



**Politecnico  
di Torino**

# **Adversarial Domain Adaptation for Real-time Semantic Segmentation**

Giulia D'Ascenzi, Patrizio de Girolamo, Carlos Rosero

Advanced Machine Learning  
2021-2022

# Introduction

**Semantic Segmentation**



**Real - Time applications**



**Domain Adaptation**



# Real - Time Semantic Segmentation

**Goal:** Assign to each pixel of an input image a category label



road	sign	truck
sidewalk	vegetation	bus
building	terrain	train
wall	sky	motorcycle
fence	person	bicycle
pole	rider	others
light	car	

**Possible Applications:** Autonomous vehicles, virtual reality, computer - aided diagnosis, etc..



# Real - Time Semantic Segmentation

## Requested characteristics:

- Time - efficient inference speed
- Low memory and resources usage
- Good segmentation performances



Not something we can get with the  
standard semantic - segmentation models



Lots of parameters and  
Floating Point Operations (FLOP)

# Real - Time Semantic Segmentation

**Additional drawback:** Semantic Segmentation models need to be trained on a large amount of densely labeled images.



It requires a large amount of human effort (time and money).



# Exploiting synthetic dataset

**Idea:** Train the networks on large photo - realistic synthetic datasets with computer-generated annotations.



**GTA5 Dataset [2]**



# Domain Shift



GTA5  
(Source)



Cityscapes  
(Target)



# Domain Shift



**GTA5**  
(Source)



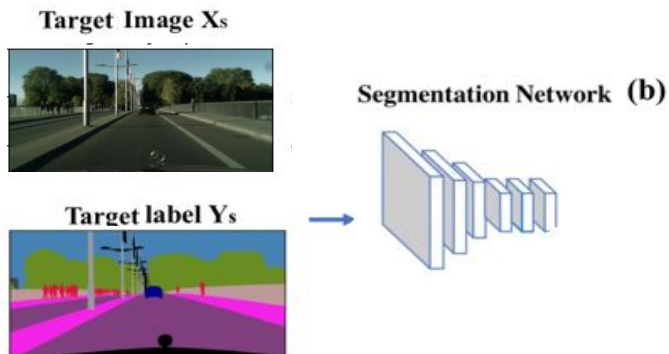
**Cityscapes**  
(Target)

