

RobustNet: Improving Domain Generalization in Urban-Scene Segmentation via Instance Selective Whitening

CVPR 2021, Oral

Sungha Choi^{*1,3} Sanghun Jung^{*2} Huiwon Yun⁴ Joanne T. Kim³
Seungryong Kim³ Jaegul Choo²

¹LG AI Research ²KAIST ³Korea University ⁴Sogang University

2023.06.30 내부미팅

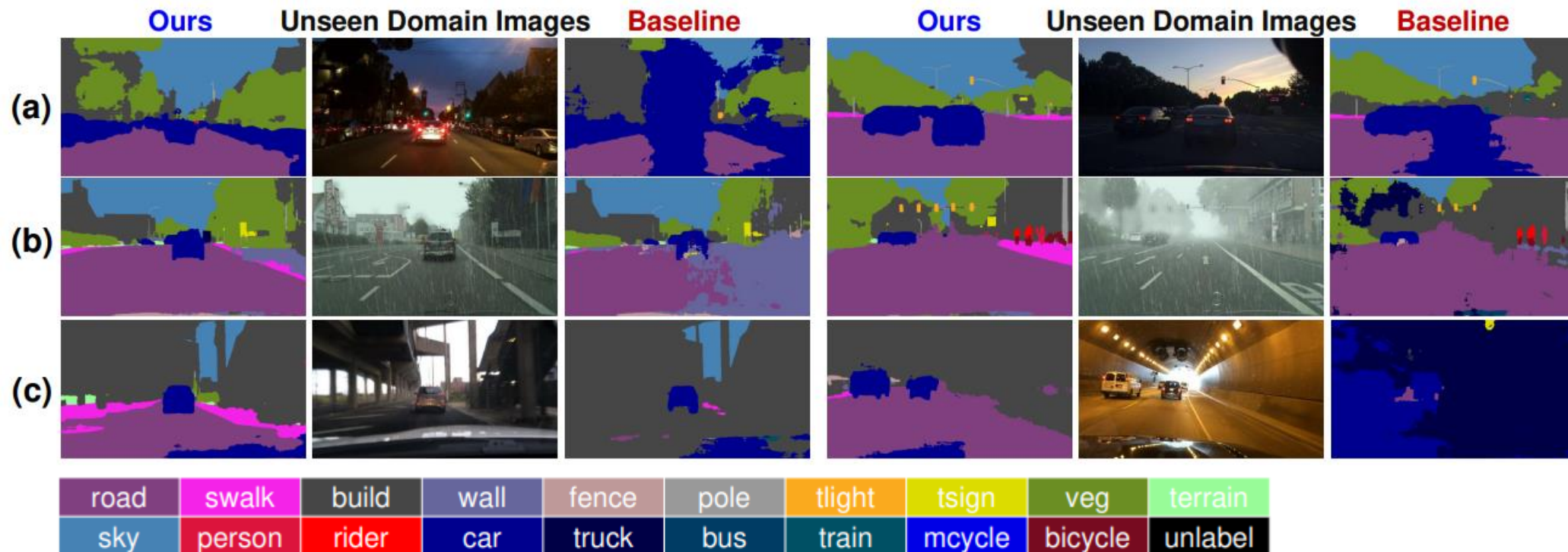
이소영

Contents

1. Abstract
2. Introduction & Related Work
3. Preliminaries
4. Proposed Method
5. Experiments
6. Discussions
7. Conclusions

Abstract

- **Instance selective whitening (ISW) loss**: robustness of the segmentation networks for unseen domains
 - Disentangles the **domain-specific style** and **domain-invariant content** : using feature covariance
 - Selectively **removes only the style** information causing domain shift
- Reasonable predictions for (a) low-illuminated, (b) rainy, and (c) unseen structures



Abstract

- **Instance selective whitening (ISW) loss**: robustness of the segmentation networks for unseen domains
 - Disentangles the **domain-specific style** and **domain-invariant content** : using feature covariance
 - Selectively **removes only the style** information causing domain shift
- Reasonable predictions for (a) low-illuminated, (b) rainy, and (c) unseen structures
- **Simple to use, effective**
- Urban-scene segmentation experiments show the superiority of our approach to existing work

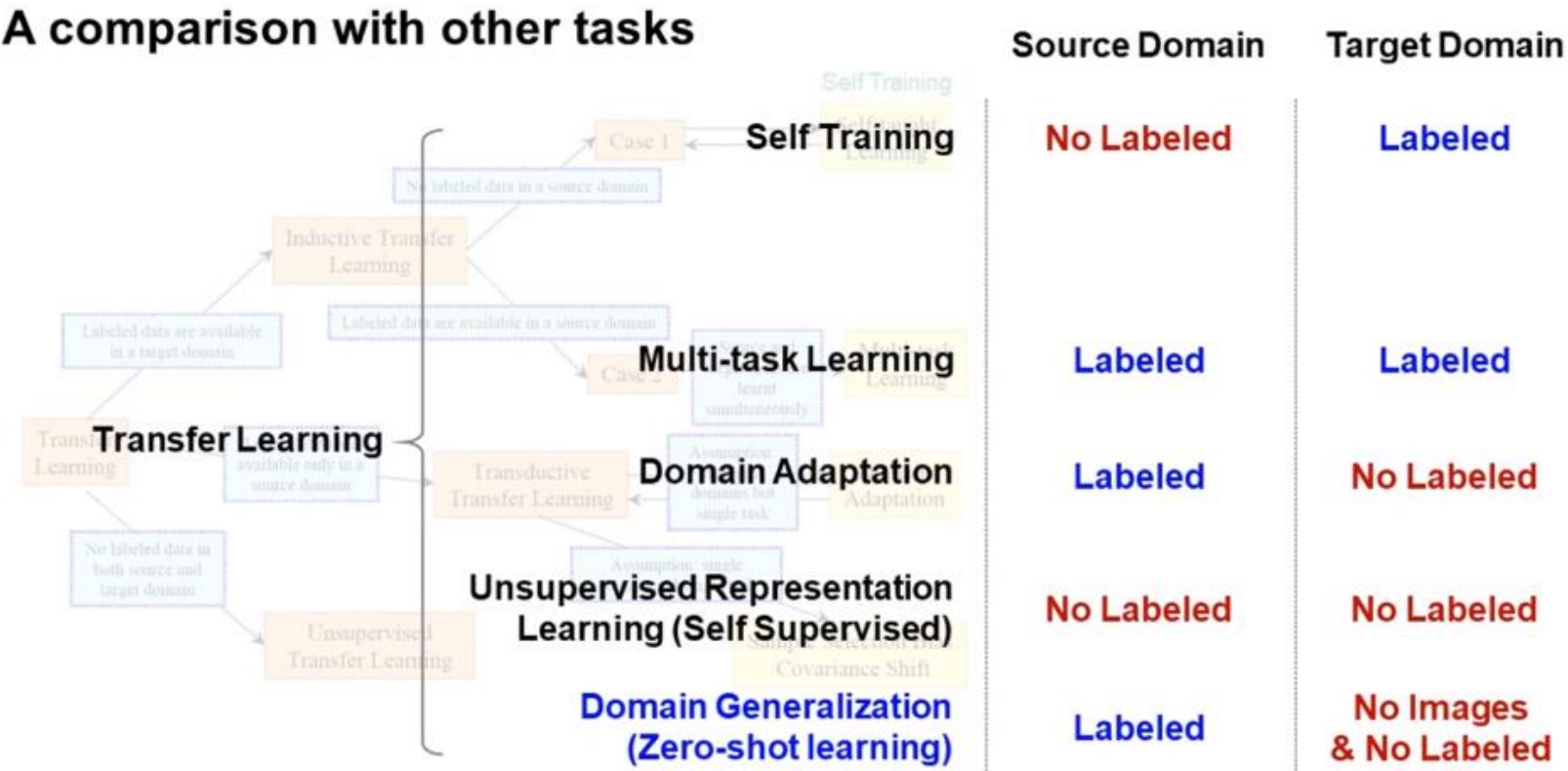
Introduction & Related Work

- When deploying unseen data, fail to perform properly due to **domain shift**
- Need to **reducing the domain gap** between source and target domain



