A photograph of a city street with a tram and people. Overlaid on the image are various colored regions representing semantic segmentation: red for people, blue for cars, green for trees, and yellow for traffic lights. A large, semi-transparent white box covers the right side of the image, containing the title and author information.

# Input or Output Space Alignment?

## A Comparative Study for Continual Unsupervised Domain Adaptation in Semantic Segmentation

---

Tim Lindenau, Bachelor Thesis  
Stuttgart, 15.09.2022

- A Motivation for Continual UDA in Semantic Segmentation
- B 5 Frameworks for Continual UDA
- C Experiment 1: Comparison of the 3 Different Style Transfer Algorithms
- D Experiment 2: Comparison of All 5 Frameworks for Continual UDA
- E Conclusion



# Motivation for Continual UDA in Semantic Segmentation

# Semantic Segmentation Provides Good Scene Understanding Required for Autonomous Driving – However: High Effort for Image Labeling



Original Image



Segmented Image: Every pixel assigned to one class



- **Semantic Segmentation:** Task of assigning every pixel in an image to one class, e.g. car, street, people
- Good results with methods based on **Convolutional Neural Networks**
- **High manual effort** for labeling training images
  - Ca. 90 minutes per image to identify & label all relevant objects
  - Large number of images required ('ballpark figure': several thousand)

# Continual UDA Helps to Achieve Good Segmentation Performance for Different Environments



The challenge: many different environments

Sunny



Fog



Night



Snow



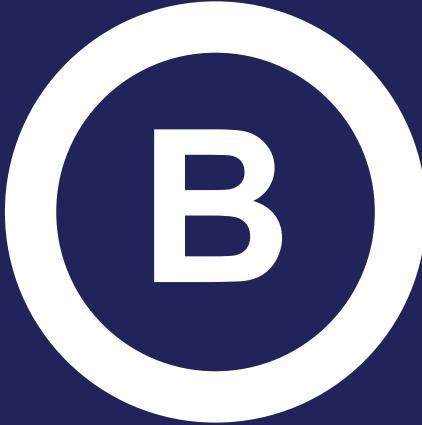
- Car experiences **many different environments** and segmentation performance must be good for all of them
- Segmentation performance decreases for environments unlike training environment. Reason: **Shift in the distribution** of input data  
→ **Labeled training data required for all environments.**  
In practice impossible due to high manual labeling effort

The technical solution: Continual UDA

**(1) Unsupervised Domain Adaptation (UDA):**

- Labeled training data for **one source environment/domain** available
- Improve performance for **one target environment** without any labels

**(2) Continual learning:** Sequentially adapt model to **many target environments** without forgetting

A large, solid blue circle with a thick white border. Inside the circle is a large, white, bold letter "B".