**Domain Adaptation**



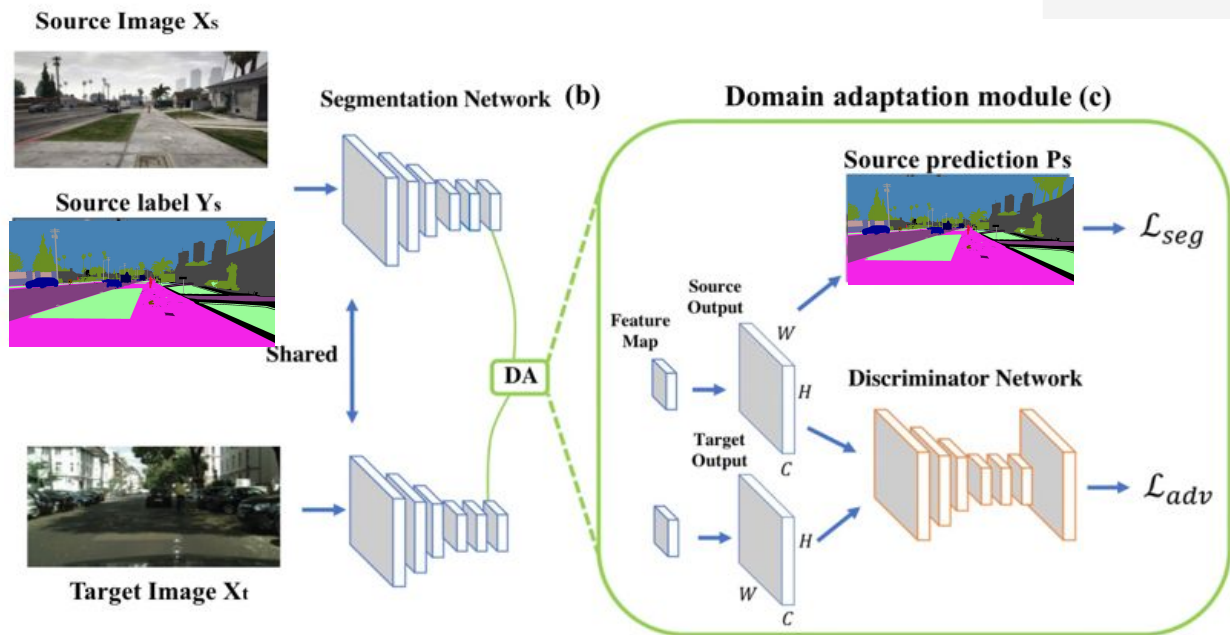


# Overview Complete Architecture



1. Real- Time semantic Segmentation Network: BiSeNet [3]. Upper bound with target only.

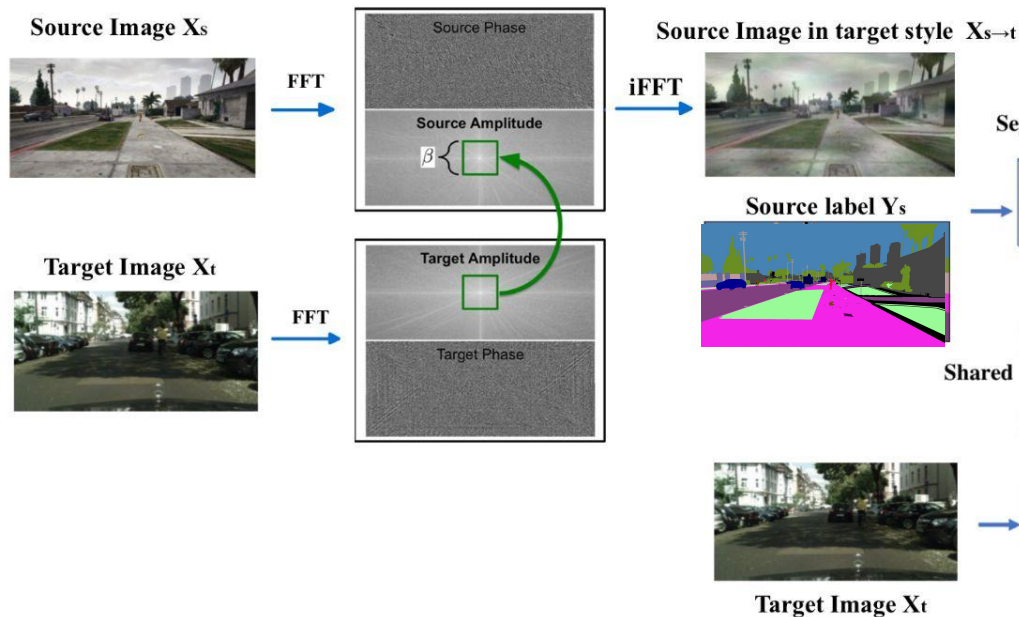
# Overview Complete Architecture



## 2. Adversarial Domain Adaptation Module [4]

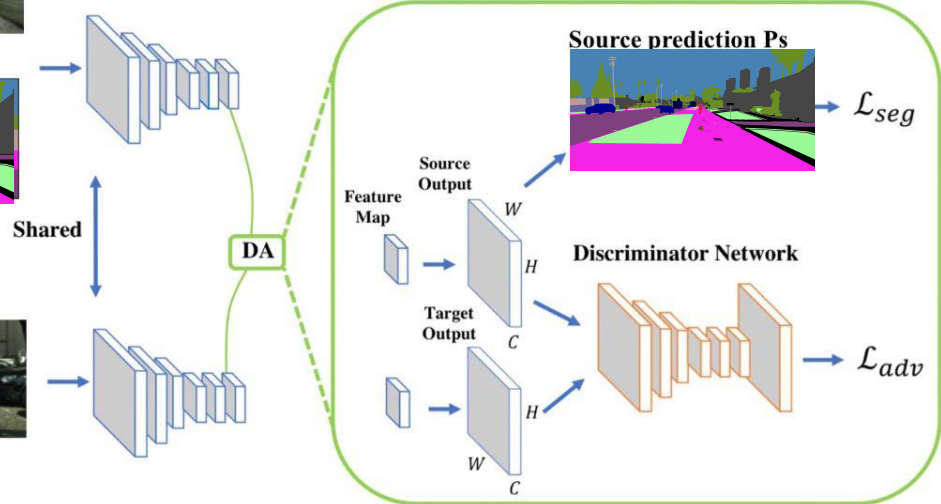
# Overview Complete Architecture

Image to Image translation (a)



