# WEDGE: Web-Image Assisted Domain Generalization for Semantic Segmentation

Namyup Kim[1], Taeyoung Son[2], Jaehyun Pahk[1], Cuiling Lan[3], Wenjun Zeng[4], and Suha Kwak[1]

*Abstract*— Domain generalization for semantic segmentation is highly demanded in real applications, where a trained model is expected to work well in previously unseen domains. One challenge lies in the lack of data which could cover the diverse distributions of the possible unseen domains for training. In this paper, we propose a WEb-image assisted Domain GEneralization (WEDGE) scheme, which is the first to exploit the diversity of web-crawled images for generalizable semantic segmentation. To explore and exploit the real-world data distributions, we collect web-crawled images which present large diversity in terms of weather conditions, sites, lighting, camera styles, etc. We also present a method which injects styles of the web-crawled images into training images on-the-fly during training, which enables the network to experience images of diverse styles with reliable labels for effective training. Moreover, we use the web-crawled images with their predicted pseudo labels for training to further enhance the capability of the network. Extensive experiments demonstrate that our method clearly outperforms existing domain generalization techniques.

## I. INTRODUCTION

Semantic segmentation has played crucial roles in many applications like autonomous vehicle and augmented reality. Recent advances in this field are mainly attributed to the development of deep neural networks, whose success depends heavily on the availability of a large-scale annotated dataset for training. However, building a large training dataset is prohibitively expensive since it demands manual annotation of pixel-level class labels. To mitigate this problem, synthetic image datasets have been introduced [1], [2]. They provide a large amount of labeled images for training at minimal cost of construction. Also, they can simulate scenes that are rarely observed in the real world yet must be considered in training (*e.g.*, accidents in autonomous driving scenarios).

When learning semantic segmentation using synthetic images, it is essential to close the gap between the synthetic and real domains caused by their appearance differences so as to avoid performance degradation of learned models on real-world images. Most of existing solutions to this issue belong to the category of domain adaptation, which aims at adapting models trained on synthetic images (*i.e.*, source domain) to real-world images (*i.e.*, target domain). In general, domain adaptation methods assume a single, particular target domain

and train models using images from both of labeled source and unlabeled target domains [3], [4], [5], [6], [7], [8], [9], [10]. Unfortunately, this setting limits applicability of learned models since, when deployed, models can face multiple and diverse target domains (*e.g.*, geolocations and weather conditions in the case of autonomous vehicle) that are latent at the training stage.

As a more realistic solution to the problem, we study *domain generalization* for semantic segmentation. The goal of this task is to learn models that generalize well to various target domains without having access to their images in training. A pioneer work [11] achieves the generalization by forcing segmentation models to be invariant to random style variations of input image. However, this method is costly since it applies an image-to-image translator [12] to every synthetic image multiple times for the style randomization. Moreover, random styles are given by a small number of images sampled from ImageNet [13], and thus often irrelevant to target applications and hard to cover a wide range of real-world image styles. Follow-up research is often limited by the knowledge of ImageNet too. For example, Chen *et al.* [14] encourage the representations learned using synthetic images to be similar with those of an ImageNet pretrained network, and Huang *et al.* [15] randomize synthetic images in a frequency space using a small subset of ImageNet as references for stylization.

In this paper, we propose a WEb-image assisted Domain GEneralization scheme, dubbed *WEDGE*, which overcomes the limitations of the previous work by using real and application-relevant images crawled from web repositories (*e.g.*, Flickr). The crawling process demands no or minimal human intervention as it only asks search keywords that are determined directly by target application (*e.g.*, "driving + road" for autonomous driving) or classes appearing in the source domain images. Moreover, unlike those of ImageNet, the retrieved images can be used for self-supervised learning as well as for stylization since they are expected to be relevant to target application.

As illustrated in Fig. 1, WEDGE utilizes images crawled from the Web in two different ways. First, it replaces neural styles of synthetic training images with those of web-crawled images on-the-fly during training. This helps enhance the generalization by giving illusions of diverse real images while exploiting groundtruth labels of synthetic images. For this purpose, we introduce a *style injection* module that conducts the style manipulation in a feature level at low cost. Since it is substantially more efficient than the image-to-image translator used in [11], it allows to perform the

[1] Namyup Kim, Jaehyun Pahk and Suha Kwak are with POSTECH, Pohang, Republic of Korea. {namyup, jhpahk, suha.kwak}@postech.ac.kr
[2] Taeyoung Son is with NALBI, Seoul, Republic of Korea. taeyoung@nalbi.ai
[3] Cuiling Lan is with Microsoft Research Asia, Beijing, China. culan@microsoft.com
[4] Wenjun Zeng is with EIT Institute for Advanced Study, Beijing, China. wenjunzeng@eias.ac.cn
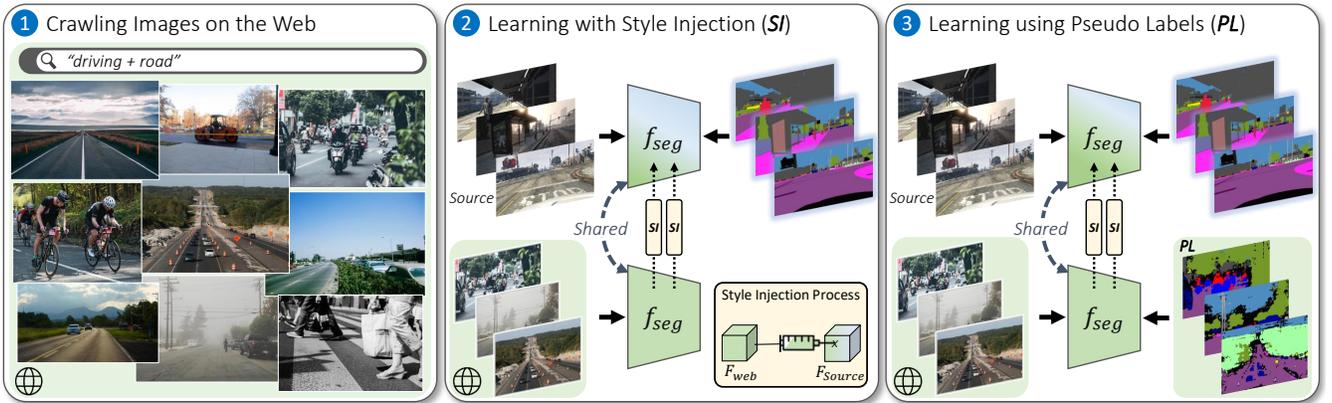
Fig. 1. Overall framework of WEDGE. (1) Crawling real and task-relevant images from the Web automatically. (2) Learning semantic segmentation while transferring feature statistics of web images to features of synthetic training images in the source domain. (3) Further training the model using both source images and web-crawled images with predicted pseudo labels.

stylization on-the-fly using a large number of web images as style references in training.

Second, the web-crawled images are used as additional training data with pseudo segmentation labels. To this end, the entire training procedure is divided into two stages. In the first stage, a segmentation model is trained with the style injection module, and the web-crawled images are used only for stylization. The learned model is then applied to the web-crawled images to estimate their pseudo labels. The second stage is identical to the former, except that it also utilizes the web-crawled images as training data by taking their pseudo labels as supervision.

To demonstrate the efficacy of WEDGE, we adopt each of the GTA5 [1] and SYNTHIA [2] datasets as the source domain for training, and evaluate trained models on three different real image datasets [16], [17], [18]. Experimental results demonstrate that WEDGE enables segmentation models to generalize well to multiple unseen real domains and clearly outperforms existing methods. In summary, the contribution of this paper is three-fold:

- To the best of our knowledge, WEDGE is the first that attempts to utilize web-crawled images for domain generalizable semantic segmentation. These images facilitate self training based on the realistic data which may better approximate unseen testing domains.
- We introduce style injection to domain generalizable semantic segmentation. Through web-crawled images, it helps achieve the generalization by giving diverse illusions of reality to the network being trained using labeled synthetic images. Also, the superiority of our particular style injection method over other potential candidates is demonstrated empirically.
- WEDGE outperforms existing domain generalization techniques [11], [14], [19], [15], [20] in every experiment.

