

# No More Discrimination: Cross City Adaptation of Road Scene Segmenters

Yi-Hsin Chen<sup>1</sup>, Wei-Yu Chen<sup>3,4</sup>, Yu-Ting Chen<sup>1\*</sup>, Bo-Cheng Tsai<sup>2\*</sup>, Yu-Chiang Frank Wang<sup>4</sup>, Min Sun<sup>1</sup>

Department of {<sup>1</sup>Electrical Engineering,<sup>2</sup>Communication Engineering}, National Tsing Hua University, Taiwan

<sup>3</sup>Department of Electrical Engineering, National Taiwan University, Taiwan

<sup>4</sup>Research Center for Information Technology Innovation, Academia Sinica, Taiwan

{yhethanchen, wyharveychen, yuting2401, vigorous0503}@gmail.com

, ycwang@citi.sinica.edu.tw, sunmin@ee.nthu.edu.tw

## Abstract

Despite the recent success of deep-learning based semantic segmentation, deploying a pre-trained road scene segmenter to a city whose images are not presented in the training set would not achieve satisfactory performance due to dataset biases. Instead of collecting a large number of annotated images of each city of interest to train or refine the segmenter, we propose an unsupervised learning approach to adapt road scene segmenters across different cities. By utilizing Google Street View and its time-machine feature, we can collect unannotated images for each road scene at different times, so that the associated static-object priors can be extracted accordingly. By advancing a joint global and class-specific domain adversarial learning framework, adaptation of pre-trained segmenters to that city can be achieved without the need of any user annotation or interaction. We show that our method improves the performance of semantic segmentation in multiple cities across continents, while it performs favorably against state-of-the-art approaches requiring annotated training data.

## 1. Introduction

Recent developments of technologies in computer vision, deep learning, and more broadly artificial intelligence, have led to the race of building advanced driver assistance systems (ADAS). From recognizing particular objects of interest toward understanding the corresponding driving environments, road scene segmentation is among the key components for a successful ADAS. With a sufficient amount of annotated training image data, existing computer vision algorithms already exhibit promising performances on the above task. However, when one applies pre-trained seg-

\*indicates equal contribution

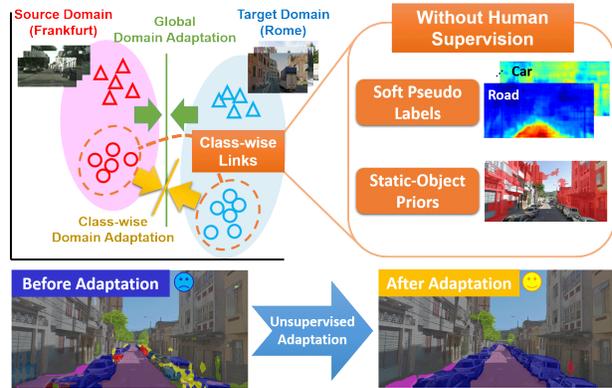


Figure 1: Illustration of our *unsupervised* domain adaptation method consisting of global and class-wise segmentation adaptations. For class-wise adaptation, we leverage “soft” pseudo labels and static object priors (obtained without human supervision) to further alleviate the domain discrimination in each class.

menters to a scene or city which is previously not seen, the resulting performance would be degraded due to dataset (domain) biases.

We conduct a pilot experiment to illustrate how severe a state-of-the-art semantic segmenter would be affected by the above dataset bias problem. We consider the segmenter of [2] which is trained on Cityscapes [5], and apply for segmenting about 400 annotated road scene images of different cities across countries: Rome, Rio, Taipei, and Tokyo. A drop in mean of intersection over union (mIoU) of 25-30% was observed (see later experiments for more details). Thus, how to suppress the dataset bias would be critical when there is a need to deploy road scene segmenters to different cities.

It is not surprising that, collecting a large number of an-

notated training image data for each city of interest would be time-consuming and expensive. For instance, pixel labeling of one Cityscapes image takes 90 minutes on average [5]. To alleviate this problem, a number of methods have been proposed to reduce human efforts in pixel-level semantic labeling. For example, researchers choose to utilize 3D information [37], rendered images [30, 31], or weakly supervised labels [32, 34, 3] for labeling. However, these existing techniques still require human annotation during data collection, and thus might not be easily scaled up to larger image datasets.

Inspired by the recent advances in domain adaptation [23, 35, 12], we propose an unsupervised learning framework for performing cross-city semantic segmentation. Our proposed model is able to adapt a pre-trained segmentation model to a new city of interest, while only the collection of *unlabeled* road scene images of that city is required. To avoid any human interaction or annotation during data collection, we utilize Google Street View with its time-machine<sup>1</sup> feature to harvest road scene images taken at the same (or nearby) locations but across different times. As detailed later in Sec. 4, this allows us to extract *static-object priors* from the city of interest. By integrating such priors with the proposed global and class-specific domain adversarial learning framework, refining/adapting the pre-trained segmenter can be easily realized.

The main contributions of this paper can be summarized as follows:

- We propose an *unsupervised* learning approach, which performs global and class-wise adaptation for deploying pre-trained road scene segmenters across cities.
- We utilize Google Street View images with time-machine features to extract static-object priors from the collected image data, without the need of user annotation or interaction.
- Along with the static-object priors, we advance adversarial learning for assigning pseudo labels to cross-city images, so that joint global and class-wise adaptation of segmenters can be achieved.

## 2. Related Work

### 2.1. CNN-based Semantic Segmentation

Semantic segmentation is among the recent breakthrough in computer vision due to the development and prevalence of Convolutional Neural Networks (CNN), which has been successfully applied to predict dense pixel-wise semantic labels [6, 18, 22, 2, 4]. For example, Long

<sup>1</sup><https://maps.googleblog.com/2014/04/go-back-in-time-with-street-view.html>

et al. [18] utilize CNN for performing pixel-level classification, which is able to produce pixel-wise outputs of arbitrary sizes. In order to achieve high resolution prediction, [22, 2] further adapt deconvolution layers into CNN with promising performances. On the other hand, Chen et al. [4] choose to add a fully-connected CRF layer at their CNN output, which refines the pixel labels with context information properly preserved. We note that, since the goal of this paper is to adapt pre-trained segmenters across cities, we do not limit the use of particular CNN-based segmentation solvers in our proposed framework.

### 2.2. Segmentation of Road Scene Images

To apply CNN-based segmenters to road scene images, there are several attempts to train segmenters on large-scale image datasets [5, 37, 30, 31]. For example, Cordts et al. [5] release a natural road scene segmentation dataset, which consists of over 5000 annotated images. Xie et al. [37] annotate 3D semantic labels in a scene, followed by transferring the 3D labels into the associated 2D video frames. [30, 31] collect semantic labels from Computer Graphic (CG) images at a large scale; however, building CG worlds for practical uses might still be computationally expensive.

On the other hand, [3] choose to relax the supervision during the data collection process, and simply require a number of point-labels per image. Moreover, [24, 26, 27] only require image-level labels during data collection and training. In addition to image-level labels, Pathak et al. [25] incorporate constraints on object sizes, [14, 34, 32] utilize weak object location knowledge, and [14] exploit object boundaries for constrained segmentation without using a large annotated dataset. Alternatively, [15, 38] apply free-form squiggles to provide partial pixel labels for data collection. Finally, [10] utilize image-level labels with co-segmentation techniques to infer semantic segmentation of foreground objects in the images of ImageNet.

### 2.3. DNN-based Domain Adaptation

Since the goal of our work is to adapt CNN-based segmenters across datasets (or cities to be more precise), we now review recent deep neural networks (DNN) based approaches for domain adaptation [23]. Based on Maximum Mean Discrepancy (MMD), Long et al. [19] minimize the mean distance between data domains, and later they incorporate the concept of residual learning [21] for further improvements. Zellinger et al. [40] consider Central Moment Discrepancy (CMD) instead of MMD, while Sener et al. [33] enforce cyclic consistency on adaptation and structured consistency on transduction in their framework.

