# Stepping Towards
# Unsupervised Scene Understanding

**Raoul de Charette**

rdecharette.github.io | raoul.de-charette@inria.fr
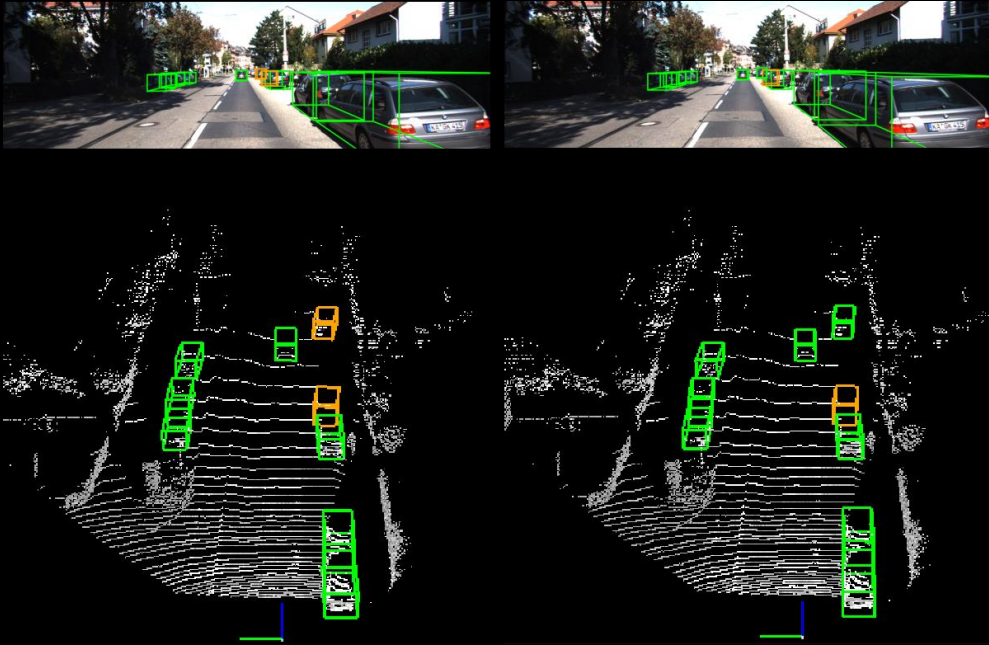
**Scene understanding**

# AD - Level 4



Waymo

Cruise

Baidu

(a quick partial/biased overview of CV pillars for AD)
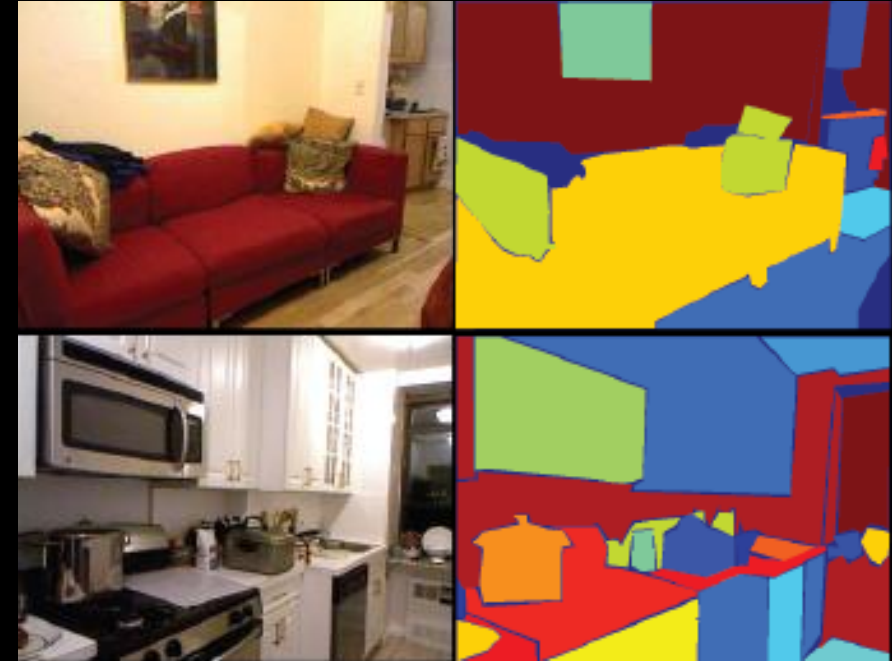
## **Object detection**



(Chen et al., ECCV 2020)

## **Semantic Segmentation**
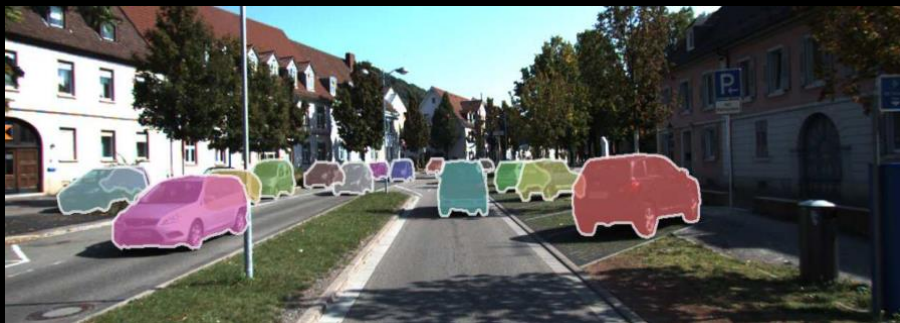


(Silberman et al., ECCV 2012)

1 color = 1 semantic class

# **Semantics** | Geometry | Motion

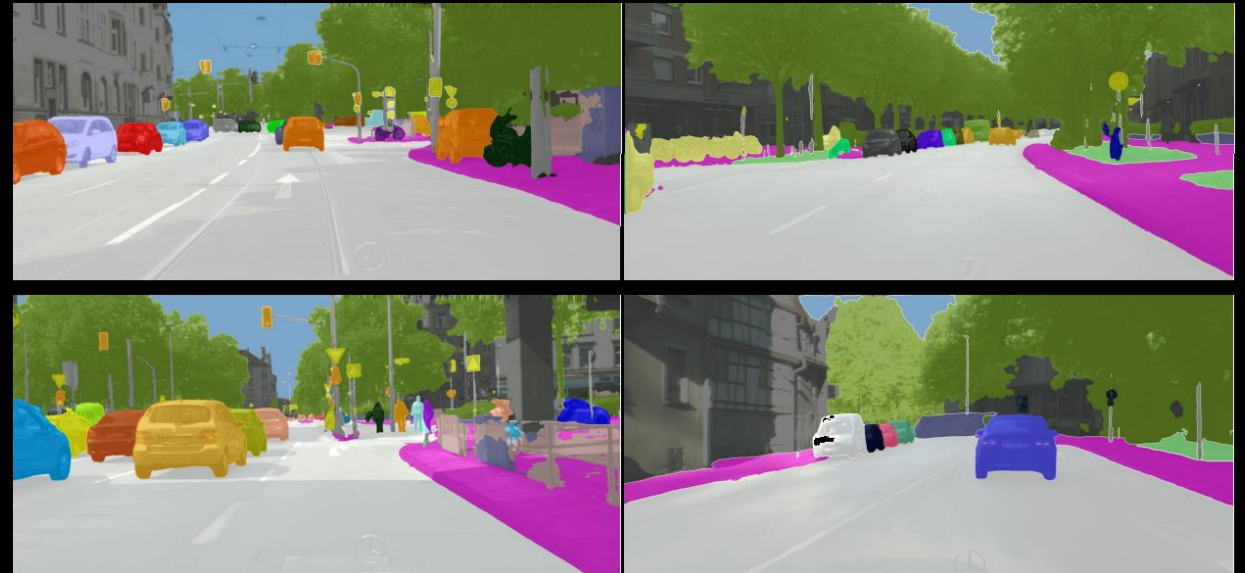### **Amodal Segmentation (2016)**



(Li and Malik, ECCV 16)



KINS (Qi et al., CVPR 19)

Segmentation of visible and invisible pixels

"Things and Stuff"

### **Panoptic segmentation (2019)**



(Kirillov et al., CVPR 2019)

1 color = 1 instance

Segmentation of class and instance at once
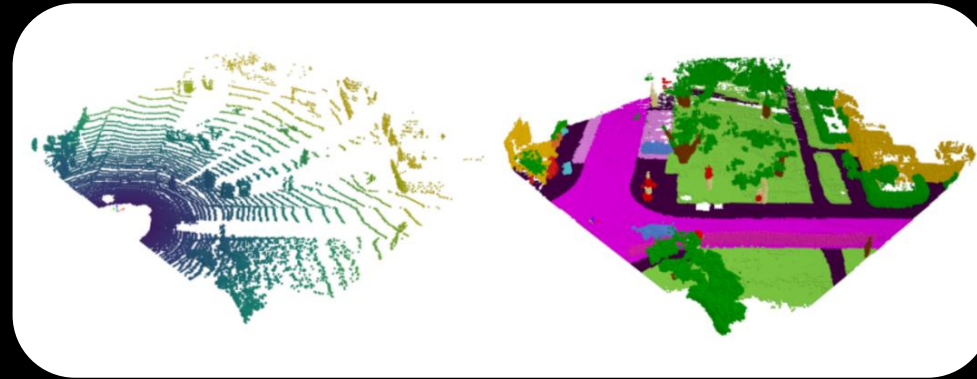
**Depth prediction**



(Bhat et al., CVPR 21)

**Reconstruction**



(Cao and de Charette, ICCV 23)

## Semantic Scene Completion



(Roldao et al., IJCV 21)



(Cao and de Charette, CVPR 22)

**Object tracking**



CAMO-MOT. (Wang et al., 22)

**Pixel tracking**



(Wang et al., ICCV 23)

**Forecasting**



(Liu et al., CVPR 21)

And many more..

# Data proficiency ($10^7$)

KITTI, 2012 / Sem.KITTI, 2019

Waymo Open Dataset, 2020

nuScenes, 2019

# Biases

(Torralba and Efros, CVPR 11)

# Cost

**For 1 annotator 8h/365d**

train/val – fine annotation – 3475 images

~8y

~1y

# World is long-tail

Amount of data

# Data proficiency ($10^7$)



KITTI, 2012 / Sem.KITTI, 2019

Waymo Open Dataset, 2020

nuScenes, 2019

# Biases



(Torralba and Efros, CVPR 11)

# Cost

For 1 annotator 8h/365d



train/val – fine annotation – 3475 images

~8y

~1y

# World is long-tail



Amount of data

**Supervised learning is doomed to Out Of Distribution**

**Multi-Task Learning**



2D

3D

**Cross-Modal Learning**



"driving at night"

VLM

**Prompt-driven Learning**

# Knowledge Distillation



**Multi-Task Learning**

**Cross-Modal Learning**

**Prompt-driven Learning**

# Knowledge Distillation



**Multi-Task Learning**
(Lopes et al., WACV 2023)

**Cross-Modal Learning**
(Jaritz et al., TPAMI 2022)

**Prompt-driven Learning**
(Fahes et al., ICCV 2023)

# Multi-task Learning

$$\{T_1, \cdots, T_n\}$$

Semantics conveys geometry cues, and vice versa.