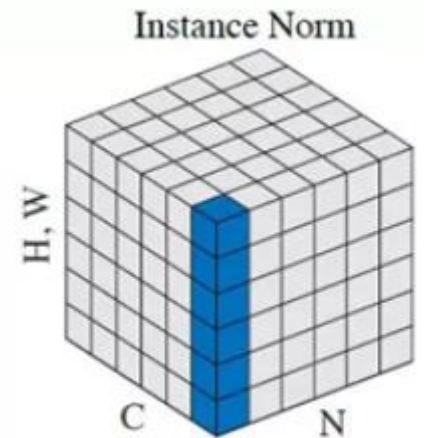
Introduction & Related Work

- Domain Adaptation (DA), Domain Generalization (DG)
- DA: access to the samples in the target domain



Introduction & Related Work

- Domain Generalization (DG): robustness to arbitrary unseen domain
 - Shared across multiple source domains: costly, labor-intensive, highly depends on the number of source datasets
- More effective?
 - IBN-Net: exploiting **instance normalization**
 - Instance normalization: standardizes features while **not considering the correlation between channels**
 - **Feature covariance** contains **domain-specific style** (texture, color) → instance norm not sufficient



Introduction & Related Work

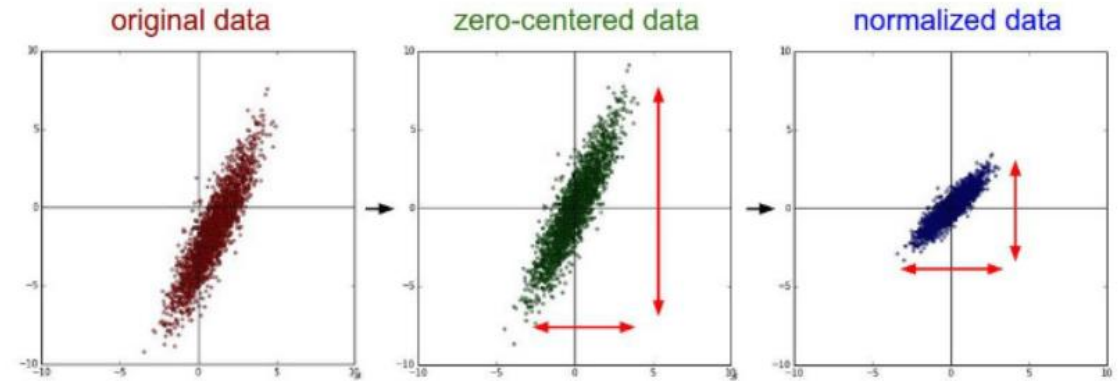
- **Whitening transformation:** removes feature correlation, each feature have unit variance
 - Feature whitening effectively eliminates domain-specific style information
 - Shown in image translation, style transfer, domain adaptation
 - **Not yet fully explored in DG**
- Semantic segmentation in DG: majority of the DG methods mainly focused on image classification
- Main scope of this paper
 - Adopting whitening transformation to DG
 - Decoupling the two factors: **domain-specific style** and **domain-invariant content**
 - **Selectively removing the domain-specific style**
 - Proposed loss can **easily be used**
 - Show **superiority to urban-scene segmentation** in DG (qualitative, quantitative)

Preliminaries

- Normalization (Standardization)

- 평균=0, 표준편차=1이 되게 변환

$$Z = \frac{X - \mu}{\sigma}$$



- Whitening

- 데이터의 평균=0, 공분산=단위행렬로 갖는 정규분포 형태로 변환 $\mathbf{x} \sim \mathcal{N}_D(\mu, \Sigma) \xrightarrow{\text{whiten}} \mathbf{z} \sim \mathcal{N}_D(0, \mathbf{I})$

- Decorrelation + Standardization

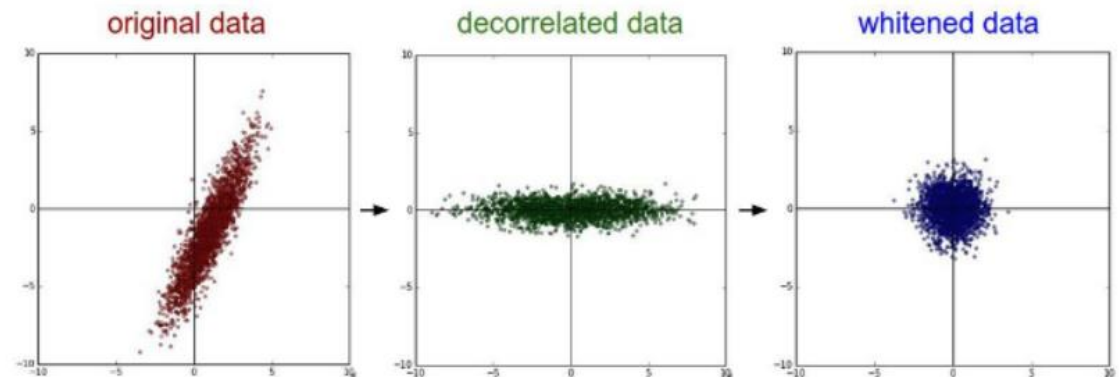
$$Z = W(X - \mu)$$

$$W^T W = \Sigma^{-1}$$

$$\Sigma = U \Lambda U^T$$

$$\Sigma^{-1/2} = U \Lambda^{-1/2} U^T$$

$$W = Q U \Lambda^{-1/2} U^T = Q \Sigma^{-1/2}$$



Preliminaries

- Covariance matrix

Definition Let X be a $K \times 1$ random vector. The covariance matrix of X , or variance-covariance matrix of X , denoted by $\text{Var}[X]$, is defined as follows:

$$\text{Var}[X] = E[(X - E[X])(X - E[X])^T]$$

Let X_1, \dots, X_K denote the K components of the vector X .

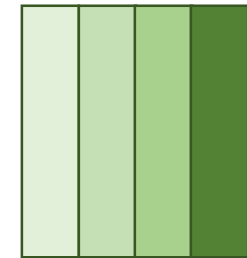
From the definition of $\text{Var}[X]$, it can easily be seen that $\text{Var}[X]$ is a $K \times K$ matrix with the following structure:

$$\begin{aligned} \text{Var}[X] &= E \begin{bmatrix} (X_1 - E[X_1])(X_1 - E[X_1]) & \dots & (X_1 - E[X_1])(X_K - E[X_K]) \\ \vdots & \ddots & \vdots \\ (X_K - E[X_K])(X_1 - E[X_1]) & \dots & (X_K - E[X_K])(X_K - E[X_K]) \end{bmatrix} \\ &= \begin{bmatrix} E[(X_1 - E[X_1])^2] & \dots & E[(X_1 - E[X_1])(X_K - E[X_K])] \\ \vdots & \ddots & \vdots \\ E[(X_K - E[X_K])(X_1 - E[X_1])] & \dots & E[(X_K - E[X_K])^2] \end{bmatrix} \\ &= \begin{bmatrix} \text{Var}[X_1] & \dots & \text{Cov}[X_1, X_K] \\ \vdots & \ddots & \vdots \\ \text{Cov}[X_K, X_1] & \dots & \text{Var}[X_K] \end{bmatrix} \end{aligned}$$

X



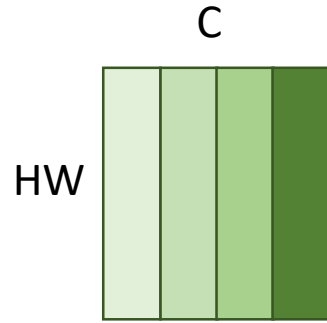
$X_1 \ X_2 \ \dots \ X_K$



Preliminaries

- Whitening transformation (WT)

$$\mathbf{X} \in \mathbb{R}^{C \times HW}$$



- WT is a linear transformation that makes
 - Variance term of each channel equal to 1
 - Covariances between each pair of channels equal to 0

$$\Rightarrow \tilde{\mathbf{X}} \cdot \tilde{\mathbf{X}}^T = (HW) \cdot \mathbf{I} \in \mathbb{R}^{C \times C}$$