B

## 5 Frameworks for Continual UDA

# I Will Discuss Five Frameworks for Continual UDA



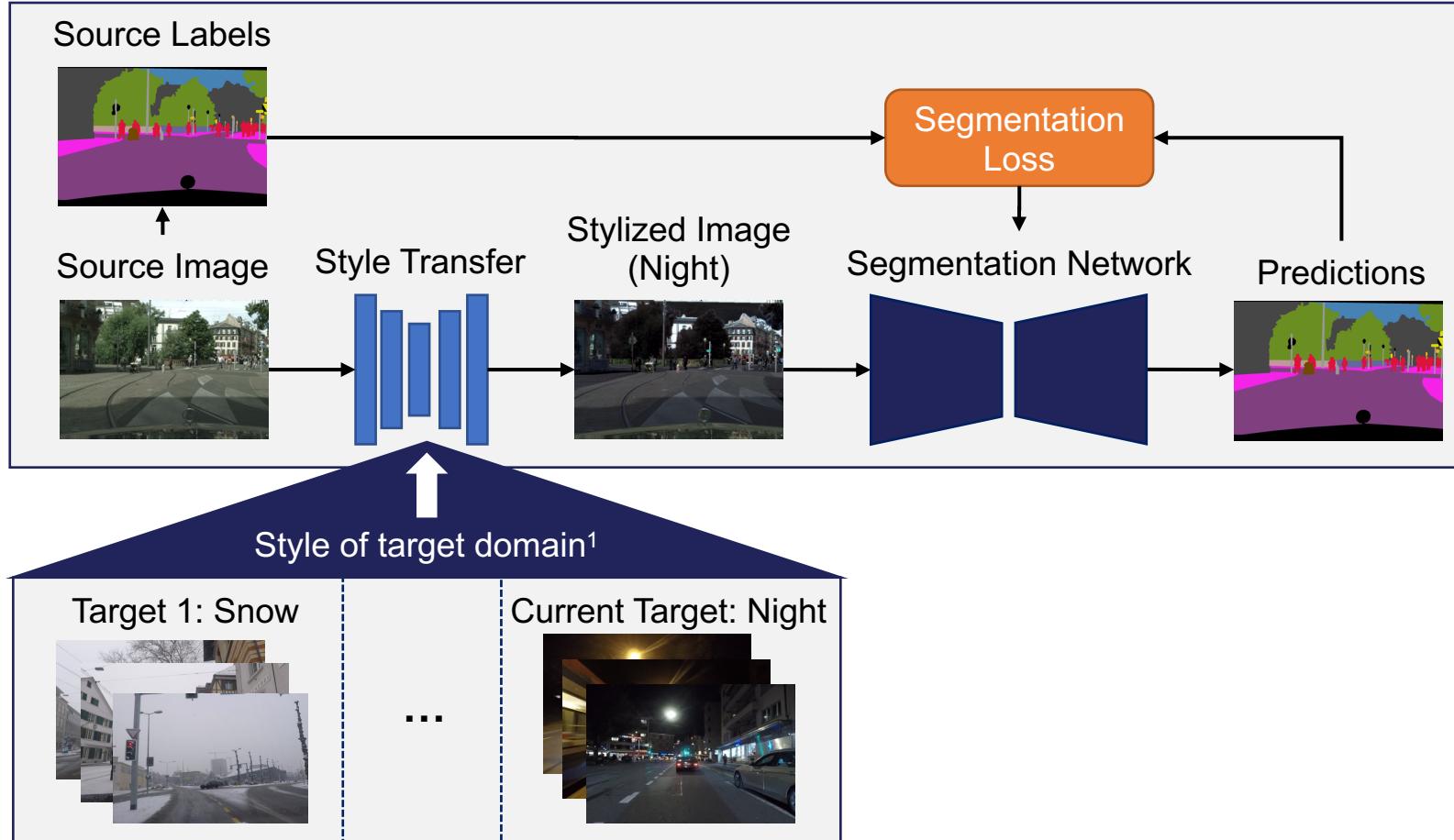
Framework	UDA	Continual Learning
a CACE      Established framework for continual UDA	Input space alignment using style transfer	Replay to prevent forgetting
b C-WCT <sup>2</sup> New framework with enhanced details		
c C-CMD      New framework with more accurate stylization		
d ETM      New framework, never compared against CACE	Output space alignment using adversarial learning	Network expansion
e MuHDI      New framework, never compared against CACE		Multiple specialized classifiers

a

# The CACE Framework for Continual UDA: Input Space Alignment to Adapt to New Domains and a Style Memory to Prevent Forgetting



Input Space Alignment



<sup>1</sup>Unlike illustrated not the images themselves, but only their style is stored which is significantly more memory efficient

## UDA: Input space aligned using style transfer

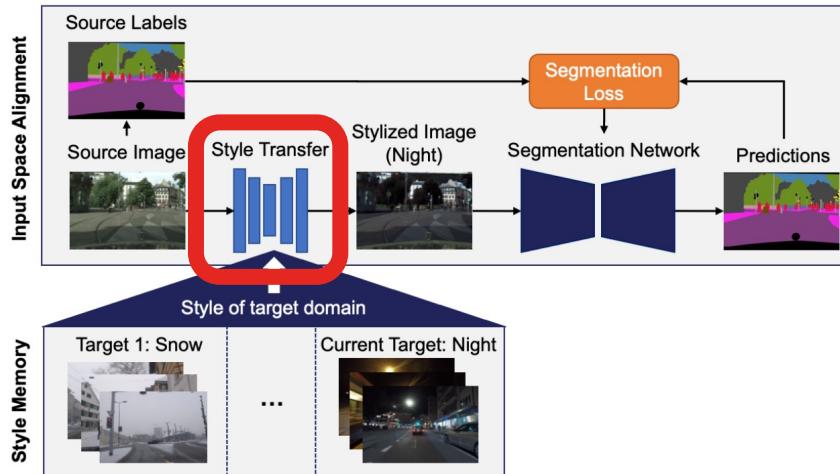
- Style transfer to stylize source images in look of target domain
- Source labels remain valid and can be used to train segmentation network

## Continual learning: Style memory

- Style memory stores styles of all experienced target domains
- When training, alternate between style from current target domain (UDA) and random previous domain (prevent forgetting)

# C-WCT<sup>2</sup> & C-CMD: Two New Frameworks Based on CACE With Enhanced Style Transfer Algorithms

CACE-like frameworks can be created by exchanging the style transfer algorithm



- Style transfer from CACE has some known limitations
- Created two new frameworks just by exchanging the style transfer algorithm (**C-WCT<sup>2</sup>, C-CMD**)

## Advanced style transfer algorithms WCT<sup>2</sup> & CMD

### CACE limitations

**Loss of fine details**  
Inevitable information loss in encoder/decoder structure (max pooling)

### Potential solution

**WCT<sup>2</sup>-algorithm**  
Increased level of detail by overcoming information loss with new pooling operation (Wavelet pooling)