TRANSFER LEARNING

**Taskonomy** [Zamir *et al.*, CVPR'18]

MULTI-TASK LEARNING

**Which tasks to learn together in MTL** [Standley *et al.*, ICML'20]

PAD-Net [Xu *et al.,* CVPR'18]

CTRL-UDA [Saha *et al.,* CVPR'21]

ATRC [Bruggemann *et al.,* ICCV'21]

GUDA [Guizilini *et al.,* ICCV'21]

3-Ways [Hoyer *et al.,* CVPR'21]

Survey: [Vandenhende et al., TPAMI 2020]

# DenseMTL:
# Multitask Learning for UDA

3 set of tasks: `S-D', `S-D-N', `S-D-N-E'

# multi-Task Exchange Block (mTEB)

# Guide feature inclusion from task *j* to task *i* and vice-versa

# Correlation-guided attention $\quad \mathcal{C}_{j\to i}$



The intuition: injecting features from *j* that will contribute to better solving task *i*.

$$Q = \text{transform}(f_i)$$ Features from task *i*

$$\left.\begin{array}{l} K^{\mathrm{T}} = \text{transform}(f_j) \\ V = \text{transform}(f_j) \end{array}\right\}$$ Features from task *j*

$$\mathcal{C}_{j\to i} = \text{softmax}\left(\frac{K^T \times Q}{\sqrt{d}}\right)$$

We account more for features of *j* which are highly spatially-correlated with features from *i*

$$\mathbf{xtask}_{j\to i} = V \times \mathcal{C}_{j\to i}$$

# Self-attention

The intuition: discover private features from *j* that help task *i*



$F_f(f_j)$     Convolution block

$F_m(f_j)$     Convolution block

We let gradient flow optimize *F_f*, *F_m* to discover relevant features in *j* for task *i*

$$\mathbf{self}_{j \to i} = F_f(f_j) \odot \sigma(F_m(f_j))$$

# Fusion

# multi-task Exchange Block (mTEB)



**mTEB can be inserted at any scale**

# multi-task Exchange Block (mTEB)



**mTEB can be inserted at any scale**

# multi-task Exchange Block (mTEB)



mTEB can be inserted at any scale

$$\mathcal{L}_{\text{tasks}} = \frac{1}{|S|} \sum_{s \in S} \sum_{t \in T} \omega_t \mathcal{L}_t^s + \sum_{t \in T} \omega_t \mathcal{L}_t^{\text{final}}$$

# The challenge of MTL metrics

- Metrics and scale differ per task: depth (RMSE), semantics (mIoU), etc.
- MTL should favor **all** metrics

0 for higher is better, 1 lower is better

$$\Delta_{\mathrm{T}}(\mathbf{f}) = 1/n \sum_{i \in T} (-1)^{g_i} (m_i - b_i)/b_i$$

MTL
performance

STL
performance

**S**egmentation

**D**epth

**N**ormal

**E**dge

# Fully supervised



| | | '*S-D*' | | | | '*S-D-N*' | | | |
|---|---|---|---|---|---|---|---|---|---|
| Methods | Semseg ↑ | Depth ↓ | Delta ↑ | Semseg ↑ | Depth ↓ | Delta ↑ | Normals ↓ | Delta ↑ |
| | mIoU % | RMSE m | $\Delta_{SD}$ % | mIoU % | RMSE m | $\Delta_{SD}$ % | mErr. ° | $\Delta_{SDN}$ % |
| STL [8] | 38.70 ±0.10 | 0.635 ±0.013 | ⤵ | *idem* | *idem* | ⤵ | 36.90 ±0.26 | ⤵ |
| MTL [8] | 39.44 ±0.34 | 0.638 ±0.004 | +1.63 ±0.37 | 39.90 ±0.41 | 0.642 ±0.003 | +1.89 ±0.67 | 36.07 ±0.09 | +1.76 ±0.53 |
| PAD-Net [135] | 35.30 ±0.84 | 0.659 ±0.004 | -5.36 ±0.83 | 36.14 ±0.30 | 0.660 ±0.006 | -4.32 ±0.68 | 36.72 ±0.08 | -2.97 ±0.43 |
| 3-ways$_{PAD-Net}$ [23] | **39.47** ±0.16 | 0.622 ±0.001 | +2.90 ±0.23 | **40.28** ±0.30 | 0.619 ±0.004 | +4.16 ±0.50 | 35.35 ±0.09 | +3.93 ±0.27 |
| Ours | 38.93 ±0.35 | **0.604** ±0.004 | +3.54 ±0.21 | **40.28** ±0.41 | **0.598** ±0.002 | +5.80 ±0.65 | **33.72** ±0.14 | **+6.49** ±0.50 |

NYUDv2

[8] MTL survey [Vandenhende *et al.*, TPAMI'2021]     [23] 3-Ways [Hoyer *et al.*, CVPR'21]     [135] PAD-Net [Xu *et al.*, CVPR'18]

# Fully supervised



| | ‘S-D’ | | | ‘S-D-N’ | | | | |
|---|---|---|---|---|---|---|---|---|
| Methods | Semseg ↑ | Depth ↓ | Delta ↑ | Semseg ↑ | Depth ↓ | Delta ↑ | Normals ↓ | Delta ↑ |
| | mIoU % | RMSE m | $\Delta_{SD}$ % | mIoU % | RMSE m | $\Delta_{SD}$ % | mErr. ° | $\Delta_{SDN}$ % |
| STL [8] | 38.70 ±0.10 | 0.635 ±0.013 | ⬏ | idem | idem | ⬏ | 36.90 ±0.26 | ⬏ |
| MTL [8] | 39.44 ±0.34 | 0.638 ±0.004 | +1.63 ±0.37 | 39.90 ±0.41 | 0.642 ±0.003 | +1.89 ±0.67 | 36.07 ±0.09 | +1.76 ±0.53 |
| PAD-Net [135] | 35.30 ±0.84 | 0.659 ±0.004 | -5.36 ±0.83 | 36.14 ±0.30 | 0.660 ±0.006 | -4.32 ±0.68 | 36.72 ±0.08 | -2.97 ±0.43 |
| 3-ways$_{PAD-Net}$ [23] | 39.47 ±0.16 | 0.622 ±0.001 | +2.90 ±0.23 | 40.28 ±0.30 | 0.619 ±0.004 | +4.16 ±0.50 | 35.35 ±0.09 | +3.93 ±0.27 |
| Ours | 38.93 ±0.35 | 0.604 ±0.004 | +3.54 ±0.21 | 40.28 ±0.41 | 0.598 ±0.002 | +5.80 ±0.65 | 33.72 ±0.14 | +6.49 ±0.50 |

NYUDv2

[8] MTL survey [Vandenhende *et al.*, TPAMI'2021]     [23] 3-Ways [Hoyer *et al.*, CVPR'21]     [135] PAD-Net [Xu *et al.*, CVPR'18]

# Fully supervised



| | 'S-D' | | | 'S-D-N' | | | | | 'S-D-N-E' | | | | | |
| Methods | Semseg ↑ mIoU % | Depth ↓ RMSE m | Delta ↑ $\Delta_{SD}$ % | Semseg ↑ mIoU % | Depth ↓ RMSE m | Delta ↑ $\Delta_{SD}$ % | Normals ↓ mErr. ° | Delta ↑ $\Delta_{SDN}$ % | Semseg ↑ mIoU % | Depth ↓ RMSE m | Normals ↓ mErr. ° | Delta ↑ $\Delta_{SDN}$ % | Edges ↑ F1 % | Delta ↑ $\Delta_{SDNE}$ % |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| STL [8] | 38.70 ±0.10 | 0.635 ±0.013 | ↘ | idem | idem | ↘ | 36.90 ±0.26 | ↘ | idem | idem | idem | ↘ | 54.90 ±0.00 | ↘ |
| MTL [8] | 39.44 ±0.34 | 0.638 ±0.004 | +1.63 ±0.37 | 39.90 ±0.41 | 0.642 ±0.003 | +1.89 ±0.67 | 36.07 ±0.09 | +1.76 ±0.53 | 39.70 ±0.35 | 0.636 ±0.001 | 36.10 ±0.12 | +1.88 ±0.33 | 55.11 ±0.15 | +1.50 ±0.20 |
| PAD-Net [135] | 35.30 ±0.84 | 0.659 ±0.004 | -5.36 ±0.83 | 36.14 ±0.30 | 0.660 ±0.006 | -4.32 ±0.68 | 36.72 ±0.08 | -2.97 ±0.43 | 36.19 ±0.24 | 0.662 ±0.005 | 36.58 ±0.06 | -2.92 ±0.37 | 54.79 ±0.07 | -2.24 ±0.26 |
| 3-ways PAD-Net [23] | 39.47 ±0.16 | 0.622 ±0.001 | +2.90 ±0.23 | 40.28 ±0.30 | 0.619 ±0.004 | +4.16 ±0.50 | 35.35 ±0.09 | +3.93 ±0.27 | 40.16 ±0.28 | 0.614 ±0.010 | 35.25 ±0.09 | +4.14 ±0.65 | 59.66 ±0.16 | +5.27 ±0.49 |
| Ours | 38.93 ±0.35 | **0.604** ±0.004 | +3.54 ±0.21 | **40.28** ±0.41 | **0.598** ±0.002 | +5.80 ±0.65 | **33.72** ±0.14 | +6.49 ±0.50 | **40.84** ±0.37 | **0.593** ±0.004 | **33.38** ±0.19 | +7.52 ±0.27 | **61.12** ±0.20 | **+8.47** ±0.12 |

NYUDv2

[8] MTL survey [Vandenhende *et al.*, TPAMI'2021]     [23] 3-Ways [Hoyer *et al.*, CVPR'21]     [135] PAD-Net [Xu *et al.*, CVPR'18]

# Fully supervised



| Methods | 'S-D' | | | 'S-D-N' | | | | | 'S-D-N-E' | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Semseg ↑ | Depth ↓ | Delta ↑ | Semseg ↑ | Depth ↓ | Delta ↑ | Normals ↓ | Delta ↑ | Semseg ↑ | Depth ↓ | Normals ↓ | Delta ↑ | Edges ↑ | Delta ↑ |
| | mIoU % | RMSE m | $\Delta_{SD}$ % | mIoU % | RMSE m | $\Delta_{SD}$ % | mErr. ° | $\Delta_{SDN}$ % | mIoU % | RMSE m | mErr. ° | $\Delta_{SDN}$ % | F1 % | $\Delta_{SDNE}$ % |
| STL [8] | 38.70 ±0.10 | 0.635 ±0.013 | ↘ | idem | idem | ↘ | 36.90 ±0.26 | ↘ | idem | idem | idem | ↘ | 54.90 ±0.00 | ↘ |
| MTL [8] | 39.44 ±0.34 | 0.638 ±0.004 | +1.63 ±0.37 | 39.90 ±0.41 | 0.642 ±0.003 | +1.89 ±0.67 | 36.07 ±0.09 | +1.76 ±0.53 | 39.70 ±0.35 | 0.636 ±0.001 | 36.10 ±0.12 | +1.88 ±0.33 | 55.11 ±0.15 | +1.50 ±0.20 |
| PAD-Net [135] | 35.30 ±0.84 | 0.659 ±0.004 | -5.36 ±0.83 | 36.14 ±0.30 | 0.660 ±0.006 | -4.32 ±0.68 | 36.72 ±0.08 | -2.97 ±0.43 | 36.19 ±0.24 | 0.662 ±0.005 | 36.58 ±0.06 | -2.92 ±0.37 | 54.79 ±0.07 | -2.24 ±0.26 |
| 3-ways$_{PAD-Net}$ [23] | **39.47** ±0.16 | 0.622 ±0.001 | +2.90 ±0.23 | **40.28** ±0.30 | 0.619 ±0.004 | +4.16 ±0.50 | 35.35 ±0.09 | +3.93 ±0.27 | 40.16 ±0.28 | 0.614 ±0.010 | 35.25 ±0.09 | +4.14 ±0.65 | 59.66 ±0.16 | +5.27 ±0.49 |
| Ours | 38.93 ±0.35 | **0.604** ±0.004 | +3.54 ±0.21 | **40.28** ±0.41 | **0.598** ±0.002 | +5.80 ±0.65 | **33.72** ±0.14 | +6.49 ±0.50 | **40.84** ±0.37 | **0.593** ±0.004 | **33.38** ±0.19 | +7.52 ±0.27 | **61.12** ±0.20 | **+8.47** ±0.12 |