Covariance matrix

$$\text{Var}[\mathbf{X}] = \mathbb{E}[(\mathbf{X} - \mathbb{E}[\mathbf{X}])(\mathbf{X} - \mathbb{E}[\mathbf{X}])^T] \quad C = \frac{1}{m} \mathbf{X} \mathbf{X}^T$$

Covariance matrix

$$\begin{bmatrix} \text{Var}[X_1] & \dots & \text{Cov}[X_1, X_K] \\ \vdots & \ddots & \vdots \\ \text{Cov}[X_K, X_1] & \dots & \text{Var}[X_K] \end{bmatrix}$$

Preliminaries

- Whitening transformation (WT)

$$\mathbf{x} \sim \mathcal{N}_D(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \xrightarrow{\text{whiten}} \mathbf{z} \sim \mathcal{N}_D(0, \mathbf{I})$$

$$\tilde{\mathbf{X}} = \boldsymbol{\Sigma}_{\mu}^{-\frac{1}{2}} (\mathbf{X} - \boldsymbol{\mu} \cdot \mathbf{1}^{\top}) \quad \boldsymbol{\mu} = \frac{1}{HW} \mathbf{X} \cdot \mathbf{1} \in \mathbb{R}^{C \times 1}$$

$$\boldsymbol{\Sigma}_{\mu} = \frac{1}{HW} (\mathbf{X} - \boldsymbol{\mu} \cdot \mathbf{1}^{\top}) (\mathbf{X} - \boldsymbol{\mu} \cdot \mathbf{1}^{\top})^{\top} \in \mathbb{R}^{C \times C}$$

- Compute $\boldsymbol{\Sigma}_{\mu}^{-\frac{1}{2}}$? \rightarrow using eigen decomposition

- The covariance matrix is real and symmetric \rightarrow can be diagonalized using eigen decomposition

$$\boldsymbol{\Sigma}_{\mu} = \mathbf{Q} \boldsymbol{\Lambda} \mathbf{Q}^{\top} \quad \left(\begin{array}{l} \mathbf{Q} \in \mathbb{R}^{C \times C} : \text{Orthogonal matrix of eigenvectors} \\ \boldsymbol{\Lambda} \in \mathbb{R}^{C \times C} : \text{Diagonal matrix each eigenvalue of the corresponding eigenvector from Q} \end{array} \right.$$

$$\Rightarrow \boldsymbol{\Sigma}_{\mu}^{-\frac{1}{2}} = \mathbf{Q} \boldsymbol{\Lambda}^{-\frac{1}{2}} \mathbf{Q}^{\top}$$

Preliminaries

- Limitations of WT
 - Eigenvalue decomposition : computationally expensive, prevents the gradient back-propagation

Universal Style Transfer via Feature Transforms

NeurIPS, 2017

Yijun Li
UC Merced
yli62@ucmerced.edu

Chen Fang
Adobe Research
cfang@adobe.com

Jimei Yang
Adobe Research
jimyang@adobe.com

Zhaowen Wang
Adobe Research
zhawang@adobe.com

Xin Lu
Adobe Research
xinl@adobe.com

Ming-Hsuan Yang
UC Merced, NVIDIA Research
mhyang@ucmerced.edu

Whitening transform. Before whitening, we first center f_c by subtracting its mean vector m_c . Then we transform f_c linearly as in (2) so that we obtain \hat{f}_c such that the feature maps are uncorrelated ($\hat{f}_c \hat{f}_c^\top = I$),

$$\hat{f}_c = E_c D_c^{-\frac{1}{2}} E_c^\top f_c, \quad (2)$$

where D_c is a diagonal matrix with the eigenvalues of the covariance matrix $f_c f_c^\top \in \mathbb{R}^{C \times C}$, and E_c is the corresponding orthogonal matrix of eigenvectors, satisfying $f_c f_c^\top = E_c D_c E_c^\top$.

Efficiency. In Table 2 (3rd row), we also compare our approach with other methods in terms of efficiency. The method by Gatys et al. [9] is slow due to loops of optimization and usually requires at least 500 iterations to generate good results. The methods [27] and [15] are efficient as the scheme is based on one feed-forward pass with a trained network. The approach [3] is feed-forward based but relatively slower as the feature swapping operation needs to be carried out for thousands of patches. Our approach is also efficient but a little bit slower than [27, 15] because we have a eigenvalue decomposition step in WCT. But note that the computational cost on this step will not increase along with the image size because the the dimension of covariance matrix only depends on filter numbers (or channels), which is at most 512 (Relu_5_1). Currently the decomposition step is implemented based on CPU. Our future work includes more efficient GPU implementations of the proposed algorithm.

Preliminaries

- **Without** the eigen-decomposition?
 1. Approximating the whitening transformation matrix using Newton's iteration (ex, IterNorm)
 2. Deep Whitening Transformation (DWT) in GDWCT : **the covariance matrix close to the identity matrix**
 - Domain-specific style and domain-invariant content are **simultaneously encoded in the covariance** of the feature map
 - Whitening all covariance elements \rightarrow diminish feature discrimination and distort the boundary of an object

$$\mathcal{L}_{\text{DWT}} = \mathbb{E}[\|\Sigma_{\mu} - \mathbf{I}\|_1],$$

Proposed Method – (1)

1. Instance Whitening Loss (IW)

- DWT loss decompose $\Sigma_{\mu} \begin{cases} \Sigma_{\mu(i,i)} \\ \Sigma_{\mu(i,j)} \end{cases}$

$$\mathcal{L}_{\text{DWT}} = \mathbb{E}[\|\Sigma_{\mu} - \mathbf{I}\|_1], \begin{cases} \|\Sigma_{\mu(i,i)} - 1\|_1 = \left\| \frac{\mathbf{x}_i^{\top} \cdot \mathbf{x}_i}{HW} - 1 \right\|_1 = \left\| \frac{|\mathbf{x}_i| |\mathbf{x}_i| \cos 0^{\circ}}{HW} - 1 \right\|_1 \\ \|\Sigma_{\mu(i,j)}\|_1 = \left\| \frac{\mathbf{x}_i^{\top} \cdot \mathbf{x}_j}{HW} \right\|_1 = \left\| \frac{|\mathbf{x}_i| |\mathbf{x}_j| \cos \theta}{HW} \right\|_1 \end{cases}$$

- **Diagonal** covariance matrix $\rightarrow \mathbf{1}$
- **Off-diagonal** of the covariance matrix $\rightarrow \mathbf{0}$
- It is **difficult to optimize** both at the same time

Proposed Method – (1)

1. Instance Whitening Loss (IW)

- To address: feature map \mathbf{X} can first be standardized using instance normalization

$$\mathbf{X}_s = (\text{diag}(\boldsymbol{\Sigma}_\mu))^{-\frac{1}{2}} \odot (\mathbf{X} - \boldsymbol{\mu} \cdot \mathbf{1}^\top)$$

$$\boldsymbol{\Sigma}_s = \frac{1}{HW} (\mathbf{X}_s)(\mathbf{X}_s)^\top \in \mathbb{R}^{C \times C}$$

- Instance Whitening Loss**

- Covariance matrix is symmetric \rightarrow **loss** can be applied only to the **strict upper triangular part**

$$\mathcal{L}_{IW} = \mathbb{E}[\|\boldsymbol{\Sigma}_s \odot \mathbf{M}\|_1]$$

$$\mathbf{M}_{i,j} = \begin{cases} 0, & \text{if } i \geq j \\ 1, & \text{otherwise} \end{cases} \quad 0 \leq i, j < C$$