## II. RELATED WORK

**Domain generalizable semantic segmentation.** The goal of domain generalization is to learn models that well generalize to unseen domains [21], [22]. Early approaches address this task mostly for classification [23], [24], [25], [26], [27], [28], but recent research demonstrates its potential for semantic

segmentation [11], [14], [26]. For example, Pan *et al.* [26] tackle this problem by feature normalization for learning domain invariant features, and Chen *et al.* [14] encourage the representation learned on a source domain to be similar with that of an ImageNet pretrained model. Also, Yue *et al.* [11] propose to learn features invariant to random style variations of input, and establish an evaluation protocol for the task. The main difference of ours from the previous work is that ours explores and exploits real images on the Web which enable models to experience a variety of real domains during training with no human intervention.

**Neural style transfer.** A pioneer work by Gatys *et al.* [29] shows that an image style can be captured by the Gram matrix of a feature map, and Johnson *et al.* [30] further enhance this idea to transfer a neural style to arbitrary images. Huang *et al.* [31] demonstrate that the channel-wise mean and standard deviation of a feature map represent image style effectively. Also, Nam and Kim [32] and Kim *et al.* [33] propose to use different normalization operations complementary to each other for style transfer. Recently, content-aware style transfer methods [34], [35], [36] are emerged to catch more details of local style patterns and to preserve content better. Huo *et al.* [36] suppose that features passing through a network form a manifold per each semantic region, and present a new style transfer technique based on manifold alignment. Style injection in WEDGE is motivated particularly by the techniques presented in [31], [36]. However, it is distinct from them in that it aims to perform feature stylization, rather than image stylization.

**Learning using data on the Web.** Modern recognition models tend to be data-hungry, yet the amount of training data is usually limited. Data on the Web have been exploited to alleviate this issue. Early studies utilize web-crawled images and videos for learning concept recognition by using their search keywords as pseudo labels [37], [38], [39], and for object localization via clustering images [37], [38] or by motion segmentation [40]. Motivated by recent advances in pseudo labeling, a large-scale web data have been used for supervised learning with their pseudo labels, which is known as webly supervised learning. For image classification, Niu *et*
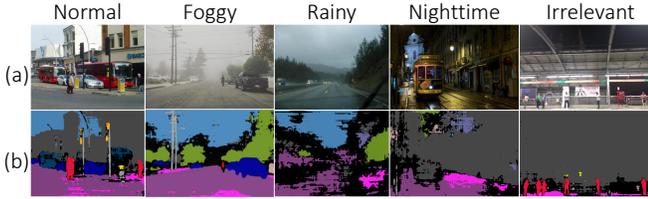
Fig. 2. Qualitative examples of web images and their pseudo labels generated by the segmentation model with ResNet101 backbone trained by the first stage of WEDGE. (a) Input images. (b) Pseudo segmentation labels. These images demonstrate large diversity of the web images, which is vital for achieving generalization to multiple latent test domains.



Fig. 3. Our style injection process. We first calculate the cross correlation matrix $\Sigma$ between features of source and web images by the cosine similarity $C(\cdot, \cdot)$. The optimal feature projection matrix for style injection is computed by $M = UV^{\top}$, where $U$ and $V$ are obtained by applying SVD to $F^s \Sigma F^{w\top}$. The source feature map is then stylized by applying $M$, and the result $F^s M^{\top}$ is fed into the next convoulution block.

*al.* [41] present a reliable way of utilizing search keywords as pseudo class labels. For semantic segmentation, Hong *et al.* [42] and Lee *et al.* [43] compute pseudo labels by segmenting web videos using attentions drawn by an image classifier. Motivated by these, we present the first that makes use of web images for domain generalization.

## III. Proposed Method

As shown in Fig. 1, WEDGE is divided into three steps:

1. **Crawling images from web repositories automatically.**
2. **The first stage of training with style injection (SI).** Learning a segmentation model on the synthetic dataset while injecting styles of the web-crawled images to its intermediate features for training.
3. **The second stage of training using pseudo labels (PL).** Further training the model using the web-crawled images and their pseudo segmentation labels as well as the synthetic dataset.

Details of each step are given in the remainder of this section.

### A. Crawling Images from the Web

We collect 4,904 images by crawling on Flickr, through the search keyword "driving + road" to find images relevant to the target application scenario, *i.e.*, autonomous driving. Examples of the collected images are presented in Fig. 2.

Using these images for domain generalization has several advantages. First, they offer a large variety of real image styles as illustrated in Fig. 2, which *potentially cover testing domains*. This is vital for achieving generalization to unseen domains. Second, they are not random but mostly relevant to target applications due to the use of search keywords and thus can be used for supervised learning given their pseudo labels. Last, they are accessible with minimal human intervention since the crawling process above is fully automated given a query. Note that previous work [11], [14], [15] also exploits external images, those of ImageNet, for style randomization; the web images are readily available like ImageNet and collected automatically, but more relevant to target application scenarios thanks to the search keyword.

The web-crawled images are often different from synthetic domain images in terms of semantic layout, and could partly contain irrelevant contents due to the ambiguity of search keywords and errors of the search engine. WEDGE is robust against these issues for the following reasons. In the first stage of training, the style injection module exploits only styles of the web images while disregarding their contents.
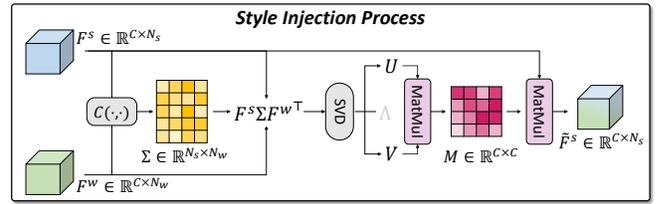
In the second stage, irrelevant parts of an image tend to be ignored in pseudo segmentation labels due to their unreliable class predictions (*i.e.*, low confidence).

### B. Stage 1: Learning with Style Injection (SI)

We propose a content-aware style injection method which transfers styles between semantically similar features of synthetic and web-crawled images. This method can inject diverse styles while better preserving the content of an image than conventional methods relying on global feature statistics (*e.g.*, AdaIN [31] or Gram matrix approximation [29]).

At each iteration of training the segmentation network, a synthetic image $I^s$ is coupled with a randomly sample web image $I^w$. Let $F^{d,l} \in \mathbb{R}^{H_d \times W_d \times C}$ be the feature map of $I^d$ from the $l^{\text{th}}$ convolution block of the network where $d \in \{s, w\}$. First, we compute a cross correlation between the $F^{s,l}$ and $F^{w,l}$ as an affinity matrix $\Sigma^l \in \mathbb{R}^{N_s \times N_w}$, where $N_d = H_d W_d$ ($d \in \{s, w\}$):

$$\Sigma^l_{i,j} = \frac{F^{s,l}_i F^{w,l}_j{}^{\top}}{\|F^{s,l}_i\| \|F^{w,l}_j\|}. \tag{1}$$

Then we find the projection matrix $M$ that minimizes the distance between features of the projected feature map $F^{s,l} M^{\top}$ and web-crawled image feature map $F^{w,l}$ weighted by the affinity matrix $\Sigma^l$; the role of the projection matrix $M$ is to project a synthetic feature onto a subspace of the semantically similar web-crawled features. The objective function is formulated by

$$\min_M J(M) = \frac{1}{N_\Sigma} \sum_{i=1}^{N_s} \sum_{j=1}^{N_w} \Sigma^l_{ij} \|F^{s,l} M^{\top} - F^{w,l}_j\|, \tag{2}$$

where $N_\Sigma$ is sum of all elements of $\Sigma^l$. The closed form solution of Eq. (2) is given by Huo *et al.* [36] as

$$M = UV^{\top}, \tag{3}$$

where $U$ and $V$ are derived from singular value decomposition of $F^{s,l} \Sigma^l F^{w,l}{}^{\top}$, *i.e.*, $F^{s,l} \Sigma^l F^{w,l}{}^{\top} = U \Lambda V^{\top}$. The synthetic feature map is projected by the estimated $M$, and the result $F^{s,l} M^{\top}$ is fed into the $l + 1^{\text{th}}$ convolutional block; Fig. 3 illustrates this process.

The overall pipeline of our style injection method follows that of the manifold alignment based style transfer (MAST) [36]. However, unlike MAST that computes a discrete affinity matrix using $k$ nearest neighbor assignment,

our method uses cosine similarity to compute the continuous affinity matrix in Eq. (1). It considers similarity of all features to produce a content-aware projection matrix and does not require hyperparameter $k$ nor $\mathcal{O}(n^2)$ time complexity to assign the nearest neighbors. Our style injection process is designed as non-parametric, which enables effective and low-cost feature adjustment.