NYUDv2

[8] MTL survey [Vandenhende *et al.*, TPAMI'2021]     [23] 3-Ways [Hoyer *et al.*, CVPR'21]     [135] PAD-Net [Xu *et al.*, CVPR'18]

# Fully supervised



| | | 'S-D' | | | 'S-D-N' | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | Methods | Semseg ↑ mIoU % | Depth ↓ RMSE m | Delta ↑ $\Delta_{SD}$ % | Semseg ↑ mIoU % | Depth ↓ RMSE m | Delta ↑ $\Delta_{SD}$ % | Normals ↓ mErr. ° | Delta ↑ $\Delta_{SDN}$ % |
| Synthia | STL [8] | 67.43 +0.15 | 5.379 +0.055 | ↘ | idem | idem | ↘ | 19.61 +0.12 | ↘ |
| | MTL [8] | 69.83 +0.25 | 5.166 +0.063 | +03.76 +0.77 | 71.27 +0.21 | 5.108 +0.076 | +05.37 +0.83 | 18.51 +0.10 | +05.45 +0.72 |
| | PAD-Net [135] | 70.87 +0.15 | 4.917 +0.014 | +06.85 +0.24 | 72.27 +0.25 | 4.949 +0.072 | +07.58 +0.56 | 19.28 +0.09 | +05.62 +0.43 |
| | 3-ways PAD-Net [23] | 77.50 +0.17 | 4.289 +0.028 | +17.60 +0.13 | 79.93 +0.5 | 4.218 +0.082 | +20.06 +0.92 | 15.54 +0.14 | +20.29 +0.84 |
| | Ours | **80.53** +0.43 | **4.161** +0.022 | **+21.04** +0.52 | **82.99** +0.38 | **4.056** +0.076 | **+23.83** +0.98 | **14.30** +0.15 | **+24.92** +0.87 |
| VKITTI2 | STL [8] | 84.53 +0.06 | 5.720 +0.027 | ↘ | idem | idem | ↘ | 23.14 +0.68 | ↘ |
| | MTL [8] | 87.73 +0.12 | 5.720 +0.029 | +01.89 +0.21 | 87.83 +0.21 | 5.714 +0.033 | +02.00 +0.27 | 22.30 +0.68 | +02.54 +0.80 |
| | PAD-Net [135] | 88.43 +0.12 | 5.571 +0.058 | +03.63 +0.45 | 88.67 +0.15 | 5.543 +0.043 | +04.09 +0.29 | 22.16 +0.70 | +04.09 +0.83 |
| | 3-ways PAD-Net [23] | 96.13 +0.15 | 4.013 +0.051 | +21.78 +0.54 | 96.87 +0.06 | 3.756 +0.013 | +24.46 +0.14 | 15.54 +0.56 | +27.25 +0.90 |
| | Ours | **97.00** +0.10 | **3.423** +0.025 | **+27.47** +0.16 | **97.53** +0.06 | **3.089** +0.006 | **+30.70** +0.05 | **14.44** +0.52 | **+33.00** +0.73 |
| Cityscapes | STL [8] | 67.93 +0.06 | 6.622 +0.020 | ↘ | idem | idem | ↘ | 44.10 +0.01 | ↘ |
| | MTL [8] | 70.43 +0.12 | 6.797 +0.520 | +00.52 +0.32 | 70.93 +0.15 | 6.736 +0.023 | +01.34 +0.28 | 43.60 +0.01 | +01.30 +0.18 |
| | PAD-Net [135] | 70.23 +0.25 | 6.777 +0.010 | +00.52 +0.27 | 70.67 +0.06 | 6.755 +0.018 | +01.00 +0.17 | 43.52 +0.00 | +01.12 +0.11 |
| | 3-ways PAD-Net [23] | **75.00** +0.10 | **6.528** +0.063 | **+05.91** +0.44 | 75.50 +0.10 | 6.491 +0.081 | +06.56 +0.61 | 41.84 +0.05 | +06.09 +0.37 |
| | Ours | 74.95 +0.10 | 6.649 +0.003 | +04.96 +0.08 | **76.08** +0.14 | **6.407** +0.013 | **+07.61** +0.04 | **40.05** +0.33 | **+08.15** +0.22 |

[8] MTL survey [Vandenhende *et al.*, TPAMI'2021]    [23] 3-Ways [Hoyer *et al.*, CVPR'21]    [135] PAD-Net [Xu *et al.*, CVPR'18]

| mTEB | | 'S-D' | | | | 'S-D-N' | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| Scales | Param. ↓ | Semseg ↑ | Depth ↓ | Delta ↑ | Param. ↓ | Semseg ↑ | Depth ↓ | Normals ↓ | Delta ↑ |
| 4 3 2 1 | #M added | mIoU % | RMSE m | $\Delta_{SD}$ % | #M added | mIoU % | RMSE m | mErr. ° | $\Delta_{SDN}$ % |
| | 0.00 | 96.88 ±0.30 | 3.604 ±0.020 | 25.81 ±0.33 | 0.00 | 97.38 ±0.02 | 3.491 ±0.041 | 14.50 ±0.57 | 30.51 ±0.98 |
| ✓ | 3.09 | 97.32 ±0.06 | 3.556 ±0.029 | 26.50 ±0.29 | 2.32 | 97.43 ±0.04 | 3.559 ±0.024 | 14.51 ±0.50 | 30.12 ±0.79 |
| ✓ | 3.09 | 97.24 ±0.03 | 3.476 ±0.018 | 27.15 ±0.16 | 2.32 | 97.49 ±0.08 | 3.353 ±0.025 | 14.45 ±0.48 | 31.43 ±0.76 |
| ✓ | 0.77 | 97.07 ±0.06 | 3.468 ±0.016 | 27.11 ±0.12 | 9.26 | 97.47 ±0.06 | 3.244 ±0.035 | 14.57 ±0.51 | 31.89 ±0.92 |
| ✓ | 0.77 | 97.00 ±0.10 | 3.423 ±0.025 | 27.47 ±0.16 | 9.26 | 97.53 ±0.06 | 3.089 ±0.006 | 14.44 ±0.52 | 33.00 ±0.73 |
| ✓ ✓ | 3.86 | 97.09 ±0.03 | 3.369 ±0.022 | 27.99 ±0.18 | 11.58 | 97.53 ±0.02 | 3.080 ±0.025 | 14.47 ±0.57 | 33.02 ±0.90 |
| ✓ ✓ ✓ | 4.63 | 97.01 ±0.02 | 3.377 ±0.008 | 27.88 ±0.06 | 13.89 | 97.39 ±0.02 | 3.136 ±0.046 | 14.81 ±0.77 | 32.13 ±1.39 |
| ✓ ✓ ✓ ✓ | 7.72 | 97.05 ±0.03 | 3.369 ±0.010 | 27.97 ±0.08 | 23.15 | 96.82 ±0.23 | 3.307 ±0.066 | 15.39 ±0.65 | 30.08 ±1.39 |

VKITTI2

**Where should tasks talk ?**

| weights | | Semseg ↑ | Depth ↓ | Delta ↑ |
|---|---|---|---|---|
| $\omega_S$ | $\omega_D$ | mIoU % | RMSE m | $\Delta_{SD}$ % |
| 1 | 1 | 83.83 ±0.15 | 5.713 ±0.060 | -0.35 ±0.47 |
| 1 | 10 | 79.87 ±0.21 | 5.708 ±0.036 | -2.66 ±0.40 |
| 10 | 1 | 86.20 ±0.71 | **5.693** ±0.055 | +1.30 ±0.22 |
| 50 | 1 | 87.73 ±0.12 | 5.720 ±0.029 | **+1.89** ±0.21 |
| 100 | 1 | 88.00 ±0.20 | 5.754 ±0.030 | +1.75 ±0.17 |
| 100 | 10 | 86.13 ±0.32 | **5.693** ±0.039 | +1.18 ±0.45 |
| 200 | 1 | 88.13 ±0.12 | 5.790 ±0.055 | +1.52 ±0.45 |
| 500 | 1 | **88.17** ±0.15 | 5.847 ±0.043 | +1.04 ±0.30 |

(a) 'S-D' gridsearch

| weights | | | Semseg ↑ | Depth ↓ | Normals ↓ | Delta ↑ |
|---|---|---|---|---|---|---|
| $\omega_S$ | $\omega_D$ | $\omega_N$ | mIoU % | RMSE m | mErr. ° | $\Delta_{SDN}$ % |
| 1 | 1 | 1 | 83.50 ±0.20 | 5.707 ±0.058 | 23.03 ±0.70 | -0.17 ±0.64 |
| 10 | 1 | 1 | 86.53 ±0.21 | 5.694 ±0.032 | 22.98 ±0.68 | +1.17 ±0.93 |
| 10 | 1 | 10 | 86.63 ±0.21 | **5.675** ±0.050 | 22.61 ±0.70 | +1.85 ±0.80 |
| 50 | 1 | 1 | 87.73 ±0.21 | 5.706 ±0.051 | 22.90 ±0.71 | +1.69 ±0.81 |
| 50 | 1 | 10 | 87.77 ±0.15 | 5.714 ±0.065 | 22.56 ±0.69 | +2.15 ±0.86 |
| 50 | 1 | 50 | 87.73 ±0.21 | 5.701 ±0.062 | 22.37 ±0.70 | +2.49 ±0.76 |
| 100 | 1 | 1 | 88.03 ±0.15 | 5.746 ±0.030 | 22.95 ±0.69 | +1.49 ±0.92 |
| 100 | 1 | 10 | 87.97 ±0.15 | 5.714 ±0.048 | 22.59 ±0.69 | +2.19 ±0.79 |
| 100 | 1 | 50 | 88.00 ±0.20 | 5.717 ±0.048 | 22.40 ±0.71 | +2.45 ±0.99 |
| 100 | 1 | 100 | 88.07 ±0.15 | 5.696 ±0.038 | **22.29** ±0.70 | **+2.75** ±1.04 |
| 150 | 1 | 10 | 88.10 ±0.20 | 5.752 ±0.059 | 22.59 ±0.70 | +2.01 ±0.86 |
| 150 | 1 | 50 | 88.10 ±0.20 | 5.738 ±0.039 | 22.41 ±0.70 | +2.35 ±0.99 |
| 150 | 1 | 100 | **88.13** ±0.15 | 5.732 ±0.037 | 22.31 ±0.71 | +2.54 ±0.94 |

(b) 'S-D-N' gridsearch

34

# MTL for Segmentation
{S,D} on Cityscapes

| Methods | road | swalk | build | wall | fence | pole | light | sign | veg | sky | person | rider | car | bus | mbike | bike | mIoU % |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 3-ways$_{PAD-Net}$ [23] | 97.21 | 79.38 | 90.50 | **47.68** | 49.68 | 51.17 | 49.41 | 64.65 | 91.40 | 93.85 | 72.41 | 46.92 | 92.66 | 80.17 | 42.43 | 66.39 | 69.74 |
| 3-ways$_{mTEB}$ | **97.62** | **82.29** | **92.44** | 46.52 | **54.76** | **59.82** | **60.94** | **73.13** | **92.22** | **94.55** | **76.40** | **58.49** | **94.26** | **85.14** | **49.41** | **71.70** | **74.36** |

MDE for depth

**Self-supervised depth can help segmentation**

# Cross-Modal Learning

# Dataset Bias or Domain Discrepancy

(slide Tuan-Hung Vu)



Cityscapes (CVPR 16)

Mapillary Vistas (ICCV 17)

...

