

# PØDA: Prompt-driven Zero-shot Domain Adaptation

## – Supplementary Material –

Mohammad Fahes<sup>1</sup> Tuan-Hung Vu<sup>1,2</sup> Andrei Bursuc<sup>1,2</sup> Patrick Pérez<sup>1,2</sup> Raoul de Charette<sup>1</sup>  
<sup>1</sup>Inria <sup>2</sup>Valeo.ai

### A. Overall pseudo-code of PØDA

Algorithm 2 presents the high-level pseudo-code of PØDA: from *source-only* training as model initialization, to prompt-driven feature augmentation, to zero-shot model adaptation.

### B. Experimental details

**Feature augmentation.** PIN operates on image features. For augmentation, we optimize  $(\boldsymbol{\mu}, \boldsymbol{\sigma})$  of source feature map  $\mathbf{f}_s$ ; it is done in batches for the sake of speed. We fix the batch size  $b = 16$  and the learning rate  $lr = 1.0$ .

**Style mixing.** In the discussion of PØDA (Sec. 4.4 and Tab. 7), we presented the performance gains that style-mixing [6] brings to our method in three settings. By randomly mixing original and augmented statistics, we introduce certain perturbations to the final augmented features. The mixed statistics  $\boldsymbol{\mu}_{\text{mix}}, \boldsymbol{\sigma}_{\text{mix}}$  are given by:

$$\boldsymbol{\mu}_{\text{mix}} = \alpha \boldsymbol{\mu}_t + (1 - \alpha) \boldsymbol{\mu}_s, \quad (4)$$

$$\boldsymbol{\sigma}_{\text{mix}} = \alpha \boldsymbol{\sigma}_t + (1 - \alpha) \boldsymbol{\sigma}_s, \quad (5)$$

where  $\alpha \in \mathbb{R}^c$  are per-channel mixing weights uniformly sampled in  $[0, 1]$ , similarly to [6]; multiplications are element-wise. Finally, the augmented features are computed as follows:

$$\mathbf{f}_{s \rightarrow t} = \text{PIN}(\mathbf{f}_s, \boldsymbol{\mu}_{\text{mix}}, \boldsymbol{\sigma}_{\text{mix}}), \quad (6)$$

with prompt-driven instance normalization PIN defined in Eq. 2.

### C. Additional experiments

**Effect of style mining initialization.** In our feature optimization step, we initialize  $(\boldsymbol{\mu}, \boldsymbol{\sigma})$  with  $(\boldsymbol{\mu}(\mathbf{f}_s), \boldsymbol{\sigma}(\mathbf{f}_s))$ . In Tab. 12, we report results using different initialization strategies. Starting from pre-defined or random initialization, instead of from original statistics, degrades badly the performance. As we do not use

---

**Algorithm 2:** Prompt-driven Zero-shot DA

---

**Input:** Source dataset  $\mathcal{D}_s = \{(\mathbf{x}_s, \mathbf{y}_s)\}$   
 CLIP encoders  $E_{\text{img}}$  and  $E_{\text{txt}}$   
 Target domain description TrgPrompt  
 Feature backbone  $M_{\text{feat}} \leftarrow E_{\text{img}}$   
 Source model:  $M = (M_{\text{feat}}, M_{\text{cls}})$   
**Result:** Target-adapted model  $M' = (M_{\text{feat}}, M'_{\text{cls}})$   
 // Initialization  
 1 TrgEmb =  $E_{\text{txt}}(\text{TrgPrompt})$   
 2  $M_{\text{cls}} \leftarrow \text{train}(M_{\text{cls}}, \mathcal{D}_s)$  ▷ source-only training  
 // Feature Augmentation  
 3  $\mathcal{F}_s \leftarrow \text{feat-ext}(M_{\text{feat}}, \{\mathbf{x}_s\})$   
 4  $\mathcal{S}_{s \rightarrow t} \leftarrow \text{augment}(\mathcal{F}_s, \text{TrgEmb})$   
 // Adaptation  
 5  $M'_{\text{cls}} \leftarrow \text{fine-tune}(M_{\text{cls}}, \mathcal{F}_s, \mathcal{S}_{s \rightarrow t}, \{\mathbf{y}_s\})$   
 ▷ fine-tuning