**Theoretical limitations**  
Matching mean and variance cannot provide for perfect distribution alignment

**CMD-algorithm**  
More exact stylizations using higher-order moments

2 New Frameworks,  
Better Results?

# Output Space Alignment: Improving Target Domain Performance Using Structural Similarity in the Output Space



Source



Large difference in appearance

Target



Output space:



Predictions good due to available labels



Structurally similar (e.g. cars driving on street)



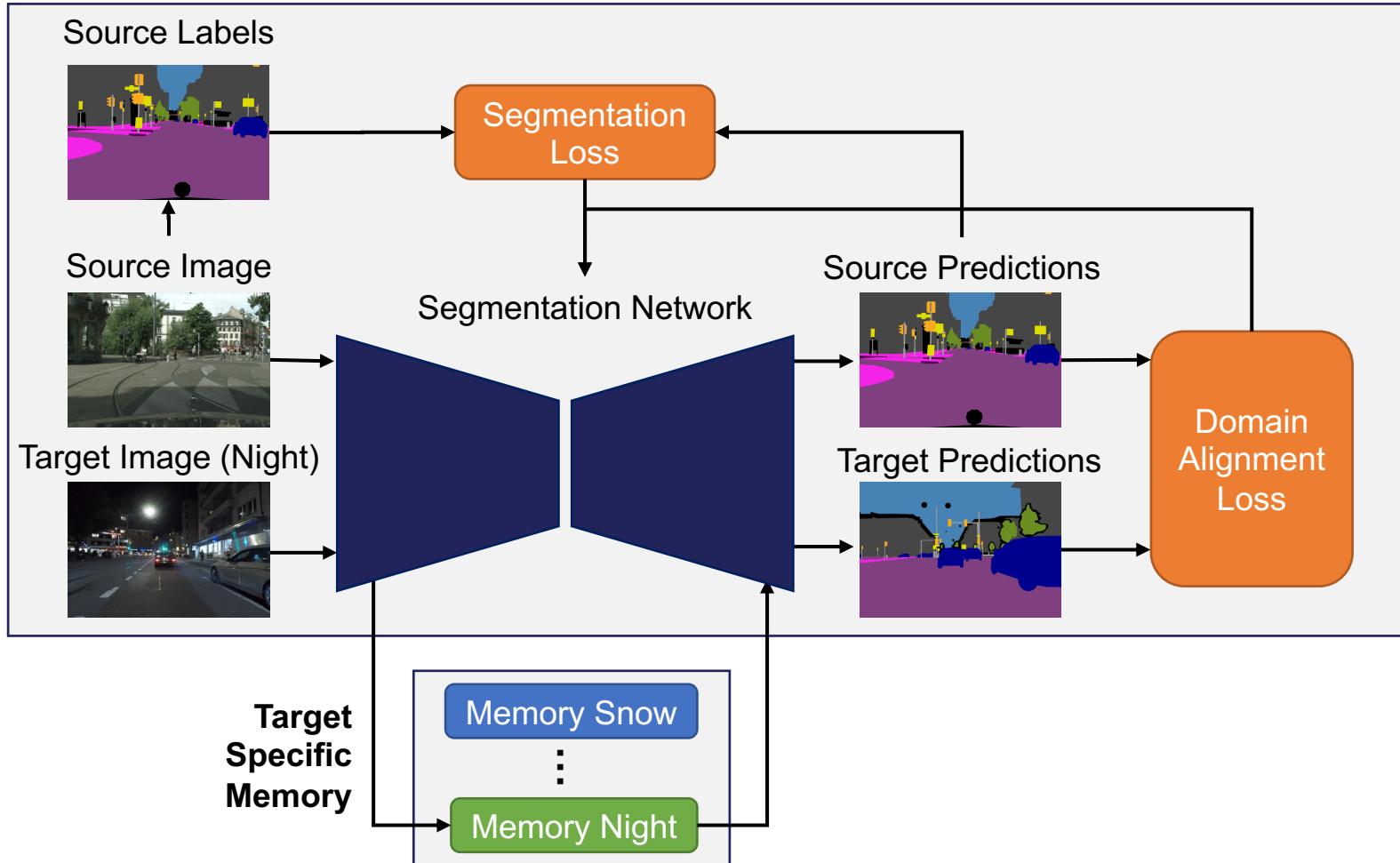
While training:  
Potential errors in predictions

**Output space alignment to improve predictions for the target domain:**

The neural network is trained to make the target domain predictions structurally similar to the predictions of the source domain

# The ETM Framework for Continual UDA: Output Space Alignment to Adapt to New Domains and Target Specific Memory Modules to Prevent Forgetting

## Output Space Alignment



## UDA: Output space alignment

- Minimize discrepancy between source domain predictions and target domain predictions