Segmentation Network (b)

Domain adaptation module (c)



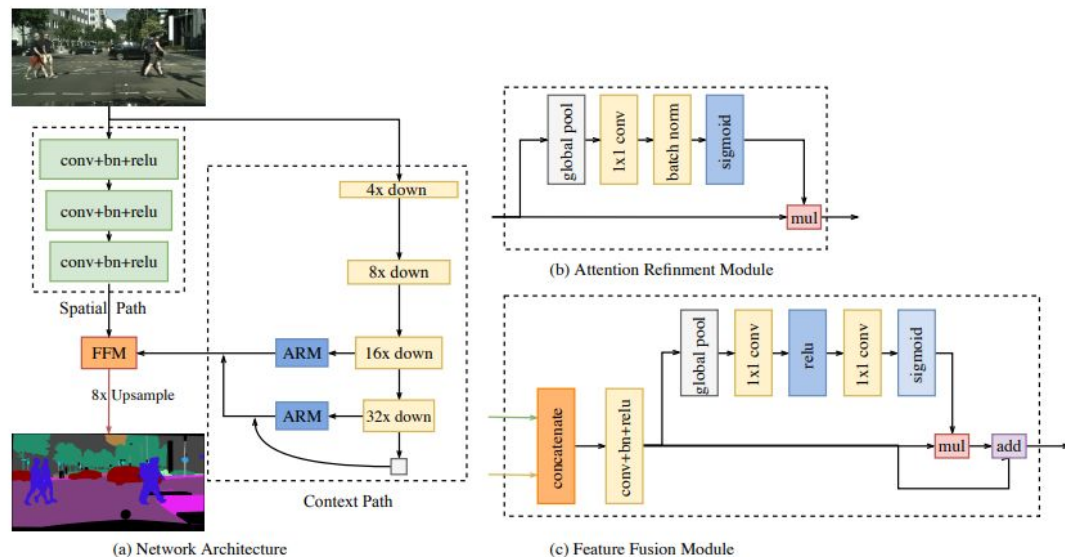
## 3. Increasing the performances: FDA [5]



# Method

# Semantic Segmentation

**Network used:** BiSeNet [3] ( backbone Resnet-18 pretrained on ImageNet)



**Capture semantics**



Context path

**Preserve spatial details**



Spatial path

# (Without) Domain adaptation

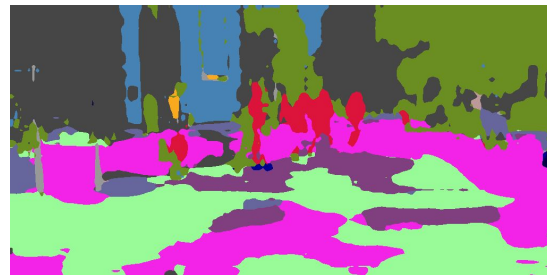
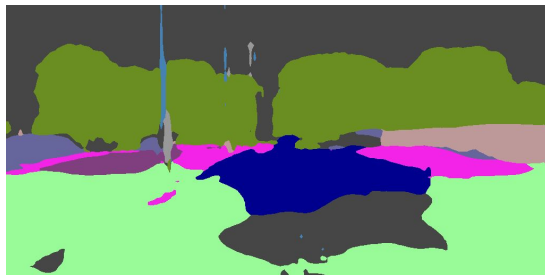
**Goal:** Transfer the knowledge obtained from the source finely annotated dataset to the target and unlabelled dataset

**First try:** Training BiSeNet using only GTA and then testing it on Cityscapes

Ground Truth:

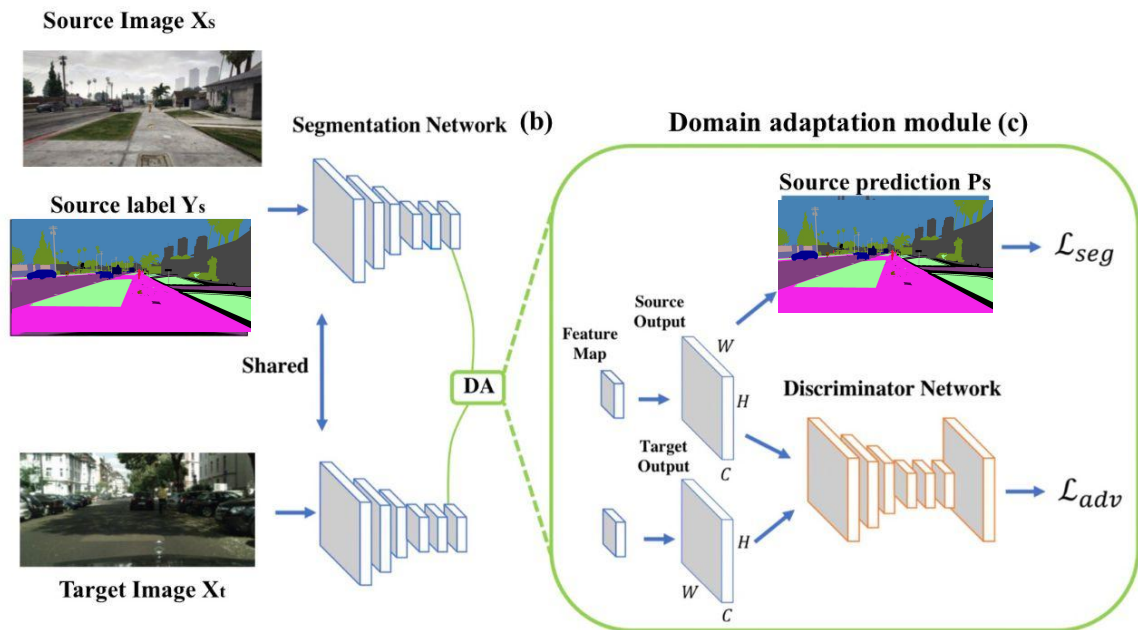


Output without DA:



# Domain adaptation

**Better Idea:** Reducing the domain shift between source and target dataset using Unsupervised Adversarial Domain Adaptation





# Unsupervised adversarial domain adaptation

**Method:** Min-Max game between the **generator** and the **discriminator**

**Generator training:**

$$\mathcal{L}(X_s, X_t) = \mathcal{L}_{seg}(X_s) + \lambda_{adv} \mathcal{L}_{adv}(X_t)$$

Where:

$$\mathcal{L}_{seg}(I_s) = - \sum_{h,w} \sum_{c \in C} Y_s^{(h,w,c)} \log \left( P_s^{(h,w,c)} \right)$$

$$\mathcal{L}_{adv}(I_t) = - \sum_{h,w} \log \left( \mathbf{D}(P_t)^{(h,w,1)} \right)$$

**Discriminator Training:**

$$\mathcal{L}_d(P) = - \sum_{h,w} (1 - z) \log \left( \mathbf{D}(P)^{(h,w,0)} \right) + z \log \left( \mathbf{D}(P)^{(h,w,1)} \right)$$



# Discriminators

Fully Convolutional Discriminator



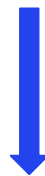
2D convolutions



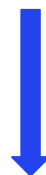
Total Params: 2.7M  
Total FLOPS: 61.8G

vs

Lighter Discriminator



Depth-wise separable convolutions



Total Params: 191K  
Total FLOPS: 4.36G

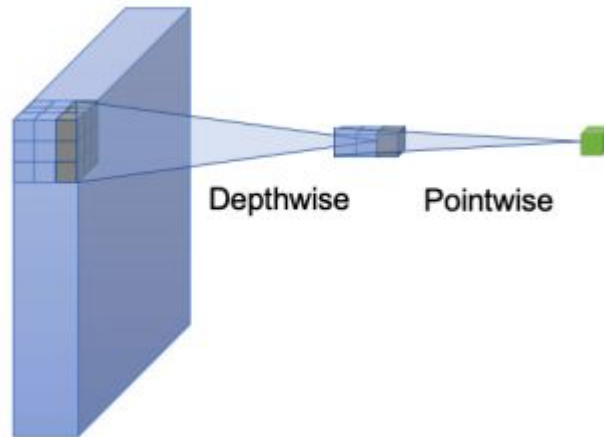
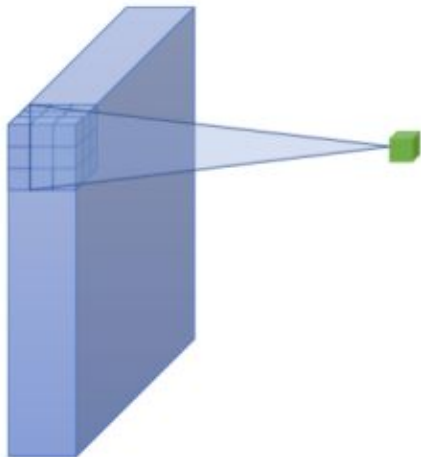


