# Curriculum Domain Adaptation for Semantic Segmentation of Urban Scenes

Yang Zhang[1], Philip David[2], and Boqing Gong[1]

[1] Center for Research in Computer Vision, University of Central Florida
[2] Computational and Information Sciences Directorate, U.S. Army Research Laboratory

yangzhang@knights.ucf.edu, philip.j.david4.civ@mail.mil, bgong@crcv.ucf.edu

## Abstract

*During the last half decade, convolutional neural networks (CNNs) have triumphed over semantic segmentation, which is a core task of various emerging industrial applications such as autonomous driving and medical imaging. However, to train CNNs requires a huge amount of data, which is difficult to collect and laborious to annotate. Recent advances in computer graphics make it possible to train CNN models on photo-realistic synthetic data with computer-generated annotations. Despite this, the domain mismatch between the real images and the synthetic data significantly decreases the models' performance. Hence we propose a curriculum-style learning approach to minimize the domain gap in semantic segmentation. The curriculum domain adaptation solves easy tasks first in order to infer some necessary properties about the target domain; in particular, the first task is to learn global label distributions over images and local distributions over landmark superpixels. These are easy to estimate because images of urban traffic scenes have strong idiosyncrasies (e.g., the size and spatial relations of buildings, streets, cars, etc.). We then train the segmentation network in such a way that the network predictions in the target domain follow those inferred properties. In experiments, our method significantly outperforms the baselines as well as the only known existing approach to the same problem.*

## 1. Introduction

This paper is concerned with domain adaptation for semantic image segmentation of urban scenes, i.e., assigning a category label to every pixel of an image or video frame [6]. Our interest in this problem is partially due to the exciting vision of autonomous driving, where understanding complex inner-city traffic scenes is an essential module and semantic segmentation is one of its key constituents [12, 19].

Machine learning methods for automatic semantic segmentation require massive amounts of high-quality annotated imagery in order to produce effective classifiers that generalize well to novel scenes. However, annotating training imagery for semantic segmentation is a very cumbersome task for humans. Cordts et al. report that the annotation and quality control take more than 1.5 hours on a single image of the Cityscapes dataset [12]. Besides, it is very difficult and time-consuming to collect imagery that depicts the large number of variabilities possible of urban scenes in different countries, seasons, and lighting conditions, etc.

To overcome both shortcomings, simulated urban environments may be used to automatically generate large amounts of annotated training imagery. This, however, introduces a new problem, that of domain mismatch between the source (simulated) domain and the target (real) domain. Figure 2 illustrates some examples drawn from the synthetic SYNTHIA [48] dataset and the real Cityscapes [12] dataset. It is readily apparent that there are significant visual differences between the two datasets. Domain adaptation techniques [48, 53, 27] may be used by machine learning methods to bridge this gap between the two domains.

In computer vision, learning domain-invariant features has been a prevalent and successful strategy to tackle the discrepancy between two domains, mainly for classification and regression problems [41, 43]. The core idea is to infer a new feature space such that the marginal distributions of the source domain (S) and the target domain (T) are about the same, i.e., $P_S(Z) \approx P_T(Z)$. Furthermore, the prediction function $P(Y|Z)$ from that space is assumed to be the same across the domains so that one can leverage the rich labeled data in the source domain to train classifiers that generalize well to the target. It is hard to verify the assumption, but the work along this line is rich and has led to impressive practical results regardless, such as the algorithms using linear transformation [22, 23, 15, 55], kernel methods [40, 20, 2, 31], and the recent deep learning meth-

---

For better reproducibility, the code is available at:
https://github.com/YangZhang4065/AdaptationSeg.

ods that directly extract domain-invariant features from raw input images [62, 36, 61, 18, 17].

In contrast to prior arts, the semantic segmentation we study in this paper is a highly structured prediction problem, for which domain adaptation is only sparsely explored in the literature [66, 27]. Under structured prediction, can we still achieve good domain adaptation results by following the above principles? Our intuition and experimental studies (cf. Section 4) tell us no. Learning a decision function for structured prediction is more involved than classification because it has to resolve the predictions in an exponentially large label space. As a result, the assumption that the source and target domains share the same prediction function becomes less likely to hold. Besides, some discriminative cues in the data would be suppressed if one matches the feature representations of the two domains without taking careful account of the structured labels. Finally, data instances are the proxy to measure the domain difference [25, 17, 18]. However, it is not immediately clear what comprises the instances in semantic segmentation [27], especially given that the top-performing segmentation methods are built upon deep neural networks [35, 44, 39, 10]. Hoffman et al. take each spatial unit in the fully convolutional network (FCN) [35] as an instance [27]. We contend that such instances are actually non-i.i.d. in either individual domain, as their receptive fields overlap with each other.

How can we avoid the assumption that the source and target domains share the same prediction function in a transformed domain-invariant feature space? Our proposed solution draws on two key observations. One is that the urban traffic scene images have strong idiosyncrasies (e.g., the size and spatial relations of buildings, streets, cars, etc.). Therefore, some tasks are "easy" and, *more importantly, suffer less because of the domain discrepancy*. Second, the structured output in semantic segmentation enables convenient posterior regularization [16], as opposed to the popular (e.g., $\ell_2$) regularization over model parameters.

Accordingly, we propose a curriculum-style [4] domain adaptation approach. Recall that, in domain adaptation, only the source domain supplies many labeled data while there are no or only scarce labels from the target. The curriculum domain adaptation begins with the easy tasks, in order to gain some high-level properties about the unknown pixel-level labels for each target image. It then learns a semantic segmentation network — the hard task, whose predictions over the target images are forced to follow those necessary properties as much as possible.

To develop the easy tasks in the curriculum, we consider label distributions over both holistic images and some landmark superpixels of the target domain. Take the former for instance. The label distribution of an image indicates the percentage of pixels that belong to each category, respectively. We argue that such tasks are easier, despite the do-

main mismatch, than assigning pixel-wise labels. Indeed, we may directly estimate the label distributions without inferring the pixel-wise labels. Moreover, the relative sizes of road, vehicle, pedestrian, etc. constrain the shape of the distributions, effectively reducing the search space. Finally, models to estimate the label distributions over superpixels may benefit from the urban scenes' canonical layout that transcends domains, e.g., buildings stand beside streets.