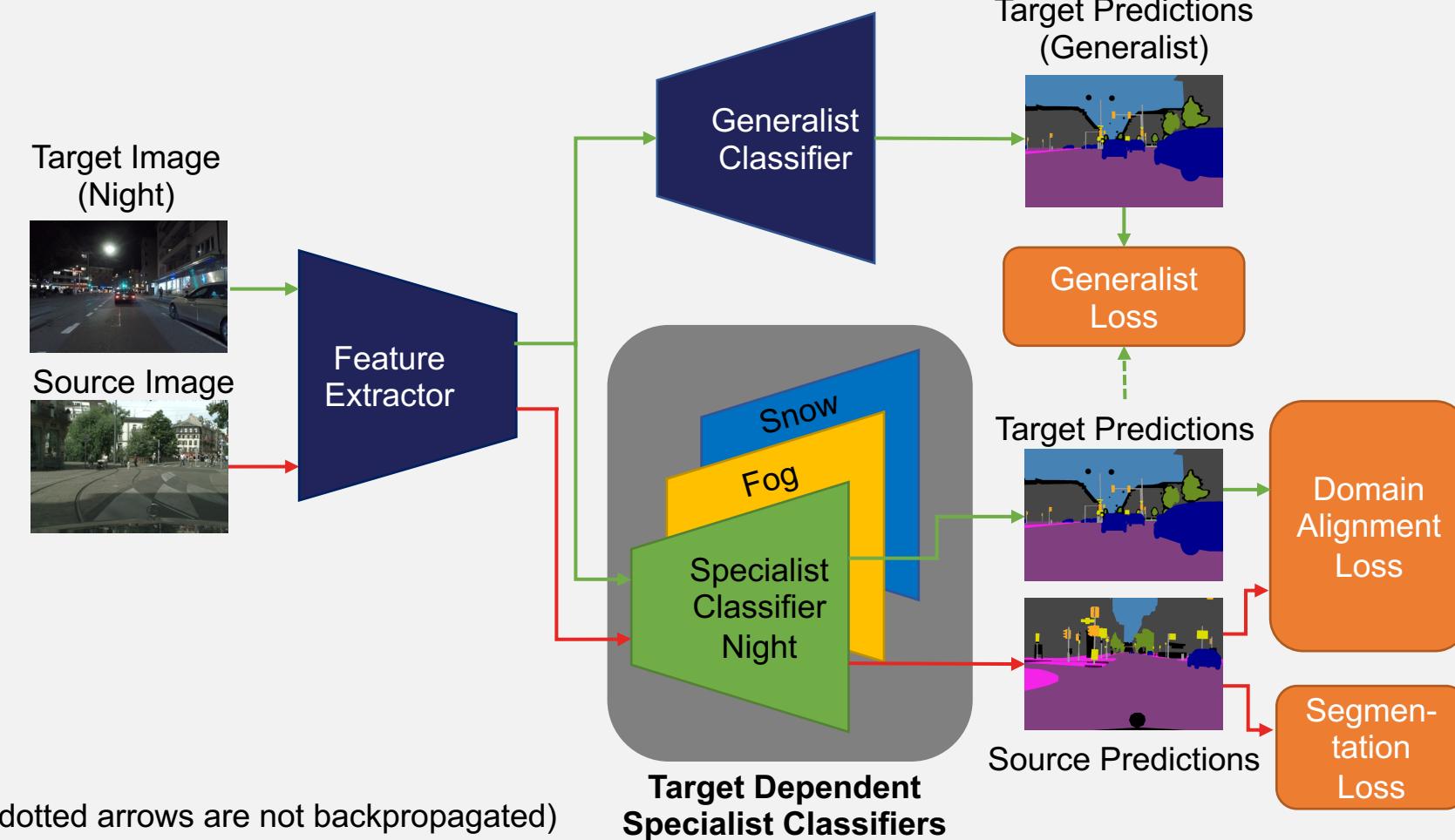
## Continual learning: Target specific memory

- For each target domain, extra memory module parallel to segmentation network
- Memory stores domain discrepancy from current target domain to source domain
- Segmentation network learns domain invariant information and can be used for every domain

# e The MuHDI Framework for Continual UDA: Output Space Alignment to Adapt to New Domains and Multiple Classifiers to Prevent Forgetting



Output Space Alignment



(dotted arrows are not backpropagated)

## UDA: Output space alignment

- Minimize discrepancy between source domain predictions and target domain predictions