---

$\boldsymbol{\mu}^0$	$\boldsymbol{\sigma}^0$	mIoU
$\boldsymbol{\mu}(\mathbf{f}_s)$	$\boldsymbol{\sigma}(\mathbf{f}_s)$	<b>25.03</b> $\pm 0.48$
$\mathbf{0}$	$\mathbf{1}$	8.59 $\pm 0.82$
$\sim \mathcal{N}(\mathbf{0}, \mathbf{I})$	$\sim \mathcal{N}(\mathbf{0}, \mathbf{I})$	6.80 $\pm 0.92$

Table 12: **Effect of style initialization.** Performance (in mIoU) of PØDA on ACDC-Night val set (Cityscapes as source), with different style statistics initializations. Starting from source images’ statistics works substantially better.

any regularization term in the CLIP cosine distance loss, we argue that initializing the optimized statistics with those of the source images is a form of regularization, favoring augmented features in a neighborhood of  $\bar{\mathbf{f}}_s$  and better preserving the semantics.

**Optimization steps.** In all our experiments, 100 iterations of optimization are performed for each batch of source features. We show in Fig. 7 the effect of the total number of iterations. We see an inflection point at around 80-100 iterations. Using few iterations is not sufficient for style alignment. Above 100, we also observe a performance drop. We refer to [3] and argue

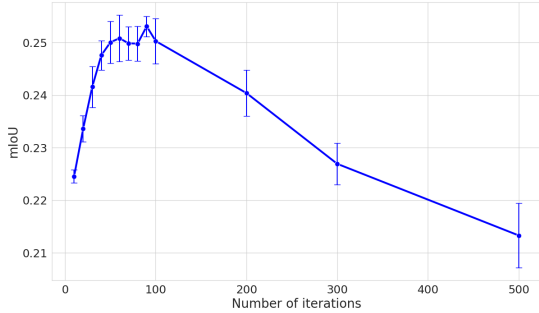


Figure 7: **Effect of the number of optimization iterations.** Performance (mIoU %) of PØDA adaptation from Cityscapes to ACDC-Night as a function of the number of statistics optimization iterations. The values are averages over 5 runs and the bars represent the standard deviation.

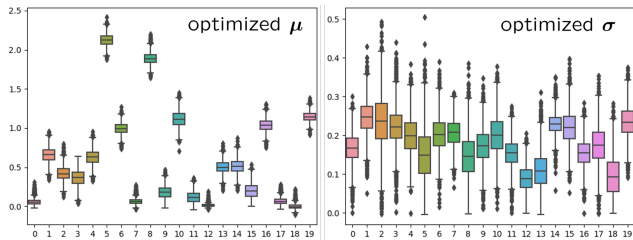


Figure 8: **Per-channel optimized statistics.** Distributions of the first 20 channels of the optimized statistics of  $\mu$  (Left) and  $\sigma$  (Right). Each boxplot shows the interquartile range (IQR) that contains 50% of the data: Its bottom and top edges delimit the first and third quartiles respectively. The horizontal line inside the box denotes the data median. The whiskers extend from the edges of the box to the furthest point within 1.5 times the IQR, in each direction. Outlier points beyond these limits are individually plotted (diamonds).

for the “over-stylization” problem in this case.

**Diversity of optimized statistics.** To verify that the global statistics — optimized for the same number of iterations with the same TrgPrompt but from different starting anchor points  $\bar{f}_s$  — are diverse, we show in Fig. 8 the boxplots of optimized parameters on the first 20 channels of  $f_{s \rightarrow t}$  (for prompt “driving at night”).

#### Training from scratch on augmented features.

In PØDA, we start with a source-only trained model (Algorithm 2, line 2) then we fine-tune it on augmented features (Algorithm 2, line 5). This is the general setting for domain adaptation. However, since our method performs domain adaptation under the assumption of label preservation, we also experimented

Method	Night	Snow	Rain	GTA5
PØDA no src pretrain	22.46	36.73	39.70	39.57
PØDA	<b>25.03</b>	<b>43.90</b>	<b>42.31</b>	<b>41.07</b>

Table 13: **Importance of source-only pre-training.** Semantic segmentation performance (mIoU %) of PØDA *vs.* its variant without source-only training, when adapting from Cityscapes to ACDC Nigt/Snow/Rain and to GTA5.