Why and when are the seemingly simple label distributions useful for the domain adaptation of semantic segmentation? In our experiments, we find that the segmentation networks trained on the source domain perform poorly on many target images, giving rise to disproportionate label assignments (e.g., many more pixels are classified to sidewalks than to streets). To rectify this, the image-level label distribution informs the segmentation network *how* to update the predictions while the label distributions of the landmark superpixels tell the network *where* to update. Jointly, they guide the adaptation of the networks to the target domain to, at least, generate proportional label predictions. Note that additional "easy tasks" can be conveniently incorporated into our framework in the future.

Our main contribution is on the proposed curriculum-style domain adaptation for the semantic segmentation of urban scenes. We select into the curriculum the easy and useful tasks of inferring label distributions for the target images and landmark superpixels, in order to gain some necessary properties about the target domain. Built upon these, we learn a pixel-wise discriminative segmentation network from the labeled source data and, meanwhile, conduct a "sanity check" to ensure the network behavior is consistent with the previously learned knowledge about the target domain. Our approach effectively eludes the assumption about the existence of a common prediction function for both domains in a transformed feature space. It readily applies to different segmentation networks, as it does not change the network architecture or tax any intermediate layers.

## 2. Related work

We discuss some related work on domain adaptation and semantic segmentation, with special focus on that transferring knowledge from virtual images to real photos.

**Domain adaptation.** Conventional machine learning algorithms rely on the assumption that the training and test data are drawn i.i.d. from the same underlying distribution. However, it is often the case that there exists some discrepancy from the training to the test stage. Domain adaptation aims to rectify this mismatch and tune the models toward better generalization at testing [60, 59, 21, 30, 25].

The existing work on domain adaptation mostly focuses on classification and regression problems [43, 41], e.g., learning from online commercial images to classify real-world objects [50, 22], and, more recently, aims to improve

the adaptability of deep neural networks [36, 18, 17, 61, 7, 37]. Among them, the most relevant work to ours is that exploring simulated data [56, 65, 48, 63, 27, 45, 53]. Sun and Saenko train generic object detectors from the synthetic images [56], while Vazquez et al. use the virtual images to improve pedestrian detections in real environment [63]. The other way around, i.e., how to improve the quality of the simulated images using the real ones, is studied in [53, 45].

**Semantic segmentation.** Semantic segmentation is the task of assigning an object label to each pixel of an image. Traditional methods [52, 58, 68] rely on local image features manually designed by domain experts. After the pioneering work [10, 35] that introduced the convolutional neural network (CNN) [32] to semantic segmentation, most recent top-performing methods are built on CNNs [64, 49, 3, 69, 39, 13].

An enormous amount of labor-intensive work is required to annotate the many images that are needed to obtain accurate segmentation models. The PASCAL VOC2012 Challenge [14] contains nearly 10,000 annotated images for the segmentation competition, and the MS COCO Challenge [34] includes over 200,000 annotated images. According to [47], it took about 60 minutes to manually segment each image in [8] and about 90 minutes for each in [12]. A plausible approach to reducing the human workload is to utilize weakly supervised information such as image labels and bounding boxes [44, 28, 42, 46].

We instead explore the use of almost effortlessly labeled virtual images for training high-quality segmentation networks. In [47], annotating a synthetic image took only 7 seconds on average through a computer game. For the urban scenes, we use the SYNTHIA [48] dataset which contains images of a virtual city.

**Domain adaptation for semantic segmentation.** Upon observing the obvious mismatch between virtual and real data [51, 47, 48], we expect domain adaptation to enhance the segmentation performance on real images by networks trained on virtual ones. To the best of our knowledge, the only attempt to algorithmically address this problem is [27]. While it regularizes the intermediate layers and constrains the output of the network, we propose a different curriculum domain adaptation strategy. We solve the easy task first and then use the learned knowledge about the target domain to regularize the network predictions.

## 3. Approach

In this section, we present the details of the proposed curriculum domain adaptation for semantic segmentation of urban scene images. Unlike previous work [43, 27] that aligns the domains via an intermediate feature space and thereby implicitly assumes the existence of the same decision function for the two domains, it is our intuition that,

for structured prediction (i.e., semantic segmentation here), the cross-domain generalization of machine learning models can be more efficiently improved if we avoid this assumption and instead train them subject to necessary properties they should retain in the target domain.

**Preliminaries.** In particular, the properties are about the pixel-wise category labels $Y_t \in \mathbb{R}^{W \times H \times C}$ of an arbitrary image $I_t \in \mathbb{R}^{W \times H}$ from the target domain, where $W$ and $H$ are the width and height of the image, respectively, and $C$ is the number of categories. We use one-hot vector encoding for the groundtruth labels, i.e., $Y_t(i, j, c)$ takes the value of 0 or 1 and the latter means that the $c$-th label is assigned by a human annotator to the pixel at $(i, j)$. Correspondingly, the prediction $\widehat{Y}_t(i, j, c) \in [0, 1]$ by a segmentation network is realized by a softmax function per pixel.

We express each target property in the form of a distribution $p_t \in \Delta$ over the $C$ categories, where $p_t(c)$ represents the occupancy proportion of the category $c$ over the $t$-th target image or a superpixel of the image. Therefore, one can immediately calculate the distribution $p_t$ given the human annotations $Y_t$ to the image. For instance, the image level label distribution is expressed by

$$p_t(c) = \frac{1}{WH} \sum_{i=1}^{W} \sum_{j=1}^{H} Y_t(i, j, c), \quad \forall c. \qquad (1)$$

Similarly, we can compute the target property/distribution from the network predictions $\widehat{Y}_t$ and denote it by $\widehat{p}_t$.

### 3.1. Domain adaptation using the target properties

Ideally, we would like to have a segmentation network to imitate human annotators on the target domain. Therefore, necessarily, the properties of their annotation results should be the same too. We capture this notion by minimizing the cross entropy $\mathcal{C}(p_t, \widehat{p}_t) = H(p_t) + \mathrm{KL}(p_t, \widehat{p}_t)$ at training, where the first term of the right-hand side is the entropy and the second is the KL-divergence.