Recently, Generative Adversarial Network (GAN) [9] has raised great attention in the fields of computer vision and machine learning. While most existing architectures are applied for synthesizing images with particular

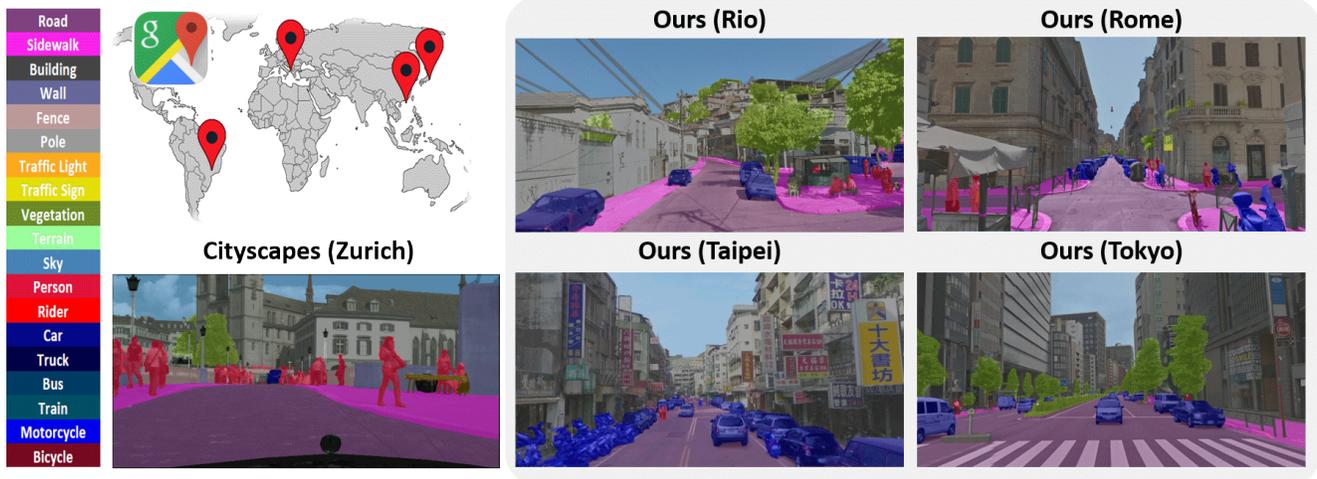


Figure 2: Example road scene images of different cities in our dataset. For evaluation purposes, we randomly select 100 images in each city to annotate pixel-level semantic labels. Color-coded labels are overlaid on each example image, where the mapping between colors and semantic classes are shown in the left panel.

styles [9, 29, 41]. Some further extend such frameworks for domain adaptation. In Coupled GAN [16], domain adaptation is achieved by first generating corresponded instances across domains, followed by performing classification.

In parallel with the appearance of GAN [9], Ganin et al. propose Domain Adversarial Neural Networks (DANN) [7, 8], which consider adversarial training for suppressing domain biases. For further extension, Variational Recurrent Adversarial Deep Domain Adaptation (VRADA) [28] utilizes Variational Auto Encoder (VAE) and RNN for time-series adaptation. Sharing a similar goal as ours, Hoffman et al. [11] extend such frameworks for semantic segmentation.

### 3. Dataset

We now detail how we collect our road scene image dataset, and explain its unique properties.

**Diverse locations and appearances.** Using Google Street View, road scene images at a global scale can be accessed across a large number of cities in the world. To address the issue of geo-location discrimination of a road scene segmenter, we download the road scene images of four cities at diverse locations, Rome, Rio, Tokyo, and Taipei, which are expected to have significant appearance differences. To ensure that we cover sufficient variations in visual appearances from each city, we randomly sample the locations in each city for image collection.

**Temporal information.** With the time-machine features of Google Street View, image pairs of the same location yet

across different times can be further obtained. As detailed later in the Sec. 4.2, this property particularly allows us to observe prior information from static objects, so that improved adaptation without any annotation can be achieved. In our work, we have collected 1600 image pairs (3200 images in total) at 1600 different locations per city with high image quality ( $647 \times 1280$  pixels).

For evaluation purposes, we select 100 image pairs from each city as the testing set, with pixel-level ground truth labels annotated by 15 image processing experts. We define 13 major classes for annotation: road, sidewalk, building, traffic light, traffic sign, vegetation, sky, person, rider, car, bus, motorcycle, and bicycle, as defined in Cityscapes [5]. Fig. 2 shows example images of our dataset. The dataset will be publicly available later for academic uses. To see more details and examples of our dataset, please refer to Appendix B or visit our website: [https://yihsinchen.github.io/segmentation\\_adaptation/](https://yihsinchen.github.io/segmentation_adaptation/).

We now summarize the uniqueness of our dataset below:

- Unlike existing datasets which typically collect images in nearby locations (e.g., road scenes of the same city), our dataset includes over 400 road scene images from four different cities around the world, with high-quality pixel-level annotations (for evaluation only).
- Our dataset include image pairs at the same location but across different times, which provide additional temporal information for further processing and learning purposes.

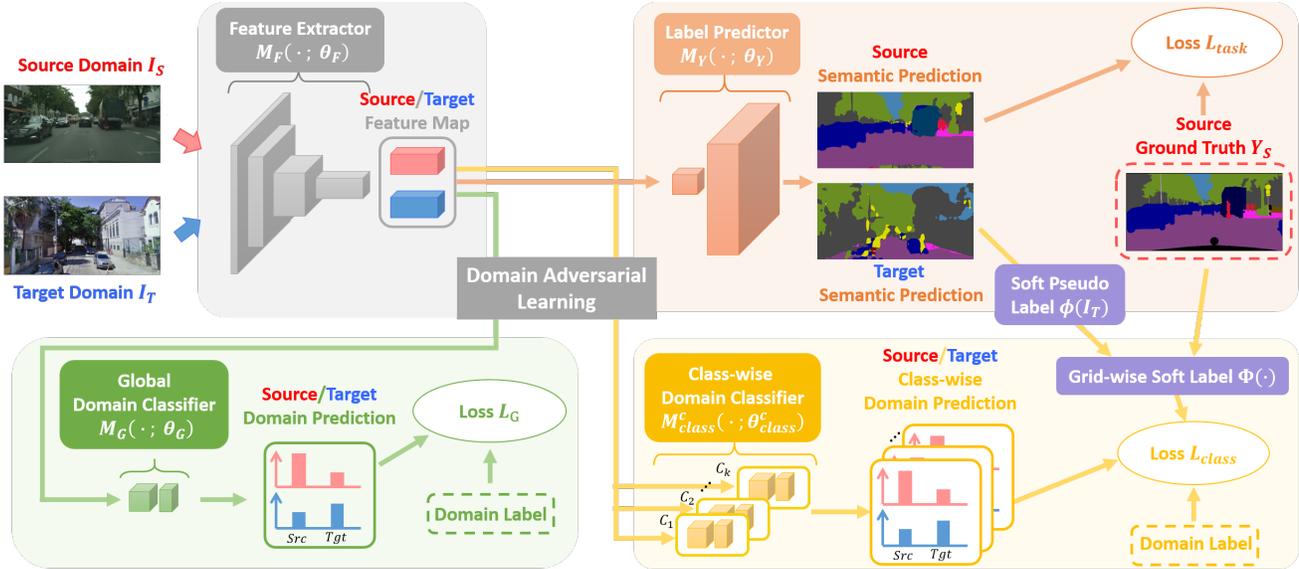


Figure 3: The overview of our proposed DNN framework. The feature extractor  $M_F$  transforms cross-domain images into a proper feature space, which is derived by performing global  $M_G$  and class-wise  $M_{class}^c$  domain alignment via adversarial learning. The label predictor  $M_Y$  regularizes the learned model by only observing the ground-truth annotation of source-domain images.

## 4. Our Method

In this section, we present the details of our proposed *unsupervised domain adaptation* framework, which is able to adapt pre-trained segmenters across different cities without using any user annotated data. In other words, while both images  $I_S$  and labels  $Y_S$  are available from the source domain  $\mathcal{S}$ , only images  $I_T$  for the target domain  $\mathcal{T}$  can be observed.

