Prototypical Pseudo Label Denoising and Target Structure Learning for Domain Adaptive Semantic Segmentation (Supplementary Material)

Pan Zhang¹ *, Bo Zhang², Ting Zhang², Dong Chen², Yong Wang¹, Fang Wen² ¹University of Science and Technology of China ²Microsoft Research Asia

1. 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 1 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 labe
nIoU	53.6	53.7

Table 1: 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 2, random crop only gives the mIoU score 52.7, whereas adding RandAugment [1] and CutOut [2] 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 2: 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 τ influences the denoising effect and requires balancing: when $\tau \to 0$, only the nearest prototype dominates whereas $\tau \to \infty$ causes that all the prototypes are accounted equally. The influence of the temperature is shown in Table 3. We empirically set $\tau = 1$ in our experiments.

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 α and β that balance the cross-entropy and the reverse cross-entropy. Table 4 shows that the final result is not sensitive to these hyper-parameters if β is not too small. Here, we follow the suggested setting as [8], *i.e.*, $\alpha = 0.1$, $\beta = 1$.

The effect of loss weight. Table 5 shows that the final result is not sensitive to the KL loss weight (γ_1) and the regularization loss weight (γ_2). In GTA5 \rightarrow Cityscapes, we set $\gamma_1 = 10$ and $\gamma_1 = 0.1$, while in SYNTHIA \rightarrow Cityscapes, we set $\gamma_1 = 10$ and $\gamma_1 = 0.1$.

^{*}This work is done during the first author's internship at Microsoft Research Asia.

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

β	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 4: The influence of α and β in the symmetric cross-entropy (SCE) loss. Here we only report the performance for the 1st training stage in the gta5 \rightarrow Cityscapes task.

γ_2 γ_1	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 5: The influence of the KL loss weight (γ_1) and the regularization loss weight (γ_2). Here we only report the performance for the 1st training stage in the gta5 \rightarrow Cityscapes task.

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

3. Detailed Ablation study

	components				road	sideway	building	wall	fence	pole	light	sign	vege.	terrace	sky	person	rider	car	truck	bus	train	motor	bike	mIoU	gain
init.		sou warn	rce 1 up		75.8 86.7	16.8 34.2	77.2 79.3	12.5 26.6	21.0 21.6	25.5 38.4	30.1 33.7	20.1 15.8	81.3 82.1	24.6 31.0	70.3 73.2	53.8 60.4	26.4 21.0	49.9 82.3	17.2 23.2	25.9 32.0	6.5 2.9	25.3 24.1	36.0 20.9	36.6 41.6	+5.0
	ST	SCE	PD	SL																				mIoU	gain
	\checkmark				87.0	39.7	77.5	31.5	25.7	41.5	38.7	20.3	84.6	38.2	74.1	63.7	21.7	86.0	29.0	37.5	0.3	34.9	26.2	45.2	+8.6
stage 1	\checkmark	\checkmark			87.7	36.8	78.2	30.9	24.8	41.5	40.0	23.2	83.0	35.3	72.9	64.1	24.6	85.9	32.9	36.5	2.0	31.0	35.0	45.6	+9.0
	\checkmark	\checkmark	\checkmark		93.2	56.7	84.1	40.4	37.5	39.5	44.1	35.1	87.1	43.2	80.3	65.8	29.8	87.7	29.6	41.9	0.0	44.4	52.6	52.3	+15.7
	\checkmark	\checkmark		\checkmark	89.0	38.6	80.7	37.1	27.2	42.8	41.5	20.7	85.8	42.4	74.8	64.8	17.8	87.6	30.8	39.4	0.0	41.0	34.6	47.6	+11.0
	\checkmark	\checkmark	\checkmark	\checkmark	91.5	52.4	82.9	42.0	35.7	40.0	44.4	43.3	87.0	43.8	79.5	66.5	31.4	86.7	41.1	52.5	0.0	45.4	53.8	53.7	+17.1
	self distill.	stage 1 init.	sup init.	self-sup init.																				mIoU	gain
stage 2				\checkmark	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 2	\checkmark	\checkmark			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
	\checkmark		\checkmark		91.0	50.2	83.1	40.1	39.8	43.5	51.9	48.1	87.9	45.9	78.5	69.6	34.3	87.9	41.3	56.6	0.0	51.7	57.1	55.7	+19.1
	\checkmark			\checkmark	89.4	56.5	81.3	46.3	42.7	45.1	52.2	51.3	88.5	47.0	82.9	69.3	36.5	87.4	46.0	57.5	0.6	45.9	54.5	56.9	+20.3
stage 3	\checkmark			\checkmark	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

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

Table 6: 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.

4. Qualitative comparison



Figure 1: Qualitative results of semantic segmentation on the Cityscapes dataset. From left to right: input, before adaptation, conventional self-training, ProDA.









References

- 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.
- [2] Terrance DeVries and Graham W Taylor. Improved regularization of convolutional neural networks with cutout. *arXiv preprint arXiv:1708.04552*, 2017. 1
- [3] 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. 6, 7
- [4] 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. 6, 7
- [5] 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. 3, 6, 7
- [6] 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. 6, 7
- [7] Haoran Wang, Tong Shen, Wei Zhang, Lingyu Duan, and Tao Mei. Classes matter: A fine-grained adversarial approach to crossdomain semantic segmentation. In *The European Conference on Computer Vision (ECCV)*, August 2020. 6, 7
- [8] 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. 1
- [9] 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. 6, 7
- [10] Zhedong Zheng and Yi Yang. Rectifying pseudo label learning via uncertainty estimation for domain adaptive semantic segmentation. *arXiv preprint arXiv:2003.03773*, 2020. 6, 7
- [11] 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. 6, 7