Given a mini-batch consisting of both source images ($S$) and target images ($T$), the overall objective function for training the cross-domain generalizing segmentation network is,

$$\min \ \frac{\gamma}{|S|} \sum_{s \in S} \mathcal{L}\left(Y_s, \widehat{Y}_s\right) + \frac{1 - \gamma}{|T|} \sum_{t \in T} \sum_{k} \mathcal{C}\left(p_t^k, \widehat{p}_t^k\right) \quad (2)$$

where $\mathcal{L}$ is the pixel-wise cross-entropy loss defined over the sufficiently labeled source domain images, enforcing the network to have the pixel level discriminative capabilities, and the second term is over the unlabeled target domain images, hinting the network what necessary properties its predictions should have in the target domain. We use $\gamma \in [0, 1]$ to balance the two strengths in training and superscript $k$ to index different types of label distributions.

Note that in the domain adaptation context, we actually cannot directly compute the label distribution $p_t^k$ from the groundtruth annotations of the target domain. Nonetheless, estimating them using the labeled source data is easier than assigning labels to every single pixel of the target images. We present the details in the next section.

**Remarks.** Mathematically, the objective function has a similar form as in model compression [9, 26]. We thus borrow some concepts to gain more intuitive understanding about our domain adaptation procedure. The "student" network follows a curriculum to learn simple knowledge about the target domain before it addresses the hard one of semantically segmenting images. The models inferring the target properties act like "teachers", as they hint what label distributions the final solution (image annotation) may have in the target domain at the image and superpixel levels.

Another perspective is to understand the target properties as a posterior regularization [16] for the network. The posterior regularization can conveniently encode a priori knowledge into the objective function. Some applications include weakly supervised segmentation [44, 49] and detection [5], and rule-regularized training of neural networks [29]. In addition to the domain adaptation setting and novel target properties, another key distinction of our work is that we decouple the label distributions from the network predictions and thus avoid the EM type of optimizations. Our approach learns the segmentation network with almost effortless changes to the popular deep learning tools.

## 3.2. Inferring the target properties

Thus far we have presented the "hard" task in the curriculum domain adaptation. In this section, we describe the "easy" ones, i.e., how to infer the target properties without accessing the image annotations of the target domain. Our contributions also include selecting the particular property of label distributions to constitute the simple tasks.

### 3.2.1 Global label distributions of images

Due to the domain disparity, a baseline segmentation network trained on the source domain (i.e., using the first term of eq. (2)) could be easily crippled given the target images. In our experiments, we find that our baseline network constantly mistakes streets for sidewalks and/or cars (cf. Figure 2). Consequently, the predicted labels for the pixels are highly disproportionate.

To rectify this, we employ the label distribution $p_t$ over the global image as our first property (cf. eq. (1)). Without access to the target labels, we have to train machine learning models from the labeled source images to estimate the label distribution $p_t$ for the target image. Nonetheless, we argue that this is less challenging than generating the per-pixel predictions despite that both tasks are influenced by the domain mismatch.

In our experiments, we examine different approaches to this task. We extract image features using the Inception-Resnet-v2 [57] as the input to the following models.

**Logistic regression.** Although multinomial logistic regression (LR) is mainly used for classification, its output is actually a valid distribution over the categories. For our purpose, we thus train it by replacing the one-hot vectors in the cross-entropy loss with the groundtruth label distribution $p_s$, which is calculated using eq. (1) and the available human labels of the source domain. Given a target image, we directly take the LR's output as the predicted label distribution.

**Mean of nearest neighbors.** We also test a nonparametric method by simply retrieving the nearest neighbors (NNs) for a target image and then transferring the mean of the NNs' label distributions to the target image. We use the $\ell_2$ distance for the NN retrieval.

Finally, we include two dumb predictions as the control experiment. One is, for any target image, to output the mean of all the label distributions in the source domain (**source mean**), and the other is to output a **uniform distribution**.

### 3.2.2 Local label distributions of landmark superpixels

The image level label distribution globally penalizes potentially disproportional segmentation output on the target domain, and yet is inadequate in providing spatial constraints. In this section, we consider the use of label distributions over some superpixels as the anchors to drive the network towards spatially desired target properties.

Note that it is not necessary, and is even harmful, to use all of the superpixels in a target image to regularize the segmentation network, because that would be too strong a force and may overrule the pixel-wise discriminativeness revealed by the labeled source images, especially when the label distributions are not inferred accurately enough.

In order to have the dual effect of both estimating the label distributions of superpixels and filtering the superpixels, we simplify the problem and employ a linear SVM in this work. In particular, we segment each image into 100 superpixels using linear spectral clustering [33]. For the superpixels of the source domain, we are able to assign a single dominant label to each of them, and then use the labels and the corresponding features extracted from the superpixels to train a multi-class SVM. Given a test superpixel of a target image, the multi-class SVM returns a class label as well as a decision value, which is interpreted as the confidence score about classifying this superpixel. We keep the top 60% superpixels, called landmark superpixels, in the target domain and calculate their label distributions as the second type of "easy" tasks. In particular, the class label of a landmark superpixel is encoded into a one-hot vector, which serves as a valid distribution about the categories in the landmark
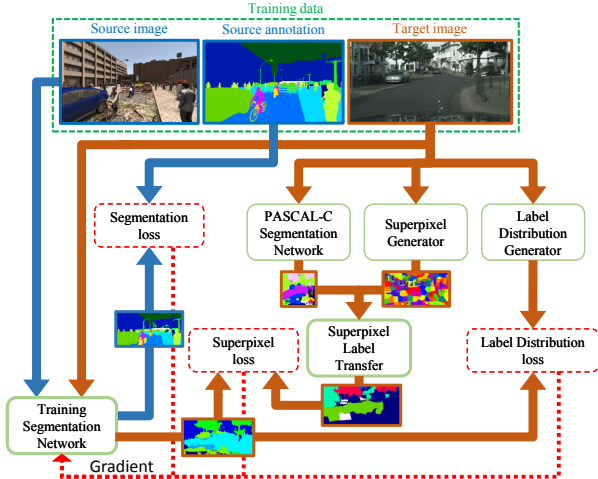
Figure 1: The overall framework of curriculum domain adaptation for semantic segmentation of urban scenes.

superpixel. Albeit simple, we find this method works very well in our experiments.

