

# Domain Adaptation on Wheels: Closing the Gap to the Open-world

**Tuan-Hung Vu**

*Research scientist*

*valeo.ai*

**Dengxin Dai**

*Director of Research*

*Huawei Zurich Research Center*

valeo.ai



ICCV23  
PARIS

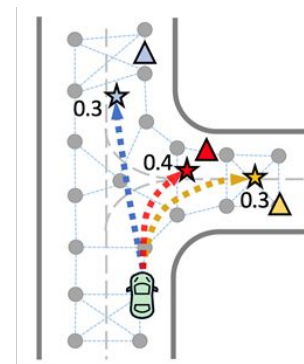
# Main Perception Tasks for Autonomous Driving



Image Semantic Segmentation



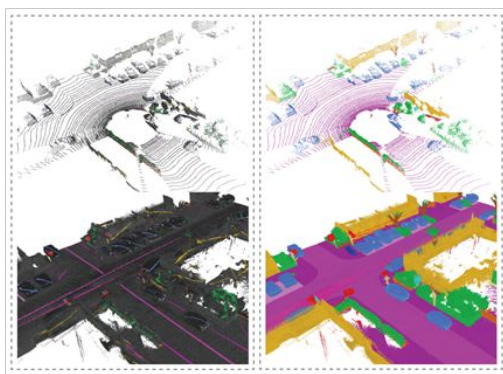
Depth Estimation



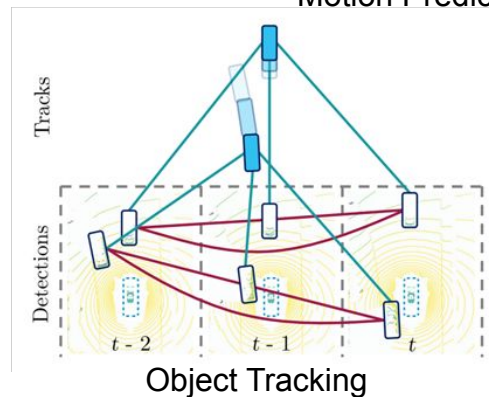
Motion Prediction



3D Object Detection



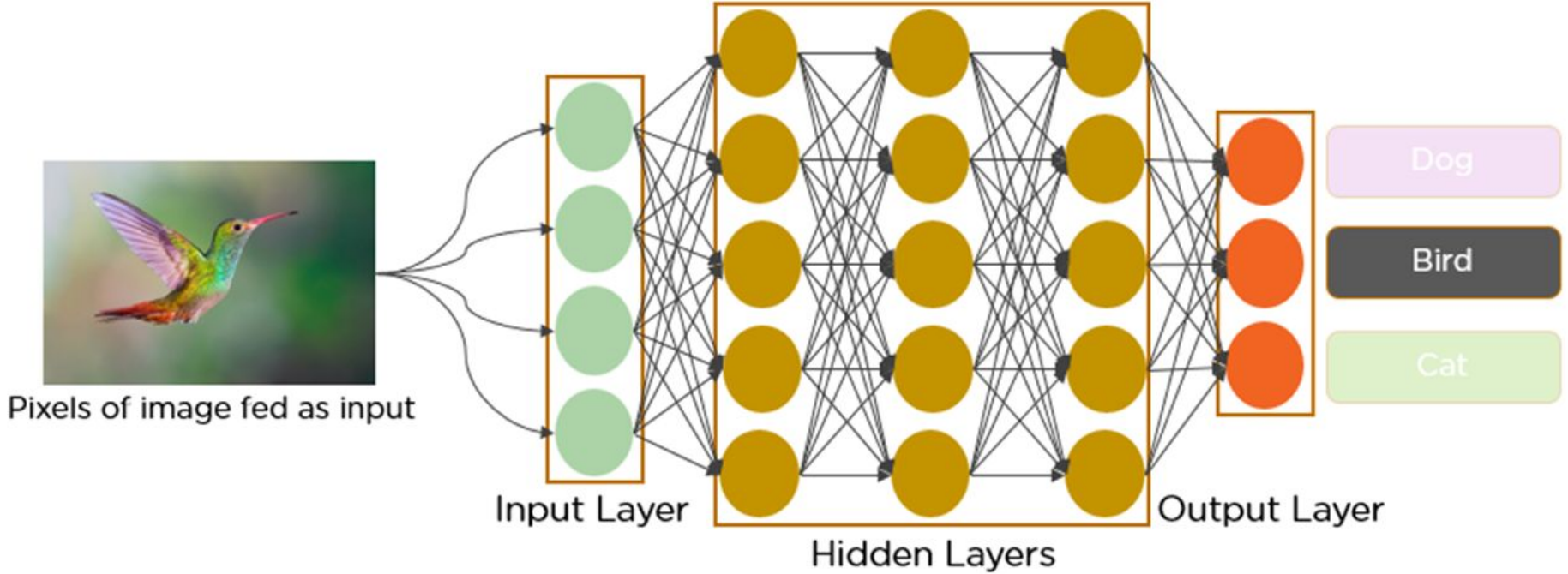
LiDAR Semantic Segmentation



Object Tracking

# Perception with Neural Networks

---



# ImageNet Classification

---



Image Classification on ImageNet

# ImageNet Classification

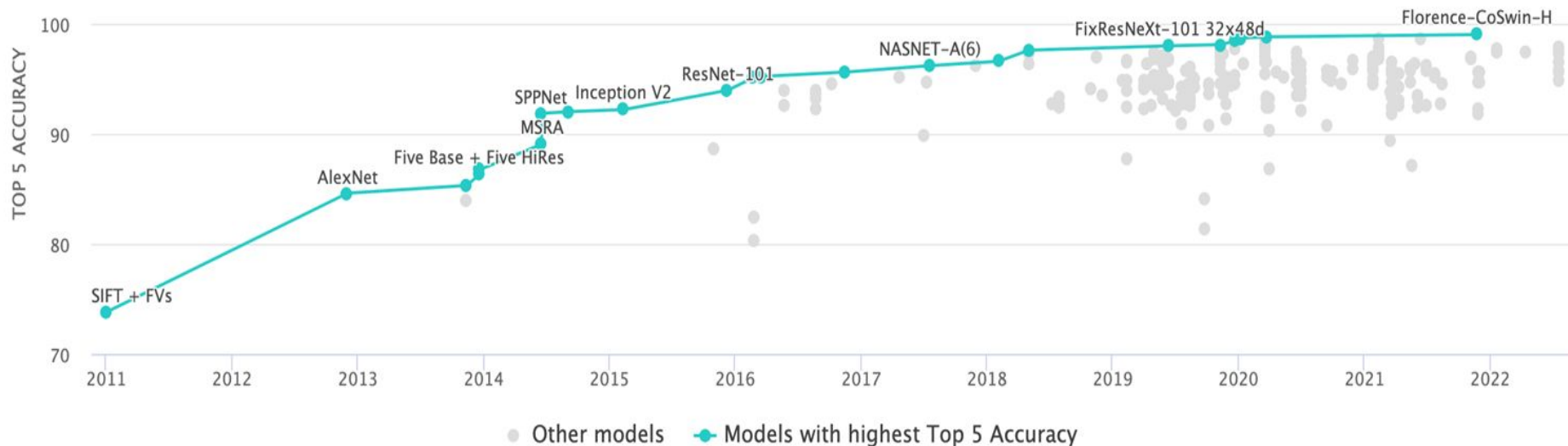
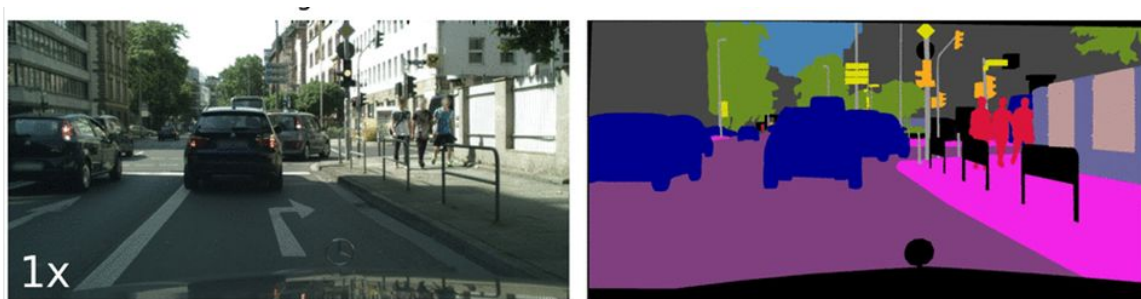


Image Classification on ImageNet

# Semantic Segmentation on Cityscapes Dataset



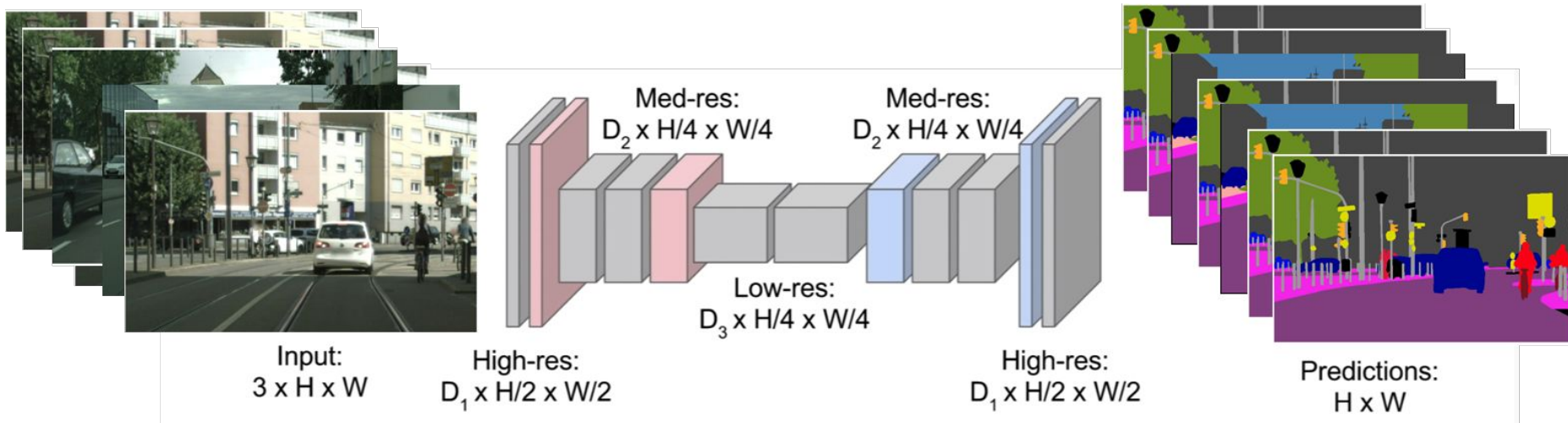
| name               | fine | coarse | 16-bit | depth | video | sub | IoU class | IoU class | IoU category | IoU category |
|--------------------|------|--------|--------|-------|-------|-----|-----------|-----------|--------------|--------------|
| LeapAI             | yes  | yes    | no     | no    | no    | no  | 86.4      | 70.9      | 93.2         | 84.2         |
| MYBank-AIoT        | yes  | yes    | no     | no    | no    | no  | 86.3      | 72.9      | 93.3         | 85.8         |
| SAIT SeeThroughNet | yes  | yes    | no     | no    | no    | no  | 86.2      | 71.5      | 93.2         | 85.7         |