The style injection is applied to multiple convolution blocks of the network, in particular lower blocks since features of deeper layers are known to be less sensitive to style variations. More details for implementation can be found in Sec. IV-A. Injecting styles of real web images to synthetic training images enlarges the training dataset by a multiple of the number of the web images, which is tremendous regarding the size of the training dataset, as well as making them look diverse and realistic in feature spaces. In addition, there are several advantages of the content-aware style injection for domain generalization of semantic segmentation over the conventional approaches [29], [31]. First, it enables style injection between semantically similar regions of web-crawled and synthetic images, which is more natural and effective for semantic segmentation. Second, since different styles are injected to different semantic regions on an image, it helps keep boundaries between the semantically different regions in style-injected features. We empirically verify the superiority of our method over other potential style injection candidates in Sec. IV-C.

Finally, the network is trained by the pixel-wise cross-entropy loss with the segmentation label of the synthetic image $I^s$. Let $P^s$ and $Y^s$ denote the segmentation prediction and the groundtruth label of $I^s$, respectively. The loss is then formulated as

$$\mathcal{L}_{\text{seg}}(P^s, Y^s) = -\frac{1}{N} \sum_{i=1}^{H} \sum_{j=1}^{W} \sum_{k=1}^{C} Y_{ijk}^s \log P_{ijk}^s, \quad (4)$$

where $N = H \times W$. Although this loss is applied only to the synthetic domain, its gradients with respect to parameters will act as if the network takes real domain images as input thanks to the style injection. Note that, in this stage, $I^w$ is used only as a style reference.

### C. Stage 2: Learning Using Pseudo Labels (PL)

Once the first stage is completed, the learned model can be used to generate pseudo labels of the web images. The pseudo labels allow us to exploit the web images for supervised learning of the segmentation network, which further enhances the generalization capability of the model by learning it directly on a variety of real-world images.

Let $P^w \in \mathbb{R}^{H \times W \times C}$ be the segmentation prediction of the network given $I^w$ as input. The pseudo segmentation label of $I^w$, denoted by $\widetilde{Y}^w \in \{0,1\}^{H \times W \times C}$, is obtained by choosing pixels with highly reliable predictions and labeling them with the classes of maximum scores:

$$\widetilde{Y}_{ijc}^w = \begin{cases} 1, & \text{if } c = \underset{k}{\operatorname{argmax}} \, P_{ijk}^w \ \text{ and } \ h(P_{ij}^w) < \tau \\ 0, & \text{otherwise} \end{cases}, \quad (5)$$

where $P_{ij}^w \in \mathbb{R}^C$ denotes the class probability distribution of the pixel $(i,j)$, $h(\cdot)$ indicates the entropy, and $\tau$ is a hyperparameter. Note that we regard the prediction $P_{ij}^w$ unreliable when its entropy is high, i.e., $h(P_{ij}^w) \geq \tau$; in this case, the pixel is assigned no label and ignored during training. Fig. 2 presents examples of the pseudo labels.

The second stage of training utilizes both of the synthetic and the web images for supervised learning. It is basically the same with the first stage including the style injection, except that the segmentation loss is now applied to $P^w$ as well as $P^s$. The total loss for the second stage is thus given by a linear combination of two segmentation losses:

$$\mathcal{L}(P^s, Y^s, P^w, \widetilde{Y}^w) = \mathcal{L}_{\text{seg}}(P^s, Y^s) + \mathcal{L}_{\text{seg}}(P^w, \widetilde{Y}^w), \quad (6)$$

where $\mathcal{L}_{\text{seg}}$ is the cross-entropy loss as given in Eq. (4).

## IV. EXPERIMENTS

In this section, we first present experimental settings in detail, then demonstrate the effectiveness of WEDGE through extensive results. Effectiveness of style injection, pseudo labeling, and other design choices of WEDGE are investigated by ablation studies.

### A. Experimental Setting

**Source datasets.** As a synthetic source domain for training, we use either the GTA5 [1] or the SYNTHIA [2] datasets. GTA5 consists of 24,966 images and shares the same set of 19 semantic classes with the real test datasets. Note that we remove 36 images of smallest file sizes from the dataset since they are non-informative, e.g., blacked-out images. Meanwhile, SYNTHIA contains 9,400 images, whose annotations cover only 16 classes of the real test datasets. Thus, we take only these 16 classes into account when evaluating models trained on SYNTHIA.

**Test datasets.** As unseen target domains for evaluation, we choose the validation splits of Cityscapes [16], BDD100k Segmentation (BDDS) [17] and Mapillary [18]. Cityscapes and BDDS have 500 and 1,000 validation images, respectively, and they are labeled for the same 19 classes. 2,000 validation images of Mapillary are annotated for 66 classes. By following the protocol of [44], we merge these classes to obtain the same 19 classes of Cityscapes.

**Web-crawled images.** From Flickr, we search for images whose widths are larger than or equal to 760 pixels, and with no copyright reserved (i.e., CC0) for their public use in future work, using the search keyword "driving + road". As a result, 4,904 web images in total are collected. Note that, given these conditions, the crawling process was done automatically, and the retrieved images are used as-is without modification.

**Networks and their training details.** Following previous work [11], we adopt DeepLab-v2 [45] with various backbone networks, ResNet50 and ResNet101 [46], as our segmentation networks. They are first pretrained on ImageNet [13], and then trained with the source dataset and our web images using SGD with momentum of 0.9 and weight decay of 5e−4. The initial learning rate is 2e−4 for the first stage

TABLE I

QUANTITATIVE RESULTS IN MIOU OF DOMAIN GENERALIZATION FROM (G)TA5 TO (C)ITYSCAPES, (B)DDS, AND (M)APILLARY.

| Methods | Backbone | G → C | G → B | G → M |
|---|---|---|---|---|
| IBN-Net [26] | ResNet50 | 29.6 | - | - |
| ASG [14] | ResNet50 | 31.9 | - | - |
| DRPC [11] | ResNet50 | 37.4 | 32.1 | 34.1 |
| RobustNet [19] | ResNet50 | 36.6 | 35.2 | 40.3 |
| WEDGE (Ours) | ResNet50 | **38.4** | **37.0** | **44.8** |
| DRPC [11] | ResNet101 | 42.5 | 38.7 | 38.1 |
| FSDR [15] | ResNet101 | 44.8 | 39.7 | 40.9 |
| PinMemory [20] | ResNet101 | 44.9 | 39.7 | 41.3 |
| WEDGE (Ours) | ResNet101 | **45.2** | **41.1** | **48.1** |

TABLE II

QUANTITATIVE RESULTS IN MIOU OF DOMAIN GENERALIZATION FROM (S)YNTHIA TO (C)ITYSCAPES, (B)DDS, AND (M)APILLARY.

| Methods | Backbone | S → C | S → B | S → M |
|---|---|---|---|---|
| DRPC [11] | ResNet50 | 35.7 | 31.5 | 32.7 |
| WEDGE (Ours) | ResNet50 | **36.1** | **32.5** | **37.2** |
| DRPC [11] | ResNet101 | 37.6 | 34.3 | 34.1 |
| FSDR [15] | ResNet101 | 40.8 | 37.4 | 39.6 |
| WEDGE (Ours) | ResNet101 | **40.9** | **38.1** | **43.1** |

TABLE III

PERFORMANCE OF WEDGE FOR DOMAIN GENERALIZATION FROM (G)TA5 AND (S)YNTHIA TO (C)ITYSCAPES, (B)DDS, AND (M)APILLARY.

| | ResNet50 | | | ResNet101 | | |
|---|---|---|---|---|---|---|
| | *Src. only* | SI (Stage 1) | PL (Stage 2) | *Src. only* | SI (Stage 1) | PL (Stage 2) |
| G → C | 28.29 | 36.25 | 38.36 | 34.28 | 43.55 | 45.18 |
| G → B | 29.16 | 36.30 | 37.00 | 32.96 | 40.35 | 41.06 |
| G → M | 40.46 | 42.75 | 44.82 | 41.31 | 47.30 | 48.06 |
| G$_{avg}$ | 32.64 | 38.43 | 40.06 | 36.18 | 43.73 | 44.77 |
| S → C | 27.06 | 35.28 | 36.09 | 29.96 | 38.22 | 40.94 |
| S → B | 23.96 | 28.62 | 32.51 | 24.28 | 30.74 | 38.07 |
| S → M | 31.67 | 36.49 | 37.18 | 36.19 | 38.61 | 43.10 |
| S$_{avg}$ | 27.56 | 33.46 | 35.26 | 30.14 | 35.86 | 40.70 |

TABLE IV

DOMAIN GENERALIZATION PERFORMANCE OF WEDGE WITH EACH VARIANT OF STYLE INJECTION METHODS AND OURS.

| SI methods | SI (Stage 1) | | | PL (Stage 2) | | |
|---|---|---|---|---|---|---|
| | G→C | G→B | G→M | G→C | G→B | G→M |
| None (source only) | 34.28 | 32.96 | 41.31 | - | - | - |
| WEDGE+AdaIN [31] | 40.54 | 40.78 | 46.33 | 40.58 | 39.74 | 46.92 |
| WEDGE+MAST [36] | 41.69 | 40.05 | 46.01 | 43.17 | 40.34 | 46.28 |
| WEDGE (Ours) | 43.55 | 40.35 | 47.30 | 45.18 | 41.06 | 48.06 |

(SI) and 1e−4 for the second stage (PL). $\tau$ in Eq. (5) is set to 5e−2 for all experiments.