## Continual Learning: Target dependent specialist classifiers

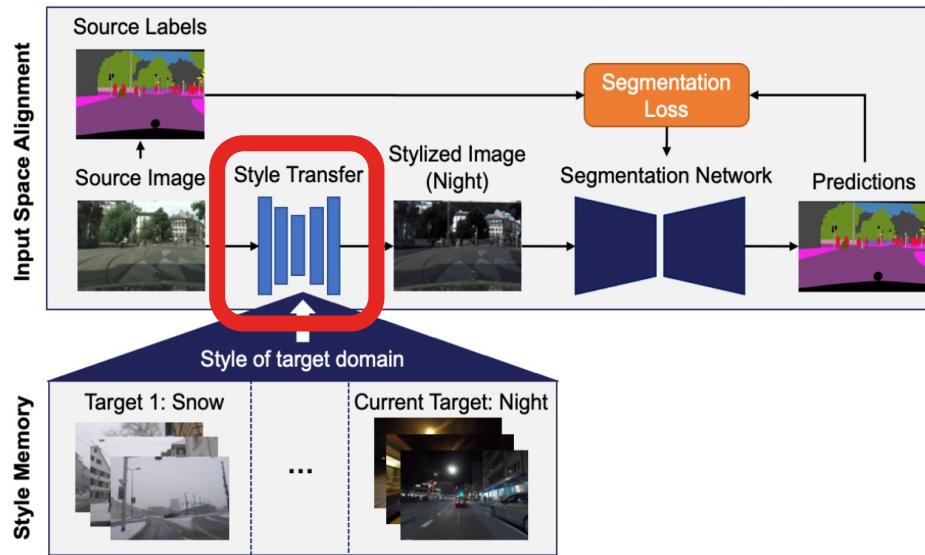
- One specialist classifier trained for each target domain
  - Segmentation performance on source domain
  - Domain alignment between source and target domain
- Knowledge of all specialist classifiers combined in one domain independent generalist classifier



# Experiment 1: Comparison of the 3 Different Style Transfer Algorithms

# Experiment 1: A Visual Analysis of the Three Style Transfer Algorithms – Are the Expected Enhancements Visible?

Experiment 1 focuses only on the style transfer algorithm



- I tested the three style transfer algorithms for transferring style in real-world street scenes
  - CACE's style transfer algorithm
  - WCT<sup>2</sup> – enhanced photorealism
  - CMD – enhanced stylization accuracy
- **Question:** Are the expected enhancements visible in the stylized images?

# All Style Transfer Algorithms Stylize Images Differently: WCT<sup>2</sup> with the Expected Higher Level of Details, CMD with Color Enhancements

Content Image

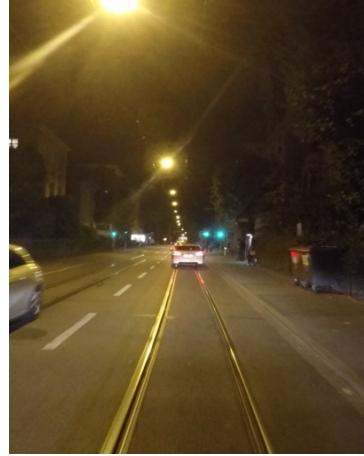


+



+

Style Reference Image



CACE



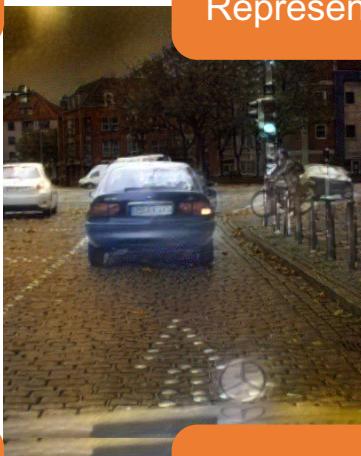
WCT<sup>2</sup>



High Details in Pavement

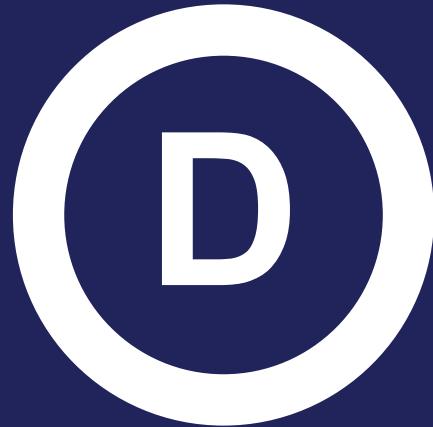
Fine Details in e.g. Bus Sign

CMD



Better Color Representation

Better Color in Street



## Experiment 2: Comparison of All 5 Frameworks for Continual UDA

# Experiment 2: Comparison of the Five Frameworks For Their Performance in Continual UDA

The Cityscapes/ACDC sequence is used to evaluate the performance of all five frameworks for continual UDA