Semantic Segmentation on Cityscapes

**Have we solved all  
perception tasks?**

# Semantic Segmentation: training and validation

---





# Dataset Bias or Domain Discrepancy

---



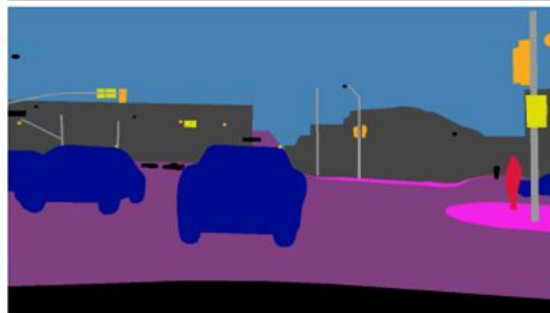
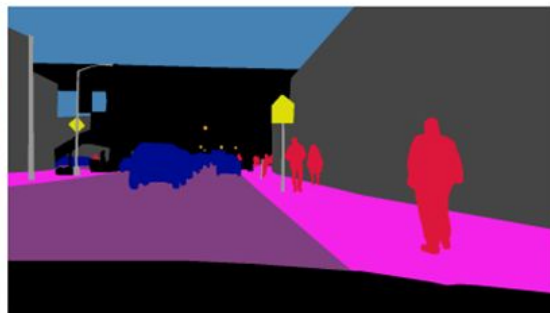
Clear weather

Rain

Physics-Based Rendering for Improving Robustness to Rain, Halder, Lalonde, and Charette, ICCV 2019

# Dataset Bias or Domain Discrepancy

---



Nighttime Image

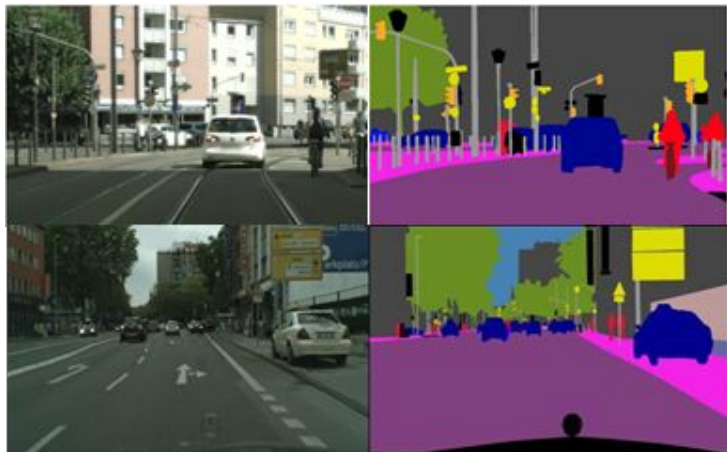
Human Annotation

Prediction

Map-Guided Curriculum Domain Adaptation and Uncertainty-Aware Evaluation for Semantic Nighttime Image Segmentation, Sakaridis, Dai, Van Gool, T-PAMI, 2020

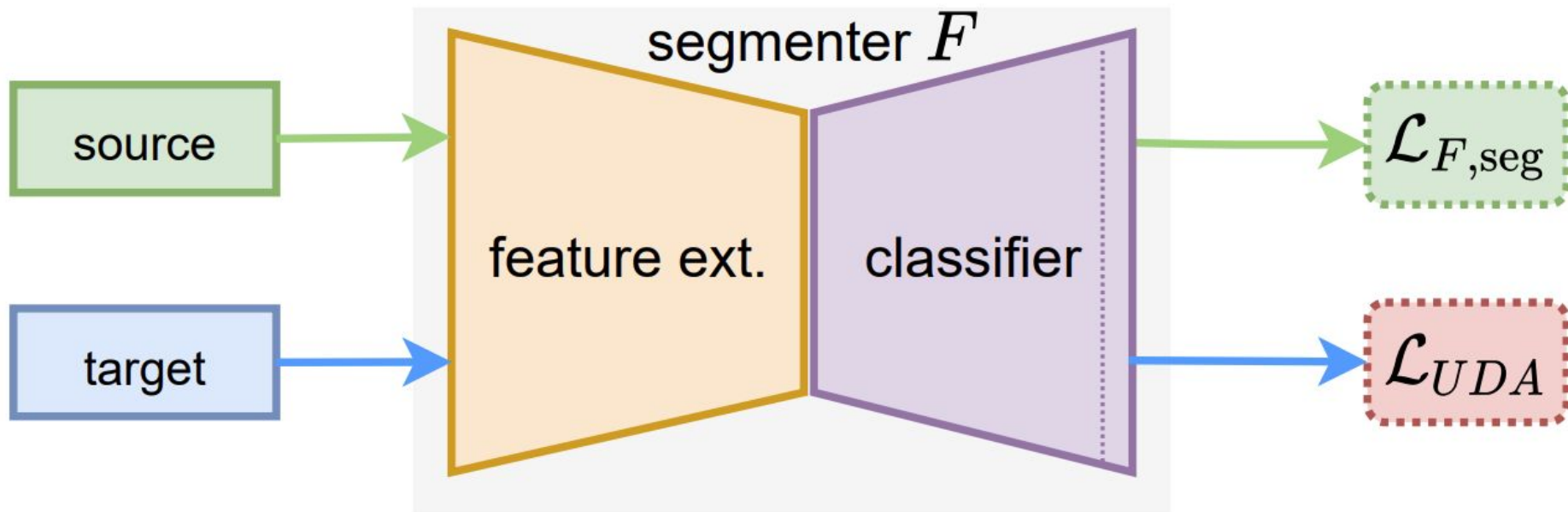
# What can we do to generalize?