Common train/test

Open-world testbed

# Out of Distribution



Physics-Based Rendering for (..) Rain. (Halder et al., ICCV 19)

# Out of Distribution

# Out of Distribution



Original($\mathcal{I}$)     Upernet [23]

$\mathcal{I} - car$     Upernet [23]

(Shetty et al. CVPR 2019)

# Out of Distribution



(a) captured by a 64-beam LiDAR   (b) captured by a 32-beam LiDAR

Complete & Label. (Yi et al. CVPR 2021)

# Domain Gap

**Train**

Sensor setup
Location (city, country)
Scenes (urban, countryside)
Conditions (weather, lighting)
Ethnicity
etc.

**Test**

# Domain Gap

Sensor setup
Location (city, country)
Scenes (urban, countryside)
Conditions (weather, lighting)
Ethnicity
etc.

Train

Test

Buolamwini and Gebru. FAccT 2018

# Domain Gap

Train

Test

Sensor setup
Location (city, country)
Scenes (urban, countryside)
Conditions (weather, lighting)
Ethnicity
etc.



Buolamwini and Gebru. FAccT 2018

# Domain Gap

Sensor setup
Location (city, country)
Scenes (urban, countryside)
Conditions (weather, lighting)
Ethnicity
etc.

Train

Test



Buolamwini and Gebru. FAccT 2018

| Source\Target | KITTI | Argoverse | nuScenes | Lyft | Waymo |
|---|---|---|---|---|---|
| KITTI | **88.0 / 82.5** | 55.8 / 27.7 | 47.4 / 13.3 | 81.7 / 51.8 | 45.2 / 11.9 |
| Argoverse | 69.5 / 33.9 | **79.2 / 57.8** | 52.5 / 21.8 | 86.9 / 67.4 | 83.8 / 40.2 |
| nuScenes | 49.7 / 13.4 | 73.2 / 21.8 | **73.4 / 38.1** | 89.0 / 38.2 | 78.8 / 36.7 |
| Lyft | 74.3 / 39.4 | 77.1 / 45.8 | 63.5 / 23.9 | **90.2 / 87.3** | 87.0 / 64.7 |
| Waymo | 51.9 / 13.1 | 76.4 / 42.6 | 55.5 / 21.6 | 87.9 / 74.5 | **90.1 / 85.3** |

Wang et al. CVPR 2020

# Domain Gap

Sensor setup
Location (city, country)
Scenes (urban, countryside)
Conditions (weather, lighting)
Ethnicity
etc.

**Train**

**Test**

Buolamwini and Gebru. FAccT 2018

| Source\Target | KITTI | Argoverse | nuScenes | Lyft | Waymo |
|---|---|---|---|---|---|
| | | -32 | -41 | -7 | -3 |
| KITTI | **88.0** / **82.5** | 55.8 / 27.7 | 47.4 / 13.3 | 81.7 / 51.8 | 45.2 / 11.9 |
| Argoverse | 69.5 / 33.9 | **79.2** / **57.8** | 52.5 / 21.8 | 86.9 / 67.4 | 83.8 / 40.2 |
| nuScenes | 49.7 / 13.4 | 73.2 / 21.8 | **73.4** / **38.1** | 89.0 / 38.2 | 78.8 / 36.7 |
| Lyft | 74.3 / 39.4 | 77.1 / 45.8 | 63.5 / 23.9 | **90.2** / **87.3** | 87.0 / 64.7 |
| Waymo | 51.9 / 13.1 | 76.4 / 42.6 | 55.5 / 21.6 | 87.9 / 74.5 | **90.1** / **85.3** |

Wang et al. CVPR 2020

**Domain Adaptation**
Source: He et al. 2022



**Domain Generalization**
Source: Frikha et al. 2022

...

VLM, DINOs, SAM, etc.

Foundation Model

Model

**Knowledge distillation**



(single mini-batch/ streaming data/ entire dataset)

**Test time adaptation**
Source: Liang et al. 2023

# xMUDA
# Cross-Modal Learning for Domain Adaptation

CVPR 20 & TPAMI 22

**xMUDA** | Jaritz et al., CVPR 20 & TPAMI 22

**3D**

sparse voxel

3D segmentation

45

**2D**

dense pixel

2D-3D lifting

**3D**

sparse voxel

3D segmentation

2D

3D

dense pixel

2D-3D lifting

sparse voxel

CROSS-MODAL LEARNING

3D segmentation

45

Camera Image $(H, W, 3)$

2D Network

dense pixel

Features $(N, F_{2D})$

Probabilities $(N, C)$

sample

classify

$P_{2D}$

project

$D_{KL}(P_{3D} \| P_{2D})$

$D_{KL}(P_{2D} \| P_{3D})$

$P_{3D}$

Lidar Point Cloud $(N, 3)$

sparse voxel

3D Network

classify

$(N, F_{3D})$

$(N, C)$

| $H, W$ | image size |
| $N$ | num. points |
| $F_{2D,3D}$ | num. feature channels |
| $C$ | num. classes |

46

$$\mathcal{L}_{\mathrm{xM}}(\boldsymbol{x}) = \boldsymbol{D}_{\mathrm{KL}}(\boldsymbol{P}_{\boldsymbol{x}}^{(n,c)}||\boldsymbol{Q}_{\boldsymbol{x}}^{(n,c)})$$

$$= -\frac{1}{N}\sum_{n=1}^{N}\sum_{c=1}^{C}\boldsymbol{P}_{\boldsymbol{x}}^{(n,c)}\log\frac{\boldsymbol{P}_{\boldsymbol{x}}^{(n,c)}}{\boldsymbol{Q}_{\boldsymbol{x}}^{(n,c)}}$$

Complete objective:

$$\min_{\theta}\left[\frac{1}{|\mathcal{S}|}\sum_{\boldsymbol{x}_s\in\mathcal{S}}\left(\mathcal{L}_{\mathrm{seg}}(\boldsymbol{x}_s,\boldsymbol{y}_s^{3\mathrm{D}})+\lambda_s\mathcal{L}_{\mathrm{xM}}(\boldsymbol{x}_s)\right)+\frac{1}{|\mathcal{T}|}\sum_{\boldsymbol{x}_t\in\mathcal{T}}\lambda_t\mathcal{L}_{\mathrm{xM}}(\boldsymbol{x}_t)\right]$$

**Single Head**

# Why does it fail ?



**2D private**    **shared**    **3D private**

$$\mathcal{L}_{\mathrm{xM}}(\boldsymbol{x}) = D_{\mathrm{KL}}(\boldsymbol{P}_{\boldsymbol{x}}^{(n,c)} || \boldsymbol{Q}_{\boldsymbol{x}}^{(n,c)})$$

$$= -\frac{1}{N} \sum_{n=1}^{N} \sum_{c=1}^{C} \boldsymbol{P}_{\boldsymbol{x}}^{(n,c)} \log \frac{\boldsymbol{P}_{\boldsymbol{x}}^{(n,c)}}{\boldsymbol{Q}_{\boldsymbol{x}}^{(n,c)}}$$

Complete objective:

$$\min_{\theta} \left[ \frac{1}{|\mathcal{S}|} \sum_{\boldsymbol{x}_s \in \mathcal{S}} \left( \mathcal{L}_{\mathrm{seg}}(\boldsymbol{x}_s, \boldsymbol{y}_s^{\mathrm{3D}}) + \lambda_s \mathcal{L}_{\mathrm{xM}}(\boldsymbol{x}_s) \right) + \frac{1}{|\mathcal{T}|} \sum_{\boldsymbol{x}_t \in \mathcal{T}} \lambda_t \mathcal{L}_{\mathrm{xM}}(\boldsymbol{x}_t) \right]$$

$(N, F_{\mathrm{2D}})$    $(N, C)$

classify    $P_{\mathrm{2D}}$

$D_{\mathrm{KL}}(P_{\mathrm{3D}} || P_{\mathrm{2D}})$    $D_{\mathrm{KL}}(P_{\mathrm{2D}} || P_{\mathrm{3D}})$

$P_{\mathrm{3D}}$

classify

$(N, F_{\mathrm{3D}})$    $(N, C)$

**Single Head**

# Disentangled objective

Single head

**Dual head**
(xMUDA)

**xMUDA / xMUDA_PL**

|  | nuScenes-Lidarseg [10]: USA/Singapore | nuScenes-Lidarseg [10]: Day/Night | Virt.KITTI [56]/ Sem.KITTI [2] | A2D2 [57]/ Sem.KITTI [2] | Waymo OD [58]: SF,PHX,MTV/KRK |
| Source | | | | | |
| Target | | | | | |
| | **Right-to-left driving** | **Illumination at night** | **Lack of Realism** | **Lidar density changes** | **Weather changes** |

**Main challenges**

| Method | nuSc-Lidarseg: USA/Singap. | | | nuSc-Lidarseg: Day/Night | | | Virt.KITTI/Sem.KITTI | | | A2D2/Sem.KITTI | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 2D | 3D | 2D+3D | 2D | 3D | 2D+3D | 2D | 3D | 2D+3D | 2D | 3D | 2D+3D |
| Baseline (src only) | 58.4 | 62.8 | 68.2 | 47.8 | 68.8 | 63.3 | 26.8 | 42.0 | 42.2 | 34.2 | 35.9 | 40.4 |
| Deep logCORAL [20] | 64.4 | 63.2 | 69.4 | 47.7 | 68.7 | 63.7 | 41.4* | 36.8 | 47.0 | 35.1* | 41.0 | 42.2 |
| MinEnt [5] | 57.6 | 61.5 | 66.0 | 47.1 | 68.8 | 63.6 | 39.2 | 43.3 | 47.1 | 37.8 | 39.6 | 42.6 |
| PL [7] | 62.0 | 64.8 | 70.4 | 47.0 | 69.6 | 63.0 | 21.5 | 44.3 | 35.6 | 34.7 | 41.7 | 45.2 |
| FDA [14] | 60.8 | - | - | 48.4 | - | - | 32.8* | - | - | 37.6* | - | - |
| xMUDA | 64.4 | 63.2 | 69.4 | 55.5 | 69.2 | 67.4 | 42.1 | 46.7 | 48.2 | 38.3 | 46.0 | 44.0 |
| xMUDA$_{PL}$ | 67.0 | 65.4 | 71.2 | 57.6 | 69.6 | 64.4 | 45.8 | 51.4 | 52.0 | 41.2 | 49.8 | 47.5 |
| Oracle | 75.4 | 76.0 | 79.6 | 61.5 | 69.8 | 69.2 | 66.3 | 78.4 | 80.1 | 59.3 | 71.9 | 73.6 |