We encode both visual and contextual information to represent a superpixel. First, we use the FCN-8s [35] pretrained on the PASCAL CONTEXT [38] dataset, which has 59 distinct classes, to obtain 59 detection scores for each pixel. We then average them within each superpixel. Finally, we represent a superpixel by the concatenation of the 59D vectors of itself, its left and right superpixels, as well as the two respectively above and below it.

### 3.3. Curriculum domain adaptation: recapitulation

We recap the proposed curriculum domain adaptation using Figure 1 before presenting the experiments in the next section. Our main idea is to execute the domain adaptation step by step, starting from the easy tasks that are less sensitive to the domain discrepancy than the semantic segmentation. We choose the labels distributions over global images and local landmark superpixels in this work; more tasks will be explored in the future. The solutions to them provide useful gradients originating from the target domain (cf. the arrows with brown color in Figure 1), while the source domain feeds the network with well-labeled images and segmentation masks (cf. the dark blue arrows in Figure 1).

## 4. Experiments

In this section, we describe the experimental setup and compare the results of our approach, its variations, and some existing baseline methods.

### 4.1. Segmentation network and optimization

In our experiments, we use FCN-8s [35] as our semantic segmentation network. We initialize its convolutional lay-

ers with VGG-19 [54], and then train it using the AdaDelta optimizer [67] with default parameters. Each mini-batch is comprised of five source images and five randomly chosen target images. When we train the baseline network with no adaptation, however, we try to use the largest possible mini-batch that includes 15 source images. The network is implemented in Keras [11] and Theano [1]. We train different versions of the network on a single Tesla K40 GPU.

Unlike the existing deep domain adaptation methods [17, 18, 36, 27] which introduce regularization to the intermediate layers, we only revise the loss function over the output. Hence, our curriculum domain adaptation can be readily applied to other segmentation networks (e.g., [39, 10]).

### 4.2. Datasets and evaluation

We use the publicly available **Cityscpaes** [12] and **SYNTHIA** [48] datasets in our experiments.

Cityscapes is a real-world, vehicle-egocentric image dataset collected in 50 cities in Germany and nearby countries. It provides four disjoint subsets: 2,993 training images, 503 validation image, 1,531 test images, and 20,021 auxiliary images. All the training, validation, and test images are accurately annotated with per pixel category labels, while the auxiliary set is coarsely labeled. There are 34 distinct categories in the dataset.

SYNTHIA [48] is a large dataset of synthetic images and provides a particular subset, called SYNTHIA-RAND-CITYSCAPES, to pair with Cityscapes. This subset contains 9,400 images that are automatically annotated with 12 object categories, one void class, and some unnamed classes. Note that the virtual city used to generate the synthetic images does not correspond to any of the real cities covered by Cityscapes. We abbreviate SYNTHIA-RAND-CITYSCAPES to SYNTHIA hereon.

**Domain idiosyncrasies.** Although both datasets depict urban scenes, and SYNTHIA is created to be as photo-realistic as possible, they are mismatched domains in several ways. The most noticeable difference is probably the coarse-grained textures in SYNTHIA; very similar texture patterns repeat in a regular manner across different images. In contrast, the Cityscapes images are captured by high-quality dash-cameras. Another major distinction is the variability in view angles. Since Cityscapes images are recorded by the dash cameras mounted on a moving car, they are viewed from almost a constant angle that is about parallel to the ground. More diverse view angles are employed by SYNTHIA — it seems like some cameras are placed on the buildings that are significantly higher than a bus. Finally, some of the SYNTHIA images are severely shadowed by extreme lighting conditions, while we find no such conditions in the Cityscapes images. These combined factors, among others, make domain adaptation from SYNTHIA to Cityscapes a very challenging problem.
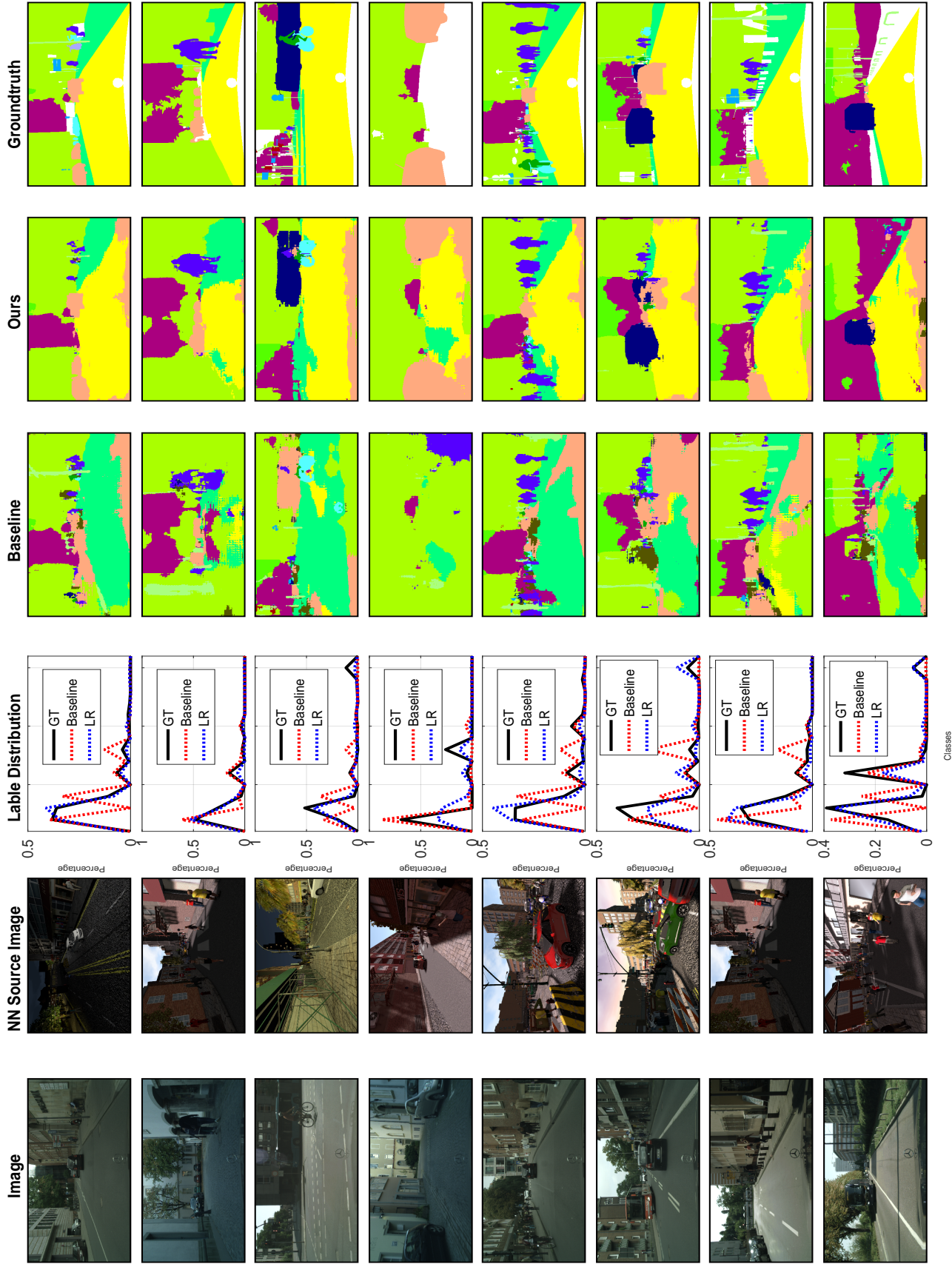
Figure 2: Qualitative semantic segmentation results on the Cityscapes dataset [48] (target domain). For each target image in the first column, we retrieve its nearest neighbor from the SYNTHIA [12] dataset (source domain). The third column plots the label distributions due to the groundtruth pixel-wise semantic annotation, the predictions by the baseline network with no adaptation, and the inferred distribution by logistic regression. The last three columns are the segmentation results by the baseline network, our domain adaptation approach, and human annotators.