**Domain shift.** When adapting image segmenters across cities, two different types of domain shifts (or dataset biases) can be expected: *global* and *class-wise domain shift*. The former comes from the overall differences in appearances between the cities, while the latter is due to distinct compositions of road scene components in each city.

To minimize the global domain shift, we follow [11] and apply the technique of adversarial learning, which introduces a domain discriminator with a loss  $\mathcal{L}_G$ . This is to distinguish the difference between source and target-domain *images*, with the goal to produce a common feature space for images across domains. To perform class-wise alignment, we extend the above idea and utilize multiple *class-wise* domain discriminators (one for each class) with the corresponding adversarial loss  $\mathcal{L}_{class}$ . Unlike the discriminator for global alignment, these class-wise discriminators are trained to suppress the difference between cross-domain images but of the same class. Since we do not have any annotation for the city of interest (i.e., target-domain images), later we will explain how our method performs unsuper-

vised learning to jointly solve the above adaptation tasks.

With the above loss terms defined, the overall loss of our approach can be written as:

$$\mathcal{L}_{total} = \mathcal{L}_{task} + \lambda_G \mathcal{L}_G + \lambda_{class} \mathcal{L}_{class}, \quad (1)$$

where  $\lambda_G$  and  $\lambda_{class}$  are weights for the global and class-wise domain adversarial loss, respectively. Note that  $\mathcal{L}_{task}$  denotes the prediction loss of source-domain images, which can be viewed as a regularization term when adapting the learned model across domains.

**Our proposed framework.** Fig. 3 illustrates our framework. Let  $\mathcal{C}$  be the set of classes, and an input image denoted as  $x$ . Our proposed architecture can be decoupled into four major components: feature extractor  $M_F(x, \theta_F)$  that transforms the input image to a high-level, semantic feature space (the gray part), label predictor  $M_Y(M_F(x, \theta_F), \theta_Y)$  that maps feature space to task label space (the orange part), and domain discriminator for global  $M_G(M_F(x, \theta_F), \theta_G)$  (the green part) and class-wise  $M_{class}^c(M_F(x, \theta_F), \theta_{class}^c)$ ,  $c \in \mathcal{C}$  alignments (the yellow part). The feature extractor and task label predictor are initialized from a pre-trained segmenter, while the domain discriminators are randomly initialized. While we utilize the front-end dilated-FCN [39] as the pre-trained segmenter in our work, it is worth noting that our framework can be generally applied to other semantic segmenters.

In Sec. 4.1 and Sec. 4.2, we will detail our unsupervised learning for global alignment and class-wise alignment, respectively. In particular, how we extract and integrate static-

object priors for the target domain images without any human annotation will be introduced in Sec. 4.3.

### 4.1. Global Domain Alignment

Previously, domain adversarial learning frameworks have been applied for solving cross-domain image classification tasks [7]. However, for cross-domain image segmentation, each image consists of multiple pixels, which can be viewed as multiple instances per observation. Thus, how to extend the idea of domain adversarial learning for adapting segmenters across image domains would be our focus.

Inspired by [11], we take each *grid* in the *fc7* feature map of the FCN-based segmenter as an instance. Let the feature maps of source and target domain images as  $M_F(I_S, \theta_F)$  and  $M_F(I_T, \theta_F)$ , each map consists of  $N$  grids. Let  $p_n(x) = \sigma(M_G(M_F(x, \theta_F)_n, \theta_G))$  be the probability that the grid  $n$  of image  $x$  belongs to the source domain, where  $\sigma$  is the sigmoid function. We note that, for cross-domain classification, Ganin et al. [7] use the same loss function plus a gradient reversal layer to update the feature extractor and domain discriminator simultaneously. If directly applying their loss function for cross-domain segmentation, we would observe:

$$\begin{aligned} \max_{\theta_F} \min_{\theta_G} \mathcal{L}_G = & - \sum_{I_S \in \mathcal{S}} \sum_{n \in N} \log(p_n(I_S)) \\ & - \sum_{I_T \in \mathcal{T}} \sum_{n \in N} \log(1 - p_n(I_T)). \end{aligned} \quad (2)$$

Unfortunately, this loss function will result in gradient vanishing as the discriminator converges to its local minimum. To alleviate the above issue, we follow [9] and decompose the above problem into two subtasks. More specifically, we have a domain discriminator  $\theta_G$  trained with  $\mathcal{L}_G^D$  for classifying these two distributions into two groups, and a feature extractor  $\theta_F$  updated by its inverse loss  $\mathcal{L}_G^{Dinv}$  which minimizes the associated distribution differences. In summary, our objective is to minimize  $\mathcal{L}_G = \mathcal{L}_G^D + \mathcal{L}_G^{Dinv}$  by iteratively update  $\theta_G$  and  $\theta_F$ :

$$\min_{\theta_G} \mathcal{L}_G^D, \min_{\theta_F} \mathcal{L}_G^{Dinv}, \quad (3)$$

where  $\mathcal{L}_G^D$  and  $\mathcal{L}_G^{Dinv}$  are defined as:

$$\begin{aligned} \mathcal{L}_G^D = & - \sum_{I_S \in \mathcal{S}} \sum_{n \in N} \log(p_n(I_S)) \\ & - \sum_{I_T \in \mathcal{T}} \sum_{n \in N} \log(1 - p_n(I_T)), \end{aligned} \quad (4)$$

$$\begin{aligned} \mathcal{L}_G^{Dinv} = & - \sum_{I_S \in \mathcal{S}} \sum_{n \in N} \log(1 - p_n(I_S)) \\ & - \sum_{I_T \in \mathcal{T}} \sum_{n \in N} \log(p_n(I_T)). \end{aligned} \quad (5)$$

### 4.2. Class-wise Domain Alignment

In addition to suppressing the global misalignment between image domains, we propose to advance the same adversarial learning architecture to perform class-wise domain adaptation.

While the idea of regularizing class-wise information during segmenter adaptation has been seen in [11], its class-wise alignment is performed based on the composition of the class components in cross-city road scene images. To be more precise, it assumes that the composition/proportion of object classes across cities would be similar. Thus, such a regularization essentially performs global instead of class-specific adaptation.

Recall that, when adapting our segmenters across cities, we only observe road scene images of the target city of interest without any label annotation. Under such unsupervised settings, we extend the idea in [20] and assign pseudo labels to pixels/grids in the images of the target domain. That is, after the global adaptation in Fig. 3, the predicted probability distribution maps  $\phi(I_T) = \text{softmax}(M_Y(M_F(I_T, \theta_F), \theta_Y))$  of target domain images can be produced. Thus,  $\phi(I_T)$  can be viewed as the ‘‘soft’’ pseudo label map for the target domain images. As a result, class-wise association across data domains can be initially estimated by relating the ground truth label in the source domain and the soft pseudo label in the target domain.

**From pixel to grid-level pseudo label assignment.** In Sec. 4.1, to train the domain discriminator, we define each grid  $n$  in the feature space as one instance, which corresponds to multiple pixels in the image space. If the (pseudo) labels of these grids can be produced, adapting class-wise information using the same adversarial learning framework can be achieved.