# Different types of convolution

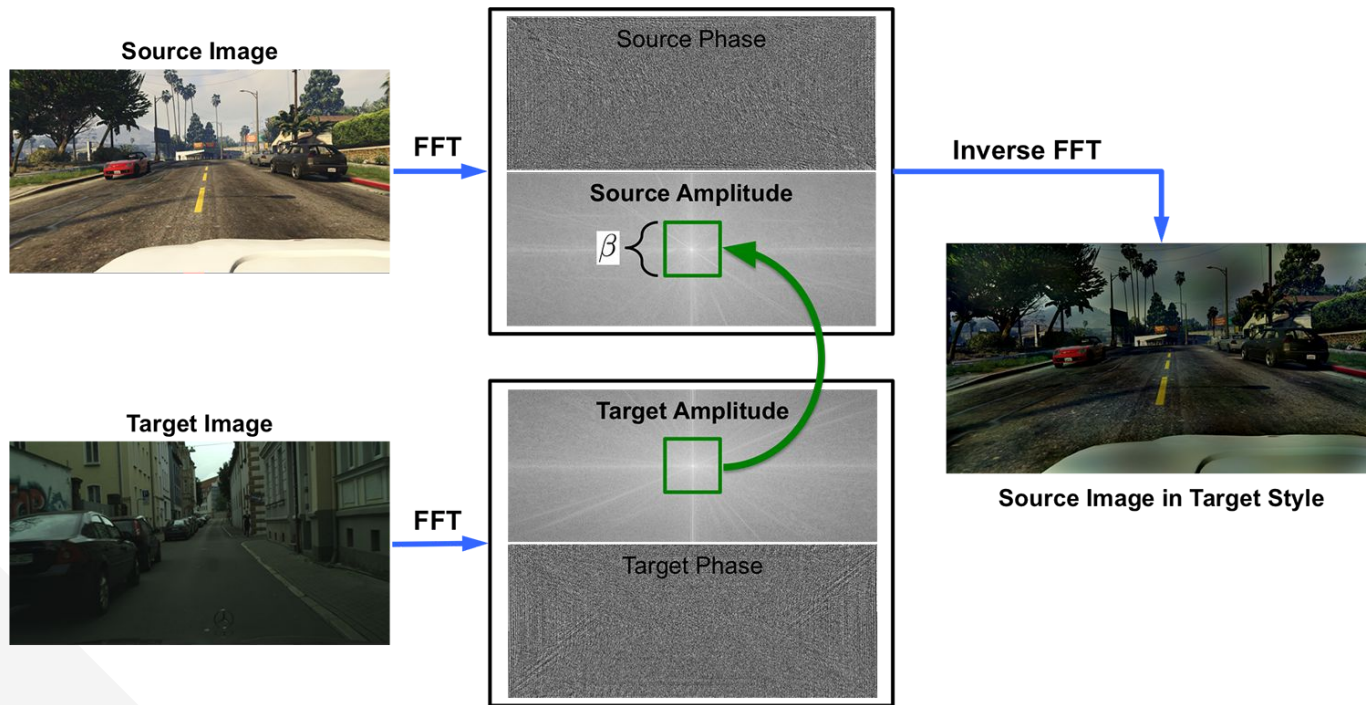
2D convolutions

VS

Depth-wise separable convolutions



# Increasing Performance: FDA



# FDA $\beta$ Ablation Study

Source (GTA)



Target (Cityscapes)