M

0	1	1	1	1	1	1
0	0	1	1	1	1	1
0	0	0	1	1	1	1
0	0	0	0	1	1	1
0	0	0	0	0	1	1
0	0	0	0	0	0	1
0	0	0	0	0	0	0

- WT is a linear transformation that makes
 - Variance term of each channel equal to 1
 - Covariances between each pair of channels equal to 0

$$\Rightarrow \tilde{\mathbf{X}} \cdot \tilde{\mathbf{X}}^\top = (HW) \cdot \mathbf{I} \in \mathbb{R}^{C \times C}$$

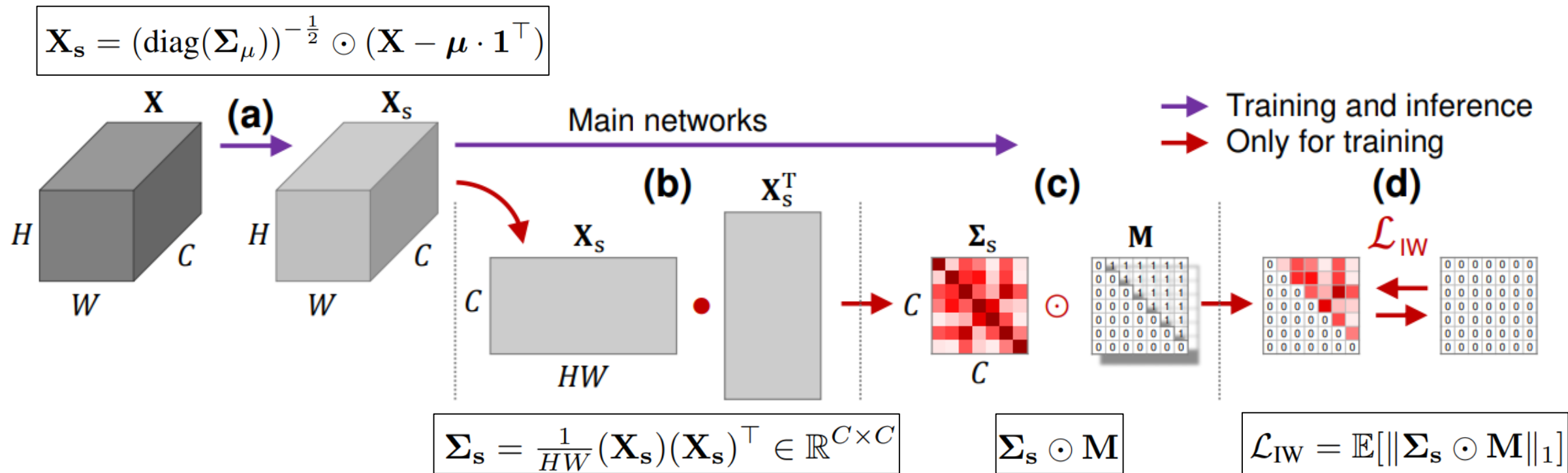
Covariance matrix

$$\text{Var}[X] = \mathbb{E}[(X - \mathbb{E}[X])(X - \mathbb{E}[X])^\top]$$

$$\begin{bmatrix} \text{Var}[X_1] & \dots & \text{Cov}[X_1, X_K] \\ \vdots & \ddots & \vdots \\ \text{Cov}[X_K, X_1] & \dots & \text{Var}[X_K] \end{bmatrix}$$

Proposed Method – (1)

1. Instance Whitening Loss (IW)



Proposed Method – (2)

2. Margin-based relaxation of whitening loss (IRW)

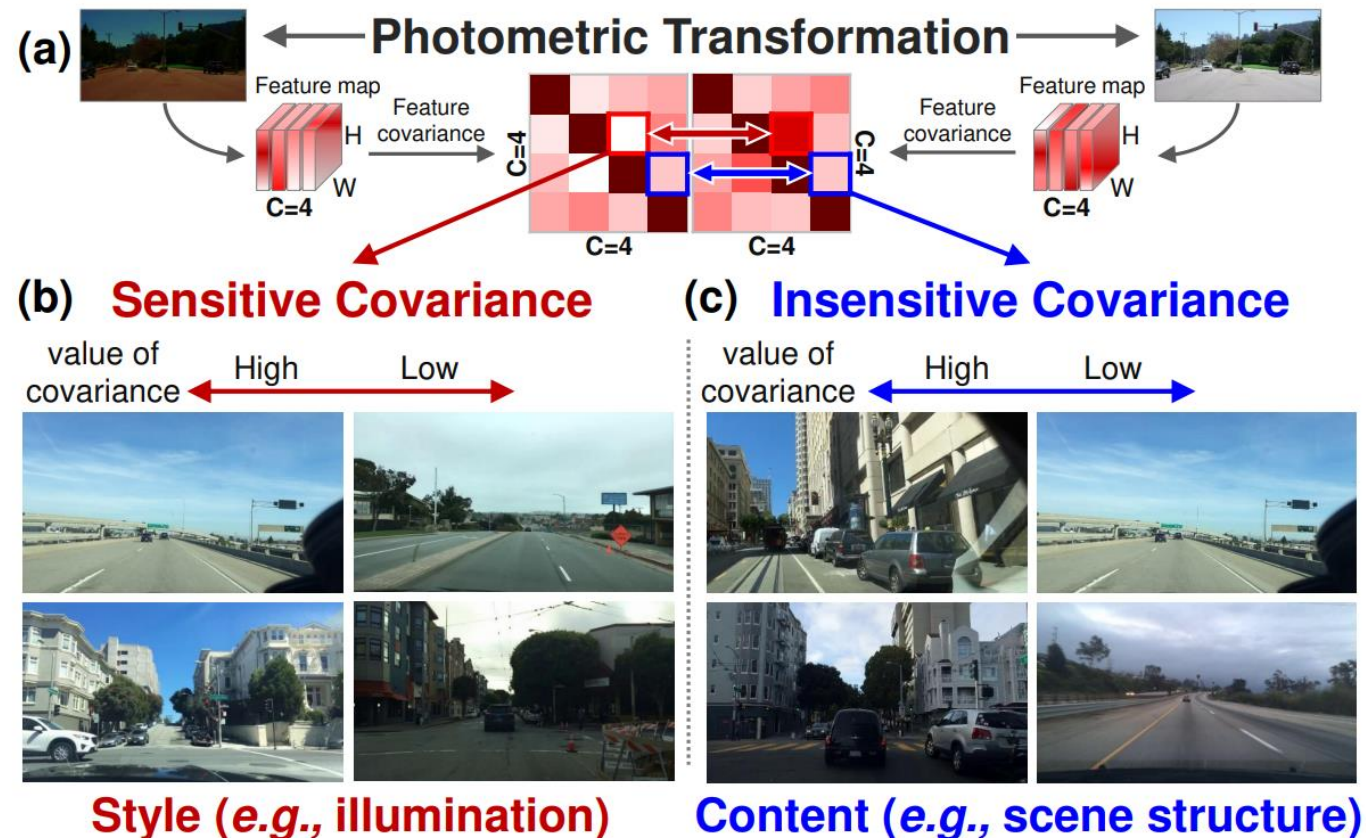
$$\mathcal{L}_{\text{IRW}} = \max(\mathbb{E}[\|\Sigma_s \odot \mathbf{M}\|_1] - \delta, 0)$$

Models (GTAV)	C	B	M	S	G
Baseline	28.95	25.14	28.18	26.23	73.45
Ours (IRW), $\delta=1/16$	32.49	32.53	37.51	27.77	72.18
Ours (IRW), $\delta=1/32$	33.30	33.17	38.03	27.43	71.96
Ours (IRW), $\delta=1/64$	33.57	33.18	38.42	27.29	71.96
Ours (IRW), $\delta=1/128$	32.85	32.40	37.36	27.43	72.21
Ours (IRW), $\delta=1/256$	32.45	32.32	37.93	27.48	72.12
Ours (IW)	33.21	32.67	37.35	27.57	72.06