## 1. Unsupervised Domain Adaptation: Learning Target Distribution with Unlabeled Samples



# UDA in Semantic Segmentation

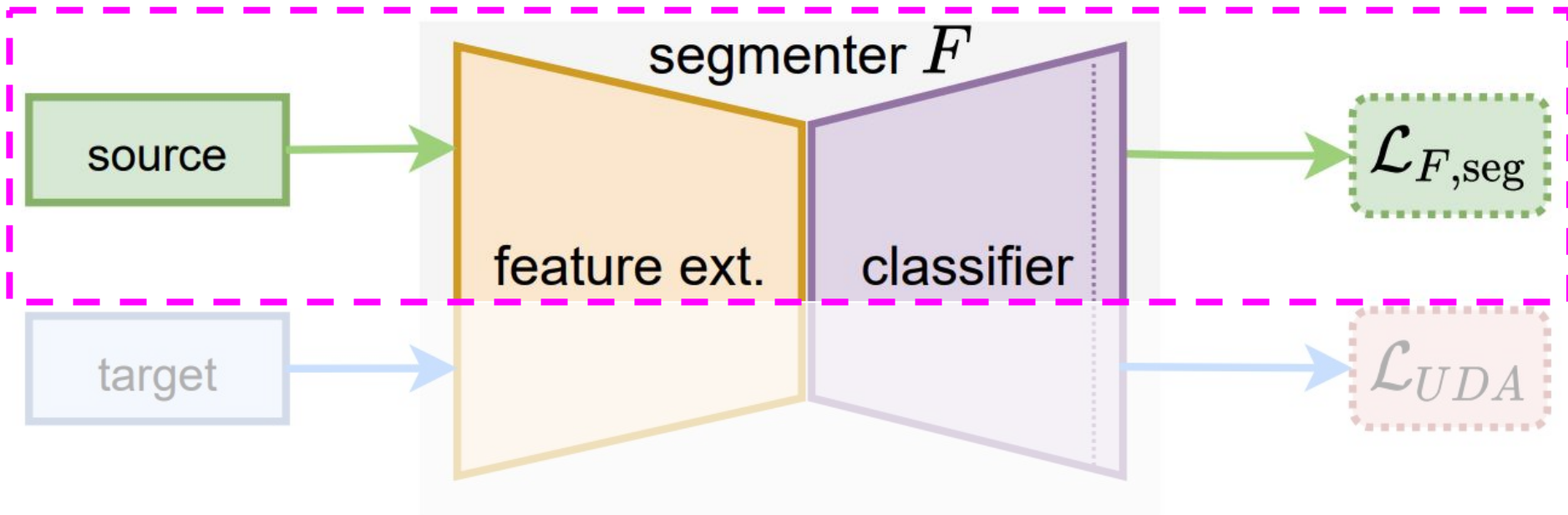
---



- A general UDA pipeline in segmentation

# UDA in Semantic Segmentation

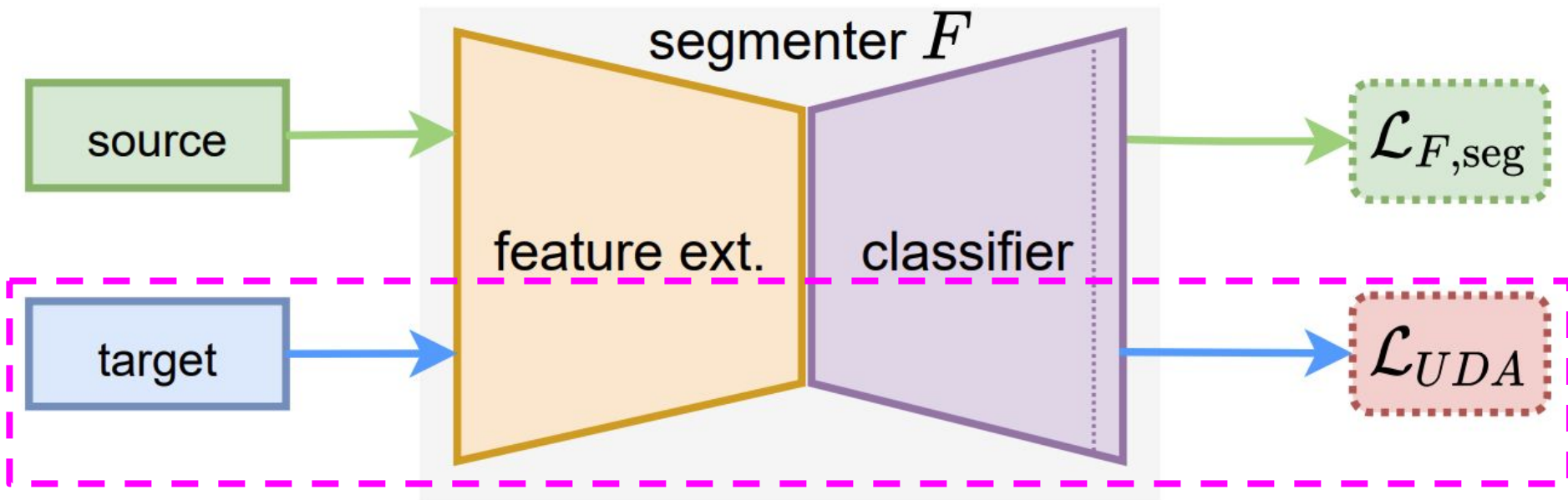
---



- Supervised training on source

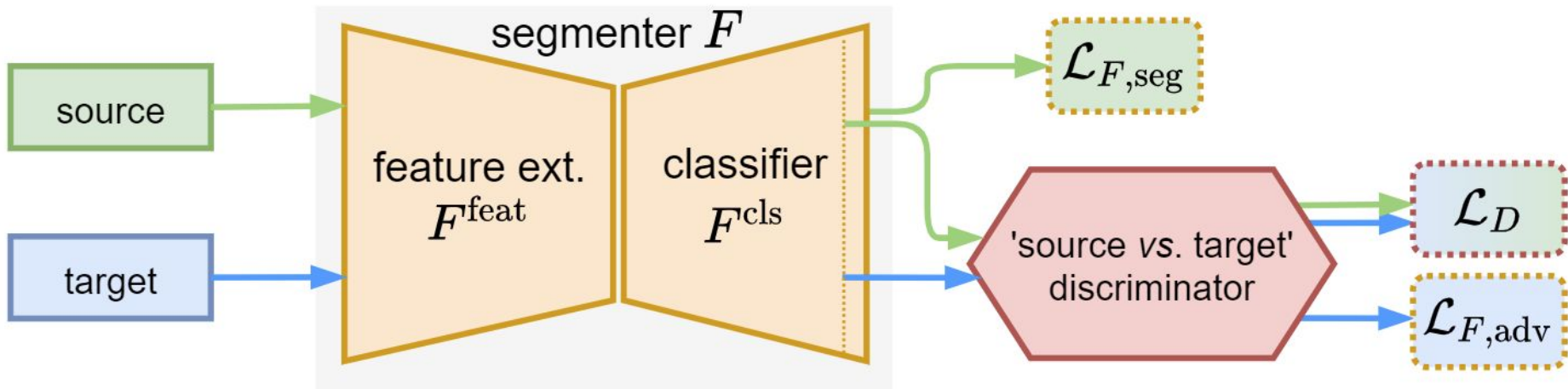
# UDA in Semantic Segmentation

---



- Different UDA techniques ~ different UDA losses

# Adversarial UDA framework in Segmentation

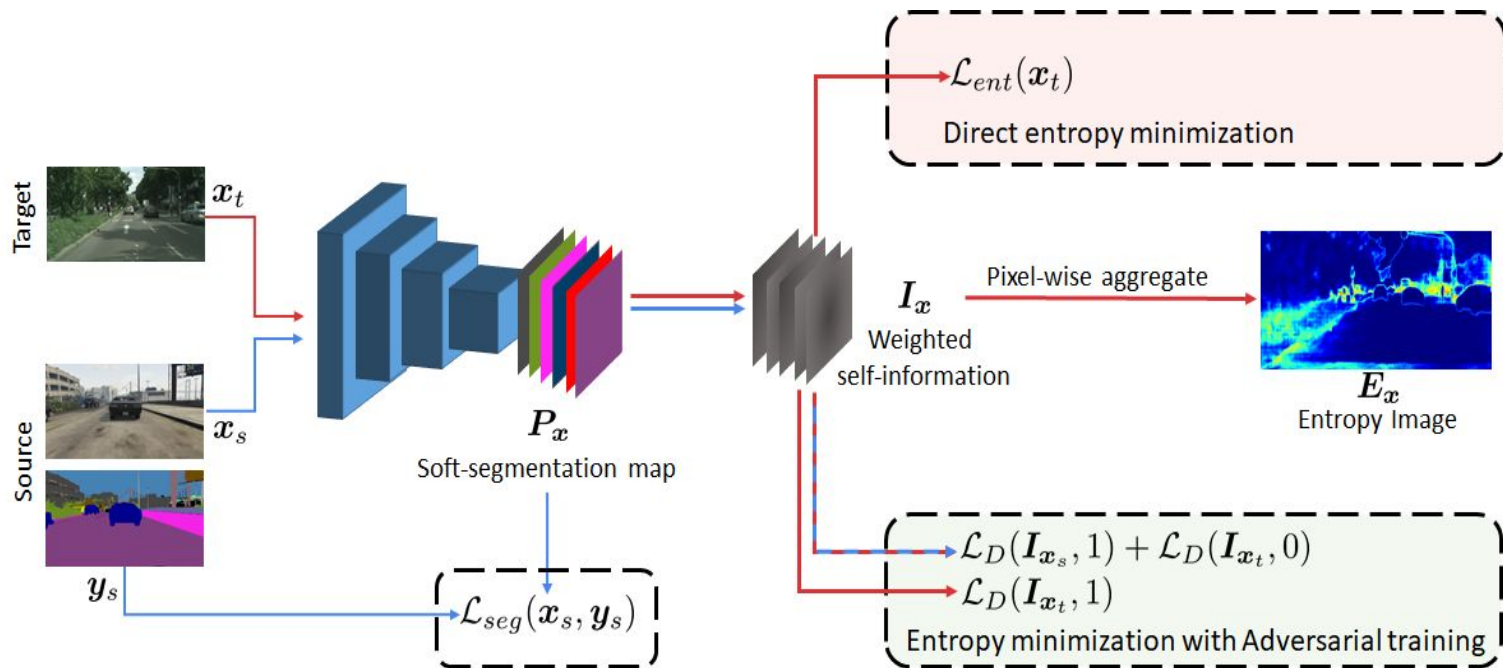


FCNs in the Wild: Pixel-level Adversarial and Constraint-based Adaptation, Hoffman et al. ICLR'17

Learning to Adapt Structured Output Space for Semantic Segmentation, Tsai et al. CVPR'18

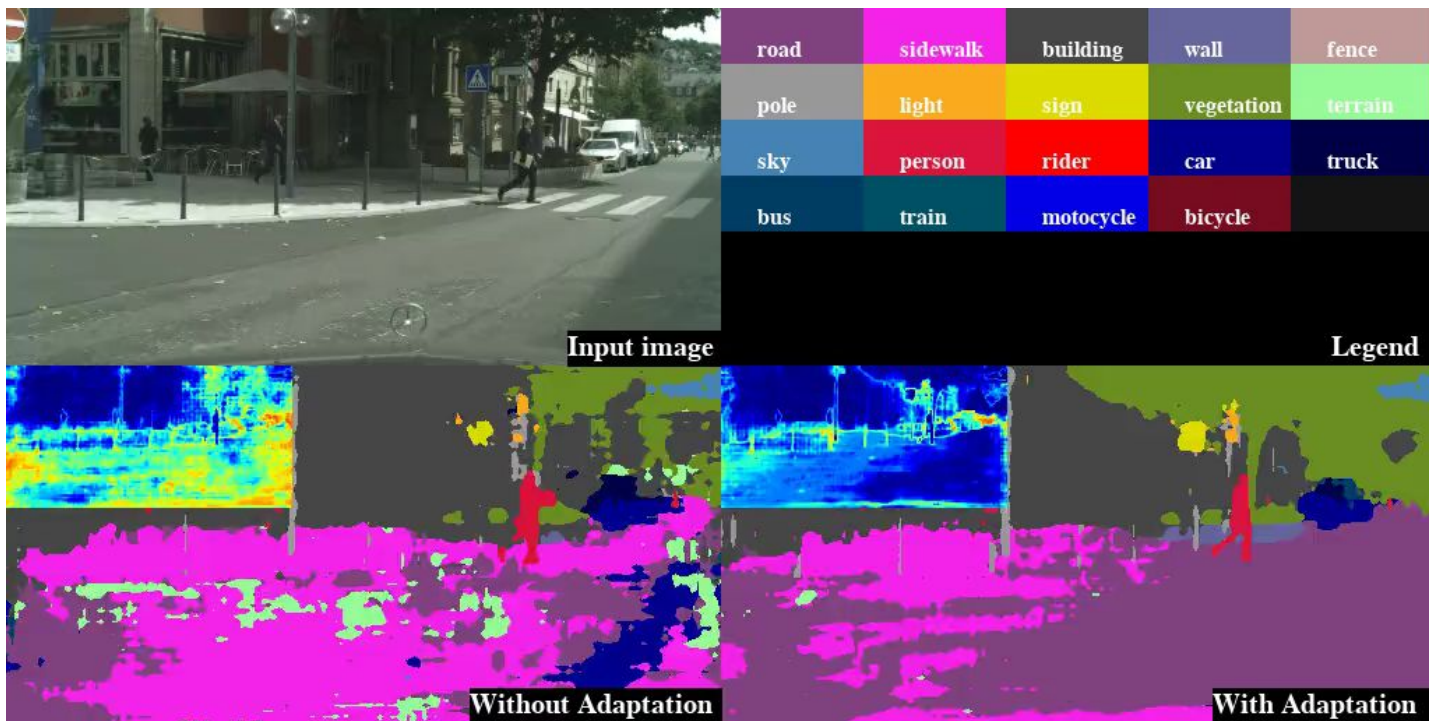
ADVENT: Adversarial Entropy Minimization for Domain Adaptation in Semantic Segmentation, Vu et al. CVPR'19

# ADVENT: adversarial UDA + entropy minimization





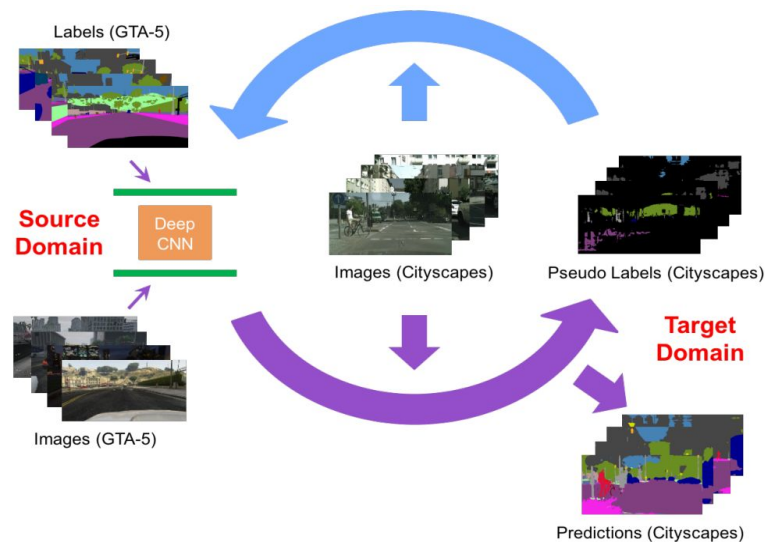
# ADVENT



ADVENT: Adversarial Entropy Minimization for Domain Adaptation in Semantic Segmentation, Vu et al. CVPR'19

# What did we learn?

- Adversarial training is great but difficult to train
- Self-training with entropy minimization works
  - ▶ Similar finding in other works
  - ▶ Self-training with pseudo-labelling
    - ▶ High-scoring predictions
    - ▶ Training with noisy labels





# Self-training for UDA

- ESL: Entropy-based criterion for pseudo-labelling
  - ▶ Low-entropy predictions as pseudo-labels
- ConDA: learnable confidence network for semantic failure detection
  - ▶ High confidence predictions as pseudo-labels

GTA5  $\triangleright$  Cityscapes

| Method           | Self-Train. | road        | sidewalk    | building    | wall        | fence       | pole        | light       | sign        | veg         | terrain     | sky         | person      | rider       | car         | truck       | bus         | train       | mbike       | bike        | mIoU        |
|------------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| AdaptSegNet [50] |             | 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         | <b>29.5</b> | 32.5        | 41.4        |
| CyCADA [49]      |             | 86.7        | 35.6        | 80.1        | 19.8        | 17.5        | <b>38.0</b> | <b>39.9</b> | <b>41.5</b> | 82.7        | 27.9        | 73.6        | <b>64.9</b> | 19.0        | 65.0        | 12.0        | 28.6        | 4.5         | 31.1        | 42.0        | 42.7        |
| DISE [64]        |             | 91.5        | 47.5        | 82.5        | 31.3        | 25.6        | 33.0        | 33.7        | 25.8        | 82.7        | 28.8        | 82.7        | 62.4        | 30.8        | 85.2        | 27.7        | 34.5        | 6.4         | 25.2        | 24.4        | 45.4        |
| AdvEnt [51]      |             | 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        |
| CBST [54]        | ✓           | 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        | <b>42.8</b> | 45.9        |
| MRKLD [55]       | ✓           | 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        | <b>26.9</b> | 26.0        | 42.3        | 47.1        |
| BDL [21]         | ✓           | 91.0        | 44.7        | 84.2        | 34.6        | <b>27.5</b> | 30.2        | 36.0        | 36.0        | 85.0        | <b>43.6</b> | 83.0        | 58.6        | <b>31.6</b> | 83.3        | 35.3        | 49.7        | 3.3         | 28.8        | 35.6        | 48.5        |
| ESL [53]         | ✓           | 90.2        | 43.9        | 84.7        | 35.9        | 28.5        | 31.2        | 37.9        | 34.0        | 84.5        | 42.2        | 83.9        | 59.0        | 32.2        | 81.8        | 36.7        | 49.4        | 1.8         | 30.6        | 34.1        | 48.6        |
| ConDA            | ✓           | <b>93.5</b> | <b>56.9</b> | <b>85.3</b> | <b>38.6</b> | 26.1        | 34.3        | 36.9        | 29.9        | <b>85.3</b> | 40.6        | <b>88.3</b> | 58.1        | 30.3        | <b>85.8</b> | <b>39.8</b> | <b>51.0</b> | 0.0         | 28.9        | 37.8        | <b>49.9</b> |