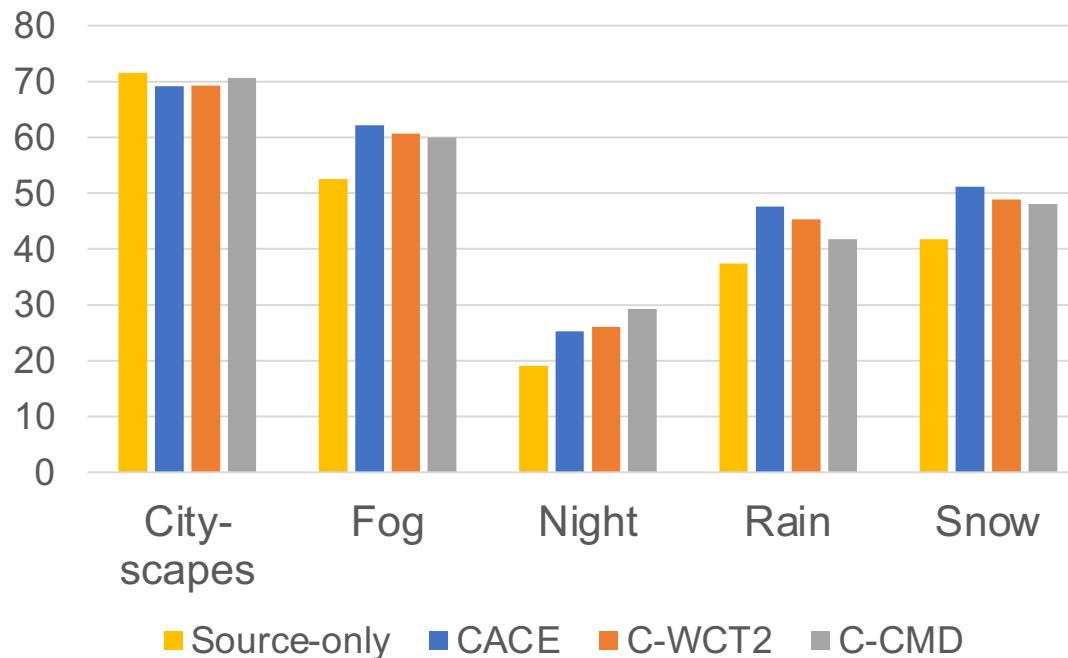
Frameworks compared for **segmentation performance** (mIoU) averaged over **all domains** (mean mIoU) after adapting to the last target domain (Snow)

- I tested the five **frameworks** for **three sequences of domains**; in today's presentation **focusing** on the **Cityscapes/ACDC sequence**<sup>1</sup>
- **2 Questions:**
  - Is there an advantage from the **enhanced style transfer** algorithms?
  - Which approach, **input space alignment or output space alignment**, delivers better performance?

# More Advanced Style Transfer Algorithms Do Not Yield a Benefit Compared to CACE, CACE Performs Best Averaged over All Domains



Domain-wise performance after adapting to the last target domain [All values measured in mIoU]



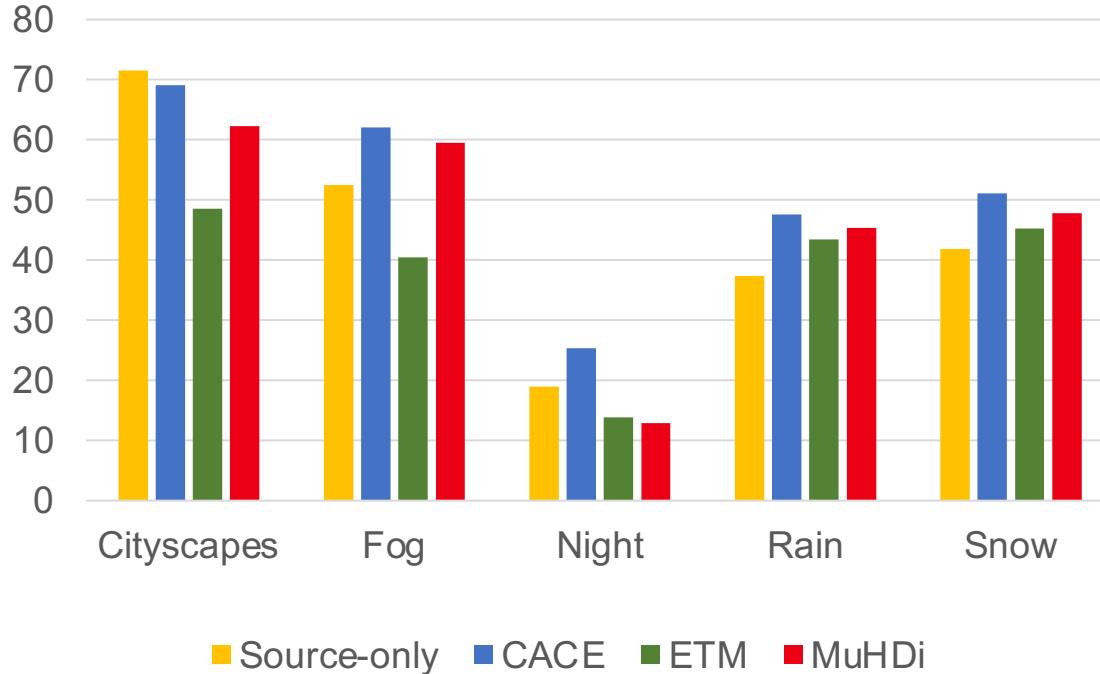
## Average performance (in mean mIoU):