To propagate and to determine the pseudo labels from pixels to each grid for the above adaptation purposes, we simply calculate the proportion of each class in each grid as the *soft* (pseudo) label. That is, let  $i$  be the pixel index in image space,  $n$  be the grid index in feature space, and  $\mathcal{R}(n)$  be the set of pixels that correspond to grid  $n$ . If  $y_i(I_S)$  denote the ground truth label of pixel  $i$  for source domain images, we then calculate source-domain grid-wise soft-label  $\Phi_n^c(I_S)$  as the probability of grid  $n$  belonging to class  $c$ :

$$\Phi_n^c(I_S) = \sum_{i \in \mathcal{R}(n)} \frac{y_i(I_S) == c}{|\mathcal{R}(n)|}. \quad (6)$$

On the other hand, due to the lack of annotated target-domain data, it is not as straightforward to assign grid-level soft pseudo labels to images in that domain. To solve this problem, we utilize  $\phi(I_T)$  derived above. Let  $\phi_i^c(I_T)$  be the pixel-wise soft pseudo label of pixel  $i$  corresponding to

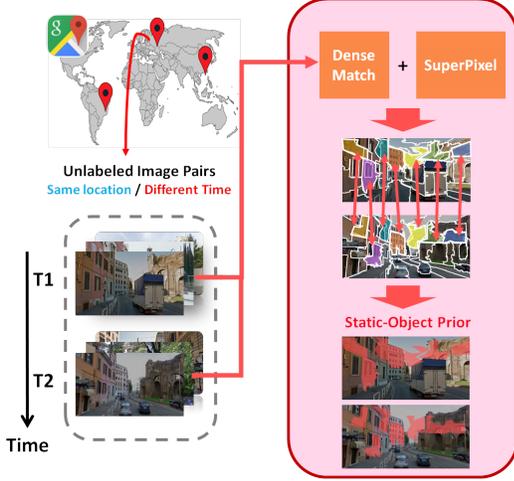


Figure 4: Illustration of static-object prior extraction. Given a pair of images at the same location but at different times, image regions belonging to static objects (e.g., the red blobs) can be identified by performing dense matching and superpixel segmentation.

class  $c$  for target-domain images, we have target grid-wise soft pseudo label  $\Phi_n^c(I_{\mathcal{T}})$  of grid  $n$ :

$$\Phi_n^c(I_{\mathcal{T}}) = \sum_{i \in \mathcal{R}(n)} \frac{\phi_i^c(I_{\mathcal{T}})}{|\mathcal{R}(n)|}. \quad (7)$$

Intuitively, grid-wise soft (pseudo) labels  $\Phi_n^c(I_{\mathcal{S}})$  and  $\Phi_n^c(I_{\mathcal{T}})$  are estimations of the probabilities that each grid  $n$  in source and target domain images belongs to object class  $c$ . To balance the appearance frequency of different classes, we normalize the estimated outputs in (6) and (7) as follows:

$$\begin{aligned} \tilde{\Phi}_n^c(I_{\mathcal{S}}) &= \frac{\Phi_n^c(I_{\mathcal{S}})}{\sum_{n \in \mathcal{N}} \Phi_n^c(I_{\mathcal{S}})} \\ \tilde{\Phi}_n^c(I_{\mathcal{T}}) &= \frac{\Phi_n^c(I_{\mathcal{T}})}{\sum_{n \in \mathcal{N}} \Phi_n^c(I_{\mathcal{T}})}. \end{aligned} \quad (8)$$

**Class-wise adversarial learning.** With the soft labels assigned to the source-domain images and the soft pseudo labels predicted for the target-domain ones, we now explain our adversarial learning for class-wise domain adaptation.

As depicted in Fig. 3, we deploy multiple *class-wise* domain discriminators  $\theta_{class}^c, c \in \mathcal{C}$  in our proposed architecture, and each discriminator is specially trained for differentiating *objects* of the corresponding class  $c$  across domains. Similar to  $p_n(x)$ , given that each object class  $c$  has a corresponded domain discriminator  $M_{class}^c$ , we define  $p_n^c(x) = \sigma(M_{class}^c(M_F(x, \theta_F)_n, \theta_{class}^c))$  as the probability predicted by  $M_{class}^c$  that the grid  $n$  of image  $x$  is from the

source domain. Combining the definition in (8), we define a pair of class-wise adversarial loss  $\mathcal{L}_{class}^D$  and  $\mathcal{L}_{class}^{Dinv}$  to guide the optimization for class-wise alignment:

$$\begin{aligned} \mathcal{L}_{class}^D &= - \sum_{I_{\mathcal{S}} \in \mathcal{S}} \sum_{c \in \mathcal{C}} \sum_{n \in \mathcal{N}} \tilde{\Phi}_n^c(I_{\mathcal{S}}) \log(p_n^c(I_{\mathcal{S}})) \\ &\quad - \sum_{I_{\mathcal{T}} \in \mathcal{T}} \sum_{c \in \mathcal{C}} \sum_{n \in \mathcal{N}} \tilde{\Phi}_n^c(I_{\mathcal{T}}) \log(1 - p_n^c(I_{\mathcal{T}})), \end{aligned} \quad (9)$$

$$\begin{aligned} \mathcal{L}_{class}^{Dinv} &= - \sum_{I_{\mathcal{S}} \in \mathcal{S}} \sum_{c \in \mathcal{C}} \sum_{n \in \mathcal{N}} \tilde{\Phi}_n^c(I_{\mathcal{S}}) \log(1 - p_n^c(I_{\mathcal{S}})) \\ &\quad - \sum_{I_{\mathcal{T}} \in \mathcal{T}} \sum_{c \in \mathcal{C}} \sum_{n \in \mathcal{N}} \tilde{\Phi}_n^c(I_{\mathcal{T}}) \log(p_n^c(I_{\mathcal{T}})). \end{aligned} \quad (10)$$

Finally, similar to (3), the class-wise alignment process is to iteratively solve the following optimization problem:

$$\min_{\theta_{class}^c} \mathcal{L}_{class}^D, \min_{\theta_F} \mathcal{L}_{class}^{Dinv}, \quad (11)$$

which minimizes the overall loss  $\mathcal{L}_{class} = \mathcal{L}_{class}^D + \mathcal{L}_{class}^{Dinv}$ .

### 4.3. Harvesting Static-Object Prior

While jointly performing global and class-wise alignment between source and target-domain images would produce promising adaptation performance, the pseudo labels are initialized by pre-trained segmenter. Under the unsupervised domain adaptation setting, since no annotation of target-domain data can be obtained, fine-tuning the segmenter by such information is not possible.

However, with the use of time-machine features from Google Street View images, we are able to leverage the temporal information for extracting the static-object priors from images in the target domain. As illustrated in Fig. 4, given an image pair of the same location but across different times, we first apply DeepMatching [36] to relate pixels within each image pair. For the regions with matched pixels across images, it implies such regions are related to static objects (e.g., building, road, etc.). Then, we additionally perform superpixel segmentation on the image pair using Entropy Rate Superpixel [17], which would group the nearby pixels into regions while the boundaries of the objects can be properly preserved. With the above derivation, we view the matched superpixels containing more than  $k$  matched pixels (we fix  $k = 3$  in this work) as the static-object prior  $\mathcal{P}_{static}(I_{\mathcal{T}})$ . Please refer to Appendix A for typical examples of mining static-object prior.

Let  $\mathcal{C}_{static}$  be the set of static-object classes. For the pixels that belong to  $\mathcal{P}_{static}(I_{\mathcal{T}})$ , we then refine their soft pseudo labels by suppressing its probabilities of being non-

Table 1: Accuracy of applying dilated-FCNs pre-trained on Cityscapes (Frankfurt) to different cities (i.e., no adaptation).

City	Dataset	mIOU (%)
Frankfurt	Cityscapes	64.6%
Rome	Ours	38.2%
Tokyo	Ours	39.2%
Rio	Ours	38.5%
Taipei	Ours	35.1%

static objects:

$$\forall i \in \mathcal{P}_{static}(I_{\mathcal{T}})$$

$$\tilde{\phi}_i^c(I_{\mathcal{T}}) = \begin{cases} \phi_i^c(I_{\mathcal{T}}) / \sum_{\hat{c} \in \mathcal{C}_{static}} \phi_i^{\hat{c}}(I_{\mathcal{T}}) & \text{if } c \in \mathcal{C}_{static} \\ 0 & \text{else} \end{cases} \quad (12)$$

## 5. Experiments

We first conduct experiments to demonstrate the issue of cross-city discrimination even using a state-of-the-art semantic segmenter. Then, we will verify the effectiveness of our proposed *unsupervised learning* method on the **Cityscapes to Our Dataset** domain adaptation task. By comparing it with a fully-supervised baseline (i.e., fine-tuning by fully annotated training data), we show that our unsupervised method would achieve comparable performances as the fully-supervised methods in most cases. Finally, we perform an extra experiment, **SYNTIA to Cityscapes**, to prove that our method could be generally applied to different datasets.