Proposed Method – (3)

3. Separation Covariance Elements (ISW)

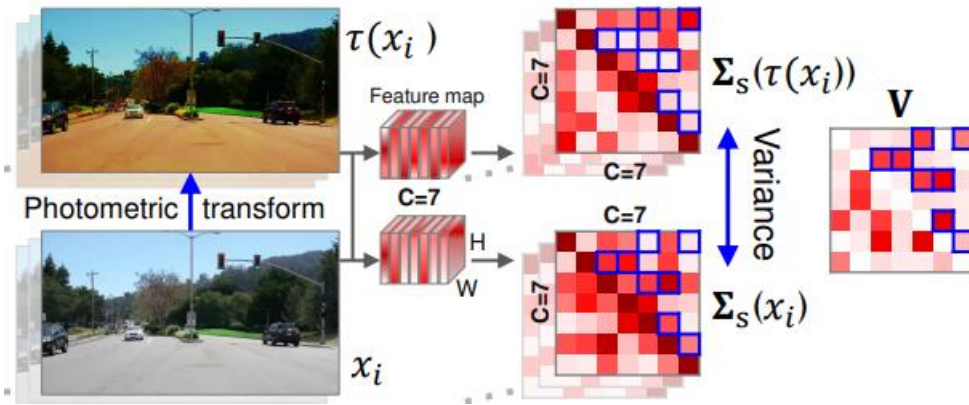
- Separate the covariance terms into two groups: domain-specific style and domain-invariant content
- Selectively **remove only the style-encoded covariances** that cause the domain shift.



Proposed Method – (3)

3. Separation Covariance Elements (ISW)

- Original / photometric transformed image \rightarrow each covariance matrices \rightarrow differences \rightarrow variance matrix (V)



$$\mathbf{V} = \frac{1}{N} \sum_{i=1}^N \sigma_i^2$$

$$\mu_{\Sigma_i} = \frac{1}{2} (\Sigma_s(x_i) + \Sigma_s(\tau(x_i)))$$

$$\sigma_i^2 = \frac{1}{2} \left((\Sigma_s(x_i) - \mu_{\Sigma_i})^2 + (\Sigma_s(\tau(x_i)) - \mu_{\Sigma_i})^2 \right)$$

- Apply k-means clustering (k = 3, m = 1)

$$C = \{c_1, c_2, \dots, c_k\} \begin{cases} G_{low} = \{c_1, \dots, c_m\} \\ G_{high} = \{c_{m+1}, \dots, c_k\} \end{cases}$$

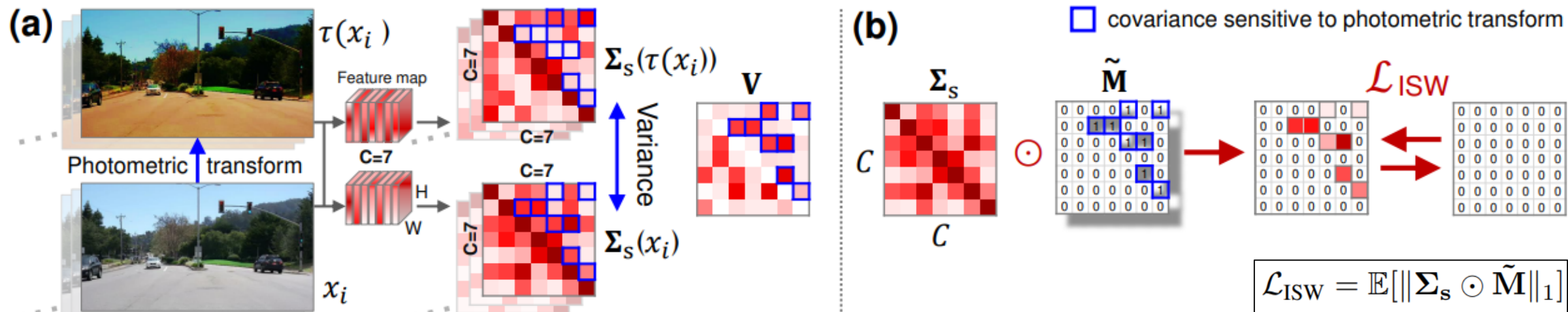
$$\tilde{M}_{i,j} = \begin{cases} 1, & \text{if } \mathbf{V}_{i,j} \in G_{high} \\ 0, & \text{otherwise} \end{cases}$$

$$\mathcal{L}_{ISW} = \mathbb{E}[\|\Sigma_s \odot \tilde{\mathbf{M}}\|_1]$$

Models (GTAV)	C	B	M	S	G
Baseline	28.95	25.14	28.18	26.23	73.45
Ours (ISW), k=2	35.46	35.00	39.38	27.70	72.08
Ours (ISW), k=3	36.58	35.20	40.33	28.30	72.10
Ours (ISW), k=5	34.84	33.58	39.25	27.52	72.31
Ours (ISW), k=10	33.58	33.76	38.96	27.68	72.24
Ours (ISW), k=20	33.66	33.29	38.70	27.47	72.10
Ours (IW)	33.21	32.67	37.35	27.57	72.06

Proposed Method – (3)

3. Separation Covariance Elements (ISW)



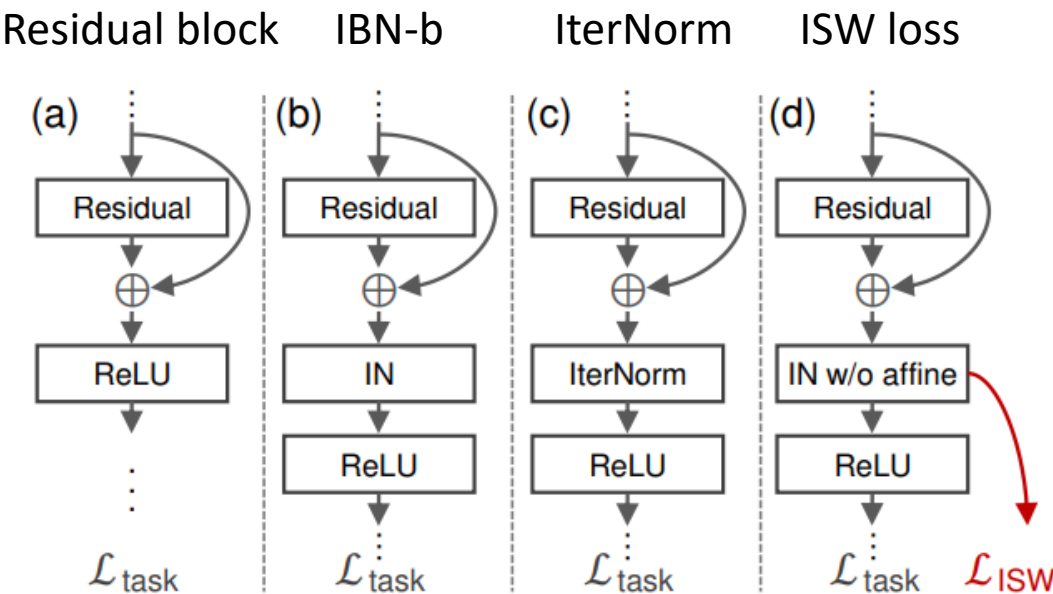
Proposed Method

Models (GTAV)	C	B	M	S	G
Ours (ISW), $\gamma=0.4$	35.60	34.07	38.98	28.10	71.96
Ours (ISW), $\gamma=0.6$	36.58	35.20	40.33	28.30	72.10
Ours (ISW), $\gamma=0.8$	35.73	34.01	39.69	27.44	71.96

- Network architecture with proposed ISW loss
- Simply add our proposed ISW loss to the instance normalization layer
 - $\gamma = 0.6$, L (layer) = 3 (IBN-Net)

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{task}} + \lambda \left(\frac{1}{L} \sum_i^L \mathcal{L}_{\text{ISW}}^i \right)$$