- Source-only: 44.4%
- C-WCT<sup>2</sup>: 50.1%
- CACE: 51.0%
- C-CMD: 49.9%

- All style-transfer-based frameworks achieved **clear performance improvements** over the source-only baseline
- On average CACE performed best, **no advantage** of using **enhanced style transfer** algorithms WCT<sup>2</sup> & CMD
- Relative performance is **dependent on domain**
  - No style transfer algorithm best over all domains
  - CACE best for three out of the four target domains (Fog, Rain, Snow)
- Additional disadvantage: C-WCT<sup>2</sup> and C-CMD suffer from **limitations in implementation**<sup>1</sup>
  - Longer runtime & increased memory consumption

# Output Space Aligned Frameworks ETM & MuHDI Show Significantly Worse Performance than the Input Space Aligned CACE Framework

Domain-wise performance after adapting to the last target domain. [All values measured in mIoU]



**Average performance** (in mean mIoU):

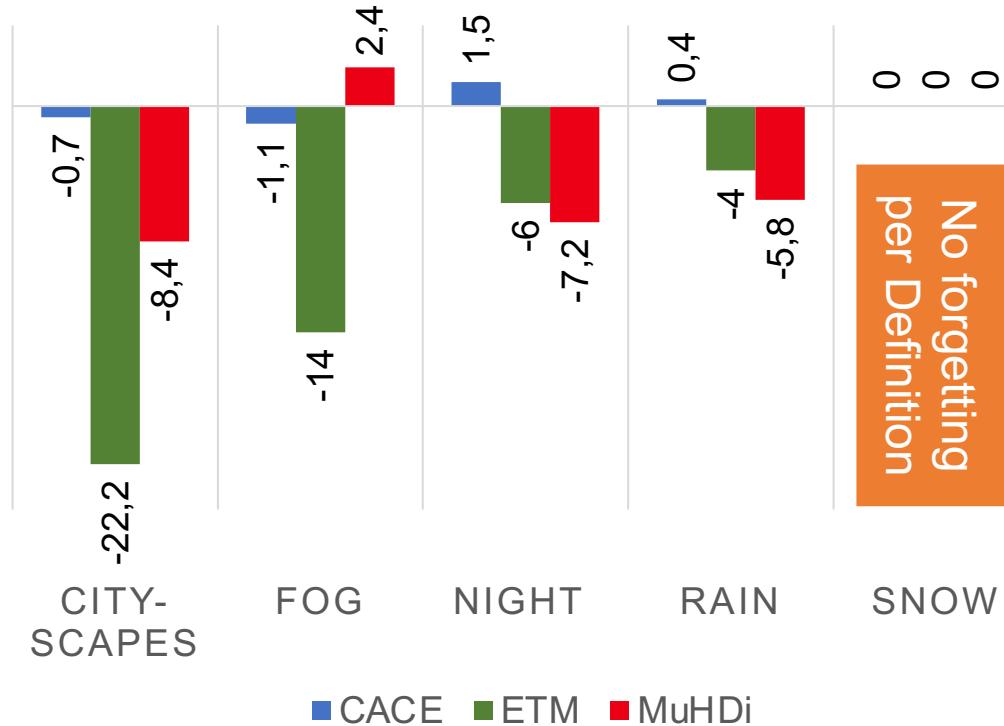
- Source-only: 44.4%
- ETM: 38.3%
- CACE: 51.0%
- MuHDI: 45.6%

- **Input space aligned frameworks (CACE) better than output space aligned frameworks (ETM, MuHDI) for every domain**
- Output space aligned frameworks with **unsatisfactory performance**
  - ETM on average 6% points below source-only
  - MuHDI better but still only in the region of source-only
- ETM & MuHDI with **adequate performance for later trained domains** (Rain, Snow), however, performance for the **early trained domains** (Cityscapes, Fog, Night) **is poor**

→ **Assumption:** Forgetting leads to poor performance on early trained domains

# ETM & MuHDI Both Suffer from Forgetting, They Cannot Maintain Performance when Adapting to New Domains

**Illustration of forgetting.** For each domain, the difference between the performance immediately after adapting to this domain and the performance after adapting to the last target domain is illustrated. [All values measure in % points  $\Delta mIoU$ ]



- **Major problem: Forgetting of old knowledge when adapting to new target domains**
  - ETM & MuHDI both with large performance drops for previous domains when adapting to a new target domain
  - CACE prevents forgetting well, for no domain significant loss in performance after adapting to additional domains



## Conclusion

# Conclusion: The CACE Framework is the Currently Best Approach for Continual UDA in Semantic Segmentation



**Result 1:**  
Comparison of different algorithms for style transfer

- **CACE performed well** and achieves clear improvements over training only on the source domain
- **No advantage** from using the more advanced style transfer algorithms WCT<sup>2</sup> and CMD
- Both more advanced algorithms additionally lead to **limitations in implementation**