### 5.1. Implementation Details

In this work, all the implementations are produced utilizing the open source TensorFlow [1] framework, and the codes will be released upon acceptance. In the following experiments, we use mini-batch size 16 and the Adam optimizer [13] with learning rate of  $5 \times 10^{-6}$ ,  $beta1 = 0.9$ , and  $beta2 = 0.999$  to optimize the network. Moreover, we set the hyper-parameters in (1):  $\lambda_G$  and  $\lambda_{class}$ , to be numbers gradually changing from 0 to 0.1 and 0 to 0.5, respectively. In addition, for the experiments using static-object priors, we use {road, sidewalk, building, wall, fence, pole, traffic light, traffic sign, vegetation, terrain, sky} as the set of static-object classes  $\mathcal{C}_{static}$  defined in Sec. 4.3.

### 5.2. Cross-City Discrimination

We apply the segmenter pre-trained on **Cityscapes** to images of different cities in **Our Dataset**. As shown in Table 1, there is a severe performance drop in the four cities

compared to its original performance on Cityscapes. Interestingly, we observe a trend that the farther the geo-distance between the target city and the pre-trained city (Frankfurt), the severer the performance degradation. This implies that different visual appearances across cities due to cultural differences would dramatically impact the accuracy of the segmenter. For example, in Taipei, as shown in Fig. 2, there are many signboards and shop signs attached to the buildings, and many scooters on the road, which are uncommon in Frankfurt. It also justifies the necessity of an effective domain adaptation method for the road scene segmenter to alleviate the discrimination.

### 5.3. Cross-City Adaptation

**Baseline.** We use a *fully-supervised* method to establish a strong baseline as the upper bound of adaptation improvement. We divide our 100 images with fine annotations to 10 subsets for each city. Each time we select one subset as the testing set, and the other 90 images as the training set and fine-tune the segmenter for 2000 steps. We repeat the procedure for 10 times and average the testing results as the baseline performance.

**Our method.** Now we apply our domain adversarial learning method to adapt the pre-segmenter in an unsupervised fashion. Meanwhile, we do the ablation study to demonstrate the contribution from each component: global alignment, class-wise alignment, and static-object prior. We summarize the experimental results in Table 2, where "Pre-trained" denotes the pre-trained model, "UB" denotes the *fully-supervised* upper bound, "GA" denotes the global alignment part of our method, "GA+CA" denotes the combination of global alignment and class-wise alignment, and finally, "Full Method" denotes our overall method that utilizes the static-object priors. On average over four cities, our global alignment method contributes 2.6% mIoU gain, our class-wise alignment method also contributes 0.9% mIoU gain, and finally, the static-object priors contributes another 0.6% mIoU improvement. Furthermore, the t-SNE visualization results in Appendix A also show that the domain shift keeps decreasing from "Pre-trained" to "GA" to "GA+CA". These results demonstrate the effectiveness of each component of our method. In Fig. 5, we show some typical examples.

### 5.4. Synthetic to Real Adaptation

We additionally apply our method to another adaptation task with a different type of domain shift: **SYNTIA to Cityscapes**. In this experiment, we take SYNTIA-RAND-CITYSCAPES [31] as the source domain, which contains 9400 synthetic road scene images with Cityscapes-compatible annotations. For the unlabeled target domain, we use the training set of Cityscapes. During evaluation, we test our adapted segmenter on the validation set of

Table 2: Segmentation performance comparisons (in mIOU), in which SW, BLDG, TL, TS, VEG, Motor stand for Sidewalk, Building, Traffic Light, Traffic Sign, Vegetation, and Motorbike, respectively. Note that GA/CA denote the components of global/class-wise adaptation in our architecture, while our method (Full Method) integrates both components with static-object priors for unsupervised domain adaptation. The performance upper bound achieved by the fully supervised baseline is noted as UB.

City	Method	Cityscapes → Our Dataset													
		Road	SW	BLDG	TL	TS	VEG	Sky	Person	Rider	Car	Bus	Motor.	Bicycle	mIOU
Rome	Pre-trained	77.7	21.9	83.5	0.1	10.7	78.9	<b>88.1</b>	21.6	10.0	67.2	30.4	6.1	0.6	38.2
	GA	79.2	25.7	84.0	<b>0.1</b>	11.8	81.0	83.3	29.3	8.9	71.8	35.9	23.7	0.9	41.2
	GA+CA	78.2	26.0	<b>84.9</b>	0.0	21.5	<b>81.7</b>	83.0	<b>31.0</b>	11.2	<b>72.0</b>	33.0	24.1	<b>1.2</b>	42.1
	Full Method	<b>79.5</b>	<b>29.3</b>	84.5	0.0	<b>22.2</b>	80.6	82.8	29.5	<b>13.0</b>	71.7	<b>37.5</b>	<b>25.9</b>	1.0	<b>42.9</b>
	UB	84.9	33.0	87.3	0.0	10.9	84.6	91.6	30.5	19.1	77.7	10.6	38.3	0.5	43.8
Rio	Pre-trained	69.0	31.8	77.0	<b>4.7</b>	3.7	71.8	<b>80.8</b>	38.2	8.0	61.2	38.9	11.5	3.4	38.5
	GA	72.8	42.2	79.0	4.4	6.1	76.2	75.3	38.9	7.1	66.5	41.2	16.9	5.5	40.9
	GA+CA	73.2	42.9	78.4	3.3	<b>7.9</b>	76.2	72.4	39.1	9.1	<b>68.3</b>	<b>43.8</b>	16.8	6.5	41.4
	Full Method	<b>74.2</b>	<b>43.9</b>	<b>79.0</b>	2.4	7.5	<b>77.8</b>	69.5	<b>39.3</b>	<b>10.3</b>	67.9	41.2	<b>27.9</b>	<b>10.9</b>	<b>42.5</b>
	UB	80.2	53.8	84.5	0.0	16.4	81.4	85.4	42.3	17.4	74.0	49.4	37.3	16.7	49.1
Tokyo	Pre-trained	81.2	26.7	71.7	8.7	5.6	73.2	<b>75.7</b>	39.3	14.9	57.6	19.0	1.6	33.8	39.2
	GA	83.5	<b>36.2</b>	72.3	10.8	7.1	77.0	66.2	44.0	18.6	61.5	<b>21.9</b>	4.9	37.5	41.7
	GA+CA	<b>83.6</b>	36.1	71.9	11.3	<b>13.0</b>	<b>77.6</b>	64.4	41.2	19.3	63.7	20.2	<b>13.9</b>	38.8	42.6
	Full Method	83.4	35.4	<b>72.8</b>	<b>12.3</b>	12.7	77.4	64.3	<b>42.7</b>	<b>21.5</b>	<b>64.1</b>	20.8	8.9	<b>40.3</b>	<b>42.8</b>
	UB	85.2	38.7	79.8	13.9	19.7	81.7	86.9	45.3	35.9	66.9	29.0	2.0	42.4	48.3
Taipei	Pre-trained	77.2	20.9	76.0	5.9	4.3	60.3	81.4	10.9	<b>11.0</b>	54.9	32.6	15.3	5.2	35.1
	GA	79.0	27.7	76.6	13.1	5.0	67.7	74.8	<b>17.5</b>	6.1	60.4	28.9	25.5	7.1	37.6
	GA+CA	<b>79.2</b>	<b>29.0</b>	<b>80.3</b>	<b>14.1</b>	<b>8.2</b>	<b>68.8</b>	81.1	16.3	10.5	<b>64.7</b>	33.8	16.2	6.5	38.8
	Full Method	78.6	28.6	80.0	13.1	7.6	68.2	<b>82.1</b>	16.8	9.4	60.4	<b>34.0</b>	<b>26.5</b>	<b>9.9</b>	<b>39.6</b>
	UB	84.0	36.6	87.7	9.9	13.7	76.2	91.9	23.4	24.1	65.1	39.4	47.8	3.2	46.4