B = 0.01



B = 0.05

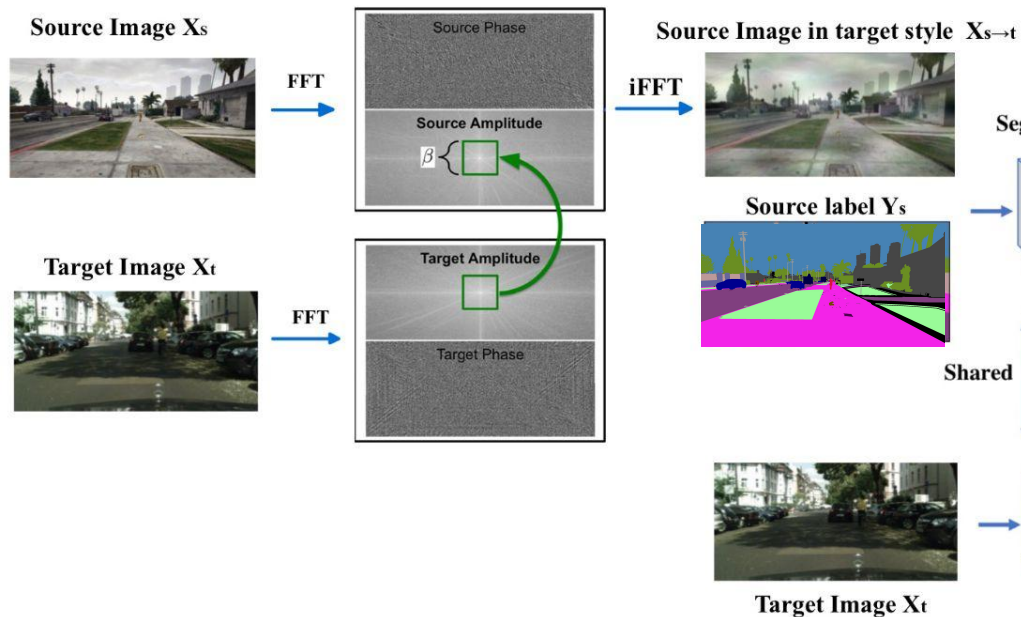


B = 0.10

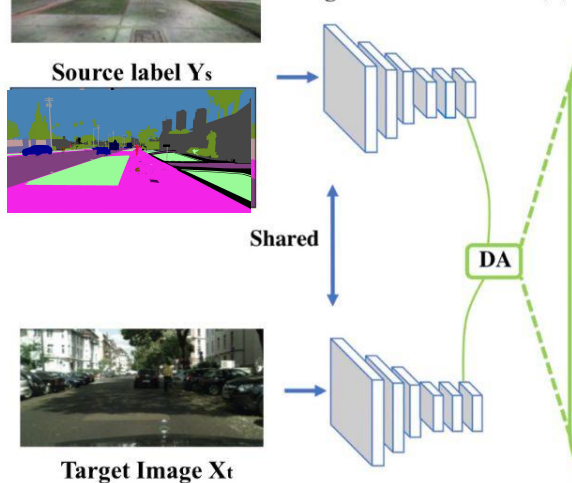


# Overview Complete Architecture

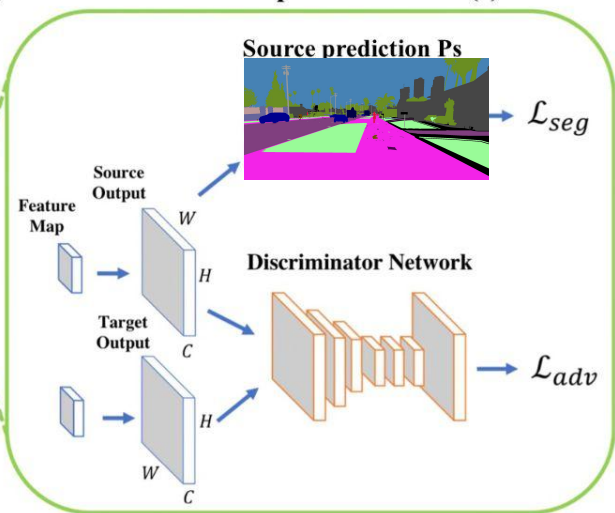
Image to Image translation (a)



Segmentation Network (b)



Domain adaptation module (c)