**Result 2:**  
Comparison of input space and output space aligned frameworks

- Output space aligned frameworks (ETM, MuHDI) **performed poorly**
- Output space alignment **effective for adapting to only one target domain (UDA)**
- Both frameworks **incapable of sufficiently preventing forgetting** when adapting to additional domains

Out of the five, the **CACE** framework is the currently **best approach** for **continual UDA** in **semantic segmentation**



Extras

# Performance Values for the Cityscapes/ACDC sequence



Domain-wise performance after adapting to the last target domain.  
[All values measure in mIoU]

	Cityscapes	Fog	Night	Rain	Snow	mean mIoU
Source-only	<b>71.5</b>	52.5	19.0	37.4	41.8	44.4
CACE	69.1	<b>62.1</b>	25.3	<b>47.6</b>	<b>51.1</b>	<b>51.0</b>
C-WCT <sup>2</sup>	69.3	60.7	26.1	45.3	48.9	50.1
C-CMD	70.6	60.0	<b>29.3</b>	41.7	48.1	49.9
ETM	48.5	40.5	13.8	43.4	45.3	38.3
MuHDI	62.3	59.5	12.9	45.5	47.8	45.6

# Performance Values for the Cityscapes/ACDC Sequence



**Illustration of pure UDA capabilities:** Performance Values immediately after adapting to each target domain are illustrated for every framework. [All values measure in mIoU]

(a) CACE

	CS	Fog	Night	Rain	Snow
CS	<b>71.5</b>	52.5	19.0	37.4	41.8
Fog	72.4	<b>61.9</b>	19.0	37.4	41.8
Night	70.6	57.6	<b>25.8</b>	44.8	45.9
Rain	69.1	59.0	27.3	<b>45.0</b>	46.1
Snow	70.8	60.8	27.3	46.0	<b>49.9</b>

(b) ETM

	CS	Fog	Night	Rain	Snow
CS	<b>70.7</b>	52.8	16.4	38.9	40.8
Fog	66.4	<b>54.4</b>	15.5	39.7	40.0
Night	63.4	48.9	<b>19.2</b>	37.9	38.4
Rain	56.5	46.1	17.7	<b>47.1</b>	39.7
Snow	48.5	40.5	13.8	43.3	<b>45.3</b>

(c) MuHDI

	CS	Fog	Night	Rain	Snow
CS	<b>70.7</b>	52.8	16.4	38.9	40.8
Fog	67.4	<b>57.1</b>	16.8	41.8	42.0
Night	66.8	56.3	<b>20.1</b>	41.6	40.1
Rain	66.6	60.8	15.3	<b>51.1</b>	47.2
Snow	62.3	59.5	12.9	45.4	<b>47.8</b>

# Performance Values for the Cityscapes/ACDC Sequence



---

**Illustration of pure UDA capabilities with adapted Implementation:** Performance Values immediately after adapting to each target domain are illustrated for adapted CACE and C-CMD. [All values measure in mIoU]

---

	Fog	Night	Rain	Snow
CACE	58.0	24.6	46.6	47.2
C-CMD	59.6	29.2	42.4	48.1

# Performance Values for the Synthetic SYNTHIA Sequence



Domain-wise performance after adapting to the last target domain.

[All values measure in mIoU]

	Dawn	Fall	Fog	Night	Rainnight	Softrain	Spring	Summer	Sunset	Winter	Winternight	mean mIoU
Source-only	70.5	56.5	58.6	53.7	26.2	37.1	56.6	47.6	66.0	40.0	40.7	50.3
Jitter	70.3	62.4	64.0	62.3	29.2	38.7	57.5	53.8	70.6	40.4	49.7	54.4
CACE	<b>71.4</b>	69.0	<b>68.5</b>	66.5	49.1	57.9	67.3	69.7	71.3	53.7	59.1	63.9
C-WCT <sup>2</sup>	71.1	69.6	67.6	68.8	42.5	51.4	68.5	71.4	<b>71.9</b>	58.2	60.1	63.7
C-CMD	69.1	66.5	66.6	65.8	<b>51.0</b>	<b>59.1</b>	66.3	68.0	69.3	60.0	57.5	63.3

# Performance Values for the Synthetic to Real GTA5/Cityscapes/IDD Sequence



**Domain-wise performance after adapting to the last target domain.**  
[All values measure in mIoU]

	GTA-5	Cityscapes	IDD	mean mIoU
Source-only	76.2	33.8	41.2	50.4
Jitter	75.0	40.6	42.7	52.8
CACE	<b>76.2</b>	<b>45.0</b> (-0.2)	44.6	<b>55.3</b>
C-WCT <sup>2</sup>	76.1	44.0 (-1.5)	43.5	54.5
ETM	66.6	38.7 (-2.5)	<b>45.5</b>	50.3
MuHDI	71.2	33.6 (-6.0)	40.0	48.3