- Architecture comparison with other methods



Experiments

- Implementation
 - Baseline architecture: **DeepLabV3+**
 - Photometric transformation in ISW: **color jittering, gaussian blur**
- Datasets
 - **Real-world** datasets: Cityscapes, BDD-100K, Mapillary
 - **Synthetic** datasets: GTAV, SYNTHIA

Experiments

- Effectiveness of ISW (1)
 - Metric: mIoU (%)
 - Cityscapes (C), BDD-100K (B), Mapillary (M), SYNTHIA (S), GTAV (G)

Models (GTAV)	C	B	M	S	G
Baseline	28.95	25.14	28.18	26.23	73.45
[†] SW [45]	29.91	27.48	29.71	27.61	73.50
[†] IBN-Net [44]	33.85	32.30	37.75	27.90	72.90
[†] IterNorm [22]	31.81	32.70	33.88	27.07	73.19
Ours (IW)	33.21	32.67	37.35	27.57	72.06
Ours (IRW)	33.57	33.18	38.42	27.29	71.96
Ours (ISW)	36.58	35.20	40.33	28.30	72.10

Models (Cityscapes)	B	M	G	S	C
Baseline	44.96	51.68	42.55	23.29	77.51
[†] SW [45]	48.49	55.82	44.87	26.10	77.30
[†] IBN-Net [44]	48.56	57.04	45.06	26.14	76.55
[†] IterNorm [22]	49.23	56.26	45.73	25.98	76.02
Ours (IW)	48.19	58.90	45.21	25.81	76.06
Ours (IRW)	48.67	59.20	45.64	26.05	76.13
Ours (ISW)	50.73	58.64	45.00	26.20	76.41

Baseline: DeeplabV3+

Backbone: ResNet-50 with an output stride of 16

[†] : own re-implemented models

Experiments

- Effectiveness of ISW (2)
 - Metric: mIoU (%)
 - Cityscapes (C), BDD-100K (B), Mapillary (M), SYNTHIA (S), GTAV (G)

Models (GTAV)	C	B	M	S	G
Baseline	25.56	22.17	28.60	23.33	66.47
[†] IBN-Net [44]	27.10	31.82	34.89	25.56	65.44
Ours (ISW)	30.98	32.06	35.31	24.31	64.99
Baseline	25.92	25.73	26.45	24.03	68.12
[†] IBN-Net [44]	30.14	27.66	27.07	24.98	67.66
Ours (ISW)	30.86	30.05	30.67	24.43	67.48

Baseline: DeeplabV3+

Backbones: ShuffleNetV2 (up), MobileNetV2 (down)

[†] : own re-implemented models

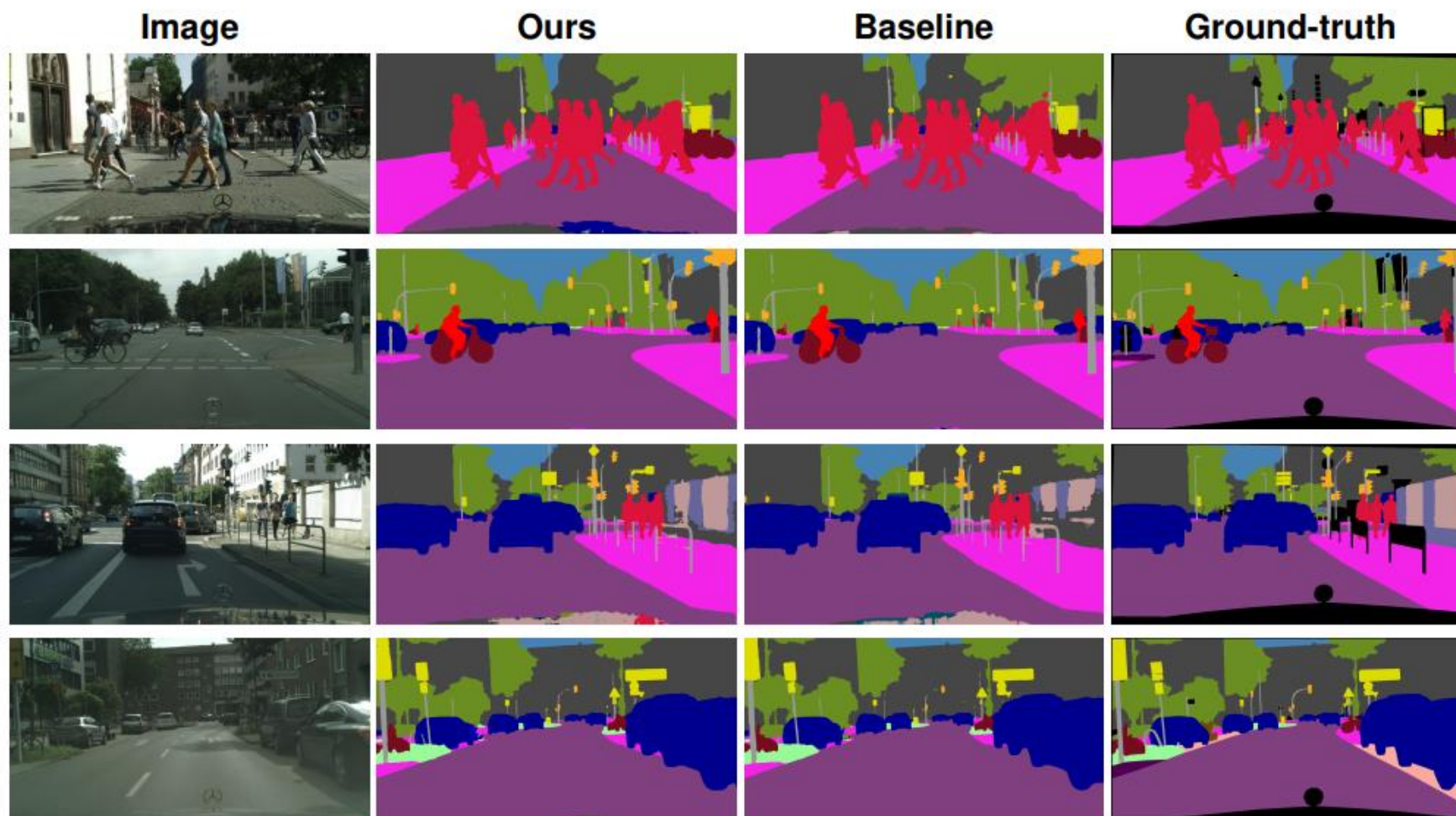
Models (G + S)	C	B	M	G	S
Baseline	35.46	25.09	31.94	68.48	67.99
IBN-Net	35.55	32.18	38.09	69.72	66.90
Ours	37.69	34.09	38.49	68.26	68.77

Backbone: ResNet-50 with an output stride of 16

Experiments

road	swalk	build	wall	fence	pole	tlght	tsign	veg	terrain
sky	person	rider	car	truck	bus	train	mcycle	bicycle	unlabel