2D+3D = ensembling

[5] MinEnt / Advent. Vu et al. CVPR 19
[7] Dong-Hyun. ICML 13
[14] Yanchao and Soatto. CVPR 20
[20] Deep LogCoral. Yifei et al. ICCV-W 17

Ground Truth

Image

Color Legend

■ Car ■ Truck ■ Bike ■ Person ■ Road ■ Sidewalk
■ Parking ■ Nature ■ Building ■ Other objects ■ Unlabeled

Baseline (no adaptation)

xMUDA

53

# Can fusion benefit from mimicking ?



Vanilla Fusion



xMUDA Fusion

| Method | Arch. | nuSc-Lidarseg: USA/Singap. | A2D2/ Sem.KITTI |
|---|---|---|---|
| Baseline (src only) | Vanilla | 66.5 | 34.2 |
| Deep logCORAL [20] | Vanilla | 64.0 | 36.2 |
| MinEnt [5] | Vanilla | 65.4 | 39.8 |
| PL [7] | Vanilla | 70.1 | 38.6 |
| xMUDA Fusion | xMUDA | 69.3 | **42.6** |
| xMUDA$_{PL}$ Fusion | xMUDA | **70.7** | 42.2 |
| Oracle | xMUDA | 80.6 | 65.7 |

# Open-vocabulary

# Vision Language Model (VLM)



Clip: (Radford et al., 2021)

# Vision Language Model (VLM)



Clip: (Radford et al., 2021)

« Teddy bears mixing sparkling chemicals as mad scientists in a steampunk style »

DALL-E 2 (Openai, 2022)

# PØDA

## Prompt-driven Zero-shot Domain Adaptation

ICCV 23

github.com/astra-vision/PODA

Fahes, Vu, Bursuc, Pérez, de Charette. ICCV 2023

segmenter
(trained on
Cityscapes

PØDA w/ prompt
"driving at night"

segmenter
(trained on Cityscapes

59

PØDA w/ prompt "driving at night"

PØDA w/ prompt "driving through fire"

segmenter (trained on Cityscapes

PØDA w/ prompt "driving in old movie"

PØDA w/ prompt "driving in sandstorm"

59

PØDA w/ prompt
"driving at night"

PØDA w/ prompt
"driving through fire"

segmenter
(trained on Cityscapes

PØDA w/ prompt
"driving in old movie"

PØDA w/ prompt
"driving in sandstorm"

**Minimal features stylization**

"driving under rain"

"driving in snow"

"driving at night" •••

+ text embedding
• image embedding
→ feat. augmentation
···▶ loss

$E_{\text{txt}}$

source domain

$E_{\text{img}}$

Cityscapes

$\bar{\mathbf{f}}_{\text{night}}$

$\bar{\mathbf{f}}_{\text{s}}$

$\bar{\mathbf{f}}_{\text{snow}}$

$\bar{\mathbf{f}}_{\text{rain}}$

+

+

+

low-level features
augmentation

$\mathbf{f}_{\text{s}}$

$\mathbf{f}_{\text{s}\to\text{t}}$

shared space

**Prompt-driven feature augmentation**

$\mathbf{f}_{\text{s}\to\text{snow}}$

$\mathbf{f}_{\text{s}\to\text{rain}}$

•••

$\mathbf{f}_{\text{s}\to\text{night}}$

source domain

$E_{\text{img}}$ $M_{\text{night}}$

Zero-shot Domain Adaptation

60

# Adaptive Instance Normalization (AdaIN)

**High-level feature statistics captures style**

AdaIN simply transfers features statistics from **y** to **x** by normalizing and rescaling

$$\mathrm{AdaIN}(x,y) = \sigma(y)\frac{x-\mu(x)}{\sigma(x)} + \mu(y)$$ No learnable parameters

# Adaptive Instance Normalization (AdaIN)

**High-level feature statistics captures style**

AdaIN simply transfers features statistics from **y** to **x** by normalizing and rescaling

$$\text{AdaIN}(x, y) = \sigma(y) \frac{x - \mu(x)}{\sigma(x)} + \mu(y)$$ No learnable parameters



Y

## 5. Adaptive Instance Normalization

If IN normalizes the input to a single style specified by the affine parameters, is it possible to adapt it to arbitrarily given styles by using adaptive affine transformations? Here, we propose a simple extension to IN, which we call adaptive instance normalization (AdaIN). AdaIN receives a content input $x$ and a style input $y$, and simply aligns the channel-wise mean and variance of $x$ to match those of $y$. Unlike BN, IN or CIN, AdaIN has no learnable affine parameters. Instead, it adaptively computes the affine parameters from the style input:

$$\text{AdaIN}(x, y) = \sigma(y) \left( \frac{x - \mu(x)}{\sigma(x)} \right) + \mu(y) \qquad (8)$$

in which we simply scale the normalized content input with $\sigma(y)$, and shift it with $\mu(y)$. Similar to IN, these statistics are computed across spatial locations.

Intuitively, let us consider a feature channel that detects brushstrokes of a certain style. A style image with this kind of strokes will produce a high average activation for this feature. The output produced by AdaIN will have the same high average activation for this feature, while preserving the spatial structure of the content image. The brushstroke feature can be inverted to the image space with a feed-forward decoder, similar to [10]. The variance of this feature channel can encoder more subtle style information, which is also transferred to the AdaIN output and the final output image.

In short, AdaIN performs style transfer in the feature space by transferring feature statistics, specifically the channel-wise mean and variance. Our AdaIN layer plays a similar role as the style swap layer proposed in [6]. While the style swap operation is very time-consuming and memory-consuming, our AdaIN layer is as simple as an IN layer, adding almost no computational cost.

# Adaptive Instance Normalization (AdaIN)

**High-level feature statistics captures style**

> AdaIN simply transfers features statistics from **y** to **x** by normalizing and rescaling
>
> $$\text{AdaIN}(x, y) = \sigma(y)\frac{x - \mu(x)}{\sigma(x)} + \mu(y)$$ No learnable parameters



Huang. and Belongie, AdaIN. ICCV 2017

# Adaptive Instance Normalization (AdaIN)

**High-level feature statistics captures style**

AdaIN simply transfers features statistics from **y** to **x** by normalizing and rescaling

$$\text{AdaIN}(x, y) = \sigma(y) \frac{x - \mu(x)}{\sigma(x)} + \mu(y)$$ No learnable parameters



x

content → AdaIN → D

y

style

**But, we don't have access to target data !**

Huang. and Belongie, AdaIN. ICCV 2017

# Prompt-Driven Instance Normalization (PIN)



"driving at night"

$E_{\mathrm{txt}}$

$+$ TrgEmb

backbone
(CLIP pre-trained, frozen)

Layer1

$E_{\mathrm{img}}$

a source image

$\bar{\mathbf{f}}_{\mathrm{s}}$

origin

$$\mathrm{PIN}(\mathbf{f}_{\mathrm{s}}, \boldsymbol{\mu}, \boldsymbol{\sigma}) = \boldsymbol{\sigma} \left( \frac{\mathbf{f}_{\mathrm{s}} - \mu(\mathbf{f}_{\mathrm{s}})}{\sigma(\mathbf{f}_{\mathrm{s}})} \right)$$

**Gradient Descent**

# Prompt-Driven Instance Normalization (PIN)



$$\text{PIN}(\mathbf{f_s}, \boldsymbol{\mu}, \boldsymbol{\sigma}) = \boldsymbol{\sigma}\left(\frac{\mathbf{f_s} - \mu(\mathbf{f_s})}{\sigma(\mathbf{f_s})}\right)$$

**Gradient Descent** $\qquad \mathcal{L}_{\mu,\sigma}(\bar{\mathbf{f}}_{s \to t}, \mathbf{TrgEmb}) = 1 - \dfrac{\bar{\mathbf{f}}_{s \to t} \cdot \mathbf{TrgEmb}}{\|\bar{\mathbf{f}}_{s \to t}\| \, \|\mathbf{TrgEmb}\|}$

# Prompt-Driven Instance Normalization (PIN)



"driving at night"

$E_{\text{txt}}$

backbone
(CLIP pre-trained, frozen)

Layer1

$E_{\text{img}}$

$\bar{\mathbf{f}}_{\text{s} \to \text{t}}^{i}$

**TrgEmb**

a source image

$\bar{\mathbf{f}}_{\text{s}}$

$\mathcal{L}_{\mu, \sigma}$

origin

$(\mu^{i}, \sigma^{i})$

$\text{PIN}(\mathbf{f}_{\text{s}}, \boldsymbol{\mu}, \boldsymbol{\sigma}) = \boldsymbol{\sigma} \left( \dfrac{\mathbf{f}_{\text{s}} - \mu(\mathbf{f}_{\text{s}})}{\sigma(\mathbf{f}_{\text{s}})} \right)$

$\mathbf{f}_{\text{s}}$  $\mathbf{f}_{\text{s} \to \text{t}}^{i}$  $\text{augment}(\cdot)$  optimized statistics

$(\mu^{i-1}, \sigma^{i-1})$

$(\mu^{i-1}, \sigma^{i-1})$

**PIN**

$(\mu^{0}, \sigma^{0})$

**Gradient Descent**

$\mathcal{L}_{\mu, \sigma}(\bar{\mathbf{f}}_{\text{s} \to \text{t}}, \mathbf{TrgEmb}) = 1 - \dfrac{\bar{\mathbf{f}}_{\text{s} \to \text{t}} \cdot \mathbf{TrgEmb}}{\|\bar{\mathbf{f}}_{\text{s} \to \text{t}}\| \, \|\mathbf{TrgEmb}\|}$

| $\boldsymbol{\mu}^0$ | $\boldsymbol{\sigma}^0$ | mIoU |
|:---:|:---:|:---:|
| $\mu(\mathbf{f_s})$ | $\sigma(\mathbf{f_s})$ | $\mathbf{25.03}_{\pm 0.48}$ |
| $\mathbf{0}$ | $\mathbf{1}$ | $8.59_{\pm 0.82}$ |
| $\sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ | $\sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ | $6.80_{\pm 0.92}$ |

**Initialization**

| $\boldsymbol{\mu}^0$ | $\boldsymbol{\sigma}^0$ | mIoU |
|:---:|:---:|:---:|
| $\mu(\mathbf{f_s})$ | $\sigma(\mathbf{f_s})$ | $\mathbf{25.03}_{\pm 0.48}$ |
| $\mathbf{0}$ | $\mathbf{1}$ | $8.59_{\pm 0.82}$ |
| $\sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ | $\sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ | $6.80_{\pm 0.92}$ |

**Initialization**

| *Layer1* | *Layer2* | *Layer3* | *Layer4* | ACDC Night |
|:---:|:---:|:---:|:---:|:---:|
| ✓ | ✗ | ✗ | ✗ | $\mathbf{25.03}_{\pm 0.48}$ |
| ✓ | ✓ | ✗ | ✗ | $23.43_{\pm 0.51}$ |
| ✓ | ✗ | ✓ | ✗ | $22.93_{\pm 0.53}$ |
| ✓ | ✗ | ✗ | ✓ | $21.05_{\pm 0.55}$ |