# Self-training for UDA

- Other self-training strategies:
  - Prototype-based pseudo-labelling: CAG\_UDA [Zheng et al. NeurIPS'19], ProDA [Zhang et al. CVPR'21]
  - Inspired by the success of prototype-based approach to deal with noisy data [Han et al. ICCV'19]
  - Prototypes treat different classes equally regardless of their occurrence frequency

|                      |      |             |             |             |             |             |             |             |             |             |             |             |             |             |             |             |             |             |             |             |
|----------------------|------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| 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        |
| 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        |
| 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        |
| 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        |
| 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        |
| 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        |
| 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        |
| <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> |

# What is missing?

---

SOTA methods still use “out-dated” network architectures and “low-res” input images

# Transformer for UDA

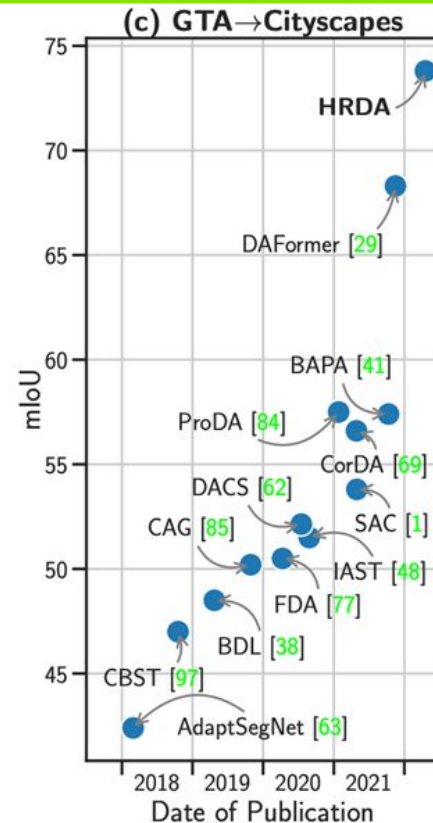
DAFormer: Improving Network Architectures and Training Strategies for Domain-Adaptive Semantic Segmentation”, Hoyer, Dai, and Van Gool, CVPR 2022

HRDA: Context-Aware High-Resolution Domain-Adaptive Semantic Segmentation”, Hoyer, Dai, and Van Gool, ECCV 2022

- Harness the robustness of SegFormer [Xie et al. NeurIPS 2021]

| Method       | Clean       | Blur        |             |             |             | Noise       |             |             |             | Digital     |             |             |             | Weather     |             |             |             |
|--------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
|              |             | Motion      | Defoc       | Glass       | Gauss       | Gauss       | Impul       | Shot        | Speck       | Bright      | Contr       | Satur       | JPEG        | Snow        | Spatt       | Fog         | Frost       |
| DLv3+ (MBv2) | 72.0        | 53.5        | 49.0        | 45.3        | 49.1        | 6.4         | 7.0         | 6.6         | 16.6        | 51.7        | 46.7        | 32.4        | 27.2        | 13.7        | 38.9        | 47.4        | 17.3        |
| DLv3+ (R50)  | 76.6        | 58.5        | 56.6        | 47.2        | 57.7        | 6.5         | 7.2         | 10.0        | 31.1        | 58.2        | 54.7        | 41.3        | 27.4        | 12.0        | 42.0        | 55.9        | 22.8        |
| DLv3+ (R101) | 77.1        | 59.1        | 56.3        | 47.7        | 57.3        | 13.2        | 13.9        | 16.3        | 36.9        | 59.2        | 54.5        | 41.5        | 37.4        | 11.9        | 47.8        | 55.1        | 22.7        |
| DLv3+ (X41)  | 77.8        | 61.6        | 54.9        | 51.0        | 54.7        | 17.0        | 17.3        | 21.6        | 43.7        | 63.6        | 56.9        | 51.7        | 38.5        | 18.2        | 46.6        | 57.6        | 20.6        |
| DLv3+ (X65)  | 78.4        | 63.9        | 59.1        | 52.8        | 59.2        | 15.0        | 10.6        | 19.8        | 42.4        | 65.9        | 59.1        | 46.1        | 31.4        | 19.3        | 50.7        | 63.6        | 23.8        |
| DLv3+ (X71)  | 78.6        | 64.1        | 60.9        | 52.0        | 60.4        | 14.9        | 10.8        | 19.4        | 41.2        | 68.0        | 58.7        | 47.1        | 40.2        | 18.8        | 50.4        | 64.1        | 20.2        |
| ICNet        | 65.9        | 45.8        | 44.6        | 47.4        | 44.7        | 8.4         | 8.4         | 10.6        | 27.9        | 41.0        | 33.1        | 27.5        | 34.0        | 6.3         | 30.5        | 27.3        | 11.0        |
| FCN8s        | 66.7        | 42.7        | 31.1        | 37.0        | 34.1        | 6.7         | 5.7         | 7.8         | 24.9        | 53.3        | 39.0        | 36.0        | 21.2        | 11.3        | 31.6        | 37.6        | 19.7        |
| DilatedNet   | 68.6        | 44.4        | 36.3        | 32.5        | 38.4        | 15.6        | 14.0        | 18.4        | 32.7        | 52.7        | 32.6        | 38.1        | 29.1        | 12.5        | 32.3        | 34.7        | 19.2        |
| ResNet-38    | 77.5        | 54.6        | 45.1        | 43.3        | 47.2        | 13.7        | 16.0        | 18.2        | 38.3        | 60.0        | 50.6        | 46.9        | 14.7        | 13.5        | 45.9        | 52.9        | 22.2        |
| PSPNet       | 78.8        | 59.8        | 53.2        | 44.4        | 53.9        | 11.0        | 15.4        | 15.4        | 34.2        | 60.4        | 51.8        | 30.6        | 21.4        | 8.4         | 42.7        | 34.4        | 16.2        |
| GSCNN        | 80.9        | 58.9        | 58.4        | 41.9        | 60.1        | 5.5         | 2.6         | 6.8         | 24.7        | 75.9        | 61.9        | 70.7        | 12.0        | 12.4        | 47.3        | 67.9        | 32.6        |
| SegFormer-B5 | <b>82.4</b> | <b>69.1</b> | <b>68.6</b> | <b>64.1</b> | <b>69.8</b> | <b>57.8</b> | <b>63.4</b> | <b>52.3</b> | <b>72.8</b> | <b>81.0</b> | <b>77.7</b> | <b>80.1</b> | <b>58.8</b> | <b>40.7</b> | <b>68.4</b> | <b>78.5</b> | <b>49.9</b> |

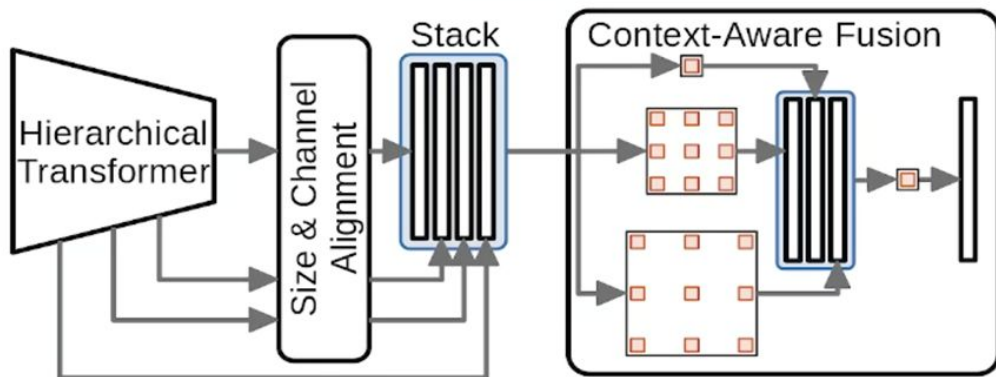
- Enable learning high-reso details and low-reso context at the same time



# Transformer for UDA - DAFormer

Design of an architecture tailored for UDA

- Hierarchical Transformer encoder [4]
- Context-aware multi-level feature fusion decoder



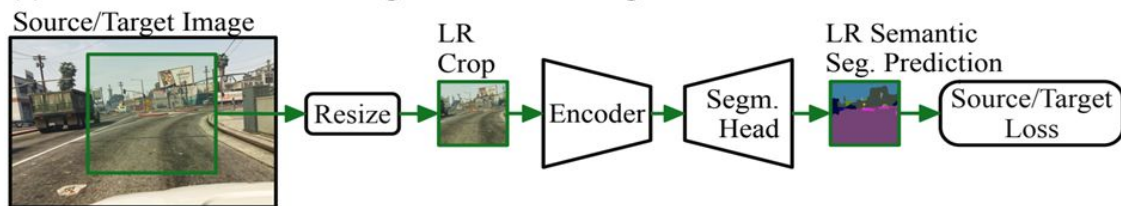
| Network Architecture | UDA  | Oracle | UDA / Oracle |
|----------------------|------|--------|--------------|
| DeepLabV2            | 56.0 | 72.1   | 77.7%        |
| DAFormer             | 68.3 | 77.6   | 88.0%        |