**Where to inject styles.** Styles of web images are injected into the feature maps output by the $1^{th}$ and $2^{nd}$ residual blocks for ResNet101 and ResNet50. The impact of injection points on performance is analyzed in the accompanying video.

### B. Comparisons with the State of the Art

WEDGE is compared with existing domain generalization techniques, IBN-Net [26], AGS [14], DRPC [11], RobustNet [19], FSDR [15] and PinMemory [20], using two source domains {(G)TA5, (S)YNTHIA}, three test domains {(C)ityscapes, (B)DDS, (M)apillary}, and two different backbone networks {ResNet50, ResNet101}. As summarized in Table I and II, WEDGE clearly outperforms all the previous arts in all the 12 experiments.

### C. In-depth Analysis on WEDGE

**Detailed performance analysis.** To investigate the contribution of each training stage in WEDGE, we measure its performance at each stage for all experiments we have conducted so far. The results in Table III show that the first stage using style injection most contributes to the performance in most experiments, which demonstrates the effectiveness of using web-crawled images and our style injection module for domain generalization. This achievement is remarkable, especially when considering that web images could be erroneous or irrelevant to the test domains. Thanks to our style injection modules, WEDGE exploits diverse and realistic styles of web images while disregarding their contents that may be irrelevant. The second stage also leads to non-trivial performance improvement, particularly in the generalization from SYNTHIA to BDDS, which imply the semantics or layouts of pseudo labels on SYNTHIA is more similar to those of BDDS than the other datasets.

**Qualitative results.** The results in Fig. 4 show that the first stage of WEDGE recovers most of the ill-classified pixels, and even finds out objects that are missing in the baseline results. Also, its second stage further improves the segmentation quality by correcting dotted errors and capturing fine details of object shapes.

**Comparison of style injection methods.** We compare our style injection method with other potential candidates based on existing style transfer techniques [31], [36] to demonstrate its advantages. Note that these techniques are also used for injecting styles of web images on-the-fly within the same framework. As summarized in Table IV, while using AdaIN [31] and MAST [36] also improves performance, our method achieves the best in both SI and PL stages except for the GTA to BDDS case in the SI stage. Moreover, our method is more efficient than MAST since it does not need $k$ nearest neighbor search, whose time complexity is $\mathcal{O}(n^2)$, that is required for MAST.

**Impact of the number of web images.** We investigate the impact of the number of web images by evaluating performance of a segmentation model trained by WEDGE with different numbers of web images. In Fig. 5, these models are compared in terms of segmentation quality on the three target datasets. As shown in the figure, the generalization capability of the model can be substantially improved by using only 1,000 web images, while using the whole web dataset further improves performance. To be specific, when using 1,000 web images, the average mIoU over the 3 test datasets is 42.4%, lacking only 2.4% compared to the average performance of our final model. The results also indicate that WEDGE consistently enhances the generalization performance when increasing the number of web-crawled images.

**Impact of using task-relevant web images.** The contribution of our crawling strategy is demonstrated by comparing WEDGE with its variants relying on other types of
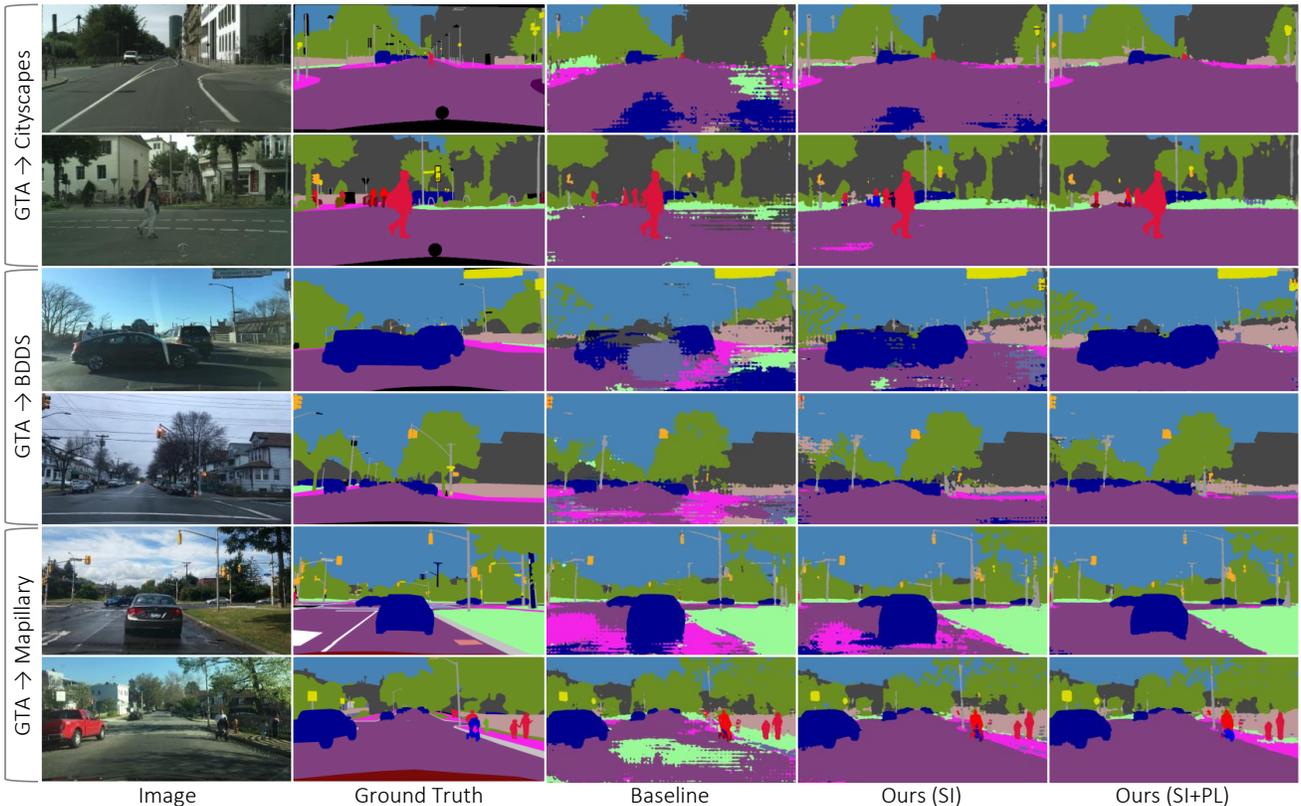
Fig. 4.    Qualitative results of WEDGE and its baseline using ResNet101 backbone and trained on the GTA5 dataset.
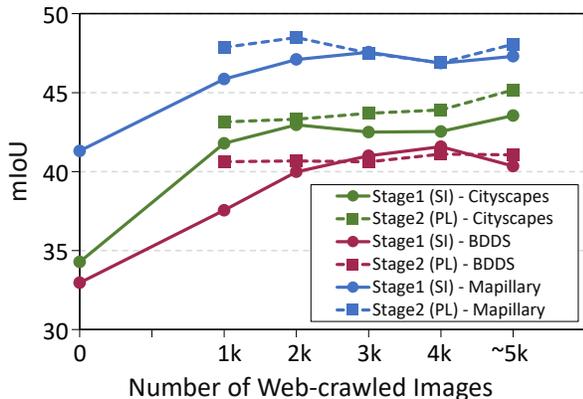


Fig. 5.    Domain generalization performance of WEDGE versus the number of web images. In all experiments ResNet101 is used as backbone.