**Minimal augmentation**

**Initialization**

| $\boldsymbol{\mu}^0$ | $\boldsymbol{\sigma}^0$ | mIoU |
|---|---|---|
| $\mu(\mathbf{f_s})$ | $\sigma(\mathbf{f_s})$ | $\mathbf{25.03}_{\pm 0.48}$ |
| $\mathbf{0}$ | $\mathbf{1}$ | $8.59_{\pm 0.82}$ |
| $\sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ | $\sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ | $6.80_{\pm 0.92}$ |

**Minimal augmentation**

| Layer1 | Layer2 | Layer3 | Layer4 | ACDC Night |
|---|---|---|---|---|
| ✓ | ✗ | ✗ | ✗ | $\mathbf{25.03}_{\pm 0.48}$ |
| ✓ | ✓ | ✗ | ✗ | $23.43_{\pm 0.51}$ |
| ✓ | ✗ | ✓ | ✗ | $22.93_{\pm 0.53}$ |
| ✓ | ✗ | ✗ | ✓ | $21.05_{\pm 0.55}$ |

**# of iterations**

63

$$\text{PIN}(\mathbf{f}_{\mathrm{s}}, \boldsymbol{\mu}, \boldsymbol{\sigma}) = \boldsymbol{\sigma} \left( \frac{\mathbf{f}_{\mathrm{s}} - \mu(\mathbf{f}_{\mathrm{s}})}{\sigma(\mathbf{f}_{\mathrm{s}})} \right)$$

64

"driving under rain"
"driving in snow"
"driving at night"

$E_{txt}$

+ text embedding
• image embedding
→ feat. augmentation
⋯► loss

source domain

$E_{img}$

Cityscapes

low-level features augmentation

$\mathbf{f}_s$    $\mathbf{f}_{s \to t}$

$\bar{\mathbf{f}}_{night}$

$\bar{\mathbf{f}}_s$

$\bar{\mathbf{f}}_{snow}$

$\bar{\mathbf{f}}_{rain}$

shared space

**Prompt-driven feature augmentation**

$\mathbf{f}_{s \to snow}$
$\mathbf{f}_{s \to rain}$
$\mathbf{f}_{s \to night}$

source domain

$E_{img}$    $M_{night}$

$\mathcal{L}_{seg}$

**Zero-shot Domain Adaptation**

65

**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   $M_{\text{cls}} \leftarrow \texttt{train}(M_{\text{cls}}, \mathcal{D}_s)$        ▷ source-only training

// Feature Augmentation

3   $\mathcal{F}_s \leftarrow \texttt{feat-ext}(M_{\text{feat}}, \{\mathbf{x}_s\})$

4   $\mathcal{S}_{s \rightarrow t} \leftarrow \texttt{augment}(\mathcal{F}_s, \text{TrgEmb})$

// Adaptation

5   $M'_{\text{cls}} \leftarrow \texttt{fine-tune}(M_{\text{cls}}, \mathcal{F}_s, \mathcal{S}_{s \rightarrow t}, \{\mathbf{y}_s\})$ ▷ fine-tuning

``driving in old movie"

Input                    Source-only                    PØDA (ours)
                                                        ¨driving at night¨

67

# Evaluation on ACDC (Sakaridis et al., ICCV'21) and GTA5 (Richter et al., ECCV'16)



**Proxies**

| Source | Target eval. | Method | mIoU[%] |
|---|---|---|---|
| CS | | TrgPrompt = "driving at night" | |
| | ACDC Night | source-only | 18.31 |
| | | CLIPstyler [1] | 21.38 ±0.36 |
| | | PØDA | **25.03** ±0.48 |
| | | TrgPrompt = "driving in snow" | |
| | ACDC Snow | source-only | 39.28 |
| | | CLIPstyler [1] | 41.09 ±0.17 |
| | | PØDA | **43.90** ±0.53 |
| | | TrgPrompt = "driving under rain" | |
| | ACDC Rain | source-only | 38.20 |
| | | CLIPstyler [1] | 37.17 ±0.10 |
| | | PØDA | **42.31** ±0.55 |
| | | TrgPrompt = "driving in a game" | |
| | GTA5 | source-only | 39.59 |
| | | CLIPstyler [1] | 38.73 ±0.16 |
| | | PØDA | **41.07** ±0.48 |
| GTA5 | CS | TrgPrompt = "driving" | |
| | | source-only | 36.38 |
| | | CLIPstyler [1] | 31.50 ±0.21 |
| | | PØDA | **40.08** ±0.52 |

[1] **Clipstyler**, CVPR 2022

68

# Comparison to CLIPStyler



CLIPStyler optimization: 65sec
PODA optimization: **0.3sec**



| Source (CS) | Rain | Night | Snow | Game |



[1] **Clipstyler**, CVPR 2022

# Comparison to CLIPStyler



Output $I_{cs}$

Crop & Augment

Threshold Rejection

"Picasso style painting"

Patchwise CLIP Loss $L_{patch}$

StyleNet $f$

"Photo"

Content Image $I_c$



Content Image — "Oil painting of flowers" — "Pop art of night city" — "Neon light"

Source (CS) — Rain — Night — Snow — Game

CLIPStyler optimization: 65sec
PODA optimization: **0.3sec**

[1] **Clipstyler**, CVPR 2022

Source (CS)

70

"Driving in snow"                    "Driving in a game"

# Prompt design

give me 5 prompts that have the
same exact meaning as "{prompt}"

Chat GPT

give me 5 random prompts of
length from 3 to 6 words
describing a random photo

| Method | ACDC Night | ACDC Snow | ACDC Rain | GTA5 |
|---|---|---|---|---|
| Source only | 18.31 | 39.28 | 38.20 | 39.59 |
| Trg | "driving at night" | "driving in snow" | "driving under rain" | "driving in a game" |
| | $25.03_{\pm 0.48}$ | $43.90_{\pm 0.53}$ | $42.31_{\pm 0.55}$ | $41.07_{\pm 0.48}$ |
| | "operating a vehicle after sunset" | "operating a vehicle in snowy conditions" | "operating a vehicle in wet conditions" | "piloting a vehicle in a virtual world" |
| | $24.38_{\pm 0.37}$ | $44.33_{\pm 0.36}$ | $42.21_{\pm 0.47}$ | $41.25_{\pm 0.40}$ |
| | "driving during the nighttime hours" | "driving on snow-covered roads" | "driving on rain-soaked roads" | "controlling a car in a digital simulation" |
| | $\mathbf{25.22}_{\pm 0.64}$ | $43.56_{\pm 0.62}$ | $\mathbf{42.51}_{\pm 0.33}$ | $41.19_{\pm 0.14}$ |
| | "navigating the roads in darkness" | "piloting a vehicle in snowy terrain" | "navigating through rainfall while driving" | "maneuvering a vehicle in a computerized racing experience" |
| | $24.73_{\pm 0.47}$ | $\mathbf{44.67}_{\pm 0.18}$ | $41.11_{\pm 0.69}$ | $40.34_{\pm 0.49}$ |
| | "driving in low-light conditions" | "driving in wintry precipitation" | "driving in inclement weather" | "operating a transport in a video game environment" |
| | $24.68_{\pm 0.34}$ | $43.11_{\pm 0.56}$ | $40.68_{\pm 0.37}$ | $41.34_{\pm 0.42}$ |
| | "travelling by car after dusk" | "travelling by car in a snowstorm" | "travelling by car during a downpour" | "navigating a machine through a digital driving simulation" |
| | $24.89_{\pm 0.24}$ | $43.83_{\pm 0.17}$ | $42.05_{\pm 0.35}$ | $\mathbf{41.86}_{\pm 0.10}$ |
| | *24.82* | *43.90* | *41.81* | *41.18* |
| | "mesmerizing northern lights display" | | | |
| | $20.05_{\pm 0.77}$ | $40.07_{\pm 0.66}$ | $38.43_{\pm 0.82}$ | $37.98_{\pm 0.31}$ |
| | "playful dolphins in the ocean" | | | |
| | $20.11_{\pm 0.31}$ | $39.87_{\pm 0.26}$ | $38.56_{\pm 0.58}$ | $37.05_{\pm 0.31}$ |
| | "breathtaking view from mountaintop" | | | |
| | $20.65_{\pm 0.33}$ | $42.08_{\pm 0.28}$ | $40.05_{\pm 0.52}$ | $40.09_{\pm 0.23}$ |
| | "cheerful sunflower field in bloom" | | | |
| | $21.10_{\pm 0.50}$ | $39.85_{\pm 0.68}$ | $40.09_{\pm 0.41}$ | $37.93_{\pm 0.55}$ |
| | "dramatic cliff overlooking the ocean" | | | |
| | $20.09_{\pm 0.98}$ | $38.20_{\pm 0.54}$ | $38.48_{\pm 0.37}$ | $37.57_{\pm 0.46}$ |
| | "majestic eagle in flight over mountains" | | | |
| | $20.70_{\pm 0.38}$ | $39.60_{\pm 0.27}$ | $40.38_{\pm 0.86}$ | $38.52_{\pm 0.21}$ |
| | *20.45* | *39.95* | *39.33* | *38.19* |

(Rows grouped vertically: "ChatGPT-generated" with "Relevant →" for the upper block; "Irrelevant ←" for the lower block)

# Prompt design

give me 5 prompts that have the
same exact meaning as "{prompt}"

 Chat GPT

give me 5 random prompts of
length from 3 to 6 words
describing a random photo

| Method | ACDC Night | ACDC Snow | ACDC Rain | GTA5 |
|---|---|---|---|---|
| Source only | 18.31 | 39.28 | 38.20 | 39.59 |
| Trg | "driving at night" | "driving in snow" | "driving under rain" | "driving in a game" |
| | $25.03_{\pm0.48}$ | $43.90_{\pm0.53}$ | $42.31_{\pm0.55}$ | $41.07_{\pm0.48}$ |
| | "operating a vehicle after sunset" | "operating a vehicle in snowy conditions" | "operating a vehicle in wet conditions" | "piloting a vehicle in a virtual world" |
| | $24.38_{\pm0.37}$ | $44.33_{\pm0.36}$ | $42.21_{\pm0.47}$ | $41.25_{\pm0.40}$ |
| | "driving during the nighttime hours" | "driving on snow-covered roads" | "driving on rain-soaked roads" | "controlling a car in a digital simulation" |
| | $\mathbf{25.22}_{\pm0.64}$ | $43.56_{\pm0.62}$ | $\mathbf{42.51}_{\pm0.33}$ | $41.19_{\pm0.14}$ |
| | "navigating the roads in darkness" | "piloting a vehicle in snowy terrain" | "navigating through rainfall while driving" | "maneuvering a vehicle in a computerized racing experience" |
| | $24.73_{\pm0.47}$ | $\mathbf{44.67}_{\pm0.18}$ | $41.11_{\pm0.69}$ | $40.34_{\pm0.49}$ |
| | "driving in low-light conditions" | "driving in wintry precipitation" | "driving in inclement weather" | "operating a transport in a video game environment" |
| | $24.68_{\pm0.34}$ | $43.11_{\pm0.56}$ | $40.68_{\pm0.37}$ | $41.34_{\pm0.42}$ |
| | "travelling by car after dusk" | "travelling by car in a snowstorm" | "travelling by car during a downpour" | "navigating a machine through a digital driving simulation" |
| | $24.89_{\pm0.24}$ | $43.83_{\pm0.17}$ | $42.05_{\pm0.35}$ | $\mathbf{41.86}_{\pm0.10}$ |
| | 24.82 | 43.90 | 41.81 | 41.18 |
| | "mesmerizing northern lights display" | | | |
| | $20.05_{\pm0.77}$ | $40.07_{\pm0.66}$ | $38.43_{\pm0.82}$ | $37.98_{\pm0.31}$ |
| | "playful dolphins in the ocean" | | | |
| | $20.11_{\pm0.31}$ | $39.87_{\pm0.26}$ | $38.56_{\pm0.58}$ | $37.05_{\pm0.31}$ |
| | "breathtaking view from mountaintop" | | | |
| | $20.65_{\pm0.33}$ | $42.08_{\pm0.28}$ | $40.05_{\pm0.52}$ | $40.09_{\pm0.23}$ |
| | "cheerful sunflower field in bloom" | | | |
| | $21.10_{\pm0.50}$ | $39.85_{\pm0.68}$ | $40.09_{\pm0.41}$ | $37.93_{\pm0.55}$ |
| | "dramatic cliff overlooking the ocean" | | | |
| | $20.09_{\pm0.98}$ | $38.20_{\pm0.54}$ | $38.48_{\pm0.37}$ | $37.57_{\pm0.46}$ |
| | "majestic eagle in flight over mountains" | | | |
| | $20.70_{\pm0.38}$ | $39.60_{\pm0.27}$ | $40.38_{\pm0.86}$ | $38.52_{\pm0.21}$ |
| | 20.45 | 39.95 | 39.33 | 38.19 |