Table 1: The $\chi^2$ distances between the groundtruth label distributions and those predicted by different methods.

| Method | Uniform | NoAdapt | Src mean | NN | **LR** |
|---|---|---|---|---|---|
| $\chi^2$ Distance | 1.13 | 0.65 | 0.44 | 0.33 | **0.27** |

Figure 2 shows some example images from both datasets. We pair each Cityscpaes image with its nearest neighbor in SYNTHIA, retrieved by the Inception-Resnet-v2 [57] features. However, the cross-dataset nearest neighbors are visually very different from the query images, verifying the dramatic disparity between the two domains.

**Experiment setup.** Since our ultimate goal is to solve the semantic segmentation problem for real images of urban scenes, we take Cityscapes as the target domain and SYNTHIA as the source domain. The Cityscapes validation set is used as our test set. We split 500 images out of the Cityscpaes training set for the validation purpose (e.g., to monitor the convergence of the networks). In training, we randomly sample mini-batches from both the images (and their labels) of SYNTHIA and the remaining images of Cityscapes yet with no labels.

As in [27], we manually find 16 common classes between the two datasets: sky, building, road, sidewalk, fence, vegetation, pole, car, traffic sign, person, bicycle, motorcycle, traffic light, bus, wall, and rider. The last four are unnamed and yet labeled in SYNTHIA.

**Evaluation.** We use the evaluation code released along with the Cityscapes dataset to evaluate our results. It calculates the PASCAL VOC intersection-over-union, i.e., IoU $= \frac{TP}{TP+FP+FN}$ [14], where TP, FP, and FN are the numbers of true positive, false positive, and false negative pixels, respectively, determined over the whole test set. Since we have to resize the images before feeding them to the segmentation network, we resize the output segmentation mask back to the original image size before running the evaluation against the groundtruth annotations.

### 4.3. Results of inferring global label distributions

Before presenting the final semantic segmentation results, we first compare the different approaches to inferring the global label distributions of the target images (cf. Section 3.2.1). We report the results on the held-out validation images of Cityscapes in this experiment, and then select the best method for the remaining experiments.

In Table 1, we compare the estimated label distributions with the groundtruth ones using the $\chi^2$ distance, the smaller the better. We see that the baseline network (NoAdapt), which is directly learned from the source domain without any adaptation methods, outperforms the dumb uniform distribution (Uniform) and yet no other methods. This confirms that the baseline network gives rise to severely disproportionate predictions over the target domain.

Another dumb prediction (Src mean), i.e., using the mean of all label distributions over the source domain as the prediction for the target images, however, performs reasonably well. To some extent, this indicates the value of the simulated source domain for the semantic segmentation task of urban scenes.

Finally, the nearest neighbors (NN) based method and the multinomial logistic regression (LR) (cf. Section 3.2.1) perform the best. We use the output of LR on the target domain in our remaining experiments.

### 4.4. Comparison results

We report the final semantic segmentation results on the test data of the target domain in this section. We compare our approach to the following competing methods.

**No adaptation (NoAdapt).** We directly train the FCN-8s model on SYNTHIA without applying any domain adaptation methods. This is the most basic baseline for our experiments.

**Superpixel classification (SP).** Recall that we have trained a multi-class SVM using the dominant labels of the superpixels in the source domain. We then use them to classify the target superpixels.

**Landmark superpixels (SP Lndmk).** Since we keep the top 60% most confidently classified superpixels as the landmarks to regularize our segmentation network during training (cf. Section 3.2.2), it is also interesting to examine the classification results of these superpixels. We run the evaluation after assigning the void class label to the other pixels of the images.

In addition to the IoU, we have also evaluated the classification results of the superpixels by accuracy. We find that the classification accuracy is 71% for all the superpixels of the target domain, while for the selected 60% landmark superpixels, the classification accuracy is more than 88%.