Reduced performance gap between UDA and supervised oracle

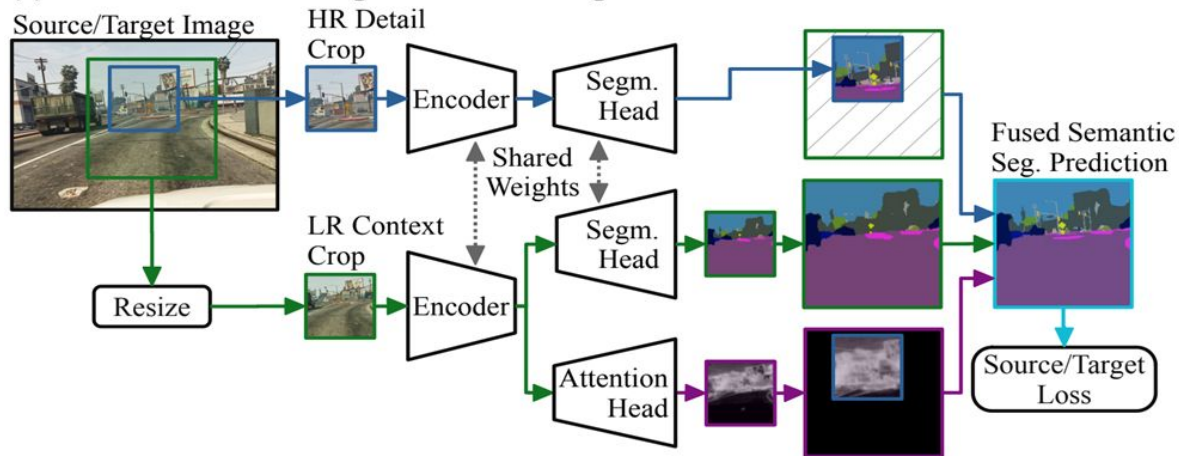


# Transformer for UDA - HRDA

(a) Previous UDA Semantic Segmentation Training



(b) Our HRDA Semantic Segmentation Training



# SoTA in 2023

| Method  | Road        | S.walk      | Build.      | Wall        | Fence       | Pole        | Tr.Light    | Sign        | Veget.      | Terrain     | Sky         | Person      | Rider       | Car         | Truck       | Bus         | Train       | M.bike      | Bike        | mIoU        |
|---|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| <b>Synthetic-to-Real: GTA→Cityscapes (Val.)</b>     |             |             |             |             |             |             |             |             |             |             |             |             |             |             |             |             |             |             |             |             |
| ADVENT [76]   | 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        |
| DACS [72]   | 89.9        | 39.7        | 87.9        | 30.7        | 39.5        | 38.5        | 46.4        | 52.8        | 88.0        | 44.0        | 88.8        | 67.2        | 35.8        | 84.5        | 45.7        | 50.2        | 0.0         | 27.3        | 34.0        | 52.1        |
| ProDA [89]  | 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        |
| DAFormer [30]                                       | 95.7        | 70.2        | 89.4        | 53.5        | 48.1        | 49.6        | 55.8        | 59.4        | 89.9        | 47.9        | 92.5        | 72.2        | 44.7        | 92.3        | 74.5        | 78.2        | 65.1        | 55.9        | 61.8        | 68.3        |
| HRDA [31]   | 96.4        | 74.4        | 91.0        | <b>61.6</b> | <u>51.5</u> | <u>57.1</u> | <u>63.9</u> | <u>69.3</u> | <u>91.3</u> | <u>48.4</u> | <u>94.2</u> | <u>79.0</u> | <u>52.9</u> | <u>93.9</u> | <u>84.1</u> | <u>85.7</u> | <u>75.9</u> | <u>63.9</u> | <u>67.5</u> | <u>73.8</u> |
| MIC (HRDA)  | <b>97.4</b> | <b>80.1</b> | <b>91.7</b> | <u>61.2</u> | <b>56.9</b> | <b>59.7</b> | <b>66.0</b> | <b>71.3</b> | <b>91.7</b> | <b>51.4</b> | <b>94.3</b> | <b>79.8</b> | <b>56.1</b> | <b>94.6</b> | <b>85.4</b> | <b>90.3</b> | <b>80.4</b> | <b>64.5</b> | <b>68.5</b> | <b>75.9</b> |
| <b>Synthetic-to-Real: Synthia→Cityscapes (Val.)</b> |             |             |             |             |             |             |             |             |             |             |             |             |             |             |             |             |             |             |             |             |
| ADVENT [76]   | 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        |
| DACS [72]   | 80.6        | 25.1        | 81.9        | 21.5        | 2.9         | 37.2        | 22.7        | 24.0        | 83.7        | –           | 90.8        | 67.6        | 38.3        | 82.9        | –           | 38.9        | –           | 28.5        | 47.6        | 48.3        |
| ProDA [89]  | <b>87.8</b> | 45.7        | 84.6        | 37.1        | 0.6         | 44.0        | 54.6        | 37.0        | <b>88.1</b> | –           | 84.4        | 74.2        | 24.3        | 88.2        | –           | 51.1        | –           | 40.5        | 45.6        | 55.5        |
| DAFormer [30]                                       | 84.5        | 40.7        | 88.4        | 41.5        | <u>6.5</u>  | 50.0        | 55.0        | 54.6        | 86.0        | –           | 89.8        | 73.2        | 48.2        | 87.2        | –           | 53.2        | –           | 53.9        | 61.7        | 60.9        |
| HRDA [31]   | 85.2        | <u>47.7</u> | <u>88.8</u> | <b>49.5</b> | 4.8         | <u>57.2</u> | <u>65.7</u> | <u>60.9</u> | 85.3        | –           | <u>92.9</u> | <u>79.4</u> | <u>52.8</u> | <u>89.0</u> | –           | <b>64.7</b> | –           | <u>63.9</u> | <b>64.9</b> | <u>65.8</u> |
| MIC (HRDA)  | <u>86.6</u> | <b>50.5</b> | <b>89.3</b> | <u>47.9</u> | <b>7.8</b>  | <b>59.4</b> | <b>66.7</b> | <b>63.4</b> | <u>87.1</u> | –           | <b>94.6</b> | <b>81.0</b> | <b>58.9</b> | <b>90.1</b> | –           | <u>61.9</u> | –           | <b>67.1</b> | <u>64.3</u> | <b>67.3</b> |

# SoTA in 2023

Day-to-Nighttime: Cityscapes→DarkZurich (Test)

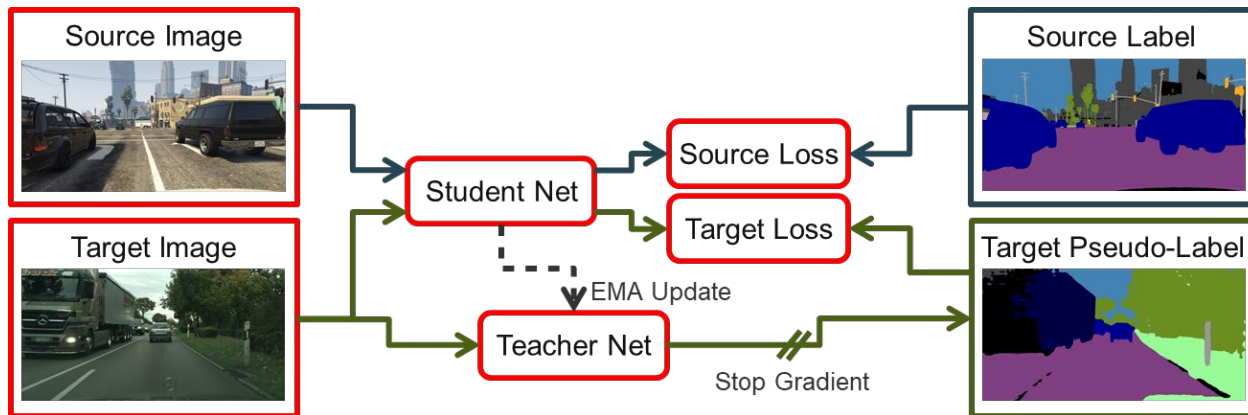
|                          |             |             |             |             |             |             |             |             |             |             |             |             |             |             |             |             |             |             |             |             |
|--------------------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| ADVENT [87]              | 85.8        | 37.9        | 55.5        | 27.7        | 14.5        | 23.1        | 14.0        | 21.1        | 32.1        | 8.7         | 2.0         | 39.9        | 16.6        | 64.0        | 13.8        | 0.0         | 58.8        | 28.5        | 20.7        | 29.7        |
| MGCDA <sup>†</sup> [76]  | 80.3        | 49.3        | 66.2        | 7.8         | 11.0        | 41.4        | 38.9        | 39.0        | <u>64.1</u> | 18.0        | 55.8        | 52.1        | 53.5        | 74.7        | <u>66.0</u> | 0.0         | 37.5        | 29.1        | 22.7        | 42.5        |
| DANNet <sup>†</sup> [92] | 90.0        | 54.0        | <u>74.8</u> | <u>41.0</u> | <u>21.1</u> | 25.0        | 26.8        | 30.2        | <b>72.0</b> | 26.2        | <b>84.0</b> | 47.0        | 33.9        | 68.2        | 19.0        | 0.3         | 66.4        | 38.3        | 23.6        | 44.3        |
| DAFormer [32]            | <u>93.5</u> | <u>65.5</u> | 73.3        | 39.4        | 19.2        | 53.3        | <u>44.1</u> | <u>44.0</u> | 59.5        | <u>34.5</u> | 66.6        | 53.4        | 52.7        | <u>82.1</u> | 52.7        | 9.5         | 89.3        | 50.5        | 38.5        | 53.8        |
| HRDA [33]                | 90.4        | 56.3        | 72.0        | 39.5        | 19.5        | <u>57.8</u> | <b>52.7</b> | 43.1        | 59.3        | 29.1        | <u>70.5</u> | <u>60.0</u> | <u>58.6</u> | <b>84.0</b> | <b>75.5</b> | <u>11.2</u> | <u>90.5</u> | <u>51.6</u> | <u>40.9</u> | <u>55.9</u> |
| MIC (HRDA)               | <b>94.8</b> | <b>75.0</b> | <b>84.0</b> | <b>55.1</b> | <b>28.4</b> | <b>62.0</b> | 35.5        | <b>52.6</b> | 59.2        | <b>46.8</b> | <u>70.0</u> | <b>65.2</b> | <b>61.7</b> | <u>82.1</u> | 64.2        | <b>18.5</b> | <b>91.3</b> | <b>52.6</b> | <b>44.0</b> | <b>60.2</b> |

Clear-to-Adverse-Weather: Cityscapes→ACDC (Test)

|                          |             |             |             |             |             |             |             |             |             |             |             |             |             |             |             |             |             |             |             |             |
|--------------------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| ADVENT [87]              | 72.9        | 14.3        | 40.5        | 16.6        | 21.2        | 9.3         | 17.4        | 21.2        | 63.8        | 23.8        | 18.3        | 32.6        | 19.5        | 69.5        | 36.2        | 34.5        | 46.2        | 26.9        | 36.1        | 32.7        |
| MGCDA <sup>†</sup> [76]  | 73.4        | 28.7        | 69.9        | 19.3        | 26.3        | 36.8        | 53.0        | 53.3        | <u>75.4</u> | 32.0        | 84.6        | 51.0        | 26.1        | 77.6        | 43.2        | 45.9        | 53.9        | 32.7        | 41.5        | 48.7        |
| DANNet <sup>†</sup> [92] | 84.3        | 54.2        | 77.6        | 38.0        | 30.0        | 18.9        | 41.6        | 35.2        | 71.3        | 39.4        | <u>86.6</u> | 48.7        | 29.2        | 76.2        | 41.6        | 43.0        | 58.6        | 32.6        | 43.9        | 50.0        |
| DAFormer [32]            | 58.4        | 51.3        | 84.0        | 42.7        | 35.1        | 50.7        | 30.0        | 57.0        | 74.8        | 52.8        | 51.3        | 58.3        | 32.6        | 82.7        | 58.3        | 54.9        | 82.4        | 44.1        | 50.7        | 55.4        |
| HRDA [33]                | <u>88.3</u> | <u>57.9</u> | <u>88.1</u> | <b>55.2</b> | <u>36.7</u> | <u>56.3</u> | <b>62.9</b> | <u>65.3</u> | 74.2        | <u>57.7</u> | 85.9        | <u>68.8</u> | <u>45.7</u> | <u>88.5</u> | <b>76.4</b> | <u>82.4</u> | <u>87.7</u> | <u>52.7</u> | <u>60.4</u> | <u>68.0</u> |
| MIC (HRDA)               | <b>90.8</b> | <b>67.1</b> | <b>89.2</b> | <u>54.5</u> | <b>40.5</b> | <b>57.2</b> | <u>62.0</u> | <b>68.4</b> | <b>76.3</b> | <b>61.8</b> | <b>87.0</b> | <b>71.3</b> | <b>49.4</b> | <b>89.7</b> | <u>75.7</u> | <b>86.8</b> | <b>89.1</b> | <b>56.9</b> | <b>63.0</b> | <b>70.4</b> |