Relevant → ChatGPT-generated ← Irrelevant

Always better

Always worse

$$y = \sigma_s^\star \frac{x - \mu_c}{\sigma_c} + \mu_s^\star,$$

$$\sigma_s^\star = \boxed{\alpha\sigma_c,} \qquad \mu_s^\star = \boxed{\beta\mu_c}$$



Content Images    Style Randomization   ⊢ - - - - - - - - - -   Our Generated Images by Inverting Latent Styles   - - - - - - - - - ⊣

**Towards robust object detection. (Fan et al., ICLR 23)**

$$y = \sigma_s^\star \frac{x - \mu_c}{\sigma_c} + \mu_s^\star,$$

$$\sigma_s^\star = \boxed{\alpha\sigma_c,} \qquad \mu_s^\star = \boxed{\beta\mu_c}$$



Content Images | Style Randomization | ←----------- Our Generated Images by Inverting Latent Styles ----------→

**Towards robust object detection. (Fan et al., ICLR 23)**

| Method | Target | S → C | C → F |
|---|---|---|---|
| Our Baseline | ✗ | 32.8 | 22.0 |
| BIN | ✗ | 44.3 | 28.4 |
| IBN | ✗ | 47.4 | 31.2 |
| SFA | ✗ | 38.4 | 25.3 |
| pAdaIN | ✗ | 43.7 | 27.6 |
| Mixstyle | ✗ | 46.4 | 30.1 |
| DSU | ✗ | 49.3 | 34.1 |
| NP (Ours) | ✗ | 54.1 | 44.0 |
| NP+ (Ours) | ✗ | 58.7 | 46.3 |

+9%  +12%

| | BDD Day → Night | | | BDD Night → Day | | | WaymoL → BDD | | | WaymoR → BDD | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | AP | AP50 | AP75 | AP | AP50 | AP75 | AP | AP50 | AP75 | AP | AP50 | AP75 |
| Faster R-CNN | 17.84 | 31.35 | 17.68 | 19.14 | 33.04 | 19.16 | 10.07 | 19.62 | 9.05 | 8.65 | 17.26 | 7.49 |
| + CycConsist | 18.35 | 32.44 | 18.07 | 18.89 | 33.50 | 18.31 | 11.55 | 23.44 | 10.00 | 9.11 | 17.92 | 7.98 |
| + CycConf | 19.09 | 33.58 | 19.14 | 19.57 | 34.34 | **19.26** | 12.27 | 26.01 | 10.24 | 9.99 | 20.58 | 8.30 |
| + NP (Ours) | 20.73 | 36.22 | 20.85 | 19.32 | 34.42 | 18.63 | 17.85 | 35.34 | 15.52 | 14.97 | 29.42 | 13.11 |
| + NP+ (Ours) | **20.97** | **36.76** | **21.10** | **19.73** | **35.30** | 19.19 | **21.18** | **42.16** | **18.67** | **19.64** | **38.69** | **17.07** |

+2%          +0.5%          +9%          +10%

72

| Method | Night | Snow | Rain | GTA5 |
|---|---|---|---|---|
| Source-only | 18.31 | 39.28 | 38.20 | 39.59 |
| Source-only-G | 21.07 | 42.84 | 42.38 | 41.54 |
| PØDA-G | $\mathbf{24.86}_{\pm 0.70}$ | $44.34_{\pm 0.36}$ | $43.17_{\pm 0.63}$ | $41.73_{\pm 0.39}$ |
| PØDA-G+style-mix | $24.18_{\pm 0.23}$ | $\mathbf{44.46}_{\pm 0.34}$ | $\mathbf{43.56}_{\pm 0.46}$ | $\mathbf{42.98}_{\pm 0.12}$ |

**Source-only**

**Source-only-G**

$(\mu, \sigma) \in \mathrm{U}(0,1)$

$\mathbf{f}_{s \to t} = \mathrm{AdaIN}(\mathbf{f}_s, \boldsymbol{\mu}, \boldsymbol{\sigma})$

Inspiration [1]

**PODA-G**

$(\mu_t, \sigma_t)$ PODA on $(\mu, \sigma)$

$\mathbf{f}_{s \to t} = \mathrm{PIN}(\mathbf{f}_s, \boldsymbol{\mu}_t, \boldsymbol{\sigma}_t)$

**PODA-G + style mixing**

$(\mu_t, \sigma_t)$ PODA on $(\mu, \sigma)$

$\mathbf{f}_{s \to t} = \mathrm{PIN}(\mathbf{f}_s, \boldsymbol{\mu}_{\mathrm{mix}}, \boldsymbol{\sigma}_{\mathrm{mix}})$

Inspiration [2]

TrgEmbed

TrgEmbed

[1] Fan et al. Towards robust object detection invariant to real-world domain shifts. ICLR 23
[2] Wu et al. Style mixing and patchwise prototypical matching for oneshot unsupervised domain adaptive semantic segmentation. AAAI 22

## Various backbones

| Backbone | Method | Night | Snow | Rain | GTA5 |
|---|---|---|---|---|---|
| Sem. FPN | src-only | 18.10 | 35.75 | 36.07 | 40.67 |
|  | PØDA | **21.48** ±0.15 | **39.55** ±0.13 | **38.34** ±0.29 | **41.59** ±0.24 |
| DLv3+ | src-only | 22.17 | 44.53 | 42.53 | 40.49 |
|  | PØDA | **26.54** ±0.12 | **46.71** ±0.43 | **46.36** ±0.20 | **43.17** ±0.13 |

## Effect of priors

| Method | Prior | ACDC Night |
|---|---|---|
| CIConv* [26] | physics | 30.60 / 34.50 (Δ=3.90) |
| SM-PPM [56] | 1 target image | 13.07 / 14.60 (Δ=1.53) |
| CLIPstyler [25] | 1 prompt | 18.31 / 21.38 (Δ=3.07) |
| PØDA | 1 prompt | 18.31 / 25.03 (Δ=6.72) |

* Results of CIConv are on DarkZurich, a subset of ACDC Night [45].

Cityscapes-Foggy. Sakaridis et al., IJCV 2018
Diverse Weather Dataset, Wu et al. ECCV'22
CUB-200. Wah et al. 2011
CUB-200-Paintings. Wang et al. CVPR'20

[8] DA-Faster. Chen et al. IJCV'21
[15] NP+. Fan et al. ICLR'23
[25] ClipStyler. Kwon and Ye, CVPR'22
[26] CIConv. Lengyel et al. ICCV'21

[42] ViSGA. Rezaeianaran et al. ICCV'21
[49] CLIP The Gap. Vidit et al. CVPR'23
[55] S-DGOD. Wu and Deng, ECCV'22
[56] SM-PPM. Wu et al., AAAI'22

| Method | Target | CS→ CS Foggy | DWD-Day Clear → | | | |
|---|---|---|---|---|---|---|
| | | | Night Clear | Dusk Rainy | Night Rainy | Day Foggy |
| DA-Faster [8] | ✓ | 32.0 | - | - | - | - |
| ViSGA [42] | ✓ | 43.3 | - | - | - | - |
| NP+ [15] | ✗ | 46.3 | - | - | - | - |
| S-DGOD [55] | ✗ | - | 36.6 | 28.2 | 16.6 | 33.5 |
| CLIP The Gap [49] | ✗ | - | 36.9 | 32.3 | 18.7 | 38.5 |
| PØDA (Faster-RCNN) | ✗ | **47.3** | **43.4** | **40.2** | **20.5** | **44.4** |

**Object Detection**

"Painting of a bird"

"Blue/Red digits"

| Method | CUB-200 paintings | Colored MNIST |
|---|---|---|
| src-only | 28.90 | 55.83 |
| PØDA | **30.91** $\pm0.69$ | **64.16** $\pm0.41$ |

**Classification**

Cityscapes-Foggy. Sakaridis et al., IJCV 2018
Diverse Weather Dataset, Wu et al. ECCV'22
CUB-200. Wah et al. 2011
CUB-200-Paintings. Wang et al. CVPR'20

[8] DA-Faster. Chen et al. IJCV'21
[15] NP+. Fan et al. ICLR'23
[25] ClipStyler. Kwon and Ye, CVPR'22
[26] CIConv. Lengyel et al. ICCV'21

[42] ViSGA. Rezaeianaran et al. ICCV'21
[49] CLIP The Gap. Vidit et al. CVPR'23
[55] S-DGOD. Wu and Deng, ECCV'22
[56] SM-PPM. Wu et al., AAAI'22