**FCNs in the wild (FCN Wld).** Hoffman et al.'s work [27] is the only existing one addressing the same problem as ours, to the best of our knowledge. They introduce a pixel-level adversarial loss to the intermediate layers of the network and impose constraints to the network output. Their experimental setup is about identical to ours except that they do not specify which part of Cityscapes is considered as the test set. Nonetheless, we include their results for comparison to put our work in a better perspective.

The comparison results are shown in Table 2. Immediately, we note that all our domain adaptation results are significantly better than those without adaptation (NoAdapt).

Table 2: Comparison results for the semantic segmentation of the Cityscapes images [12] by adapting from SYNTHIA [48].

| Method    % | IoU | Class-wise IoU | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | bike | fence | wall | t-sign | pole | mbike | t-light | sky | bus | rider | veg | bldg | car | person | sidewalk | road |
| NoAdapt [27] | 17.4 | 0.0 | 0.0 | 1.2 | 7.2 | 15.1 | 0.1 | 0.0 | 66.8 | _3.9_ | 1.5 | 30.3 | 29.7 | 47.3 | 51.1 | 17.7 | 6.4 |
| FCN Wld [27] | _20.2_ | _0.6_ | _0.0_ | **4.4** | **11.7** | _20.3_ | _0.2_ | _0.1_ | _68.7_ | 3.2 | _3.8_ | _42.3_ | _30.8_ | **54.0** | **51.2** | _19.6_ | _11.5_ |
| NoAdapt | 22.0 | **18.0** | 0.5 | 0.8 | 5.3 | 21.5 | 0.5 | 8.0 | _75.6_ | 4.5 | **9.0** | _72.4_ | 59.6 | 23.6 | _35.1_ | 11.2 | 5.6 |
| **Ours (I)** | _25.5_ | 16.7 | **0.8** | _2.3_ | _6.4_ | 21.7 | 1.0 | 9.9 | 59.6 | _12.1_ | 7.9 | 70.2 | _67.5_ | _32.0_ | 29.3 | _18.1_ | _51.9_ |
| SP Lndmk | 23.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | **83.1** | **26.1** | 0.0 | 73.1 | 67.7 | 41.1 | 5.8 | 10.6 | 60.8 |
| SP | 25.6 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 80.5 | 22.1 | 0.0 | 71.9 | 69.3 | _45.9_ | 24.6 | 19.8 | **75.0** |
| **Ours (SP)** | _28.1_ | 10.2 | _0.4_ | _0.1_ | 2.7 | _8.1_ | _0.8_ | _3.7_ | 68.7 | 21.4 | _7.9_ | _75.5_ | _74.6_ | 42.9 | _47.3_ | _23.9_ | 61.8 |
| **Ours (I+SP)** | **29.0** | 13.1 | 0.5 | 0.1 | 3.0 | 10.7 | 0.7 | 3.7 | 70.6 | 20.7 | 8.2 | **76.1** | 74.9 | 43.2 | 47.1 | **26.1** | 65.2 |

We denote by (**Ours (I)**) the network trained using the global label distributions over the target images (and the labeled source images). Although one may wonder that the image-wise label distributions are too high-level to supervise the pixel-wise discriminative network, the gain is actually significant. They are able to correct some obvious errors of the baseline network, such as the disproportional predictions about road and sidewalk (cf. the results of **Ours (I)** vs. NoAdapt in the last two columns).

It is interesting to see that both superpixel classification-based segmentation results (SP and SP Lndmk) are also better than the baseline network (NoAdapt). The label distributions obtained over the landmark superpixels boost the segmentation network (**Ours (SP)**) to the mean IoU of 28.1%, which is better than those by either superpixel classification or the baseline network individually. We have also tried to use the label distributions over all the superpixels to train the network, and observe little improvement over NoAdapt. This is probably because it is too forceful to regularize the network output at every single superpixel especially when the estimated label distributions are not accurate enough.

The superpixel-based methods, including **Ours (SP)**, miss small objects such as fences, traffic lights (t-light), and traffic signs (t-sign), and instead are very accurate for categories like the sky, road, and building, that typically occupy larger image regions. On the contrary, the label distributions on the images give rise to a network (**Ours (I)**) that performs better on the small objects than **Ours (SP)**. In other words, they mutually complement to some extent. Re-training the network by using the label distributions over both global images and local landmark superpixels (**Ours (I+SP)**), we achieve the best semantic segmentation results on the target domain. In the future work, it is worth exploring other target properties, perhaps still in the form of label distributions, that handle the small objects well, in order to further complement the superpixel-level label distributions.

**Comparison with FCNs in the wild [27].** Although we use the same segmentation network (FCN-8s) as [27], our baseline results (NoAdapt) are better than those reported in [27]. This may be due to subtle differences in terms of implementation or experimental setup. Although our own baseline results are superior, we gain larger improvements (7%) over them than the performance gain of [27] (3%) over the seemingly underperforming baseline network there.

**Comparison with learning domain-invariant features.** At our first attempt to solve the domain adaptation problem for the semantic segmentation of urban scenes, we tried to learn domain invariant features following the deep domain adaptation methods [36] for classification. In particular, we impose the maximum mean discrepancy [24] over the layer before the output. We name such network layer the feature layer. Since there are virtually three output layers in FCN-8s, we experiment with all the three feature layers correspondingly. We have also tested the domain adaptation by reversing the gradients of a domain classifier [17]. However, none of these efforts lead to any noticeable gain over the baseline network so the results are omitted.

## 5. Conclusion

In this paper, we address domain adaptation for the semantic segmentation of urban scenes. We propose a curriculum style approach to this problem. We learn to estimate the global label distributions of the images and local label distributions of the landmark superpixels of the target domain. Such tasks are easier to solve than the pixel-wise label assignment. Therefore, we use their results to effectively regularize our training of the semantic segmentation network such that its predictions meet the inferred label distributions over the target domain. Our method outperforms several competing methods that do domain adaptation from simulated images to real photos of urban traffic scenes. In future work, we will explore more target properties that can be conveniently inferred to enrich our curriculum domain adaptation framework.