TABLE V

| Style reference | G → C | G → B | G → M |
|---|---|---|---|
| None (source only) | 34.28 | 32.96 | 41.31 |
| ImageNet | 42.42 | 38.24 | 47.37 |
| Web images: "indoor" | 42.28 | 40.47 | 47.69 |
| Web images: "driving+road" | 45.18 | 41.06 | 48.06 |

## V. CONCLUSION

We have introduced WEDGE, the first web-image assisted domain generalization scheme for learning semantic segmentation. It explores and exploits web images that depict large diversity of real world scenes, which potentially cover latent test domains and thus help improve generalization capability of trained models. It utilizes the web-crawled images in two effective ways, namely style injection and pseudo labeling, which lead to consistent performance improvement on various test domains. WEDGE clearly outperformed existing domain generalization techniques in all experiments. Also, extensive ablation studies demonstrated that WEDGE is able to utilize noisy and irrelevant web-crawled images reliably and is not sensitive to their number in training.

style references instead of the task-relevant web images. Specifically, we utilize images sampled from the ImageNet dataset and web images crawled by the search keyword "indoor", both of which are irrelevant to the target task. Also, the number of style references is set to 5,000 for fair comparisons to WEDGE. Note that since these images are totally irrelevant to the target task, they are not suitable for pseudo labeling thus are used only for style injection. As shown in Table V, our method using task-relevant web images (*i.e.*, "driving+road") clearly outperforms the others. Using the real yet irrelevant images improves performance, suggesting the robustness of our method, but the results are still inferior to those of our method, meaning that our crawling strategy is useful and using relevant images matters.

# Appendix

This material presents implementation details and experimental results omitted from the main paper due to the space limit. First, Sec. VI demonstrates advantages of the feature-level style injection compared to the image-level style transfer. Sec. VII examines how much sensitive WEDGE is to the hyper-parameter $\tau$ and Sec. VIII investigates the impact of style injection points by an ablation study. Sec. IX discusses using domain-specific web images and Sec. X provides more qualitative examples of the web-crawled images we use. Then Sec. XI presents more qualitative results of WEDGE.

## VI. ADVANTAGES OF FEATURE-LEVEL STYLE INJECTION

Since WEDGE injects style representations in feature levels, one may wonder its advantages over image-level style transfer. This section demonstrates the effectiveness of WEDGE, especially its style injection (SI) module, compared to image-level style transfer. To this end, we adopt AdaIN [31], exploiting feature statistics as style representation like WEDGE. We generate 100,000 stylized GTA5 [1] images by AdaIN using web-crawled images as style references; a few examples are shown in Fig. 6. We then train a segmentation model on the stylized GTA5 dataset.

As summarized in Table VI, we compare the model of the $1^{\text{st}}$ stage of WEDGE (SI only) with the model trained on the stylized GTA5 images generated by AdaIN. The results show WEDGE using SI outperforms AdaIN on all experimental settings except G→C and G→B with VGG16 [46]. Moreover, our feature-level approach has another benefit over the image-level counterpart in terms of efficiency. AdaIN requires an additional network for style transfer. On the other hand, SI in WEDGE is non-parametric and adjusting feature statistics of source images by those of web-crawled images, thus demands a much lower computational cost than AdaIN.

## VII. SENSITIVITY TO HYPER-PARAMETER $\tau$

This section demonstrates the impact of the thresholding parameter s$\tau$ on the quality of pseudo labels in terms of semantic segmentation performance. Specifically, pseudo segmentation labels are generated using $\tau$, which is a hyper-parameter that filters out unreliable predictions. To this end, we design multiple variants of our model that are trained from different pseudo segmentation labels generated from various $\tau$. The pseudo segmentation label of $I^w$, denoted by $\widetilde{Y}^w \in \{0,1\}^{H \times W \times C}$, is obtained by choosing pixels with highly reliable predictions and labeling them with the classes of maximum scores:

$$\widetilde{Y}_{ijc}^w = \begin{cases} 1, & \text{if } c = \underset{k}{\arg\max}\ P_{ijk}^w \ \text{ and } \ h(P_{ij}^w) < \tau \\ 0, & \text{otherwise} \end{cases}, \quad (7)$$

where $h(\cdot)$ indicates the entropy, $P_{ij}^w \in \mathbb{R}^C$ denotes the class probability distribution of the pixel $(i,j)$, and $\tau$ is a hyper-parameter. We sample $\tau$ from $\{0.1, 0.05, 0.01, 0.005\}$, where $\tau = 0.05$ means our model in the main paper.
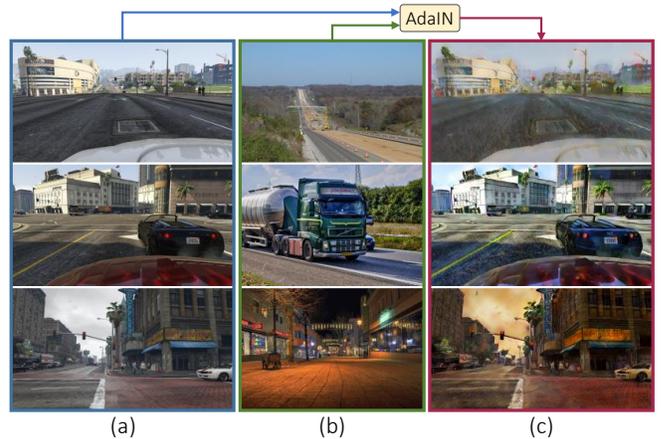


Fig. 6.  Examples of GTA5 images stylized by AdaIN. (a) Content images. (b) Style images. (c) Stylized images.

TABLE VI

QUANTITATIVE RESULTS IN MIOU AND PARAMETERS OF DOMAIN GENERALIZATION FROM (G)TA5 TO (C)ITYSCAPES , (B)DDS AND (M)APPILLARY.

| Methods | Backbone | Params | G → C | G → B | G → M |
|---|---|---|---|---|---|
| Deeplab-v2 [45] +AdaIN | VGG16 | 53.1M | 35.33 | 34.49 | 40.17 |
| | ResNet50 | 48.6M | 33.31 | 34.02 | 38.55 |
| | ResNet101 | 67.6M | 39.41 | 36.20 | 41.50 |
| WEDGE (SI) | VGG16 | 29.6M | 35.33 | 34.48 | 40.54 |
| | ResNet50 | 25.1M | 36.25 | 36.30 | 42.75 |
| | ResNet101 | 44.0M | 43.55 | 40.35 | 47.30 |

As summarized in Fig. 7, the results are marginally different across the variation of the hyper-parameters except 0.1, but the setting we adopt in the paper is slightly better than the others. Examples of the pseudo labels of web-crawled images are presented in Fig. 8, which demonstrates both pros and cons of different threshold values. With a moderate thresholding (*e.g.*, 0.1), the pseudo labels cover more real texture or parts of an object but have more noisy semantic labels. With a strict thresholding, on the other hand, the pseudo labels have more accurate semantic information but cover smaller regions of web-crawled images. The thresholding hyper-parameter we choose is in the middle, and leads to the best performance.

## VIII. DETAILS OF STYLE INJECTION

Style representations of web-crawled images are injected into the feature maps output by $1^{\text{th}}$ and $2^{\text{nd}}$ residual blocks for both ResNet101 and ResNet50 [46], and those of $2^{\text{nd}}$ and $3^{\text{rd}}$ blocks for VGG16 [47]. To verify the effectiveness of our method, this section presents ablation studies with various combinations of injection points. We present experimental results with ResNet101 combined with six different combinations in Table VII. The results show that semantic segmentation performance is degraded when $4^{\text{th}}$ residual block is included. We suspect this is because deeper features
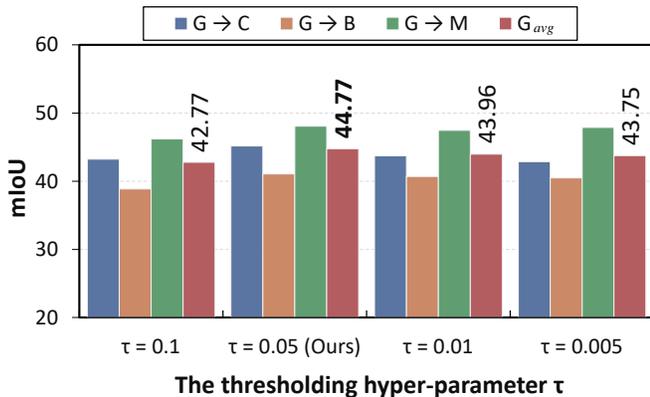
Fig. 7. Performance of our ResNet101 backbone model trained with different thresholded pseudo labels. $G_{avg}$ is the average performance of the three test domains. The performances across the set of hyper-parameters $\tau$ except 0.1 are marginally different.