# Limitations ?

- Global stylization: large structural classes benefit more from it

# Limitations ?

- Global stylization: large structural classes benefit more from it



| 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 | **43.83** | 6.11 | 17.08 | 23.52 | **24.51** | 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 | **35.49** | 23.34 | 19.01 | **45.41** | 8.81 | **27.87** | 21.06 | 8.48 | 38.17 | 1.84 | 0.00 | 11.54 | **10.38** | 4.89 | 21.38 ±0.36 |
| CS | | PØDA | **77.54** | **26.90** | 42.71 | **13.51** | **21.36** | 33.52 | 23.70 | **21.73** | 39.91 | **9.51** | 19.40 | **28.80** | **11.85** | **50.89** | **10.14** | 0.00 | **20.76** | 8.76 | **14.50** | **25.03** ±0.48 |
| | | | | | | | | | TrgPrompt = "driving in snow" | | | | | | | | | | | | | |
| | ACDC snow | source-only | 70.47 | 23.50 | 63.80 | 17.96 | **27.36** | **38.52** | **56.26** | **45.00** | **83.00** | 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 | **53.23** | 0.13 | 67.66 | 23.88 | **9.14** | 36.48 | **42.67** | 28.76 | 41.09 ±0.17 |
| | | PØDA | **75.40** | **34.61** | **75.22** | **26.77** | 27.34 | 35.20 | 52.68 | 44.37 | 82.01 | **14.16** | **93.72** | 50.51 | **0.99** | **69.11** | **26.64** | 2.72 | **46.98** | 42.64 | **33.09** | **43.90** ±0.53 |
| | | | | | | | | | TrgPrompt = "driving under rain" | | | | | | | | | | | | | |
| | ACDC rain | source-only | 74.10 | 31.98 | 63.07 | **15.08** | **23.92** | **41.31** | **50.12** | **44.43** | 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 | **9.65** | **74.54** | 13.16 | **20.33** | 12.73 | 14.06 | 18.26 | 37.17 ±0.10 |
| | | PØDA | **76.60** | **38.52** | **78.01** | 15.02 | 22.53 | 40.33 | 45.39 | 41.40 | **86.85** | **37.97** | **96.46** | 50.39 | 6.35 | 74.19 | **19.19** | 7.98 | **22.06** | 21.04 | 23.65 | **42.31** ±0.55 |
| | | | | | | | | | TrgPrompt = "driving in a game" | | | | | | | | | | | | | |
| | GTA5 | source-only | 68.72 | 22.65 | 78.79 | 36.81 | **17.31** | 39.66 | **39.33** | 14.84 | **72.61** | 22.53 | 87.31 | 57.50 | 26.14 | 74.29 | **44.57** | **20.45** | 0.00 | 18.30 | 10.35 | 39.59 |
| | | CLIPstyler | 73.06 | **29.89** | 77.86 | 25.50 | 11.69 | **39.72** | 35.88 | **24.04** | 67.38 | 12.75 | 88.77 | 46.58 | 33.38 | 72.03 | 42.79 | 11.12 | 0.00 | 28.84 | **14.61** | 38.73 ±0.16 |
| | | PØDA | **73.93** | 22.69 | **78.82** | **37.52** | 14.17 | 36.97 | 33.14 | 17.34 | 72.44 | **26.22** | **88.85** | **62.69** | **37.04** | **74.33** | 43.03 | 11.91 | 0.00 | **35.33** | 13.91 | **41.07** ±0.48 |
| | | | | | | | | | TrgPrompt = "driving" | | | | | | | | | | | | | |
| GTA5 | CS | source-only | 58.97 | 20.92 | 72.84 | 16.53 | **24.58** | 31.37 | 34.77 | 23.62 | 82.12 | 17.04 | 66.28 | **63.46** | 14.72 | **81.27** | **20.83** | 17.19 | 4.68 | **20.57** | 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 | **81.04** | 54.64 | 0.52 | 73.47 | 20.65 | **22.30** | 4.03 | 15.79 | 10.73 | 31.50 ±0.21 |
| | | PØDA | **84.34** | **36.73** | **79.43** | **18.33** | 16.54 | **36.93** | **38.45** | **33.81** | **82.44** | **19.14** | 75.90 | 62.65 | **16.47** | 75.48 | 15.68 | 19.57 | **11.28** | 16.53 | **21.76** | **40.08** ±0.52 |

# Limitations ?

- Global stylization: large structural classes benefit more from it

| 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 | **43.83** | 6.11 | 17.08 | 23.52 | **24.51** | 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 | **35.49** | 23.34 | 19.01 | **45.41** | 8.81 | **27.87** | 21.06 | 8.48 | 38.17 | 1.84 | 0.00 | 11.54 | **10.38** | 4.89 | 21.38 ±0.36 |
| | | PØDA | **77.54** | **26.90** | 42.71 | **13.51** | **21.36** | 33.52 | 23.70 | **21.73** | 39.91 | **9.51** | 19.40 | **28.80** | **11.85** | **50.89** | **10.14** | 0.00 | **20.76** | 8.76 | **14.50** | **25.03** ±0.48 |
| | | | | | | | | | | TrgPrompt = "driving in snow" | | | | | | | | | | | |
| | ACDC snow | source-only | 70.47 | 23.50 | 63.80 | 17.96 | **27.36** | **38.52** | **56.26** | **45.00** | **83.00** | 10.75 | 83.65 | 47.73 | 0.72 | 61.42 | 21.87 | 5.90 | 21.58 | 35.83 | 31.01 | 39.28 |
| CS | | CLIPstyler | 74.29 | 31.25 | 69.17 | 15.21 | 25.21 | 36.83 | 44.79 | 42.56 | 76.87 | 11.07 | 91.48 | **53.23** | 0.13 | 67.66 | 23.88 | **9.14** | 36.48 | **42.67** | 28.76 | 41.09 ±0.17 |
| | | PØDA | **75.40** | **34.61** | **75.22** | **26.77** | 27.34 | 35.20 | 52.68 | 44.37 | 82.01 | **14.16** | **93.72** | 50.51 | **0.99** | **69.11** | **26.64** | 2.72 | **46.98** | 42.64 | **33.09** | **43.90** ±0.53 |
| | | | | | | | | | | TrgPrompt = "driving under rain" | | | | | | | | | | | |
| | ACDC rain | source-only | 74.10 | 31.98 | 63.07 | **15.08** | **23.92** | **41.31** | **50.12** | **44.43** | 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 | **9.65** | **74.54** | 13.16 | **20.33** | 12.73 | 14.06 | 18.26 | 37.17 ±0.10 |
| | | PØDA | **76.60** | **38.52** | **78.01** | 15.02 | 22.53 | 40.33 | 45.39 | 41.40 | **86.85** | **37.97** | **96.46** | **50.39** | 6.35 | 74.19 | **19.19** | 7.98 | **22.06** | **21.04** | **23.65** | **42.31** ±0.55 |
| | | | | | | | | | | TrgPrompt = "driving in a game" | | | | | | | | | | | |
| | GTA5 | source-only | 68.72 | 22.65 | 78.79 | 36.81 | **17.31** | 39.66 | **39.33** | 14.84 | **72.61** | 22.53 | 87.31 | 57.50 | 26.14 | 74.29 | **44.57** | **20.45** | 0.00 | 18.30 | 10.35 | 39.59 |
| | | CLIPstyler | 73.06 | **29.89** | 77.86 | 25.50 | 11.69 | **39.72** | 35.88 | **24.04** | 67.38 | 12.75 | 88.77 | 46.58 | 33.38 | 72.03 | 42.79 | 11.12 | 0.00 | 28.84 | **14.61** | 38.73 ±0.16 |
| | | PØDA | **73.93** | 22.69 | **78.82** | **37.52** | 14.17 | 36.97 | 33.14 | 17.34 | 72.44 | **26.22** | **88.85** | **62.69** | **37.04** | **74.33** | 43.03 | 11.91 | 0.00 | **35.33** | 13.91 | **41.07** ±0.48 |
| | | | | | | | | | | TrgPrompt = "driving" | | | | | | | | | | | |
| | | source-only | 58.97 | 20.92 | 72.84 | 16.53 | **24.58** | 31.37 | 34.77 | 23.62 | 82.12 | 17.04 | 66.28 | **63.46** | 14.72 | **81.27** | **20.83** | 17.19 | 4.68 | **20.57** | 19.56 | 36.38 |
| GTA5 | CS | CLIPstyler | 66.70 | 23.63 | 64.12 | 5.08 | 3.66 | 20.67 | 19.31 | 18.10 | 81.68 | 12.36 | **81.04** | 54.64 | 0.52 | 73.47 | 20.65 | **22.30** | 4.03 | 15.79 | 10.73 | 31.50 ±0.21 |
| | | PØDA | **84.34** | **36.73** | **79.43** | **18.33** | 16.54 | **36.93** | **38.45** | **33.81** | **82.44** | **19.14** | 75.90 | 62.65 | **16.47** | 75.48 | 15.68 | 19.57 | **11.28** | 16.53 | **21.76** | **40.08** ±0.52 |

76

# Limitations ?

- Global stylization: large structural classes benefit more from it
- Spatial-/Semantic- invariant transformation

|  | fog |
|---|---|
| Src only | 44.73 |
| **Foggy driving** | 43.50(+-0.17) |

Relevant

|  | fog |
|---|---|
| Src only | 44.73 |
| **Foggy driving** | 43.50(+-0.17) |
| Driving in low visibility | 44.40(+-0.17) |
| Operating a vehicle in thick fog | 44.13(+-0.14) |
| Navigating through a foggy environment | 43.68(+-0.07) |
| Driving in a dense fog | 44.37(+-0.12) |
| Piloting a car when visibility is limited due to fog | 44.16(+-0.18) |

Relevant

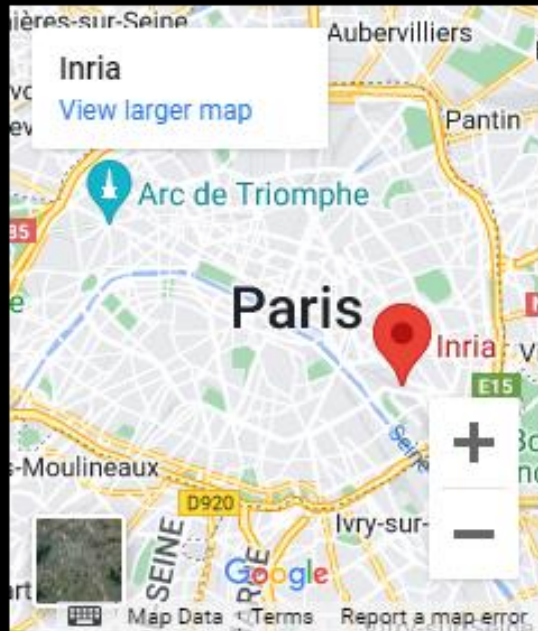|  | fog |
|---|---|
| Src only | 44.73 |
| **Foggy driving** | 43.50(+-0.17) |
| Driving in low visibility | 44.40(+-0.17) |
| Operating a vehicle in thick fog | 44.13(+-0.14) |
| Navigating through a foggy environment | 43.68(+-0.07) |
| Driving in a dense fog | 44.37(+-0.12) |
| Piloting a car when visibility is limited due to fog | 44.16(+-0.18) |

Relevant

| "Mesmerizing Northern Lights display" | 41.95(+-0.30) |
|---|---|
| "Adorable baby's first steps" | **45.08(**+-0.15) |
| "Intense athlete mid-competition" | **45.79**(+-0.16) |
| "Playful dolphins in the ocean" | 41.72(+-0.23) |
| "Breathtaking view from mountaintop" | 42.88(+-0.22) |

Irrelevant

78

# Astra-Vision

Inria valeo.ai

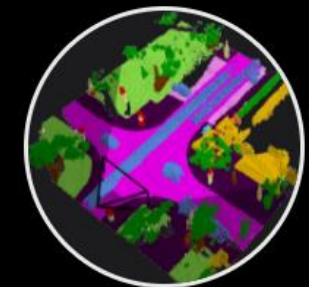## 2D/3D Robust Scene Understanding

**Inria Paris**

Learning with less supervision

Vision in complex conditions

3D scene understanding

Regular job openings.

astra-vision.github.io

**DenseMTL**, in *WACV 23.* Lopes, Vu, de Charette
github.com/astra-vision/DenseMTL

**xMUDA**, in *TPAMI 22 & CVPR 20.* Jaritz, Vu, de Charette, Wirbel, Pérez.
https://github.com/valeoai/xmuda_journal

**PØDA**, in ICCV 23. Fahes, Vu, Bursuc, Pérez, de Charette
github.com/astra-vision/PODA

# Any Questions ?