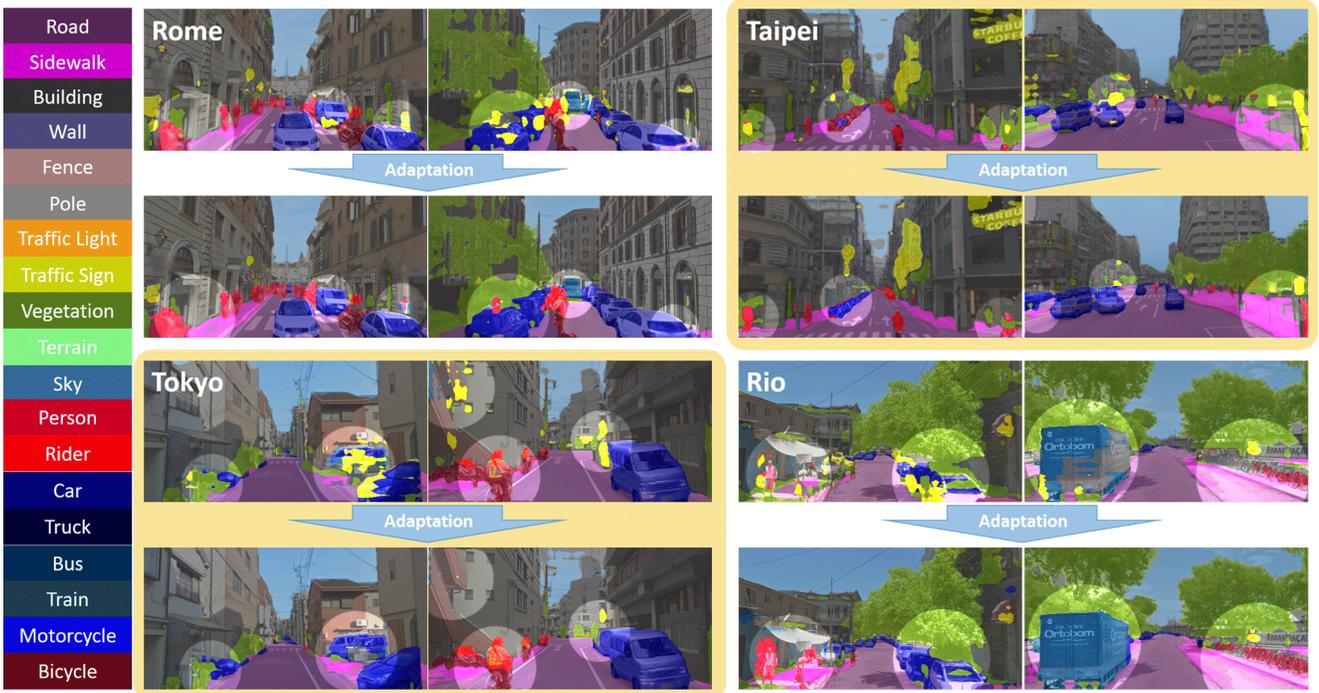


Figure 5: Examples of cross-city adaptation. The first/third and second/fourth rows show the results before and after adaptation, respectively. The regions with improved segmentation adaptation are highlighted for better visualization.

Cityscapes. We note that, since there are no paired images with temporal information in Cityscapes (as those in our dataset), we cannot extract static-object priors in this ex-

periment. Nevertheless, from the results shown in Table 3, performing global and class-wise alignment using our proposed method still achieves 3.1% and 1.9% mIOU gain, re-

Table 3: Experimental results for the SYNTHIA-to-Cityscapes segmentation adaptation task.

Method	SYNTHIA → Cityscapes													
	Road	SW	BLDG	TL	TS	VEG	Sky	Person	Rider	Car	Bus	Motor.	Bicycle	mIOU
Pre-trained	24.3	19.5	48.3	<b>1.5</b>	5.4	77.4	76.1	<b>42.8</b>	<b>9.7</b>	62.5	9.8	0.5	<b>20.9</b>	30.7
GA	56.5	24.0	<b>78.9</b>	1.1	<b>5.9</b>	77.8	77.3	35.8	5.4	61.7	5.2	0.9	8.4	33.8
GA+CA	<b>62.7</b>	<b>25.6</b>	78.3	1.2	5.4	<b>81.3</b>	<b>81.0</b>	37.4	6.4	<b>63.5</b>	<b>16.1</b>	<b>1.2</b>	4.6	<b>35.7</b>

spectively. These results again demonstrate the robustness of our proposed method. For typical examples of this adaptation task, please refer to Appendix C.

## 6. Conclusion

In this paper, we present an *unsupervised* domain adaptation method for semantic segmentation, which alleviates cross-domain discrimination on road scene images across different cities. We propose a unified framework utilizing domain adversarial learning, which performs joint global and class-wise alignment by leveraging soft labels from source and target-domain data. In addition, our method uniquely identifies and introduce static-object priors to our method, which are retrieved from images via natural synchronization of static objects over time. Finally, we provide a new dataset containing road scene images of four cities across countries, good-quality annotations and paired images with temporal information are also included. We demonstrate the effectiveness of each component of our method on tasks with different levels of domain shift.

## References

- [1] M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. S. Corrado, A. Davis, J. Dean, M. Devin, et al. Tensorflow: Large-scale machine learning on heterogeneous distributed systems. *arXiv preprint arXiv:1603.04467*, 2016. [7](#)
- [2] V. Badrinarayanan, A. Kendall, and R. Cipolla. Segnet: A deep convolutional encoder-decoder architecture for image segmentation. *arXiv preprint arXiv:1511.00561*, 2015. [1](#), [2](#)
- [3] A. Bearman, O. Russakovsky, V. Ferrari, and L. Fei-Fei. Whats the point: Semantic segmentation with point supervision. In *ECCV*. Springer, 2016. [2](#)
- [4] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille. Semantic image segmentation with deep convolutional nets and fully connected crfs. In *ICLR*, 2015. [2](#)
- [5] 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 *CVPR*. IEEE, 2016. [1](#), [2](#), [3](#)
- [6] C. Farabet, C. Couprie, L. Najman, and Y. LeCun. Learning hierarchical features for scene labeling. *IEEE transactions on pattern analysis and machine intelligence*, 35(8):1915–1929, 2013. [2](#)
- [7] Y. Ganin and V. Lempitsky. Unsupervised domain adaptation by backpropagation. In *ICML*, 2015. [3](#), [5](#)
- [8] Y. Ganin, E. Ustinova, H. Ajakan, P. Germain, H. Larochelle, F. Laviolette, M. Marchand, and V. Lempitsky. Domain-adversarial training of neural networks. *Journal of Machine Learning Research*, 17(59):1–35, 2016. [3](#)
- [9] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In *NIPS*, 2014. [2](#), [3](#), [5](#)
- [10] M. Guillaumin, D. Küttel, and V. Ferrari. Imagenet auto-annotation with segmentation propagation. In *IJCV*. Springer, 2014. [2](#)
- [11] 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. [3](#), [4](#), [5](#)
- [12] A. Khosla, T. Zhou, T. Malisiewicz, A. A. Efros, and A. Torralba. Undoing the damage of dataset bias. In *ECCV*. Springer, 2012. [2](#)
- [13] D. Kingma and J. Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015. [7](#)
- [14] A. Kolesnikov and C. H. Lampert. Seed, expand and constrain: Three principles for weakly-supervised image segmentation. In *ECCV*. Springer, 2016. [2](#)
- [15] D. Lin, J. Dai, J. Jia, K. He, and J. Sun. Scribble-sup: Scribble-supervised convolutional networks for semantic segmentation. In *CVPR*. IEEE, 2016. [2](#)
- [16] M.-Y. Liu and O. Tuzel. Coupled generative adversarial networks. In *NIPS*, 2016. [3](#)
- [17] M.-Y. Liu, O. Tuzel, S. Ramalingam, and R. Chellappa. Entropy rate superpixel segmentation. In *CVPR*. IEEE, 2011. [6](#)
- [18] J. Long, E. Shelhamer, and T. Darrell. Fully convolutional networks for semantic segmentation. In *CVPR*, 2015. [2](#)
- [19] M. Long, Y. Cao, J. Wang, and M. I. Jordan. Learning transferable features with deep adaptation networks. In *ICML*, 2015. [2](#)
- [20] M. Long, J. Wang, G. Ding, J. Sun, and P. S. Yu. Transfer feature learning with joint distribution adaptation. In *ICCV*, 2013. [5](#)
- [21] M. Long, H. Zhu, J. Wang, and M. I. Jordan. Unsupervised domain adaptation with residual transfer networks. In *NIPS*, 2016. [2](#)
- [22] H. Noh, S. Hong, and B. Han. Learning deconvolution network for semantic segmentation. In *ICCV*. IEEE, 2015. [2](#)
- [23] S. J. Pan and Q. Yang. A survey on transfer learning. *IEEE Transactions on knowledge and data engineering*, 22(10):1345–1359, 2010. [2](#)
- [24] G. Papandreou, L.-C. Chen, K. P. Murphy, and A. L. Yuille. Weakly-and semi-supervised learning of a deep convolutional network for semantic image segmentation. In *ICCV*. IEEE, 2015. [2](#)