# Results





## Results - Upper Bound / Augmentation

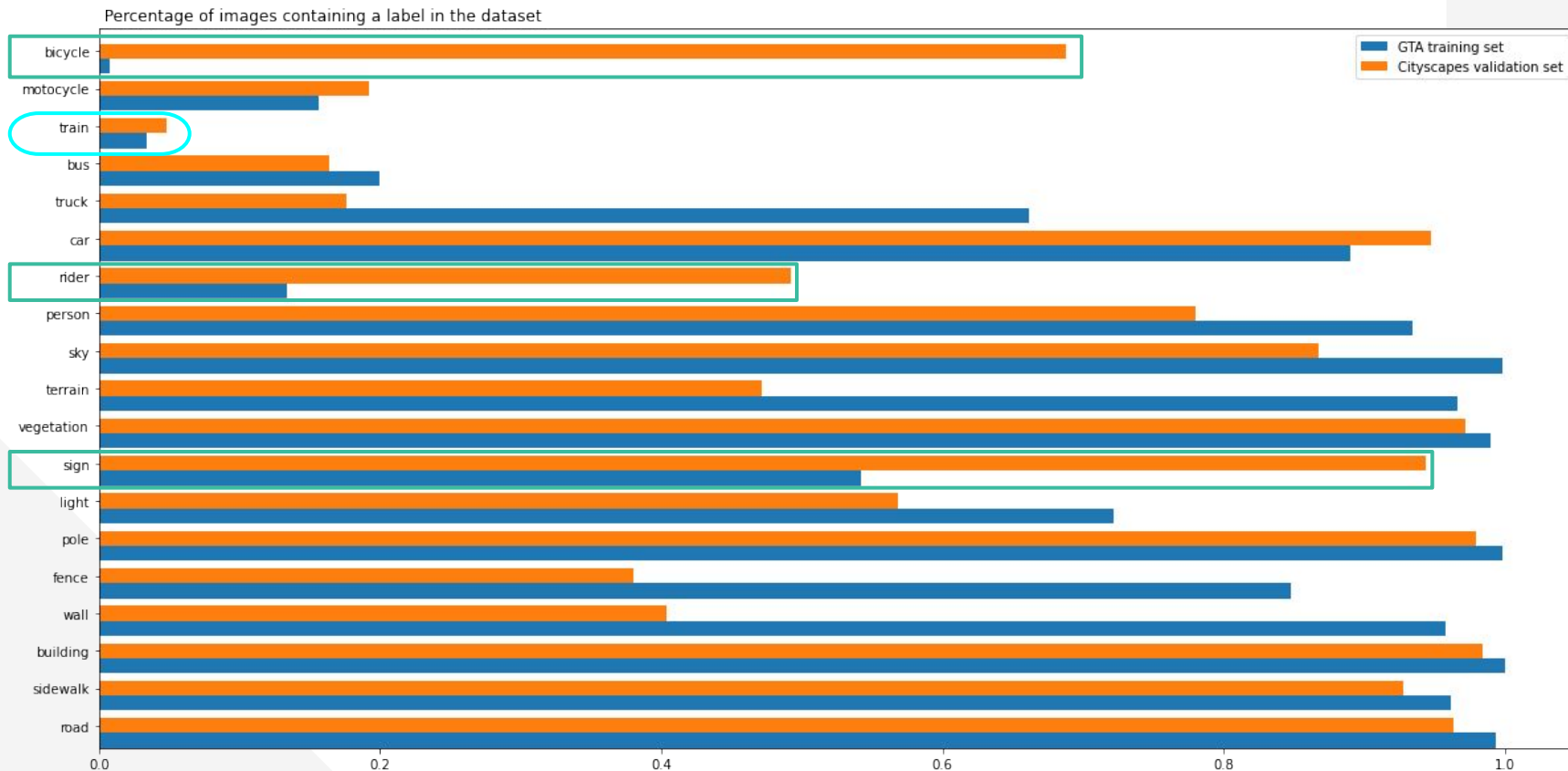
Augmentation	Accuracy(%)	mIoU(%)
None	79.3	47.3
Blur	79.1	47.0
Horizontal Flipping	79.8	49.7
Horizontal Flipping and Blur	79.7	<b>49.9</b>