# What did we learn?

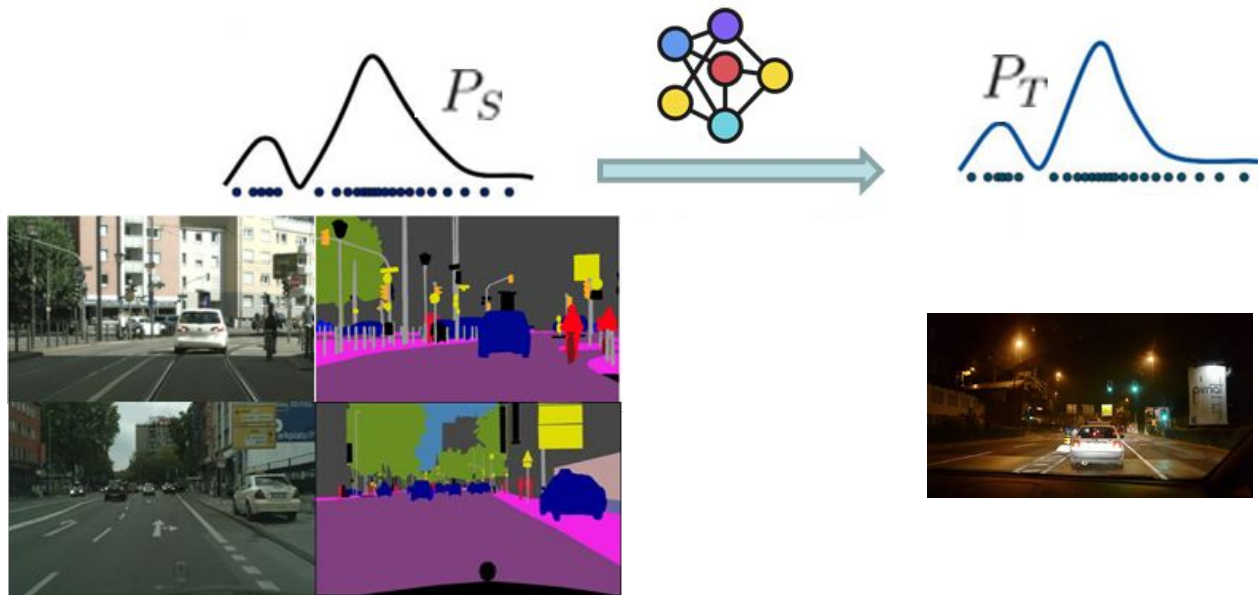
- Self-training



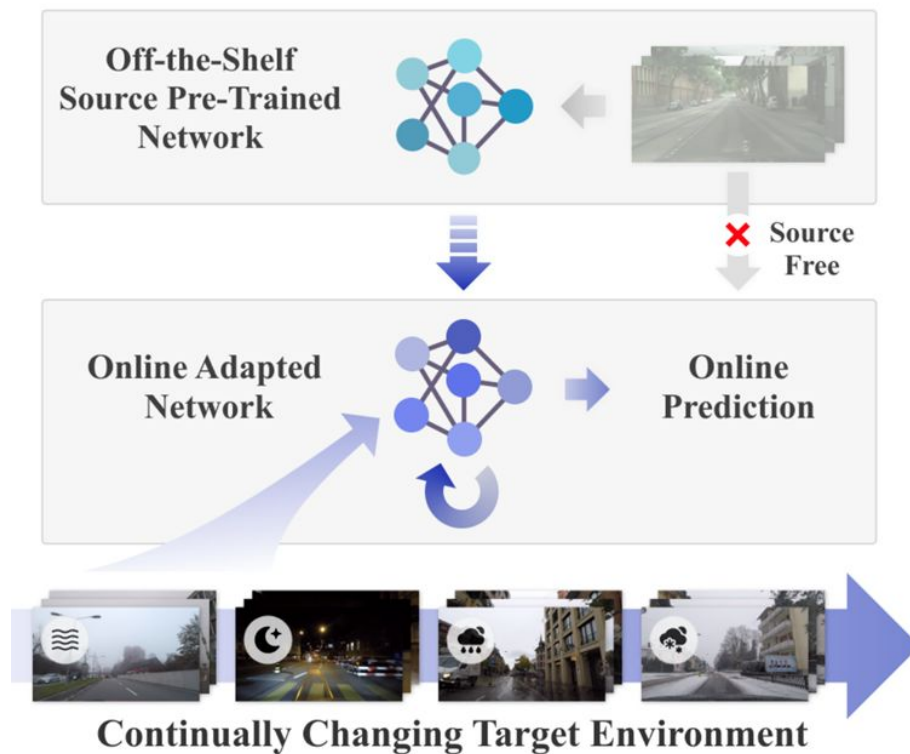
- Robust encoder architecture, e.g. SegFormer
- High-resolution recognition, e.g. HRDA

# What can we do to generalize?

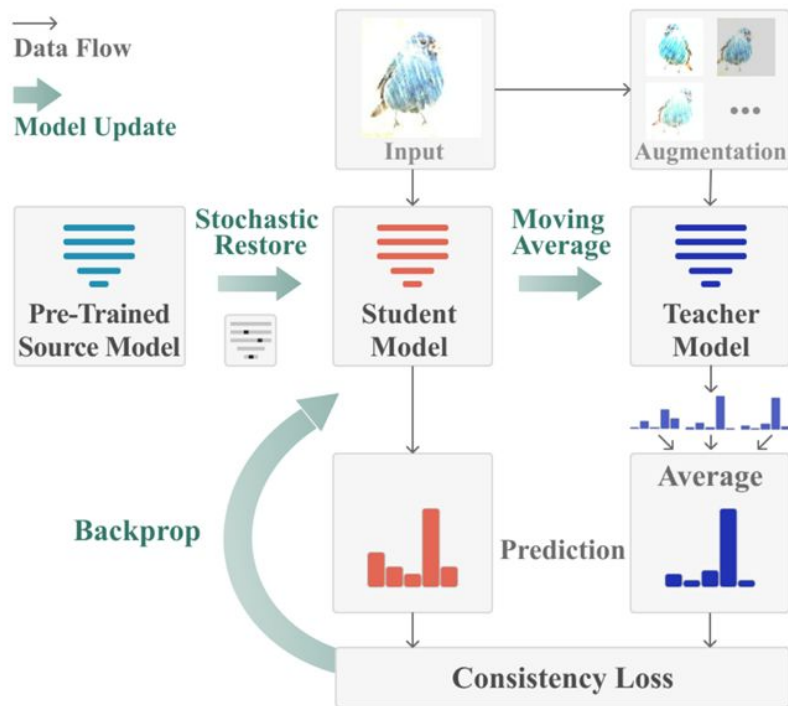
1. Unsupervised Domain Adaptation: Learning Target Distribution with Unlabeled Samples
2. Test-time Adaptation: Learning Target Distribution at Test Time from a Single Sample



# Continual Test-Time Domain Adaptation



# Continual Test-Time Domain Adaptation



- Self-training with (better) predictions by a **teacher network**
- Self-training with (better) **Augmentation-Averaged Pseudo-Labels**
- **Stochastic Weights Restoration** to avoid catastrophic forgetting

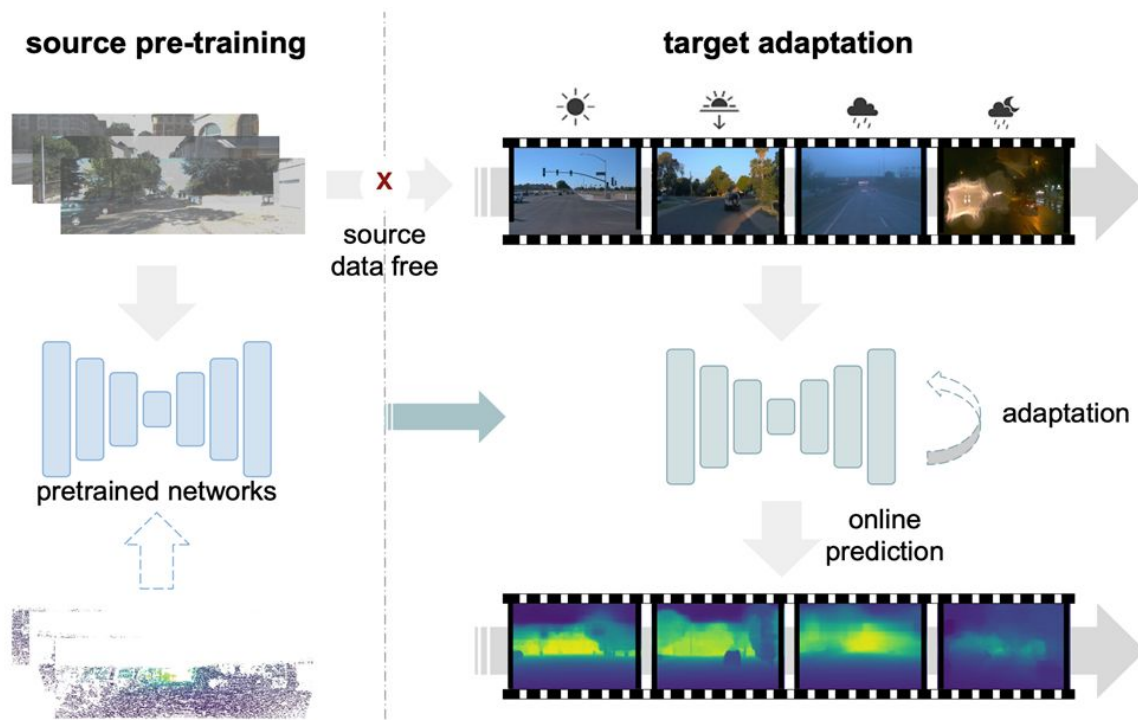
# Continual Test-Time Domain Adaptation

Table 2. Classification error rate (%) for the standard CIFAR10-to-CIFAR10C online continual test-time adaptation task. Results are evaluated on WideResNet-28 with the largest corruption severity level 5. \* denotes the requirement on additional domain information.

| Method              | Weight-avg. | Aug-avg. | Stochastic Restore | $t$         |             |             |             |             |             |             |             |             |             |            |             |               |             |             |                   | Mean |
|---------------------|-------------|----------|--------------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|------------|-------------|---------------|-------------|-------------|-------------------|------|
|                     |             |          |                    | Gaussian    | shot        | impulse     | defocus     | glass       | motion      | zoom        | snow        | frost       | fog         | brightness | contrast    | elastic_trans | pixelate    | jpeg        |                   |      |
| Source              |             |          |                    | 72.3        | 65.7        | 72.9        | 46.9        | 54.3        | 34.8        | 42.0        | 25.1        | 41.3        | 26.0        | 9.3        | 46.7        | 26.6          | 58.5        | 30.3        | 43.5              |      |
| BN Stats Adapt      |             |          |                    | 28.1        | 26.1        | 36.3        | 12.8        | 35.3        | 14.2        | 12.1        | 17.3        | 17.4        | 15.3        | 8.4        | 12.6        | 23.8          | 19.7        | 27.3        | 20.4              |      |
| Pseudo-label        |             |          |                    | 26.7        | 22.1        | 32.0        | 13.8        | 32.2        | 15.3        | 12.7        | 17.3        | 17.3        | 16.5        | 10.1       | 13.4        | 22.4          | 18.9        | 25.9        | 19.8              |      |
| TENT-online* [61]   |             |          |                    | 24.8        | 23.5        | 33.0        | 12.0        | 31.8        | 13.7        | 10.8        | 15.9        | 16.2        | 13.7        | 7.9        | 12.1        | 22.0          | 17.3        | 24.2        | 18.6              |      |
| TENT-continual [61] |             |          |                    | 24.8        | <b>20.6</b> | 28.6        | 14.4        | 31.1        | 16.5        | 14.1        | 19.1        | 18.6        | 18.6        | 12.2       | 20.3        | 25.7          | 20.8        | 24.9        | 20.7              |      |
| CoTTA (Ours)        | ✓           |          |                    | 27.2        | 22.8        | 30.8        | 12.1        | 30.1        | 13.9        | 11.9        | 17.2        | 16.0        | 14.3        | 9.4        | 13.1        | 19.9          | 15.4        | 19.9        | 18.3              |      |
| CoTTA (Ours)        | ✓           | ✓        |                    | 24.5        | 21.0        | <b>26.0</b> | 12.3        | 27.9        | 13.9        | 12.0        | 16.6        | 15.9        | 14.7        | 9.4        | 13.6        | 19.8          | 14.7        | 18.7        | 17.4              |      |
| CoTTA (Ours)        | ✓           | ✓        | ✓                  | <b>24.3</b> | 21.3        | 26.6        | <b>11.6</b> | <b>27.6</b> | <b>12.2</b> | <b>10.3</b> | <b>14.8</b> | <b>14.1</b> | <b>12.4</b> | <b>7.5</b> | <b>10.6</b> | <b>18.3</b>   | <b>13.4</b> | <b>17.3</b> | <b>16.2</b> (0.1) |      |