	ACDC Night	Nighttime Driving [2]	NightCity [5]
Source-only	18.31	29.61	25.63
PØDA	<b>25.03<math>\pm</math>0.48</b>	<b>33.98<math>\pm</math>0.61</b>	<b>28.90<math>\pm</math>0.61</b>

Table 14: **PØDA at night.** Segmentation performance (mIoU %) of PØDA adapted from Cityscapes to nighttime with TrgPrompt = “driving at night”, on three different night-time driving datasets.

training the model from scratch on augmented features. The results (Tab. 13) show the importance of the first, source-only training step.

**Testing PØDA on other datasets.** PØDA does not use target datasets at any point in training. Although there is no reason for the improvements observed to be specific for the datasets we test on, we show in Tab. 14 the performance of the model adapted using “driving at night” on two additional night-time driving scenes datasets:

- **Nighttime Driving [2]** test set, which consists of 50 annotated images of night driving scenes, with resolution of  $1920 \times 1080$ .
- **NightCity [5]**, which is a large dataset of 4297 night-time driving scenes collected from many cities around the world; We tested on the validation/testing set, which consists of 1299 images of resolution  $1024 \times 512$ .

## D. Class-wise performance

We report class-wise IoUs in Tab. 15.

## E. PØDA for Object Detection

Here, we share the implementation details for our object detection experiments (Sec. 5 and Tab. 10). We used the implementation of Faster R-CNN [4] from the MMDetection library.<sup>1</sup> With Cityscapes as source

<sup>1</sup><https://github.com/open-mmlab/mmdetection>

Source	Target eval.	Method	road	sidewalk	building	wall	fence	pole	traffic light	traffic sign	vegetation	terrain	sky	person	rider	car	truck	bus	train	motorcycle	bicycle	mIoU%
TrgPrompt = "driving at night"																						
ACDC	Night	source-only	70.42	18.32	<b>43.83</b>	6.11	17.08	23.52	<b>24.51</b>	19.76	39.74	6.11	0.78	21.62	8.96	23.08	2.53	0.00	3.27	8.42	9.87	18.31
		CLIPstyler	73.96	23.26	42.16	3.31	7.21	<b>35.49</b>	23.34	19.01	<b>45.41</b>	8.81	<b>27.87</b>	21.06	8.48	38.17	1.84	0.00	11.54	<b>10.38</b>	4.89	21.38±0.36
		PØDA	<b>77.54</b>	<b>26.90</b>	42.71	<b>13.51</b>	<b>21.36</b>	33.52	23.70	<b>21.73</b>	39.91	<b>9.51</b>	19.40	<b>28.80</b>	<b>11.85</b>	<b>50.89</b>	<b>10.14</b>	0.00	<b>20.76</b>	8.76	<b>14.50</b>	<b>25.03</b> ±0.48
TrgPrompt = "driving in snow"																						
ACDC	snow	source-only	70.47	23.50	63.80	17.96	<b>27.36</b>	<b>38.52</b>	<b>56.26</b>	<b>45.00</b>	<b>83.00</b>	10.75	83.65	47.73	0.72	61.42	21.87	5.90	21.58	35.83	31.01	39.28
		CLIPstyler	74.29	31.25	69.17	15.21	25.21	36.83	44.79	42.56	76.87	11.07	91.48	<b>53.23</b>	0.13	67.66	23.88	<b>9.14</b>	36.48	<b>42.67</b>	28.76	41.09±0.17
		PØDA	<b>75.40</b>	<b>34.61</b>	<b>75.22</b>	<b>26.77</b>	27.34	35.20	52.68	44.37	82.01	<b>14.16</b>	<b>93.72</b>	50.51	<b>0.99</b>	<b>69.11</b>	<b>26.64</b>	2.72	<b>46.98</b>	42.64	<b>33.09</b>	<b>43.90</b> ±0.53
TrgPrompt = "driving under rain"																						
ACDC	rain	source-only	74.10	31.98	63.07	<b>15.08</b>	<b>23.92</b>	<b>41.31</b>	<b>50.12</b>	<b>44.43</b>	79.93	22.07	87.45	47.99	4.39	68.92	10.35	18.52	13.64	7.03	21.58	38.20
		CLIPstyler	73.71	36.09	68.91	3.77	16.99	36.94	39.75	36.44	78.21	20.64	91.79	40.34	<b>9.65</b>	<b>74.54</b>	13.16	<b>20.33</b>	12.73	14.06	18.26	37.17±0.10
		PØDA	<b>76.60</b>	<b>38.52</b>	<b>78.01</b>	15.02	22.53	40.33	45.39	41.40	<b>86.85</b>	<b>37.97</b>	<b>96.46</b>	<b>50.39</b>	6.35	74.19	<b>19.19</b>	7.98	<b>22.06</b>	<b>21.04</b>	<b>23.65</b>	<b>43.31</b> ±0.55
TrgPrompt = "driving in a game"																						
GTA5	CS	source-only	68.72	22.65	78.79	36.81	<b>17.31</b>	39.66	<b>39.33</b>	14.84	<b>72.61</b>	22.53	87.31	57.50	26.14	74.29	<b>44.57</b>	<b>20.45</b>	0.00	18.30	10.35	39.59
		CLIPstyler	73.06	<b>29.89</b>	77.86	25.50	11.69	<b>39.72</b>	35.88	<b>24.04</b>	67.38	12.75	88.77	46.58	33.38	72.03	42.79	11.12	0.00	28.84	<b>14.61</b>	38.73±0.16
		PØDA	<b>73.93</b>	22.69	<b>78.82</b>	<b>37.52</b>	14.17	36.97	33.14	17.34	72.44	<b>26.22</b>	<b>88.85</b>	<b>62.69</b>	<b>37.04</b>	<b>74.33</b>	43.03	11.91	0.00	<b>35.33</b>	13.91	<b>41.07</b> ±0.48
TrgPrompt = "driving"																						
GTA5	CS	source-only	58.97	20.92	72.84	16.53	<b>24.58</b>	31.37	34.77	23.62	82.12	17.04	66.28	<b>63.46</b>	14.72	<b>81.27</b>	<b>20.83</b>	17.19	4.68	<b>20.57</b>	19.56	36.38
		CLIPstyler	66.70	23.63	64.12	5.08	3.66	20.67	19.31	18.10	81.68	12.36	<b>81.04</b>	54.64	0.52	73.47	20.65	<b>22.30</b>	4.03	15.79	10.73	31.50±0.21
		PØDA	<b>84.34</b>	<b>36.73</b>	<b>79.43</b>	<b>18.33</b>	16.54	<b>36.93</b>	<b>38.45</b>	<b>33.81</b>	<b>82.44</b>	<b>19.14</b>	75.90	62.65	<b>16.47</b>	75.48	15.68	19.57	<b>11.28</b>	16.53	<b>21.76</b>	<b>40.08</b> ±0.52