TABLE VII

PERFORMANCE OF THE MODELS WITH RESNET101 BACKBONE ON THE SETTING FROM (G)TA5 TO (C)ITYSCAPES, (B)DDS AND (M)APPILLARY.

| Style injection points 1 2 3 4 | | | | $G \rightarrow C$ | $G \rightarrow B$ | $G \rightarrow M$ | Average |
|---|---|---|---|---|---|---|---|
| ✓ | ✓ | | | 43.55 | 40.35 | 47.30 | 43.73 |
| ✓ | ✓ | ✓ | | 44.03 | 39.30 | 47.30 | 43.54 |
| ✓ | ✓ | ✓ | ✓ | 41.01 | 38.49 | 46.46 | 41.99 |
| | ✓ | ✓ | | 42.00 | 39.03 | 44.93 | 41.99 |
| | | ✓ | ✓ | 37.92 | 35.02 | 38.47 | 37.20 |
| | ✓ | ✓ | ✓ | 38.64 | 35.19 | 41.39 | 38.44 |

TABLE VIII

DOMAIN GENERALIZATION PERFORMANCE OF WEDGE WITH DIFFERENT TYPES OF STYLE REFERENCE. IN ALL EXPERIMENTS RESNET101 IS USED AS BACKBONE.

| Method | Keywords | $G \rightarrow C$ | $G \rightarrow B$ | $G \rightarrow M$ | Average |
|---|---|---|---|---|---|
| WEDGE (SI) | driving + snow | 42.44 | 39.78 | 45.13 | 42.45 |
| | driving + rain | 42.87 | 38.13 | 45.95 | 42.32 |
| | driving + fog | 43.40 | 40.22 | 46.54 | 43.39 |
| | driving + road | **43.55** | **40.35** | **47.30** | **43.73** |

are known to contain semantic information rather than styles, which makes them inappropriate for style injection. As a result, using the output feature maps from the $\{1^{st}, 2^{nd}\}$ residual blocks turn out to be the most effective combination for ResNet101. Therefore, we choose our injection points based on these observations when applying style injection to other backbone networks.

## IX. COMPARISON WITH USING DOMAIN-SPECIFIC WEB IMAGES

Since our task at hand is domain generalization that assumes arbitrary target domains, we employ the keyword that does not indicate any specific domains. Nevertheless, we experiment with the keywords "driving + {rain, show, fog}". As summarized in Table VIII, these specific keywords are not as useful as the general one "driving + road" in our framework.

## X. EXAMPLES OF WEB-CRAWLED IMAGES

This section exhibits a part of our web-crawled dataset. Qualitative examples of the web-crawled images are presented in Fig. 11, which demonstrates the diversity of the images in terms of time, geolocation, weather condition, and so on. Such diversity enables WEDGE to achieve the generalization to latent real domains. Note that these images often depict entities and semantic layouts that are irrelevant to those of source (and target) domains. However, they are used as-is with no manual filtering process since the style injection and pseudo labeling of WEDGE offer reliable and effective ways to utilize them.
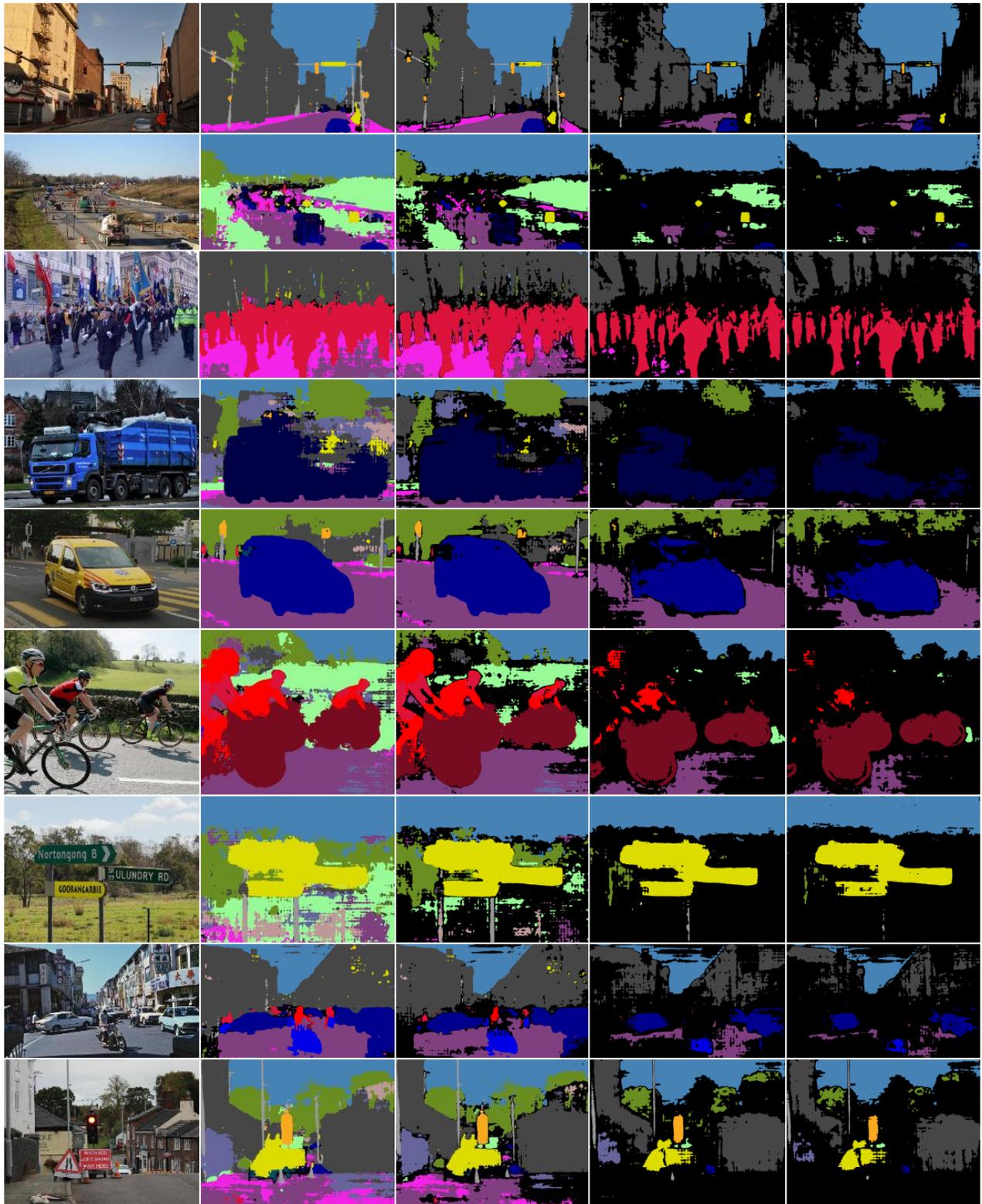
## XI. MORE QUALITATIVE RESULTS

We present qualitative examples of semantic segmentation results by WEDGE for both of Stage 1 (SI) and Stage 2 (SI+PL) in Fig. 9 and Fig. 10. In these figures, the semantic segmentation results are color-coded by following the standard Cityscapes color map [16]; the colors associated to the classes are exhibited in Table IX.



TABLE IX

THE COLOR CODE OF CLASSES ON THE TEST DATASETS.