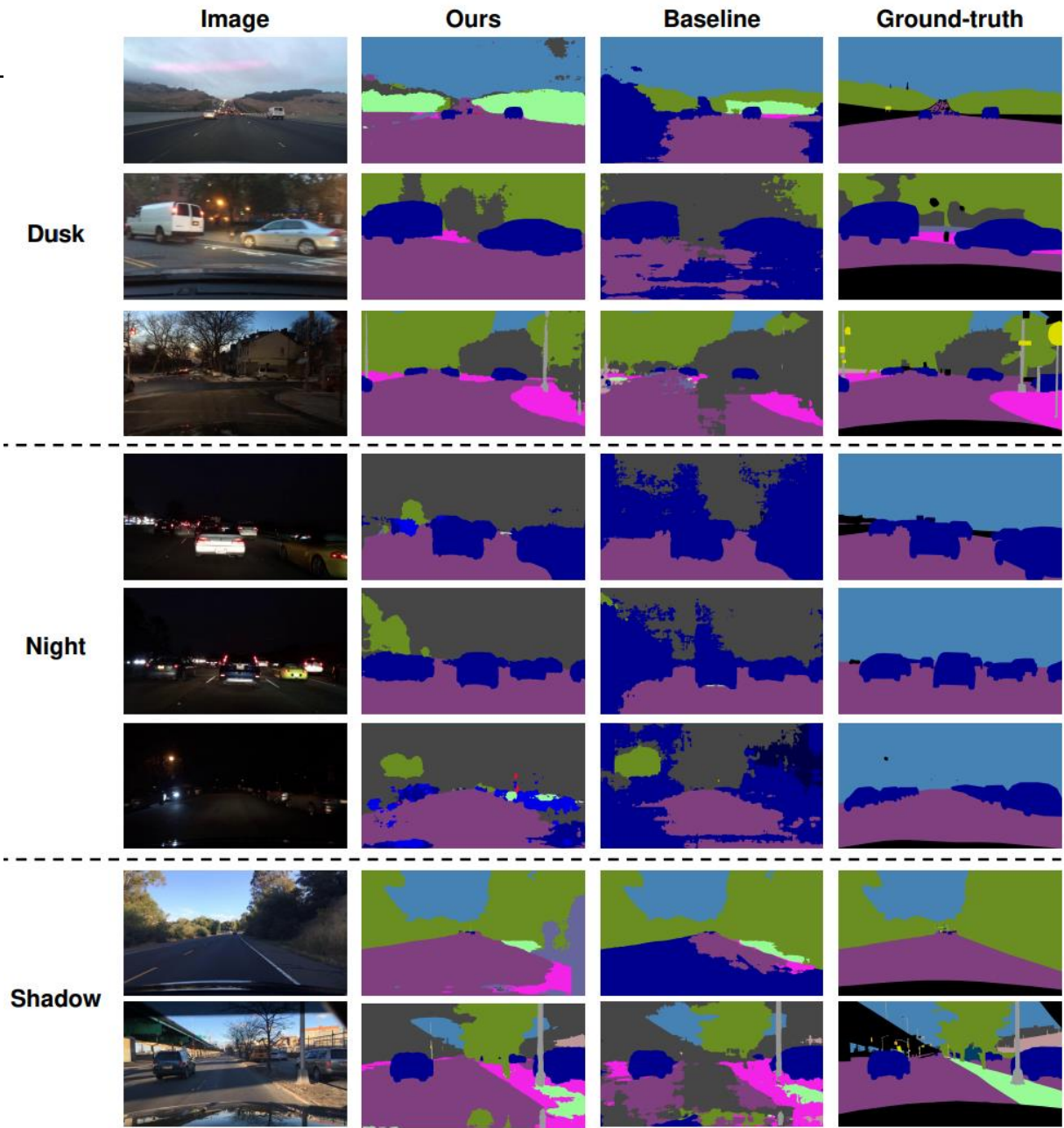
- Effectiveness of ISW (2)
 - Segmentation results on **seen domain** images (Cityscape)



Experiments

- Effectiveness of ISW (2)
 - Segmentation results under **illumination changes**
 - Train C, inference B

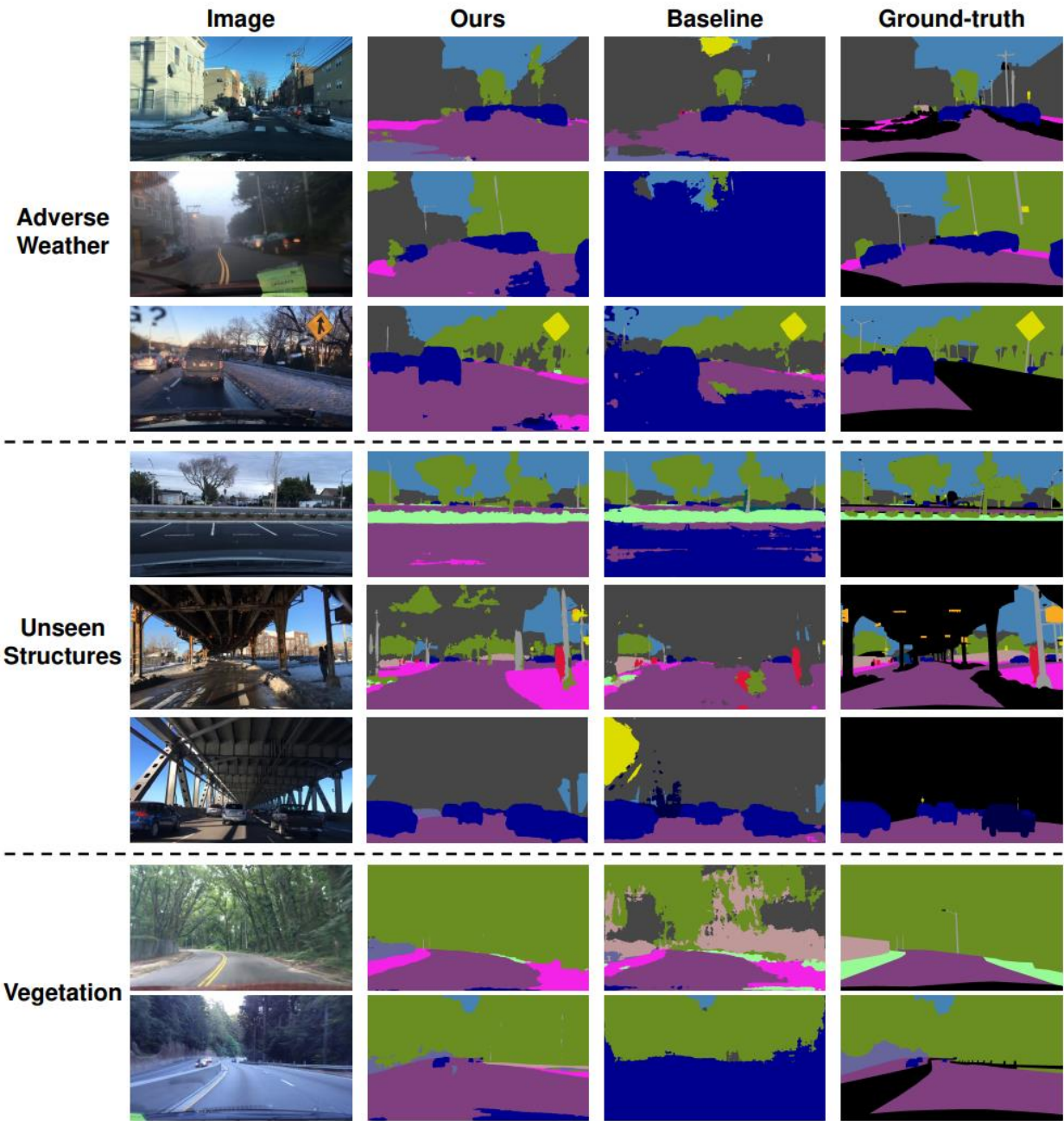
road	swalk	build	wall	fence	pole	tlight	tsign	veg	terrain
sky	person	rider	car	truck	bus	train	mcycle	bicycle	unlabel



Experiments

- Effectiveness of ISW (2)
 - Segmentation results under **various circumstances**
 - Train C, inference B