# References

[1] R. Al-Rfou, G. Alain, A. Almahairi, C. Angermueller, D. Bahdanau, N. Ballas, F. Bastien, J. Bayer, A. Belikov, A. Belopolsky, and others. Theano: A Python framework for fast computation of mathematical expressions. *arXiv preprint arXiv:1605.02688*, 2016. 5

[2] R. Aljundi, R. Emonet, D. Muselet, and M. Sebban. Landmarks-based kernelized subspace alignment for unsupervised domain adaptation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 56–63, 2015. 1

[3] V. Badrinarayanan, A. Kendall, and R. Cipolla. Segnet: A deep convolutional encoder-decoder architecture for image segmentation. *arXiv preprint arXiv:1511.00561*, 2015. 3

[4] Y. Bengio, J. Louradour, R. Collobert, and J. Weston. Curriculum learning. In *Proceedings of the 26th annual international conference on machine learning*, pages 41–48. ACM, 2009. 2

[5] H. Bilen, M. Pedersoli, and T. Tuytelaars. Weakly supervised object detection with posterior regularization. In *BMVA British Machine Vision Conference (BMVC)*, 2014. 4

[6] E. Borenstein and S. Ullman. Class-specific, top-down segmentation. In *European conference on computer vision*, pages 109–122. Springer, 2002. 1

[7] K. Bousmalis, N. Silberman, D. Dohan, D. Erhan, and D. Krishnan. Unsupervised pixel-level domain adaptation with generative adversarial networks. *arXiv preprint arXiv:1612.05424*, 2016. 3

[8] G. J. Brostow, J. Fauqueur, and R. Cipolla. Semantic object classes in video: A high-definition ground truth database. *Pattern Recognition Letters*, 30(2):88–97, 2009. 3

[9] C. Bucilu, R. Caruana, and A. Niculescu-Mizil. Model compression. In *Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 535–541. ACM, 2006. 4

[10] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille. Semantic image segmentation with deep convolutional nets and fully connected crfs. *arXiv preprint arXiv:1412.7062*, 2014. 2, 3, 5

[11] F. Chollet. keras. https://github.com/fchollet/keras, 2015. 5

[12] 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 *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3213–3223, 2016. 1, 3, 5, 6, 8, 12

[13] J. Dai, K. He, and J. Sun. Instance-aware semantic segmentation via multi-task network cascades. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3150–3158, 2016. 3

[14] M. Everingham, S. A. Eslami, L. Van Gool, C. K. Williams, J. Winn, and A. Zisserman. The pascal visual object classes challenge: A retrospective. *International Journal of Computer Vision*, 111(1):98–136, 2015. 3, 7

[15] B. Fernando, A. Habrard, M. Sebban, and T. Tuytelaars. Unsupervised visual domain adaptation using subspace align-

[16] K. Ganchev, J. Gillenwater, B. Taskar, et al. Posterior regularization for structured latent variable models. *Journal of Machine Learning Research*, 11(Jul):2001–2049, 2010. 2, 4

[17] Y. Ganin and V. Lempitsky. Unsupervised domain adaptation by backpropagation. In *International Conference on Machine Learning (ICML)*, pages 1180–1189, 2015. 2, 3, 5, 8

[18] Y. Ganin, E. Ustinova, H. Ajakan, P. Germain, H. Larochelle, F. Laviolette, M. Marchand, and V. S. Lempitsky. Domain-adversarial training of neural networks. *CoRR*, arXiv:1505.07818, 2015. 2, 3, 5

[19] A. Geiger, P. Lenz, C. Stiller, and R. Urtasun. Vision meets robotics: The kitti dataset. *The International Journal of Robotics Research*, 32(11):1231–1237, 2013. 1

[20] B. Gong, K. Grauman, and F. Sha. Connecting the dots with landmarks: Discriminatively learning domain invariant features for unsupervised domain adaptation. In *International Conference on Machine Learning (ICML)*, pages 222–230, 2013. 1

[21] B. Gong, F. Sha, and K. Grauman. Overcoming dataset bias: An unsupervised domain adaptation approach. In *NIPS Workshop on Large Scale Visual Recognition and Retrieval (LSVRR)*, 2012. 2

[22] B. Gong, Y. Shi, F. Sha, and K. Grauman. Geodesic flow kernel for unsupervised domain adaptation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2066–2073, 2012. 1, 2

[23] R. Gopalan, R. Li, and R. Chellappa. Domain adaptation for object recognition: An unsupervised approach. In *IEEE International Conference on Computer Vision (ICCV)*, pages 999–1006, 2011. 1

[24] A. Gretton, K. M. Borgwardt, M. J. Rasch, B. Schölkopf, and A. Smola. A kernel two-sample test. *Journal of Machine Learning Research*, 13(Mar):723–773, 2012. 8

[25] A. Gretton, A. Smola, J. Huang, M. Schmittfull, K. Borgwardt, and B. Schölkopf. Covariate shift by kernel mean matching. In J. Quiñonero-Candela, M. Sugiyama, A. Schwaighofer, and N. D. Lawrence, editors, *Dataset Shift in Machine Learning*. The MIT Press, 2008. 2

[26] G. Hinton, O. Vinyals, and J. Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015. 4

[27] 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. 1, 2, 3, 5, 7, 8, 12

[28] S. Hong, J. Oh, H. Lee, and B. Han. Learning transferrable knowledge for semantic segmentation with deep convolutional neural network. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3204–3212, 2016. 3

[29] Z. Hu, X. Ma, Z. Liu, E. Hovy, and E. Xing. Harnessing deep neural networks with logic rules. *arXiv preprint arXiv:1603.06318*, 2016. 4

ment. In *IEEE International Conference on Computer Vision (ICCV)*, pages 2960–2967, 2013. 1

[30] A. Khosla, T. Zhou, T. Malisiewicz, A. A. Efros, and A. Torralba. Undoing the damage of dataset bias. In *European Conference on Computer Vision (ECCV)*, pages 158–171, 2012. 2