- [25] D. Pathak, P. Krahenbuhl, and T. Darrell. Constrained convolutional neural networks for weakly supervised segmentation. In ICCV. IEEE, 2015. 2
- [26] D. Pathak, E. Shelhamer, J. Long, and T. Darrell. Fully convolutional multi-class multiple instance learning. 2015. 2
- [27] P. O. Pinheiro and R. Collobert. From image-level to pixel-level labeling with convolutional networks. In CVPR. IEEE, 2015. 2
- [28] S. Purushotham, W. Carvalho, T. Nilanon, and Y. Liu. Variational recurrent adversarial deep domain adaptation. In ICLR, 2017. 3
- [29] A. Radford, L. Metz, and S. Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. In ICLR, 2016. 3
- [30] S. R. Richter, V. Vineet, S. Roth, and V. Koltun. Playing for data: Ground truth from computer games. In ECCV. Springer, 2016. 2
- [31] 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 CVPR. IEEE, 2016. 2, 7
- [32] F. Saleh, M. S. A. Akbarian, M. Salzmann, L. Petersson, S. Gould, and J. M. Alvarez. Built-in foreground/background prior for weakly-supervised semantic segmentation. In ECCV. Springer, 2016. 2
- [33] O. Sener, H. O. Song, A. Saxena, and S. Savarese. Learning transferrable representations for unsupervised domain adaptation. In NIPS, 2016. 2
- [34] W. Shimoda and K. Yanai. Distinct class-specific saliency maps for weakly supervised semantic segmentation. In ECCV. Springer, 2016. 2
- [35] A. Torralba and A. A. Efros. Unbiased look at dataset bias. In CVPR. IEEE, 2011. 2
- [36] P. Weinzaepfel, J. Revaud, Z. Harchaoui, and C. Schmid. Deepflow: Large displacement optical flow with deep matching. In ICCV. IEEE, 2013. 6
- [37] J. Xie, M. Kiefel, M.-T. Sun, and A. Geiger. Semantic instance annotation of street scenes by 3d to 2d label transfer. In CVPR. IEEE, 2016. 2
- [38] J. Xu, A. G. Schwing, and R. Urtasun. Learning to segment under various forms of weak supervision. In CVPR. IEEE, 2015. 2
- [39] F. Yu and V. Koltun. Multi-scale context aggregation by dilated convolutions. In ICLR, 2016. 4
- [40] W. Zellingner, T. Grubinger, E. Lughofer, T. Natschläger, and S. Saminger-Platz. Central moment discrepancy (CMD) for domain-invariant representation learning. In ICLR, 2017. 2
- [41] J.-Y. Zhu, P. Krähenbühl, E. Shechtman, and A. A. Efros. Generative visual manipulation on the natural image manifold. In ECCV. Springer, 2016. 3

# Appendix

## A. Visualize GA, CA and Static-Object prior

In Sec. 4.1-4.3 of the main paper, we explain how each component in our structure enhance the performance of segmentation, and also show quantitative results in experiment. Here we'll further illustrate effects of these components:

**T-SNE Visualization** To visualize the adaptation results on common feature space with t-SNE, we randomly select 100 images from each domain, and for each image we extracted its average  $f_{c7}$  feature from each class, so for both source and target we have 100 feature points from each class.

As shown in Fig. 6, with pre-trained model only, there is an obvious shift between source and target domain. After applying the global alignment (GA), the distance between clusters with same labels becomes closer, while we could still observe a gap between domains. Once we further apply the class-wise alignment (CA), the gap between domains nearly vanishes. This result again demonstrates the effectiveness of each component of our proposed method.

**Harvesting Static-Object Prior** In Sec. 4.3, we propose a novel pipeline to extract the static-object prior using the natural synchronization of static objects over time. For better understanding, we show some typical results of our proposed pipeline in Fig. 7. Clearly, most of the regions iden-

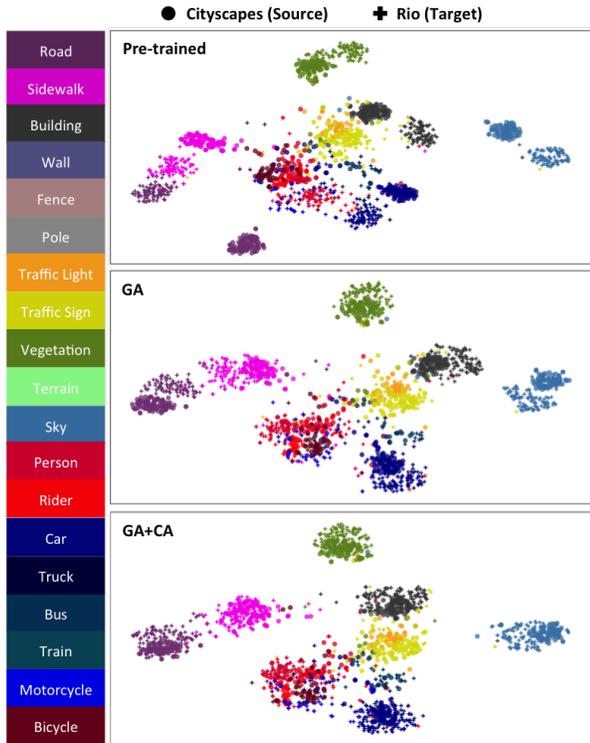


Figure 6: t-SNE visualization results. For simplicity, we only show the results of the task *Cityscapes*  $\rightarrow$  *Rio*. We could clearly observe that the alignment between domains becomes better from *pre-trained* to *GA+CA*.

tified by our method truly belong to static-objects. This demonstrates the effectiveness of our method.

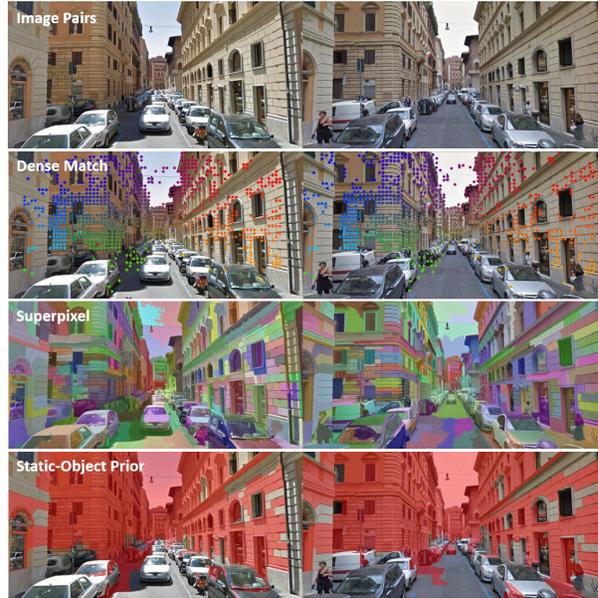


Figure 7: Typical results of our static-object prior pipeline. The first row is the original unlabeled image pair of same place across time. The second row is the result of dense matching, noted by points of same color. The third row is the result of superpixel segmentation marked by different colors. Combining the results from the above two rows, we could extract static-object prior of this image pair, as shown by the red regions in the last row.