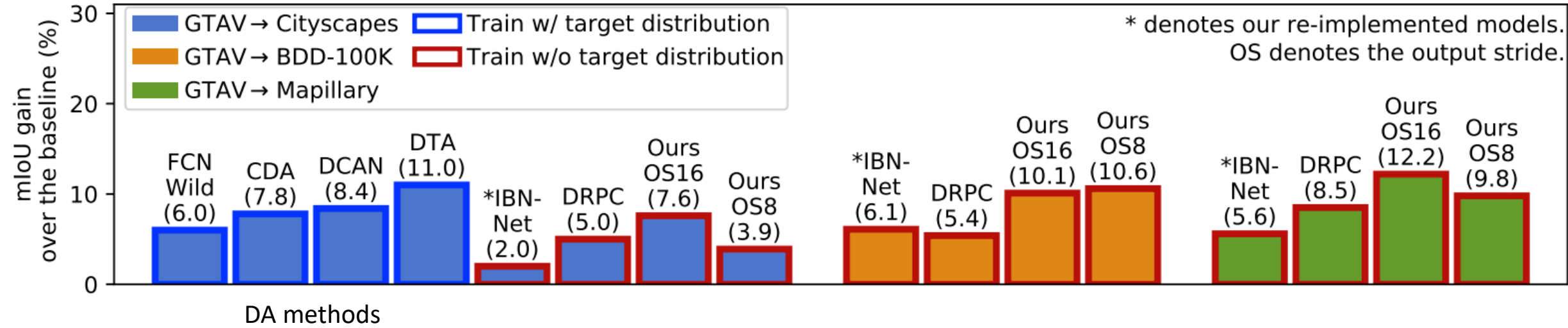
road	swalk	build	wall	fence	pole	tlght	tsign	veg	terrain
sky	person	rider	car	truck	bus	train	mcycle	bicycle	unlabel



Experiments

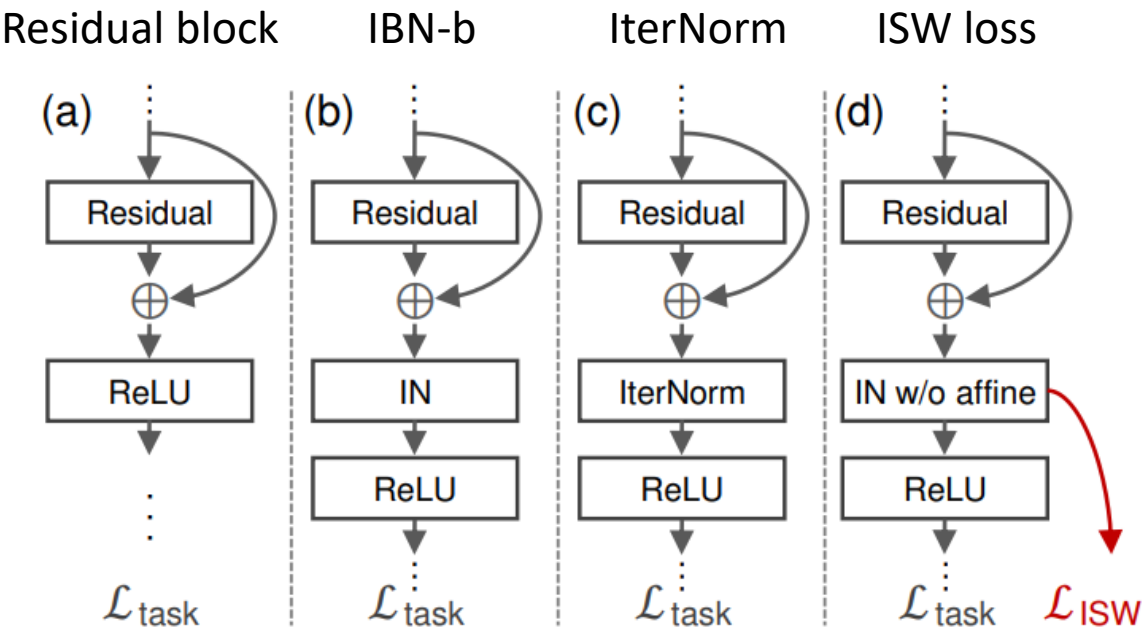
- Comparison with other DG and DA methods

Models (GTAV)		C		B		M	
DG methods	Baseline	22.20	7.40 ↑	N/A		N/A	
	IBN-Net [44]	29.60					
	Baseline	32.45	4.97↑	26.73	5.41↑	25.66	8.46↑
	DRPC [64]	37.42		32.14		34.12	
	Baseline	28.95	7.63↑	25.14	10.06↑	28.18	12.15↑
	Ours (ISW)	36.58		35.20		40.33	



Experiments

- Computational cost analysis



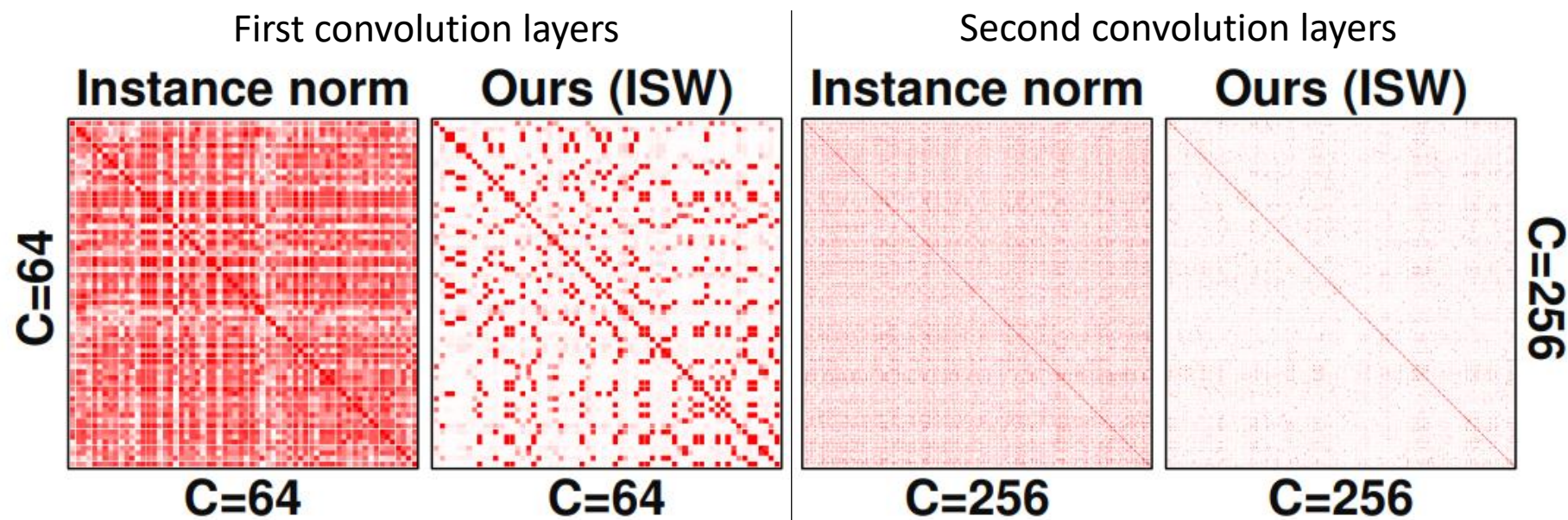
Share the same network architecture
Different **normalization methods**

Models	# of Params	GFLOPS	Inference Time (ms)
Baseline	45.082M	554.31	10.48
[†] IBN-Net [44]	45.083M	554.31	10.51
[†] IterNorm [22]	45.081M	554.31	40.31
Ours	45.081M	554.31	10.43

→ whitening transformation
without additional computational cost

Experiments

- Comparison of **covariance matrices**



ISW selectively eliminates the covariance

Experiments

- **Reconstructing** images with whitened features
 - Using U-Net



Image contents are properly maintained
Style (such as illumination and colors) vanish

Discussions

- Affine parameters
 - Adding affine parameter or 1x1 convolution layer after the normalization layer: not improve
 - Conjecture: above approach do not have sufficient complexity in recovering the original distribution
- Photometric transformation
 - We found that applying color transform and Gaussian blur does not harm the content information
 - Expect various photometric augmentation techniques

Conclusions

- Focused on solving the **domain generalization** problem in **urban-scene segmentation**
- A novel **instance selective whitening (ISW) loss**
 - Disentangling the covariances (of the intermediate features): the style- and content-related ones
 - **Suppressing only the style-related covariances** → learn the domain-invariant feature representation

END