Fig. 8. Examples of the pseudo labels with the different thresholding hyper-parameter $\tau$. (a) Input images. Pseudo labels with (b) $\tau = 0.1$. (c) $\tau = 0.05$ (Ours). (d) $\tau = 0.01$. (e) $\tau = 0.005$.
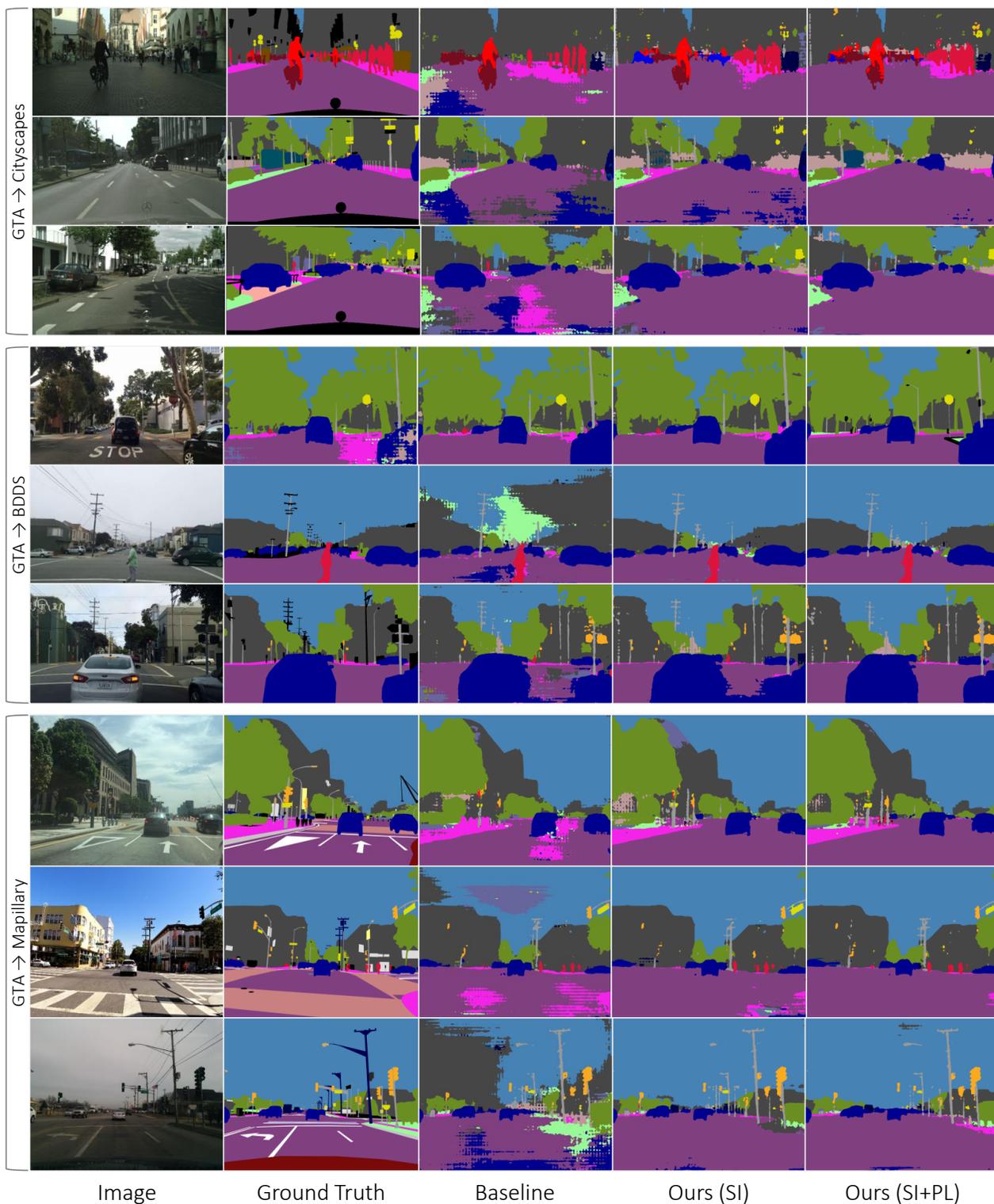
Image      Ground Truth      Baseline      Ours (SI)      Ours (SI+PL)

Fig. 9.    Qualitative results of WEDGE and its baseline using ResNet101 backbone and trained on the GTA5 dataset.

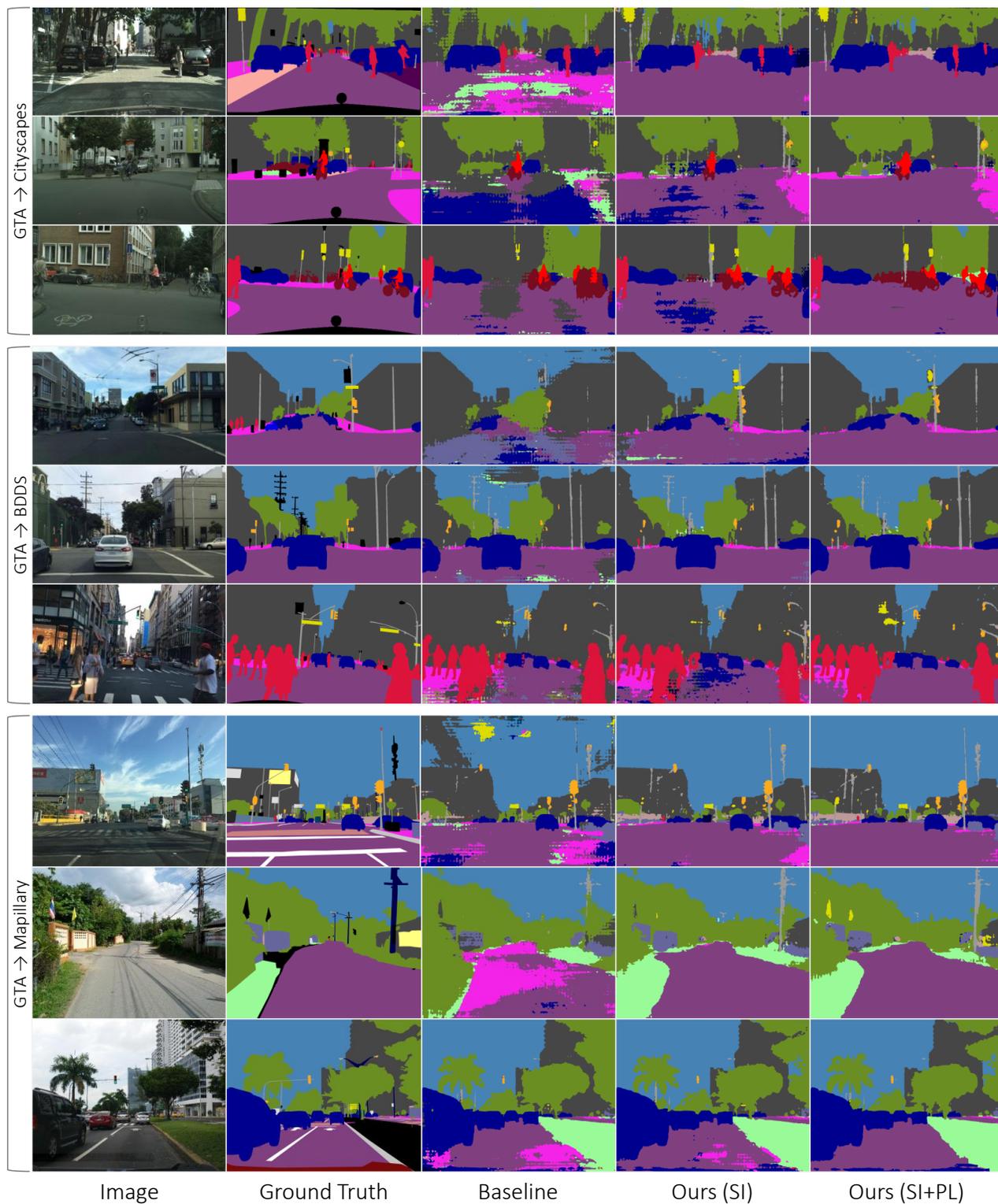| Image | Ground Truth | Baseline | Ours (SI) | Ours (SI+PL) |
|---|---|---|---|---|

Fig. 10.    Qualitative results of WEDGE and its baseline using ResNet101 backbone and trained on the GTA5 dataset.

Fig. 11.    500 random samples of the web-crawled images used in WEDGE.

# REFERENCES

[1] S. R. Richter, V. Vineet, S. Roth, and V. Koltun, "Playing for data: Ground truth from computer games," in *Proc. European Conference on Computer Vision (ECCV)*, 2016.

[2] G. Ros, L. Sellart, J. Materzynska, D. Vazquez, and A. M. Lopez, "The synthia dataset: A large collection of synthetic images for semantic segmentation of urban scenes," in *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.