# Results

Method	mIoU %
<b>Target only</b>	<b>49.9</b>
No domain adaptation	13.4
Fully Conv	24.2
LW Conv	24.4
LW Conv, FDA ( $\beta = 0.01$ )	27.0
LW Conv , FDA ( $\beta = 0.05$ )	26.6
<b>LW Conv, FDA (<math>\beta = 0.10</math>)</b>	<b>27.4</b>
LW Conv , FDA ( $\beta = 0.10$ , no blur)	26.9

# Results - Analysis



# Results - Analysis

Method	Sign	Rider	Train	Bicycle	mIoU %
Target only	<b>45.9</b>	<b>30.2</b>	<b>25.9</b>	<b>53.3</b>	<b>49.9</b>
No domain adaptation	0.0	0.0	0.0	0.0	13.4
Fully Conv	0.1	0.1	0.0	0.0	24.2
LW Conv	0.2	1.0	0.0	0.0	24.4
LW Conv, FDA ( $\beta = 0.01$ )	1.4	1.4	0.0	0.0	27.0
LW Conv , FDA ( $\beta = 0.05$ )	0.6	8.2	0.0	0.0	26.6
LW Conv, FDA ( $\beta = 0.10$ )	<b>1.4</b>	<b>3.5</b>	<b>0.0</b>	<b>0.0</b>	<b>27.4</b>
LW Conv , FDA ( $\beta = 0.10$ , no blur )	1.3	4.5	0.0	0.0	26.9

# Conclusion

**Semantic  
Segmentation**



**Real-Time  
Applications**



**Fourier  
Transform**



**Narrower Domain Gap**



# **Thank You**

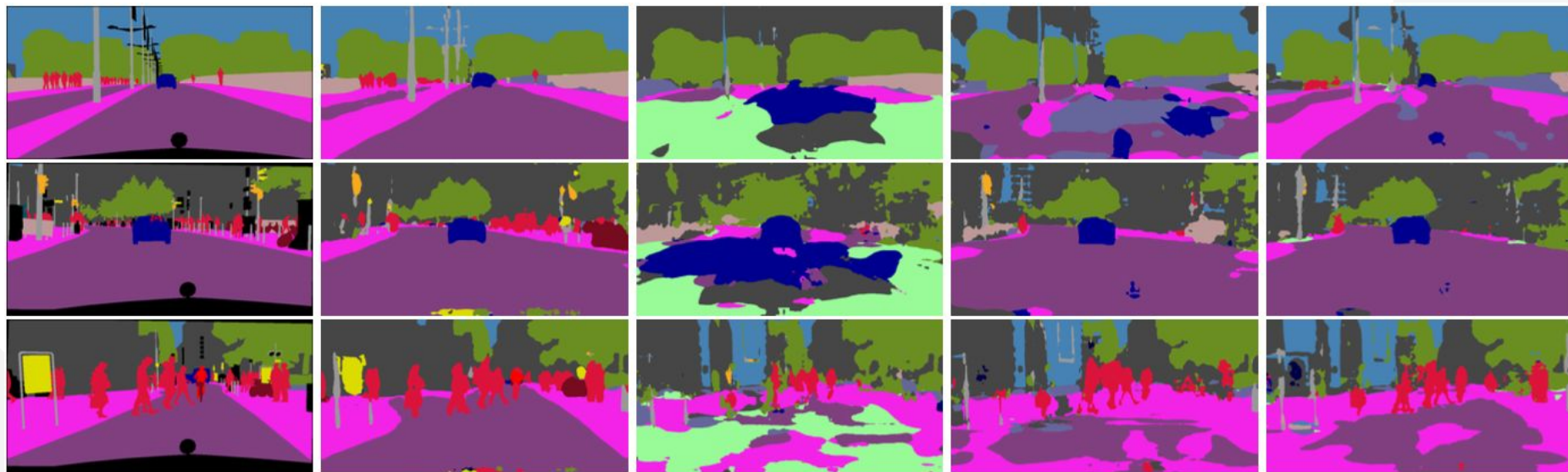
Any Questions?



# **Appendix**



# Visual Comparison of Results



((a)) GT

((b)) Only Target

((c)) No DA

((d)) LW no FDA

((e)) LW FDA

road	building	fence	light	vegetation	sky	rider	truck	train	bicycle
sidewalk	wall	pole	sign	terrain	person	car	bus	motorcycle	others