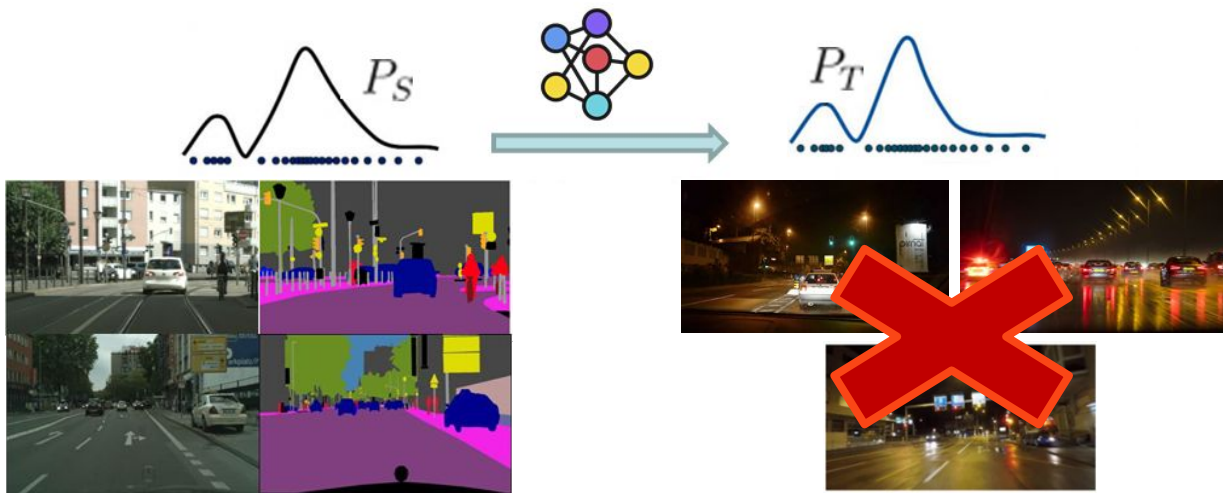


# Continual Test-Time Domain Adaptation



# What can we do to generalize?

1. Unsupervised Domain Adaptation: Learning Target Distribution with Unlabeled Samples
2. Test-time Adaptation: Learning Target Distribution at Test Time from a Single Sample
3. Zero-shot Adaptation: Learning Target Distribution with Text Prompt



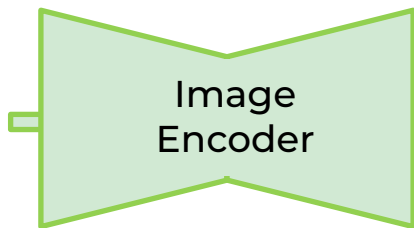
# 2022 - Foundation Models

---

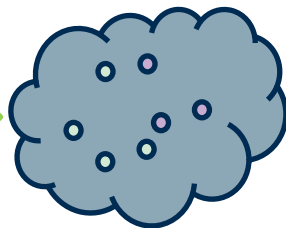
## Multimodal Foundation Models

- ┆ Vision-Language Models - VLM: **CLIP / BLIP / ALIGN**

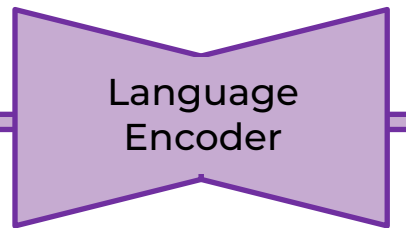
Image space



Multi Modal  
Space



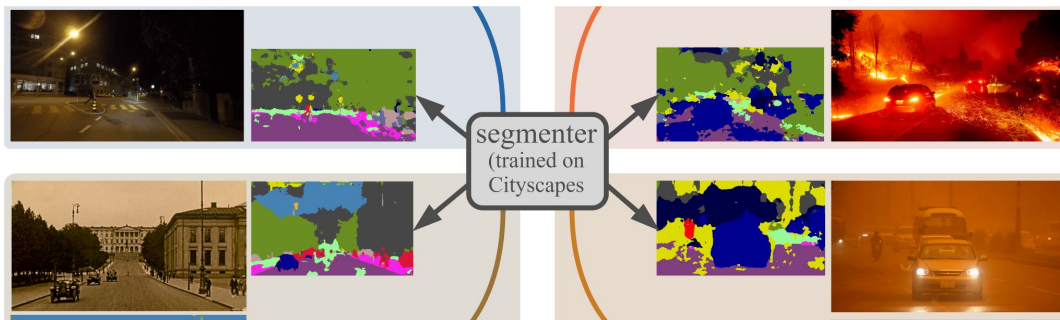
Language space



Driver  
stopping at  
pedestrian  
crossing

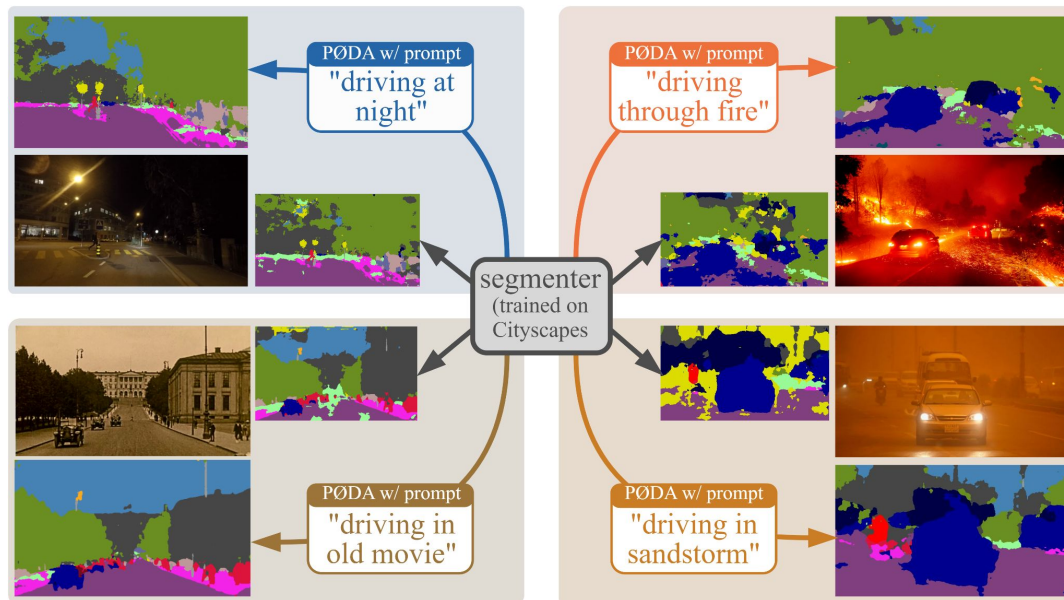
# Prompt-driven Zero-shot Domain Adaptation

Harness foundation models for DA?



# Prompt-driven Zero-shot Domain Adaptation

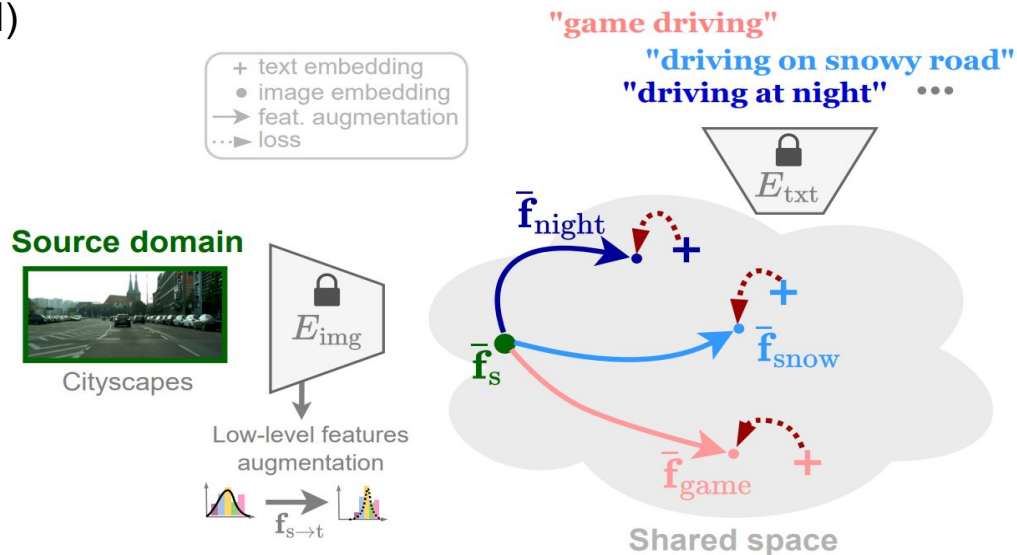
Harness foundation models for DA?



# Prompt-driven Zero-shot Domain Adaptation

## Prompt-driven Instance Normalization (PIN)

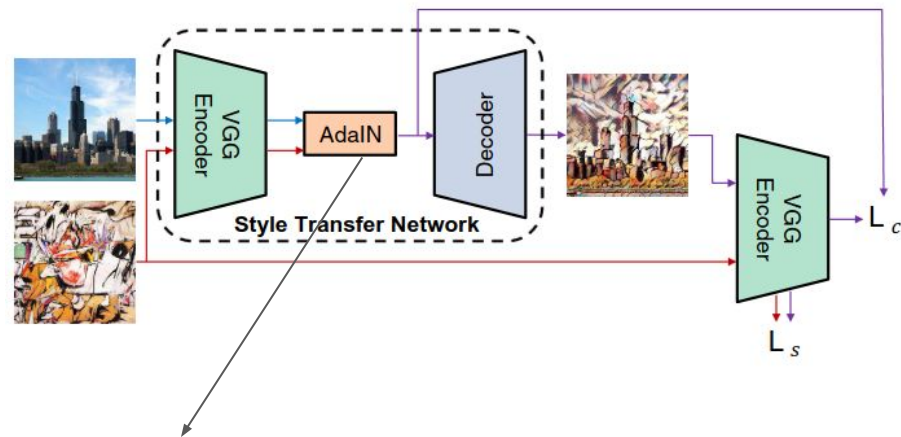
- Stylize features using prompts
- Preserve semantics