Table 15: **Zero-shot domain adaptation in semantic segmentation.** Performance (mIoU%) of PØDA compared against CLIPstyler [3] and source-only baseline. Results are grouped by source domain and target domain with associated TrgPrompt. CS stands for Cityscapes [1]. This table provides details of the main results in Tab. 2.

dataset, we trained all models for 8 epochs using the SGD optimizer with 0.9 momentum and  $1e-4$  weight decay. The initial learning rate  $lr$  is set as  $1e-2$  and is dropped by a factor of 10 after the 7th epoch; the same  $lr$  scheme is used in source only and PØDA trainings. With Day-Sunny split of the DWD dataset as source, models are trained for 20 epochs using a similar SGD optimizer. When training on source, the learning rate starts at  $1e-3$  and drops at the 9th epoch to  $1e-4$ ; in PØDA training, the learning rate is ten times less.

## References

- [1] 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 *CVPR*, 2016. 3
- [2] Dengxin Dai and Luc Van Gool. Dark model adaptation: Semantic image segmentation from daytime to nighttime. In *ITSC*, 2018. 2
- [3] Gihyun Kwon and Jong Chul Ye. Clipstyler: Image style transfer with a single text condition. In *CVPR*, 2022. 1, 3
- [4] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *NeurIPS*, 2015. 2
- [5] Xin Tan, Ke Xu, Ying Cao, Yiheng Zhang, Lizhuang Ma, and Rynson WH Lau. Night-time scene parsing with a large real dataset. *IEEE T-IP*, 2021. 2
- [6] Xinyi Wu, Zhenyao Wu, Yuhang Lu, Lili Ju, and Song Wang. Style mixing and patchwise prototypical matching for one-shot unsupervised domain adaptive semantic segmentation. In *AAAI*, 2022. 1