[3] J. Hoffman, D. Wang, F. Yu, and T. Darrell, "Fcns in the wild: Pixel-level adversarial and constraint-based adaptation," *arXiv preprint arXiv:1612.02649*, 2016.

[4] A. Valada, J. Vertens, A. Dhall, and W. Burgard, "Adapnet: Adaptive semantic segmentation in adverse environmental conditions," in *Proc. IEEE International Conference on Robotics and Automation (ICRA)*, 2017.

[5] Y.-H. Tsai, W.-C. Hung, S. Schulter, K. Sohn, M.-H. Yang, and M. Chandraker, "Learning to adapt structured output space for semantic segmentation," in *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.

[6] Y. Li, L. Yuan, and N. Vasconcelos, "Bidirectional learning for domain adaptation of semantic segmentation," in *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.

[7] Y. Zou, Z. Yu, X. Liu, B. Kumar, and J. Wang, "Confidence regularized self-training," in *Proc. IEEE International Conference on Computer Vision (ICCV)*, 2019.

[8] F. Pan, I. Shin, F. Rameau, S. Lee, and I. S. Kweon, "Unsupervised intra-domain adaptation for semantic segmentation through self-supervision," in *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.

[9] G. Kang, Y. Wei, Y. Yang, Y. Zhuang, and A. Hauptmann, "Pixel-level cycle association: A new perspective for domain adaptive semantic segmentation," in *Proc. Neural Information Processing Systems (NeurIPS)*, 2020.

[10] X. Zhang, H. Zhang, J. Lu, L. Shao, and J. Yang, "Target-targeted domain adaptation for unsupervised semantic segmentation," in *Proc. IEEE International Conference on Robotics and Automation (ICRA)*, 2021.

[11] X. Yue, Y. Zhang, S. Zhao, A. Sangiovanni-Vincentelli, K. Keutzer, and B. Gong, "Domain randomization and pyramid consistency: Simulation-to-real generalization without accessing target domain data," in *Proc. IEEE International Conference on Computer Vision (ICCV)*, 2019.

[12] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros, "Unpaired image-to-image translation using cycle-consistent adversarial networks," in *Proc. IEEE International Conference on Computer Vision (ICCV)*, 2017.

[13] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "ImageNet: a large-scale hierarchical image database," in *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2009.

[14] W. Chen, Z. Yu, Z. Wang, and A. Anandkumar, "Automated synthetic-to-real generalization," in *Proc. International Conference on Machine Learning (ICML)*, 2020.

[15] J. Huang, D. Guan, A. Xiao, and S. Lu, "Fsdr: Frequency space domain randomization for domain generalization," in *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021.

[16] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, and B. Schiele, "The cityscapes dataset for semantic urban scene understanding," in *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.

[17] F. Yu, H. Chen, X. Wang, W. Xian, Y. Chen, F. Liu, V. Madhavan, and T. Darrell, "Bdd100k: A diverse driving dataset for heterogeneous multitask learning," in *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.

[18] G. Neuhold, T. Ollmann, S. Rota Bulo, and P. Kontschieder, "The mapillary vistas dataset for semantic understanding of street scenes," in *Proc. IEEE International Conference on Computer Vision (ICCV)*, 2017.

[19] S. Choi, S. Jung, H. Yun, J. T. Kim, S. Kim, and J. Choo, "Robustnet: Improving domain generalization in urban-scene segmentation via instance selective whitening," in *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021.

[20] J. Kim, J. Lee, J. Park, D. Min, and K. Sohn, "Pin the memory: Learning to generalize semantic segmentation," in *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022.

[21] K. Muandet, D. Balduzzi, and B. Schölkopf, "Domain generalization via invariant feature representation," in *Proc. International Conference on Machine Learning (ICML)*, 2013.

[22] C. Gan, T. Yang, and B. Gong, "Learning attributes equals multi-source domain generalization," in *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.

[23] H. Li, S. J. Pan, S. Wang, and A. C. Kot, "Domain generalization with adversarial feature learning," in *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.

[24] D. Li, Y. Yang, Y.-Z. Song, and T. M. Hospedales, "Deeper, broader and artier domain generalization," in *Proc. IEEE International Conference on Computer Vision (ICCV)*, 2017.

[25] D. Li, Y. Yang, Y.-Z. Song, and T. Hospedales, "Learning to generalize: Meta-learning for domain generalization," in *Proc. AAAI Conference on Artificial Intelligence (AAAI)*, 2018.

[26] X. Pan, P. Luo, J. Shi, and X. Tang, "Two at once: Enhancing learning and generalization capacities via ibn-net," in *Proc. European Conference on Computer Vision (ECCV)*, 2018.

[27] H. Nam, H. Lee, J. Park, W. Yoon, and D. Yoo, "Reducing domain gap via style-agnostic networks," *arXiv preprint arXiv:1910.11645*, vol. 2, no. 7, p. 8, 2019.

[28] K. Zhou, Y. Yang, Y. Qiao, and T. Xiang, "Domain generalization with mixstyle," in *Proc. International Conference on Learning Representations (ICLR)*, 2021.

[29] L. A. Gatys, A. S. Ecker, and M. Bethge, "Image style transfer using convolutional neural networks," in *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.

[30] J. Johnson, A. Alahi, and L. Fei-Fei, "Perceptual losses for real-time style transfer and super-resolution," in *European conference on computer vision*, 2016.

[31] X. Huang and S. Belongie, "Arbitrary style transfer in real-time with adaptive instance normalization," in *Proc. IEEE International Conference on Computer Vision (ICCV)*, 2017.

[32] H. Nam and H.-E. Kim, "Batch-instance normalization for adaptively style-invariant neural networks," in *Proc. Neural Information Processing Systems (NeurIPS)*, 2018.

[33] J. Kim, M. Kim, H. Kang, and K. H. Lee, "U-gat-it: Unsupervised generative attentional networks with adaptive layer-instance normalization for image-to-image translation," in *Proc. International Conference on Learning Representations (ICLR)*, 2020.

[34] D. Y. Park and K. H. Lee, "Arbitrary style transfer with style-attentional networks," in *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.

[35] S. Liu, T. Lin, D. He, F. Li, M. Wang, X. Li, Z. Sun, Q. Li, and E. Ding, "Adaattn: Revisit attention mechanism in arbitrary neural style transfer," in *Proc. IEEE International Conference on Computer Vision (ICCV)*, 2021.

[36] J. Huo, S. Jin, W. Li, J. Wu, Y.-K. Lai, Y. Shi, and Y. Gao, "Manifold alignment for semantically aligned style transfer," in *Proc. IEEE International Conference on Computer Vision (ICCV)*, 2021.

[37] X. Chen, A. Shrivastava, and A. Gupta, "Neil: Extracting visual knowledge from web data," in *Proc. IEEE International Conference on Computer Vision (ICCV)*, 2013.

[38] S. K. Divvala, A. Farhadi, and C. Guestrin, "Learning everything about anything: Webly-supervised visual concept learning," in *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2014.

[39] X. Chen and A. Gupta, "Webly supervised learning of convolutional networks," in *Proc. IEEE International Conference on Computer Vision (ICCV)*, 2015.

[40] A. Prest, C. Leistner, J. Civera, C. Schmid, and V. Ferrari, "Learning object class detectors from weakly annotated video," in *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2012.

[41] L. Niu, A. Veeraraghavan, and A. Sabharwal, "Webly supervised learning meets zero-shot learning: A hybrid approach for fine-grained classification," in *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.

[42] S. Hong, D. Yeo, S. Kwak, H. Lee, and B. Han, "Weakly supervised semantic segmentation using web-crawled videos," in *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 2224–2232.

[43] J. Lee, E. Kim, S. Lee, J. Lee, and S. Yoon, "Frame-to-frame aggregation of active regions in web videos for weakly supervised

semantic segmentation," in *Proc. IEEE International Conference on Computer Vision (ICCV)*, 2019.

[44] Y. He, S. Rahimian, B. Schiele, and M. Fritz, "Segmentations-leak: Membership inference attacks and defenses in semantic image segmentation," in *European Conference on Computer Vision*. Springer, 2020, pp. 519–535.

[45] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille, "Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs," *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 2017.

[46] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016.

[47] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," in *Proc. International Conference on Learning Representations (ICLR)*, 2015.