## B. Dataset

To demonstrate the uniqueness of our dataset for road scene semantic segmenter adaptation, here we show more examples of it.

**Unlabeled Image Pairs** There are more examples collected at different cities with diverse appearances in Fig. 8. Valuable temporal information which facilitates *unsupervised* adaptation is contained in these image pairs.

**Labeled Image** We also show more annotated images in Fig. 9 to demonstrate the label-quality of our dataset.

## C. Synthetic to Real Adaptation

In Sec. 5.4 of the main paper, we have shown the quantitative results of this adaptation task in Table 3. We conclude that our method could perform well even under this challenging setting. To better support our conclusion, here we show some typical examples of this task in Fig. 10.

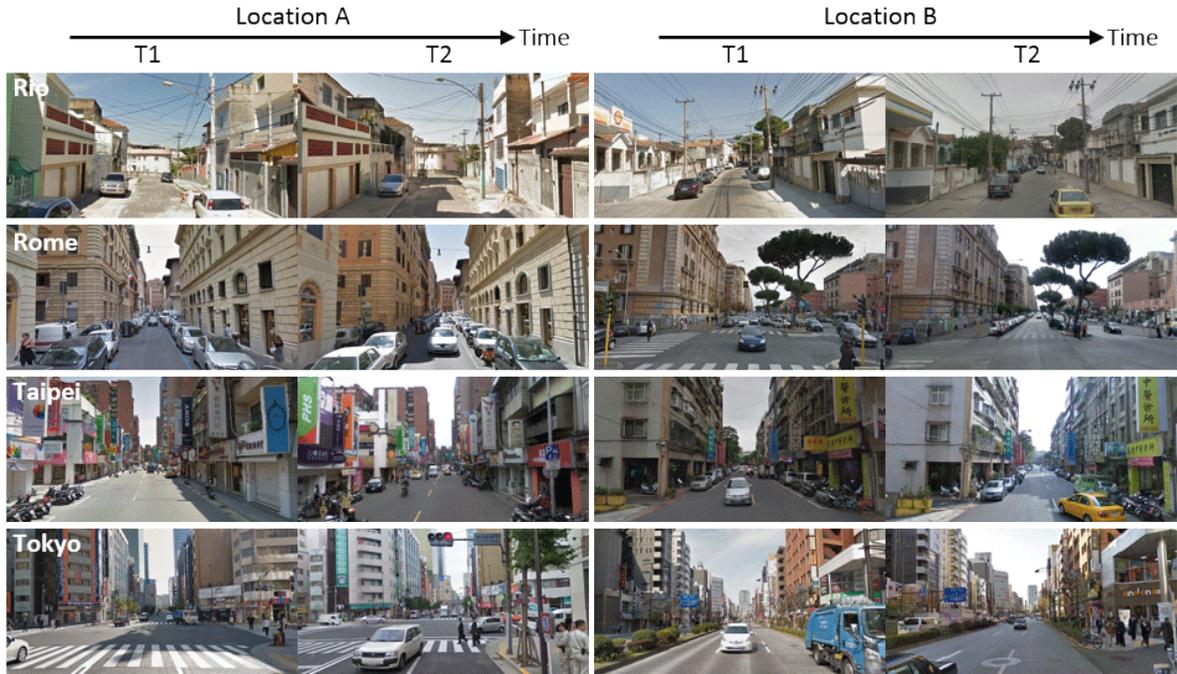


Figure 8: Examples of the unlabeled image pairs of different cities in our dataset. In each row, we show two image pairs at different locations in one city.

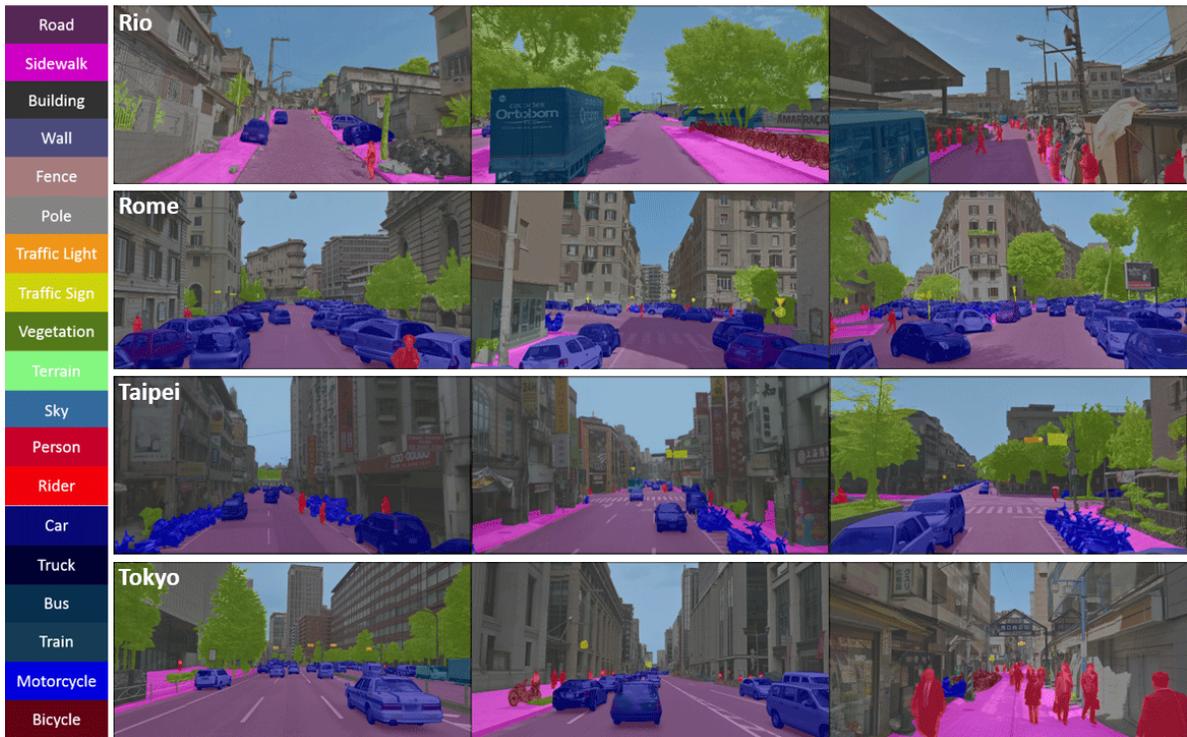


Figure 9: Examples of the labeled images of different cities in our dataset. Each image is annotated in good quality.

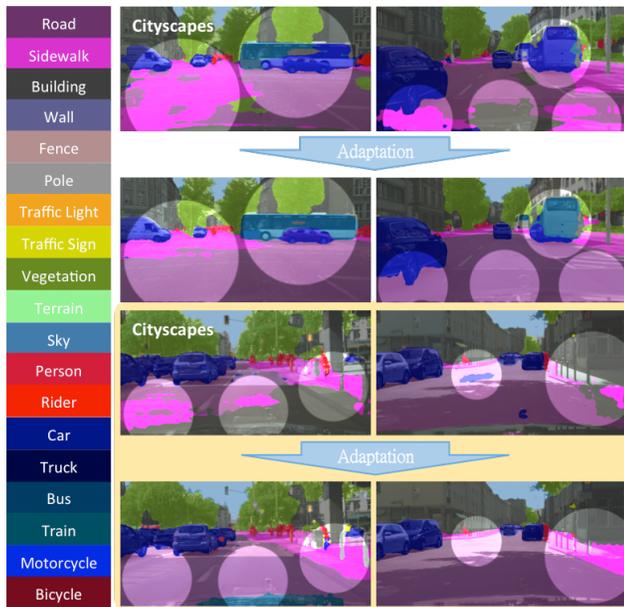


Figure 10: adaptation task: STNTHIA to Cityscapes. The first row and second show the results before and after adaptation, respectively.