# Prompt-driven Zero-shot Domain Adaptation

## Prompt-driven Instance Normalization (PIN)

- Stylize features using prompts
- Preserve semantics

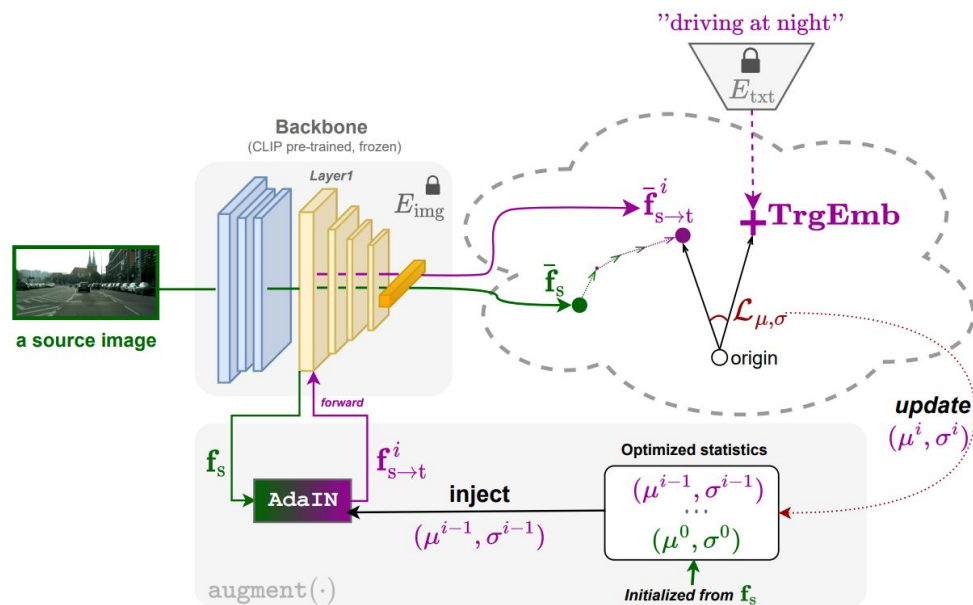


$$\text{AdaIN}(x, y) = \sigma(y) \frac{x - \mu(x)}{\sigma(x)} + \mu(y)$$

# Prompt-driven Zero-shot Domain Adaptation

## Prompt-driven Instance Normalization (PIN)

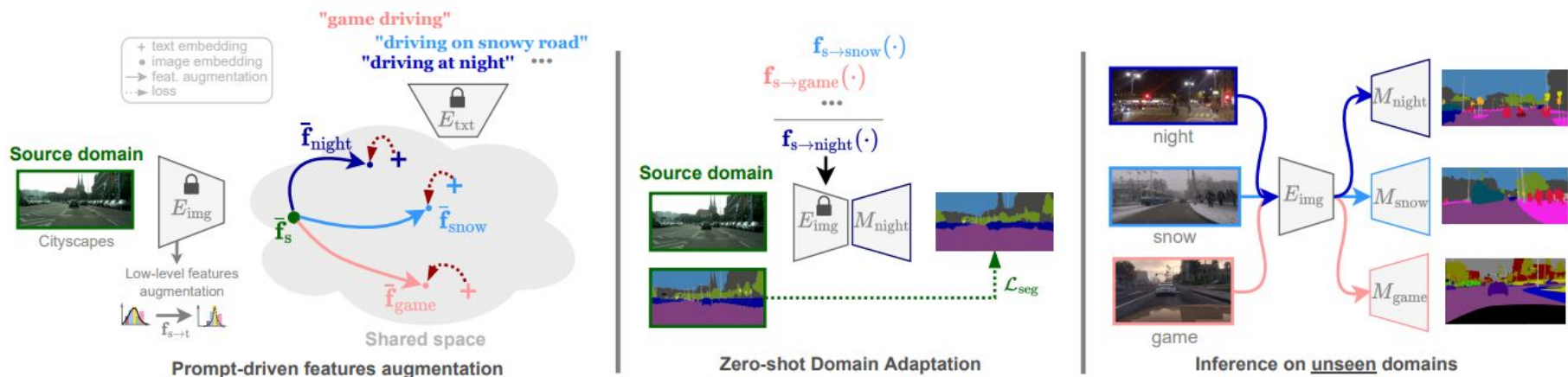
- Stylize features using prompts
- Preserve semantics



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



# Prompt-driven Zero-shot Domain Adaptation



# Prompt-driven Zero-shot Domain Adaptation

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



Figure 5. **CLIPstyler [21] stylization.** A sample Cityscapes image stylized using adhoc target prompts. Translated images exhibit visible artifacts, potentially harming adaptation *e.g.* rain in Tab. 1

# Prompt-driven Zero-shot Domain Adaptation

---

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

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

**Table 8. Effect of different priors for zero-shot/one-shot adaptation.** We report mIoU% for source-only / adapted models, and gain brought by adaptation ( $\Delta$  in mIoU). Note that [26, 56] use a deeper backbone making results not directly comparable.

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

↑ Relevant

ChatGPT-generated

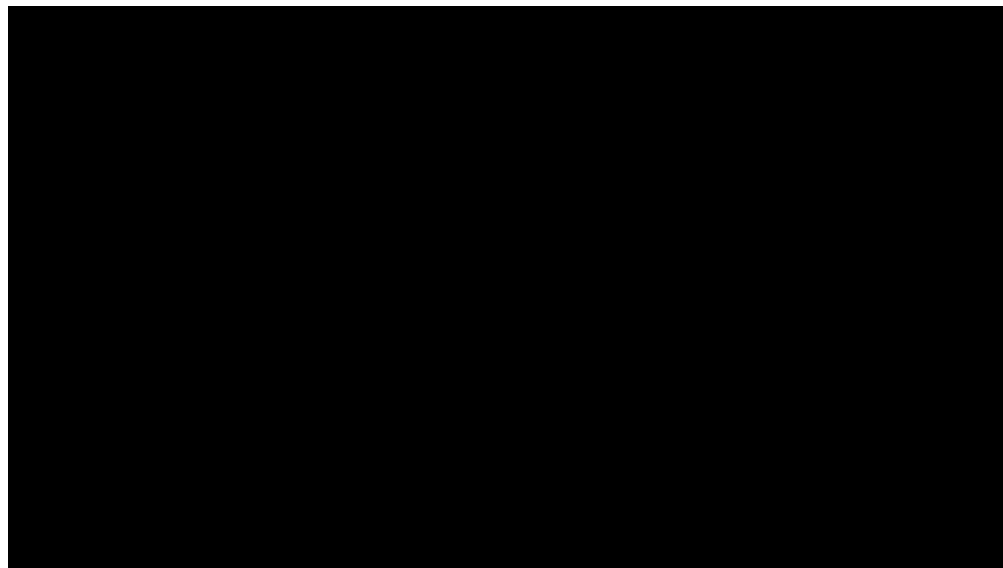
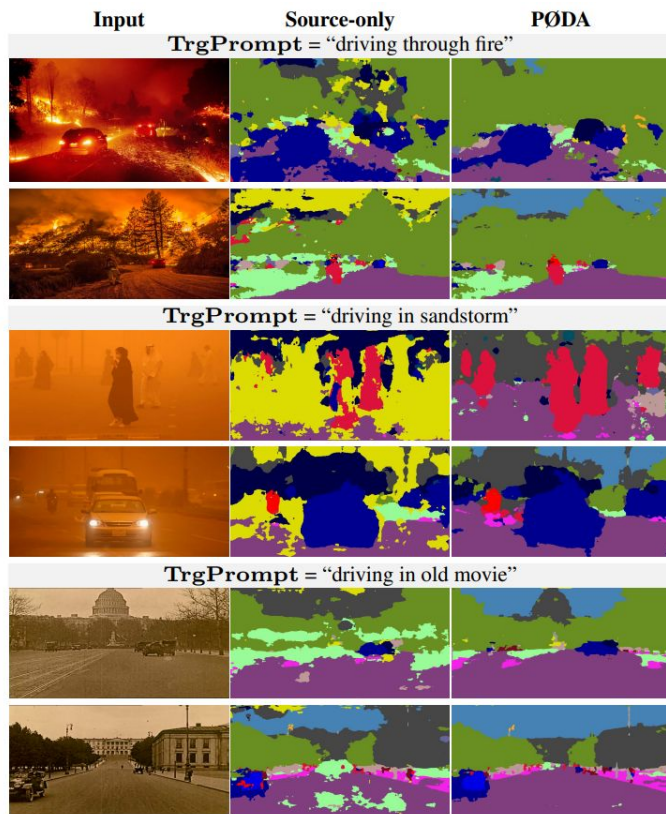
↓ Irrelevant

# Prompt-driven Zero-shot Domain Adaptation

| Method            | Target | CS→Foggy    | Night<br>Clear | Dusk<br>Rainy | Night<br>Rainy | Day<br>Foggy |
|-------------------|--------|-------------|----------------|---------------|----------------|--------------|
| <i>Backbone</i>   |        | ResNet-50   | ResNet-101     |               |                |              |
| DA-Faster [6]     | ✓      | 32.0        | -              | -             | -              | -            |
| ViSGA [38]        | ✓      | 43.3        | -              | -             | -              | -            |
| NP+ [12]          | ✗      | 46.3        | -              | -             | -              | -            |
| S-DGOD [48]       | ✗      | -           | 36.6           | 28.2          | 16.6           | 33.5         |
| CLIP The Gap [44] | ✗      | -           | 36.9           | 32.3          | 18.7           | 38.5         |
| PØDA              | ✗      | <b>47.3</b> | <b>40.3</b>    | <b>37.4</b>   | <b>19.0</b>    | <b>41.7</b>  |

Table 7. **PØDA** for object detection (mAP%).

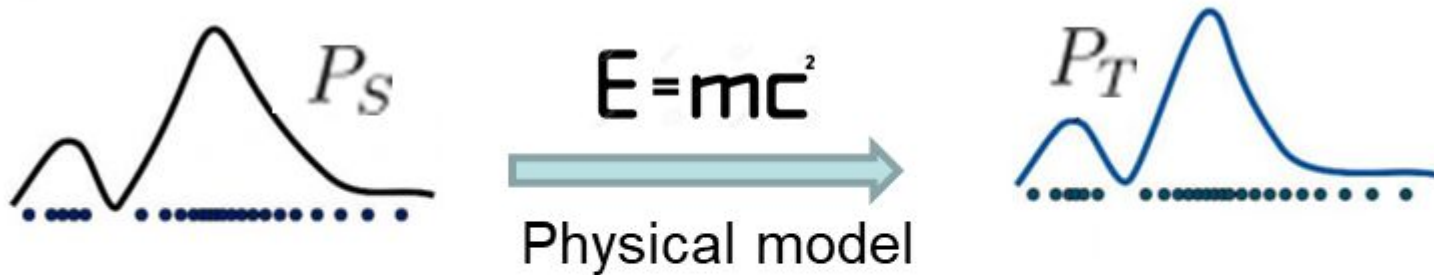
# Prompt-driven Zero-shot Domain Adaptation



Poster on Friday

# What **else** can we do to advance?

1. Unsupervised Domain Adaptation: Learning Target Distribution with Unlabeled Samples
2. Test-time Adaptation: Learning Target Distribution at Test Time from a Single Sample
3. Zero-shot Adaptation: Learning Target Distribution with Text Prompt
4. Data Synthesis: Simulate Target Distribution via Physics-Based Model



# Data synthesis



Semantic Foggy Scene Understanding with Synthetic Data, Sakaridis, Dai, and Van Gool, IJCV, 2018



Flare7K: A Phenomenological Nighttime Flare Removal Dataset. Dai, Li, Zhou, Feng, and Loy, NeurIPS, 2022



Physics-Based Rendering for Improving Robustness to Rain. Halder, Lalonde, and Charette, ICCV 2019



# What **else** can we do to advance?

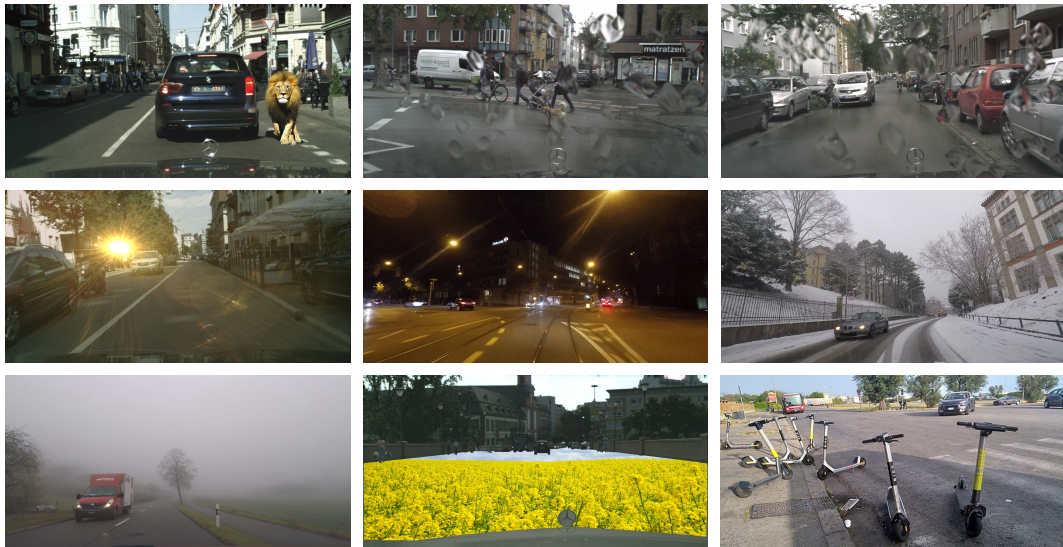
1. Unsupervised Domain Adaptation: Learning Target Distribution with Unlabeled Samples
2. Test-time Adaptation: Learning Target Distribution at Test Time from a Single Sample
3. Zero-shot Adaptation: Learning Target Distribution with Text Prompt
4. Data Synthesis: Simulate Target Distribution via Physics-Based Model
5. Robustness Benchmark

# BRAVO Challenge

---

A unified robustness benchmark for vision perception in autonomous driving

- Semantic segmentation
- Two tracks: single- and multi-domain training
- 3,901 images
- 7 metrics for a comprehensive assessment
- 6 assessment modalities on the test datasets





Thank you!