[31] B. Kulis, K. Saenko, and T. Darrell. What you saw is not what you get: Domain adaptation using asymmetric kernel transforms. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2011. 1

[32] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998. 3

[33] Z. Li and J. Chen. Superpixel segmentation using linear spectral clustering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1356–1363, 2015. 4

[34] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollr, and C. L. Zitnick. Microsoft coco: Common objects in context. In *European Conference on Computer Vision*, pages 740–755. Springer, 2014. 3

[35] J. Long, E. Shelhamer, and T. Darrell. Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3431–3440, 2015. 2, 3, 5

[36] M. Long, Y. Cao, J. Wang, and M. I. Jordan. Learning transferable features with deep adaptation networks. In *International Conference on Machine Learning (ICML)*, 2015. 2, 3, 5, 8

[37] A. M. López, J. Xu, J. L. Gomez, D. Vázquez, and G. Ros. From virtual to real world visual perception using domain adaptation–the dpm as example. *arXiv preprint arXiv:1612.09134*, 2016. 3

[38] R. Mottaghi, X. Chen, X. Liu, N.-G. Cho, S.-W. Lee, S. Fidler, R. Urtasun, and A. Yuille. The role of context for object detection and semantic segmentation in the wild. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 891–898, 2014. 5

[39] H. Noh, S. Hong, and B. Han. Learning deconvolution network for semantic segmentation. In *IEEE International Conference on Computer Vision (ICCV)*, pages 1520–1528, 2015. 2, 3, 5

[40] S. J. Pan, J. T. Tsang, Ivor W.and Kwok, and Q. Yang. Domain adaptation via transfer component analysis. *Transactions on Neural Networks*, 22(2):199 – 210, 2011. 1

[41] S. J. Pan and Q. Yang. A survey on transfer learning. *Transactions on Knowledge and Data Engineering*, 22(10):1345–1359, 2010. 1, 2

[42] 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 *Proceedings of the IEEE International Conference on Computer Vision*, pages 1742–1750, 2015. 3

[43] V. M. Patel, R. Gopalan, R. Li, and R. Chellappa. Visual domain adaptation: A survey of recent advances. *Signal Processing Magazine*, 32(3):53–69, 2015. 1, 2, 3

[44] D. Pathak, P. Krahenbuhl, and T. Darrell. Constrained convolutional neural networks for weakly supervised segmentation. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1796–1804, 2015. 2, 3, 4

[45] X. Peng and K. Saenko. Synthetic to real adaptation with deep generative correlation alignment networks. *arXiv preprint arXiv:1701.05524*, 2017. 3

[46] P. O. Pinheiro and R. Collobert. From image-level to pixel-level labeling with convolutional networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1713–1721, 2015. 3

[47] S. R. Richter, V. Vineet, S. Roth, and V. Koltun. Playing for data: Ground truth from computer games. In *European Conference on Computer Vision*, pages 102–118. Springer, 2016. 3, 11, 12

[48] 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 *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3234–3243, 2016. 1, 3, 5, 6, 8

[49] G. Ros, S. Stent, P. F. Alcantarilla, and T. Watanabe. Training constrained deconvolutional networks for road scene semantic segmentation. *arXiv preprint arXiv:1604.01545*, 2016. 3, 4

[50] K. Saenko, B. Kulis, M. Fritz, and T. Darrell. Adapting visual category models to new domains. In *European Conference on Computer Vision (ECCV)*, pages 213–226, 2010. 2

[51] A. Shafaei, J. J. Little, and M. Schmidt. Play and learn: using video Games to train computer vision models. *arXiv preprint arXiv:1608.01745*, 2016. 3

[52] J. Shotton, M. Johnson, and R. Cipolla. Semantic texton forests for image categorization and segmentation. In *2008 IEEE Conference on Computer Vision and Pattern Recognition*, pages 1–8, June 2008. 3

[53] A. Shrivastava, T. Pfister, O. Tuzel, J. Susskind, W. Wang, and R. Webb. Learning from simulated and unsupervised images through adversarial training. *arXiv preprint arXiv:1612.07828*, 2016. 1, 3

[54] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *CoRR*, arXiv:1409.1556, 2014. 5

[55] B. Sun, J. Feng, and K. Saenko. Return of frustratingly easy domain adaptation. In *AAAI Conference on Artificial Intelligence (AAAI)*, 2016. 1

[56] B. Sun and K. Saenko. From virtual to reality: Fast adaptation of virtual object detectors to real domains. In *BMVA British Machine Vision Conference (BMVC)*, 2014. 3

[57] C. Szegedy, S. Ioffe, V. Vanhoucke, and A. Alemi. Inception-v4, inception-resnet and the impact of residual connections on learning. *arXiv preprint arXiv:1602.07261*, 2016. 4, 7

[58] J. Tighe and S. Lazebnik. Superparsing: scalable nonparametric image parsing with superpixels. In *European conference on computer vision*, pages 352–365. Springer, 2010. 3

[59] T. Tommasi, N. Patricia, B. Caputo, and T. Tuytelaars. A deeper look at dataset bias. *CoRR*, arXiv:1505.01257, 2015. 2

[60] A. Torralba and A. A. Efros. Unbiased look at dataset bias. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2011. 2

[61] E. Tzeng, J. Hoffman, T. Darrell, and K. Saenko. Simultaneous deep transfer across domains and tasks. In *IEEE*

*International Conference on Computer Vision (ICCV)*, pages 4068–4076, 2015. 2, 3

[62] E. Tzeng, J. Hoffman, N. Zhang, K. Saenko, and T. Darrell. Deep domain confusion: Maximizing for domain invariance. *CoRR*, arXiv:1412.3474, 2014. 2

[63] D. Vazquez, A. M. López, J. Marín, D. Ponsa, and D. Gerónimo. Virtual and real world adaptation for pedestrian detection. *Transactions of Pattern Recognition and Machine Analyses (PAMI)*, 36(4):797 – 809, 2014. 3

[64] Z. Wu, C. Shen, and A. v. d. Hengel. Wider or Deeper: Revisiting the ResNet Model for Visual Recognition. *arXiv preprint arXiv:1611.10080*, 2016. 3

[65] J. