update the pseudo labels on-the-fly and the training with the denoised pseudo labels significantly boosts the performance to 52.3. The model with this component alone sets the state-of-the-art, outperforming the

prior best score by 2.0.

Figure 2 illustrates how the pseudo labels progress as the training proceeds. In the early phase, false pseudo labels are likely to occur and fail to recognize the tiny objects, *e.g.*, the poles and traffic signs. As the training goes on, pseudo labels get denoised by virtue of the prototypes, and tiny objects can be gradually identified in the refined pseudo labels. The training after 40k iterations can already correct most of the incorrect labels. In Figure 3, we further quantify the mIoU score of the pseudo labels during the whole training process. While convention self-training only updates the pseudo labels after each stage and leads to step-wise mIoU increase, the pseudo labels given by ProDA can quickly attain a high quality, showing a distinct advantage over the

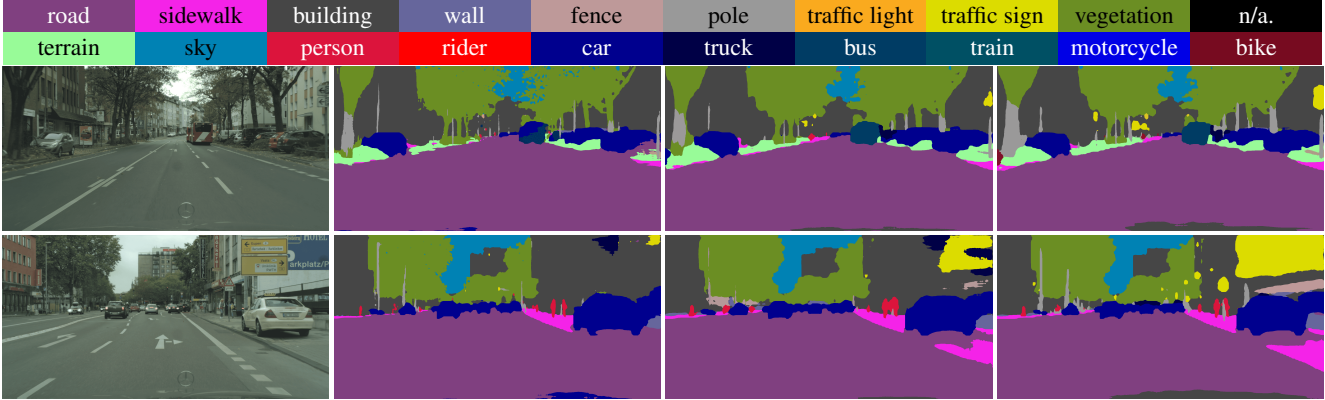


Figure 2: ProDA online refines the pseudo labels throughout the training. The first column is the segmentation input. The 2nd to 4th columns illustrate the pseudo labels after 1k, 10k, and 40k iterations.

		components				mIoU gain
initialization		source				36.6
		warm up				41.6 +5.0
stage 1		self training	sce	prototypical denoising	structure learning	mIoU gain
	✓					45.2 +8.6
	✓	✓				45.6 +9.0
	✓	✓	✓			52.3 +15.7
	✓	✓	✓	✓		47.6 +11.0
✓	✓	✓	✓	✓	53.7 +17.1	
stage 2		self distill.	stage 1 init.	supervised init.	self-supervised init.	mIoU gain
	✓				✓	55.8 +19.2
	✓	✓				56.3 +19.7
	✓		✓			55.7 +19.1
✓				✓	56.9 +20.3	
stage 3	✓			✓		57.5 +20.9

Table 3: Ablation study of each proposed component. The whole training involves three stages, where knowledge distillation can be applied in the last two stages.

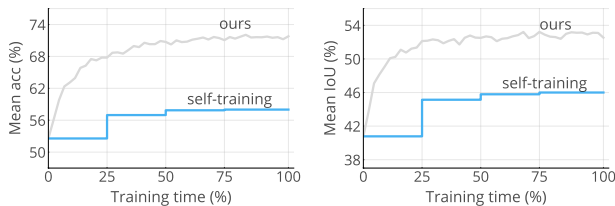


Figure 3: The mean accuracy and mean IoU score of the pseudo labels throughout the training. Comparing to the conventional self-training that updates pseudo labels only after the training stage, the pseudo labels in our method steadily improves as the training proceeds.

vanilla self-training of multiple training stages.

**How to prevent degenerate solution?** Simultaneously

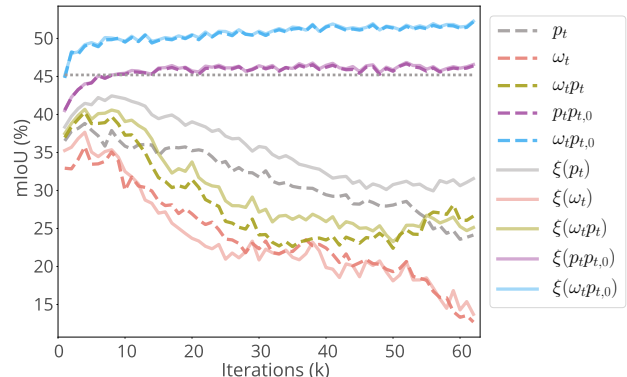


Figure 4: The training curves of different label reweighting schemes when online updating the labels. Adopting fixed  $p_{t,0}$  labels avoid the trivial solutions that plague the training with non-fixed pseudo labels  $p_t$ . The dotted line denotes the performance of conventional self-training.

learning the feature and updating the labels during self-training derives degenerate solutions. The key to avoiding such degeneration is to adopt a fixed soft label  $p_{t,0}$  as the boilerplate upon which we apply the weight  $w_t$  for rectification (Equation 3). To explain this, we investigate different variants for the online label updating in Figure 4, where the labels for self-training can be dynamic  $p_t$  or fixed  $p_{t,0}$ , and the modulation weight can be the network prediction  $p_t$  or the confidence  $w_t$  estimated according to prototypes. In addition, we investigate the benefit of using hard labels ( $\xi(\cdot)$ ). Figure 4 shows that the performance climbs up for a while and then starts to drop significantly when using non-fixed soft predictions ( $p_t$  and its variants), whereas the learning with the fixed ones ( $p_{t,0}$  and its variants) steadily improves throughout the training and surpasses the conventional self-training. We conjecture that fixing  $p_{t,0}$  makes the refined pseudo labels hardly deviate from this initial prediction, thus avoiding the trivial solution. Moreover, in contrast to using  $p_t$  for reweighting, prototypical reweighting improves

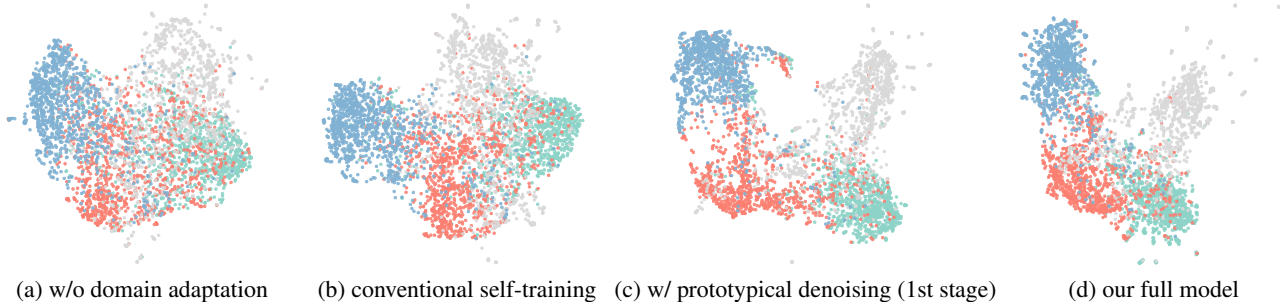


Figure 5: The visualization of feature space, where we map features to 2D space with UMAP [38]. For a clear illustration, we only show four categories, *i.e.*, blue for building, gray for traffic sign, orange for pole, and green for vegetation.

the mIoU by more than 5.0, corroborating the importance of using prototypes for label denoising. Besides, we observe slight improvement ( $\sim 0.2$ ) using hard labels over the soft ones.

**The effectiveness of target structure learning.** We propose to cultivate the intrinsic knowledge for the target domain and learn its underlying structure. As shown in Table 11, without the label denoising, the structure learning improves the performance from 45.6 to 47.6. Note that this is also competitive among self-training approaches, but we do not need to carefully choose the threshold for selecting the confident pseudo labels. The target structure learning assists the pseudo label denoising by forming compact feature clusters and brings the performance gain by 1.4.

**The effectiveness of distilling to self-supervised model.** At the second training stage, we apply knowledge distillation and transfer the dark knowledge of the 1st stage model to the current phase. Table 11 also compares different initialization strategies for this stage. Compared to resuming the last stage training, the initialization from a self-supervised pretrained model (*i.e.*, SimCLRv2) improves the mIoU by 0.6, whereas the initialization with the supervised pretraining degrades the performance. This is because the self-supervised pretraining possesses stronger transferability and can benefit a broad of downstream tasks. The initialization in this way helps the model escape from the local optima in the last stage. Table 11 also proves the effectiveness of the knowledge distillation: ablating this component drops the mIoU by 1.1. It is surprising to see that the third stage with the knowledge distillation further improves the performance by 0.6, attaining the 20.9 mIoU improvement relative to the model without domain adaptation.

**The UMAP visualization of target features.** To better develop intuition, we visualize the learned features for ProDA in Figure 5. The model before the domain adaptation mixes the features of the same class. The conventional self-training could produce more separated feature space, yet it is still hard for linear classification. When we apply prototypical pseudo label denoising, features among differ-

threshold	0.0	0.2	0.4	0.6	0.8	0.9	0.95
mIoU	53.7	53.8	53.7	53.7	53.8	53.7	53.3

Table 4: The effect of threshold in pseudo label selection during the prototypical label denoising.

momentum	0.99	0.999	0.9999	0.99999
mIoU	53.5	53.6	53.7	52.3

Table 5: The effect of prototype momentum during the prototype online update.

ent classes are better separated, though the distribution is still dispersed. In comparison, the full model gives the most compact feature clusters that are amenable to classification.

#### 5.4. Parameter sensitivity analysis.

To showcase that ProDA is robust to the hyper-parameter selection, we analyze the impact of parameters. In Table 4, we use different threshold values to select the pseudo labels and the performance is not sensitive to this threshold as opposed to the conventional self-training. Hence, ProDA does not apply thresholding for convenient usage. We also study the effect of momentum value when online computing the prototypes, and ProDA is robust to a wide numerical range as shown in Table 5.

## 6. Conclusions

In this paper, we propose ProDA which resorts to prototypes to online denoise the pseudo labels and learn the compact feature space for the target domain. Knowledge distillation to a self-supervised pretrained model further boosts the performance. The proposed method outperforms state-of-the-art methods by a large margin, greatly reducing the gap with supervised learning. We will make the code publicly available. We believe the proposed method is a general amelioration to the self-training and we hope to explore its capability in other tasks in the future.



## References

- [1] Alexey Abramov, Christopher Bayer, and Claudio Heller. Keep it simple: Image statistics matching for domain adaptation. *arXiv preprint arXiv:2005.12551*, 2020. 2
- [2] Yuki Markus Asano, Christian Rupprecht, and Andrea Vedaldi. Self-labelling via simultaneous clustering and representation learning. *arXiv preprint arXiv:1911.05371*, 2019. 2
- [3] Shai Ben-David, John Blitzer, Koby Crammer, Alex Kulesza, Fernando Pereira, and Jennifer Wortman Vaughan. A theory of learning from different domains. *Machine learning*, 79(1-2):151–175, 2010. 2
- [4] Mathilde Caron, Piotr Bojanowski, Armand Joulin, and Matthijs Douze. Deep clustering for unsupervised learning of visual features. In *European Conference on Computer Vision*, 2018. 1, 2
- [5] Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin. Unsupervised learning of visual features by contrasting cluster assignments. *arXiv preprint arXiv:2006.09882*, 2020. 2
- [6] Wei-Lun Chang, Hui-Po Wang, Wen-Hsiao Peng, and Wei-Chen Chiu. All about structure: Adapting structural information across domains for boosting semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1900–1909, 2019. 1, 2
- [7] Chaoqi Chen, Weiping Xie, Wenbing Huang, Yu Rong, Xinghao Ding, Yue Huang, Tingyang Xu, and Junzhou Huang. Progressive feature alignment for unsupervised domain adaptation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 627–636, 2019. 2
- [8] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *IEEE transactions on pattern analysis and machine intelligence*, 40(4):834–848, 2017. 2, 5
- [9] Minghao Chen, Hongyang Xue, and Deng Cai. Domain adaptation for semantic segmentation with maximum squares loss. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2090–2099, 2019. 2
- [10] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. *arXiv preprint arXiv:2002.05709*, 2020. 2
- [11] Ting Chen, Simon Kornblith, Kevin Swersky, Mohammad Norouzi, and Geoffrey Hinton. Big self-supervised models are strong semi-supervised learners. *arXiv preprint arXiv:2006.10029*, 2020. 5
- [12] Yun-Chun Chen, Yen-Yu Lin, Ming-Hsuan Yang, and Jia-Bin Huang. Crdoco: Pixel-level domain transfer with cross-domain consistency. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1791–1800, 2019. 2
- [13] Jaehoon Choi, Taekyung Kim, and Changick Kim. Self-ensembling with gan-based data augmentation for domain adaptation in semantic segmentation. In *Proceedings of the IEEE international conference on computer vision*, pages 6830–6840, 2019. 2
- [14] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban scene understanding. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3213–3223, 2016. 2, 5
- [15] Ekin D Cubuk, Barret Zoph, Jonathon Shlens, and Quoc V Le. Randaugment: Practical automated data augmentation with a reduced search space. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 702–703, 2020. 5, 12
- [16] Terrance DeVries and Graham W Taylor. Improved regularization of convolutional neural networks with cutout. *arXiv preprint arXiv:1708.04552*, 2017. 5, 12
- [17] Jeff Donahue and Karen Simonyan. Large scale adversarial representation learning. In *Advances in Neural Information Processing Systems*, pages 10542–10552, 2019. 2
- [18] Liang Du, Jingang Tan, Hongye Yang, Jianfeng Feng, Xiangyang Xue, Qibao Zheng, Xiaoqing Ye, and Xiaolin Zhang. Ssf-dan: Separated semantic feature based domain adaptation network for semantic segmentation. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 982–991, 2019. 2
- [19] Spyros Gidaris, Praveer Singh, and Nikos Komodakis. Unsupervised representation learning by predicting image rotations. *arXiv preprint arXiv:1803.07728*, 2018. 2
- [20] Rui Gong, Wen Li, Yuhua Chen, and Luc Van Gool. Dlow: Domain flow for adaptation and generalization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2477–2486, 2019. 2
- [21] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in neural information processing systems*, pages 2672–2680, 2014. 2
- [22] Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre H Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Daniel Guo, Mohammad Gheshlaghi Azar, et al. Bootstrap your own latent: A new approach to self-supervised learning. *arXiv preprint arXiv:2006.07733*, 2020. 2
- [23] Jiangfan Han, Ping Luo, and Xiaogang Wang. Deep self-learning from noisy labels. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 5138–5147, 2019. 2
- [24] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9729–9738, 2020. 2, 4
- [25] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 5
- [26] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015. 2
- [27] Judy Hoffman, Eric Tzeng, Taesung Park, Jun-Yan Zhu,

- Phillip Isola, Kate Saenko, Alexei Efros, and Trevor Darrell. Cycada: Cycle-consistent adversarial domain adaptation. In *International conference on machine learning*, pages 1989–1998. PMLR, 2018. 2, 5, 6
- [28] Weixiang Hong, Zhenzhen Wang, Ming Yang, and Junsong Yuan. Conditional generative adversarial network for structured domain adaptation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1335–1344, 2018. 1, 2
- [29] Abhishek Kumar, Prasanna Sattigeri, Kahini Wadhawan, Leonid Karlinsky, Rogerio Feris, Bill Freeman, and Gregory Wornell. Co-regularized alignment for unsupervised domain adaptation. In *Advances in Neural Information Processing Systems*, pages 9345–9356, 2018. 2
- [30] Gustav Larsson, Michael Maire, and Gregory Shakhnarovich. Learning representations for automatic colorization. In *European conference on computer vision*, pages 577–593. Springer, 2016. 2
- [31] Chen-Yu Lee, Tanmay Batra, Mohammad Haris Baig, and Daniel Ulbricht. Sliced wasserstein discrepancy for unsupervised domain adaptation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 10285–10295, 2019. 2
- [32] Kuang-Huei Lee, Xiaodong He, Lei Zhang, and Linjun Yang. Cleannet: Transfer learning for scalable image classifier training with label noise. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5447–5456, 2018. 2
- [33] Junnan Li, Yongkang Wong, Qi Zhao, and Mohan S Kankanhalli. Learning to learn from noisy labeled data. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5051–5059, 2019. 2
- [34] Yuncheng Li, Jianchao Yang, Yale Song, Liangliang Cao, Jiebo Luo, and Li-Jia Li. Learning from noisy labels with distillation. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1910–1918, 2017. 2
- [35] Yunsheng Li, Lu Yuan, and Nuno Vasconcelos. Bidirectional learning for domain adaptation of semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6936–6945, 2019. 1, 2, 5, 6, 17, 18
- [36] Mingsheng Long, Yue Cao, Jianmin Wang, and Michael Jordan. Learning transferable features with deep adaptation networks. In *International conference on machine learning*, pages 97–105. PMLR, 2015. 2
- [37] Yawei Luo, Liang Zheng, Tao Guan, Junqing Yu, and Yi Yang. Taking a closer look at domain shift: Category-level adversaries for semantics consistent domain adaptation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2507–2516, 2019. 2, 5, 6, 17, 18
- [38] Leland McInnes, John Healy, and James Melville. Umap: Uniform manifold approximation and projection for dimension reduction. *arXiv preprint arXiv:1802.03426*, 2018. 8
- [39] Robert Mendel, Luis Antonio de Souza, David Rauber, João Paulo Papa, and Christoph Palm. Semi-supervised segmentation based on error-correcting supervision. In *European Conference on Computer Vision*, pages 141–157. Springer, 2020. 2
- [40] Ishan Misra and Laurens van der Maaten. Self-supervised learning of pretext-invariant representations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6707–6717, 2020. 2
- [41] Sebastian Nowozin, Botond Cseke, and Ryota Tomioka. f-gan: Training generative neural samplers using variational divergence minimization. In *Advances in neural information processing systems*, pages 271–279, 2016. 2
- [42] Avital Oliver, Augustus Odena, Colin A Raffel, Ekin Dogus Cubuk, and Ian Goodfellow. Realistic evaluation of deep semi-supervised learning algorithms. In *Advances in neural information processing systems*, pages 3235–3246, 2018. 4
- [43] Guo-Jun Qi and Jiebo Luo. Small data challenges in big data era: A survey of recent progress on unsupervised and semi-supervised methods. *arXiv preprint arXiv:1903.11260*, 2019. 2
- [44] Hang Qi, Matthew Brown, and David G Lowe. Low-shot learning with imprinted weights. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5822–5830, 2018. 4
- [45] Stephan R Richter, Vibhav Vineet, Stefan Roth, and Vladlen Koltun. Playing for data: Ground truth from computer games. In *European conference on computer vision*, pages 102–118. Springer, 2016. 1, 2, 5
- [46] German Ros, Laura Sellart, Joanna Materzynska, David Vazquez, and Antonio M Lopez. The synthia dataset: A large collection of synthetic images for semantic segmentation of urban scenes. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3234–3243, 2016. 1, 2, 5
- [47] Kuniaki Saito, Donghyun Kim, Stan Sclaroff, Trevor Darrell, and Kate Saenko. Semi-supervised domain adaptation via minimax entropy. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 8050–8058, 2019. 2
- [48] Kuniaki Saito, Donghyun Kim, Stan Sclaroff, and Kate Saenko. Universal domain adaptation through self supervision. *arXiv preprint arXiv:2002.07953*, 2020. 2
- [49] Kuniaki Saito, Kohei Watanabe, Yoshitaka Ushiku, and Tatsuya Harada. Maximum classifier discrepancy for unsupervised domain adaptation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3723–3732, 2018. 1, 2
- [50] Swami Sankaranarayanan, Yogesh Balaji, Arpit Jain, Ser Nam Lim, and Rama Chellappa. Learning from synthetic data: Addressing domain shift for semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3752–3761, 2018. 2
- [51] Jake Snell, Kevin Swersky, and Richard Zemel. Prototypical networks for few-shot learning. In *Advances in neural information processing systems*, pages 4077–4087, 2017. 4
- [52] Hwanjun Song, Minseok Kim, and Jae-Gil Lee. Selfie: Refurbishing unclean samples for robust deep learning. In *International Conference on Machine Learning*, pages 5907–5915, 2019. 2
- [53] Hwanjun Song, Minseok Kim, Dongmin Park, and Jae-Gil Lee. Learning from noisy labels with deep neural networks: A survey. *arXiv preprint arXiv:2007.08199*, 2020. 2
- [54] Yu Sun, Eric Tzeng, Trevor Darrell, and Alexei A Efros. Unsupervised domain adaptation through self-supervision.

- arXiv preprint arXiv:1909.11825*, 2019. 2
- [55] Yi-Hsuan Tsai, Wei-Chih Hung, Samuel Schulter, Kihyuk Sohn, Ming-Hsuan Yang, and Manmohan Chandraker. Learning to adapt structured output space for semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7472–7481, 2018. 1, 2, 5, 6
- [56] Y.-H. Tsai, W.-C. Hung, S. Schulter, K. Sohn, M.-H. Yang, and M. Chandraker. Learning to adapt structured output space for semantic segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 5, 14, 17, 18
- [57] Yi-Hsuan Tsai, Kihyuk Sohn, Samuel Schulter, and Manmohan Chandraker. Domain adaptation for structured output via discriminative patch representations. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1456–1465, 2019. 5, 6
- [58] Tuan-Hung Vu, Himalaya Jain, Maxime Bucher, Matthieu Cord, and Patrick Pérez. Advent: Adversarial entropy minimization for domain adaptation in semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2517–2526, 2019. 2, 5, 6, 17, 18
- [59] Ziyu Wan, Bo Zhang, Dongdong Chen, Pan Zhang, Dong Chen, Jing Liao, and Fang Wen. Bringing old photos back to life. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2747–2757, 2020. 1, 2
- [60] Haoran Wang, Tong Shen, Wei Zhang, Lingyu Duan, and Tao Mei. Classes matter: A fine-grained adversarial approach to cross-domain semantic segmentation. *arXiv preprint arXiv:2007.09222*, 2020. 2
- [61] Haoran Wang, Tong Shen, Wei Zhang, Lingyu Duan, and Tao Mei. Classes matter: A fine-grained adversarial approach to cross-domain semantic segmentation. In *The European Conference on Computer Vision (ECCV)*, August 2020. 5, 6, 17, 18
- [62] Xinshao Wang, Yang Hua, Elyor Kodirov, and Neil M Robertson. Proselfc: Progressive self label correction for training robust deep neural networks. *arXiv preprint arXiv:2005.03788*, 2020. 2
- [63] Yisen Wang, Xingjun Ma, Zaiyi Chen, Yuan Luo, Jinfeng Yi, and James Bailey. Symmetric cross entropy for robust learning with noisy labels. In *IEEE International Conference on Computer Vision*, 2019. 2, 4, 12
- [64] Zuxuan Wu, Xin Wang, Joseph E Gonzalez, Tom Goldstein, and Larry S Davis. Ace: adapting to changing environments for semantic segmentation. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2121–2130, 2019. 2
- [65] Zhirong Wu, Yuanjun Xiong, Stella X Yu, and Dahua Lin. Unsupervised feature learning via non-parametric instance discrimination. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3733–3742, 2018. 2
- [66] Qizhe Xie, Minh-Thang Luong, Eduard Hovy, and Quoc V Le. Self-training with noisy student improves imagenet classification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10687–10698, 2020. 2
- [67] Jiaolong Xu, Liang Xiao, and Antonio M López. Self-supervised domain adaptation for computer vision tasks. *IEEE Access*, 7:156694–156706, 2019. 2
- [68] Jihan Yang, Ruijia Xu, Ruiyu Li, Xiaojuan Qi, Xiaoyong Shen, Guanbin Li, and Liang Lin. An adversarial perturbation oriented domain adaptation approach for semantic segmentation. In *AAAI*, pages 12613–12620, 2020. 5, 6
- [69] Qiming Zhang, Jing Zhang, Wei Liu, and Dacheng Tao. Category anchor-guided unsupervised domain adaptation for semantic segmentation. In *Advances in Neural Information Processing Systems*, pages 433–443, 2019. 1, 2, 5, 6, 17, 18
- [70] Richard Zhang, Phillip Isola, and Alexei A Efros. Split-brain autoencoders: Unsupervised learning by cross-channel prediction. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1058–1067, 2017. 2
- [71] Zhilu Zhang and Mert Sabuncu. Generalized cross entropy loss for training deep neural networks with noisy labels. In *Advances in neural information processing systems*, pages 8778–8788, 2018. 2
- [72] Han Zhao, Remi Tachet des Combes, Kun Zhang, and Geoffrey J Gordon. On learning invariant representation for domain adaptation. *arXiv preprint arXiv:1901.09453*, 2019. 2
- [73] Zhedong Zheng and Yi Yang. Rectifying pseudo label learning via uncertainty estimation for domain adaptive semantic segmentation. *arXiv preprint arXiv:2003.03773*, 2020. 2, 5, 6, 17, 18
- [74] Zhedong Zheng and Yi Yang. Unsupervised scene adaptation with memory regularization in vivo. In *IJCAI*, 2020. 2
- [75] Yang Zou, Zhiding Yu, BVK Vijaya Kumar, and Jinsong Wang. Unsupervised domain adaptation for semantic segmentation via class-balanced self-training. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 289–305, 2018. 1, 2, 3, 5, 6, 17, 18
- [76] Yang Zou, Zhiding Yu, Xiaofeng Liu, BVK Kumar, and Jinsong Wang. Confidence regularized self-training. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 5982–5991, 2019. 1, 2, 3, 4, 5, 6

## Appendix A. Influences of Design Choices

**Prototype initialization strategy.** In our implementation, the proposed ProDA initializes the prototypes of the target domain according to the pseudo predictions for the target images. Alternatively, the target prototypes can also be initialized according to the ground truth labeling in the source domain. However, both choices have their pros and cons: the former suffers from the noises in the pseudo labels whereas the latter suffers from domain gap as the prototypes of the two domains may not accurately align. Table 6 shows that the two initialization strategies induce comparable results, as the prototypes are online updated and can rapidly converge to the true cluster centroids. The quantitative performance is measured on the dataset GTA5  $\rightarrow$  Cityscapes, whereas the other dataset shows similar results.

	source ground truth	target pseudo label
mIoU	53.6	53.7

Table 6: The performance of different target prototype initialization strategies. Here we only report the performance for the 1st training stage in the gta5  $\rightarrow$  Cityscapes task.

**Strong augmentation.** In the target structure learning, we take weak and strong augmentation views for the target image. We employ random crop for weak augmentation and explore the effects of different augmentation types for the strong augmented view. As shown in Table 7, random crop only gives the mIoU score 52.7, whereas adding RandAugment [15] and CutOut [16] respectively improve the mIoU by 0.78 and 0.5. The strongest augmentation gives the best performance, indicating the importance of data augmentation when learning the compact feature space for the target domain.

	crop	crop & RandAug	crop & Cutout	crop & RandAug & Cutout
mIoU	52.7	53.5	53.2	53.7

Table 7: The influence of various strong augmentations. Here we only report the performance for the 1st training stage in the gta5  $\rightarrow$  Cityscapes task.

**Effect of temperature during prototypical denoising.** We rely on the prototypical context to rectify the pseudo labels. We compute the softmax over feature distance to all the prototypes, and the softmax temperature  $\tau$  influences the denoising effect and requires balancing: when  $\tau \rightarrow 0$ , only the nearest prototype dominates whereas  $\tau \rightarrow \infty$  causes that all the prototypes are accounted equally. The influence of the temperature is shown in Table 8. We empirically set  $\tau = 1$  in our experiments.

	0.1	0.5	1	2	3	5	10
mIoU	48.8	52.1	53.7	51.9	47.5	44.9	40.9

Table 8: The effects of temperature during the prototypical denoising. Here we only report the performance for the 1st training stage in the gta5  $\rightarrow$  Cityscapes task.

**Symmetric cross-entropy loss.** We employ the symmetric cross-entropy loss (SCE) for robust learning to stabilize the early training phase. The SCE has coefficients  $\alpha$  and  $\beta$  that balance the cross-entropy and the reverse cross-entropy. Table 9 shows that the final result is not sensitive to these hyper-parameters if  $\beta$  is not too small. Here, we follow the suggested setting as [63], *i.e.*,  $\alpha = 0.1, \beta = 1$ .

**The effect of loss weight.** Table 10 shows that the final result is not sensitive to the KL loss weight ( $\gamma_1$ ) and the regularization loss weight ( $\gamma_2$ ). In GTA5  $\rightarrow$  Cityscapes, we set  $\gamma_1 = 10$  and  $\gamma_2 = 0.1$ , while in SYNTHIA  $\rightarrow$  Cityscapes, we set  $\gamma_1 = 10$  and  $\gamma_2 = 0$ .

$\alpha \backslash \beta$	0.1	0.5	1	5
0.01	46.4	52.7	53.8	53.6
0.1	47.6	52.9	53.7	53.5
0.5	50.4	53.1	53.3	53.5
1	51.1	52.7	53.1	53.6

Table 9: The influence of  $\alpha$  and  $\beta$  in the symmetric cross-entropy (SCE) loss. Here we only report the performance for the 1st training stage in the gta5  $\rightarrow$  Cityscapes task.

$\gamma_1 \backslash \gamma_2$	0.02	0.1	0.2
2	52.9	53.7	53.5
10	53.2	53.7	53.4
20	53.4	53.6	52.1
50	53.6	52.0	52.1

Table 10: The influence of the KL loss weight ( $\gamma_1$ ) and the regularization loss weight ( $\gamma_2$ ). Here we only report the performance for the 1st training stage in the gta5  $\rightarrow$  Cityscapes task.

## Appendix B. Algorithm

The training procedure of our ProDA is summarized in Algorithm 1, which is composed of three stages. The first stage consists of prototypical pseudo label denoising and target structure learning. In the second and third stages, we apply knowledge distillation to a self-supervised model. For detailed equations and loss functions, please refer to our main paper.

---

### Algorithm 1: ProDA

---

**Input:** training dataset:  $(\mathcal{X}_s, \mathcal{Y}_s, \mathcal{X}_t)$ ; prototype momentum:  $\lambda$ ; weak, strong augmentations:  $\mathcal{T}, \mathcal{T}'$ ; the pretrained SimCLRv2 model:  $h'_\theta$ ; pseudo label selection threshold:  $T$ ;

**Output:** the output model  $h_\theta$ .

- 1 Warmup:  $h_\theta = g_\theta \circ f_\theta \leftarrow (\mathcal{X}_s, \mathcal{Y}_s, \mathcal{X}_t)$  according to [56];
- 2 Generate soft pseudo label:  $p_{t,0} \leftarrow h_\theta(\mathcal{X}_t)$ ;
- 3 Prototype initialization:  $\eta_c \leftarrow (f_\theta, \mathcal{X}_t)$ ;
- 4 EMA model initialization:  $\tilde{h}_\theta \leftarrow h_\theta$ ;
- 5 **for**  $m \leftarrow 0$  **to**  $epochs$  **do**
- 6     **for**  $i \leftarrow 0$  **to**  $len(\mathcal{X}_t)$  **do**
- 7         Get source images  $x_s^{(i)}$ ;
- 8         Train the model  $h_\theta$  using loss  $\ell_{ce}^s$ ;
- 9
- 10        Get target images  $x_t^{(i)}$ ;
- 11        Calculate the denoising weight  $\omega_t^{(i,k)}$ ;
- 12        Update the pseudo label  $\hat{y}_t^{(i,k)}$ ;
- 13        Train model  $h_\theta$  using loss  $\ell_{sce}^t$ ;
- 14
- 15        Calculate the soft label  $z_{\mathcal{T}}, z_{\mathcal{T}'}$ ;
- 16        Train the model  $h_\theta$  using loss  $\ell_{kl}^t$  and  $\ell_{reg}^t$ ;
- 17
- 18        Calculate the batch prototype  $\eta'_c$ ;
- 19         $\eta_c \leftarrow \lambda \eta_c + (1 - \lambda) \eta'_c$ ;
- 20        Update the EMA model  $\tilde{h}_\theta$ ;
- 21
- 22 **for**  $stage \leftarrow 1$  **to**  $2$  **do**
- 23     Generate the pseudo label:  $\hat{y}_t \leftarrow \xi(h_\theta(\mathcal{X}_t), T)$ ;
- 24     Student model initialization:  $h_\theta^\dagger \leftarrow h'_\theta$ ;
- 25     **for**  $m \leftarrow 0$  **to**  $epochs$  **do**
- 26         **for**  $i \leftarrow 0$  **to**  $len(\mathcal{X}_t)$  **do**
- 27             Get source images  $x_s^{(i)}$ ;
- 28             Tune the model  $h_\theta^\dagger$  using loss  $\ell_{ce}^s$ ;
- 29
- 30             Get target images  $x_t^{(i)}$ ;
- 31             Calculate the teacher probability  $h_\theta(x_t^{(i)})$ ;
- 32             Calculate the student probability  $h_\theta^\dagger(x_t^{(i)})$ ;
- 33             Tune the model  $h_\theta^\dagger$  using loss  $\ell_{ce}^t$  and KL loss;
- 34
- 34      $h_\theta \leftarrow h_\theta^\dagger$ ;

---

### Appendix C. Detailed Ablation study

Here we show a detailed ablation study for all the 19 classes on GTA5 → Cityscapes.

components		road	sideway	building	wall	fence	pole	light	sign	vege.	terrace	sky	person	rider	car	truck	bus	train	motor	bike	mIoU	gain		
init.	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			
	warm up	86.7	34.2	79.3	26.6	21.6	38.4	33.7	15.8	82.1	31.0	73.2	60.4	21.0	82.3	23.2	32.0	2.9	24.1	20.9	41.6	+5.0		
stage 1	ST	✓																				mIoU	gain	
	SCE	✓																						
	PD		✓																					
	SL			✓																				
		✓	✓	✓	✓																			
stage 2	self distill.	90.0	57.4	81.8	42.0	40.2	43.8	50.3	50.9	87.6	42.6	80.0	69.2	32.9	87.8	45.5	56.9	0.0	46.0	55.4	55.8	+19.2		
	stage 1 init.	✓	✓																					
	sup init.			✓																				
	self-sup init.	✓																						
stage 3	self distill.	91.4	53.3	83.4	41.3	37.8	43.9	53.0	47.9	88.3	46.1	79.9	70.5	33.2	89.0	48.4	54.6	0.0	50.5	56.7	56.3	+19.7		
	stage 1 init.	✓																						
	sup init.	✓		✓																				
	self-sup init.	✓																						
stage 3	✓			✓																				
		87.8	56.0	79.7	46.3	44.8	45.6	53.5	53.5	88.6	45.2	82.1	70.7	39.2	88.8	45.5	59.4	1.0	48.9	56.4	57.5	+20.9		

Table 11: Ablation study of each proposed component. The whole training involves three stages, where knowledge distillation can be applied in the last two stages. Here, ST stands for self-training, PD for prototypical denoising, and SL for structure learning.

## Appendix D. Qualitative comparison

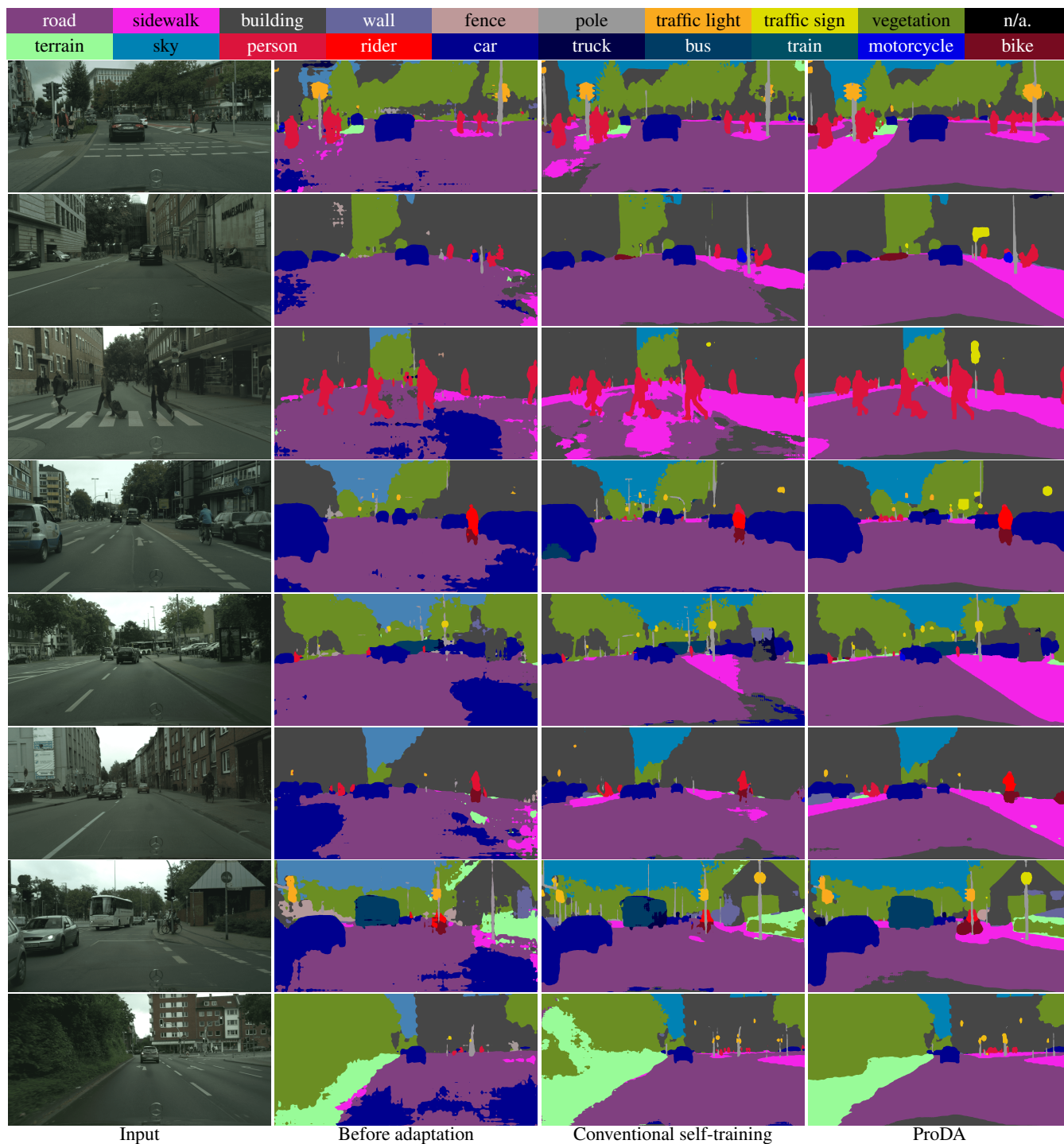


Figure 6: Qualitative results of semantic segmentation on the Cityscapes dataset. From left to right: input, before adaptation, conventional self-training, ProDA.



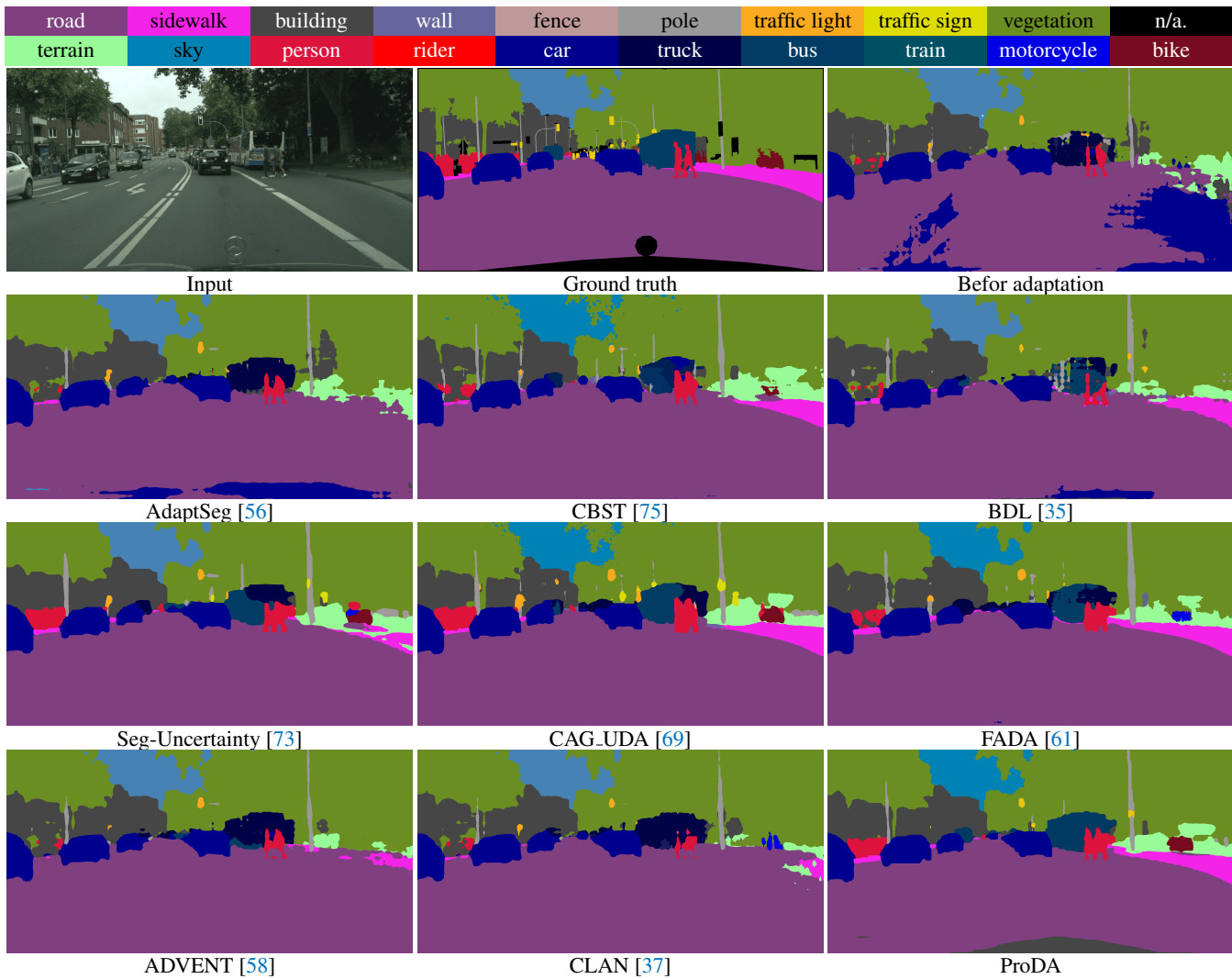


Figure 7: Qualitative comparisons of different methods.

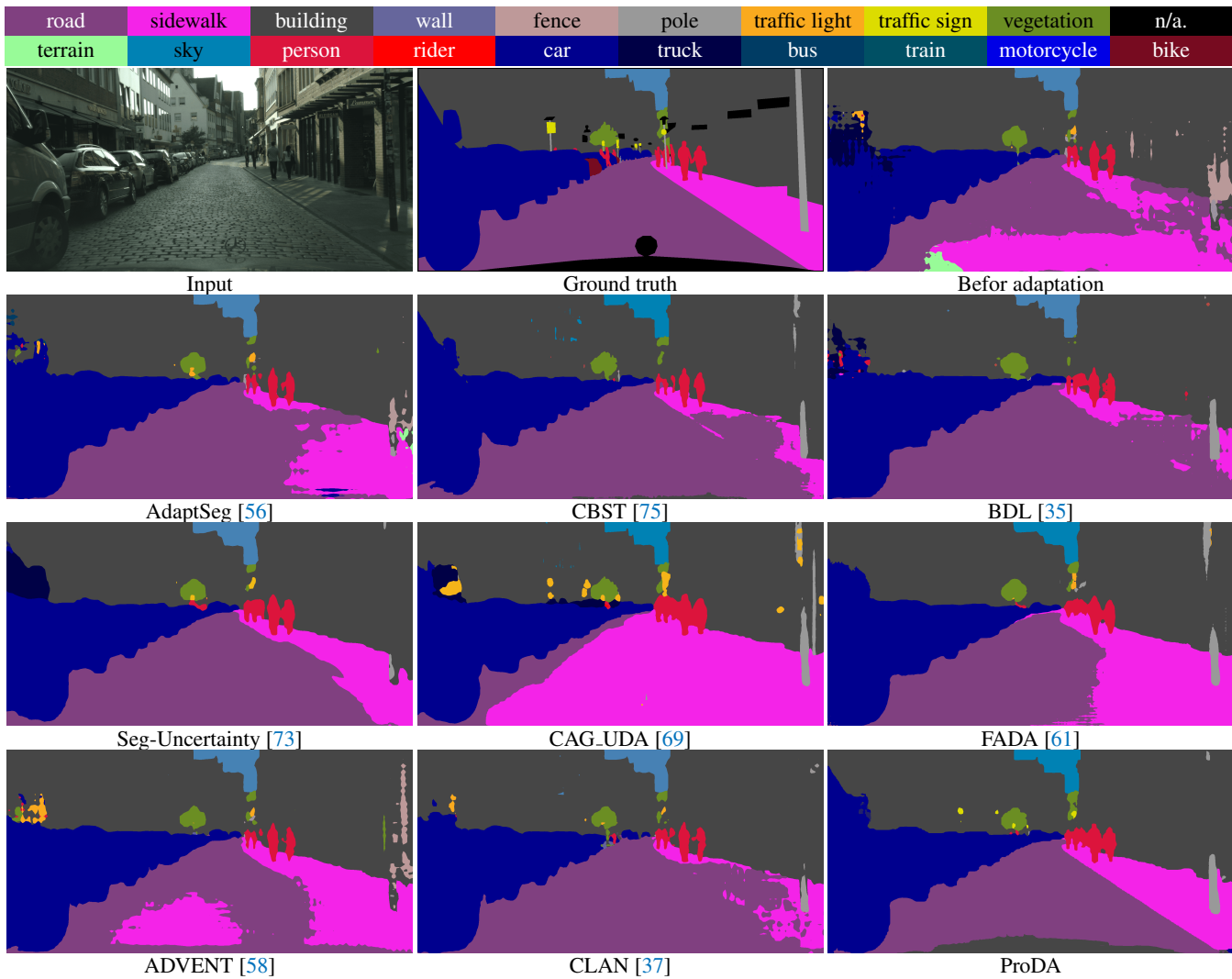


Figure 8: Qualitative comparisons of different methods.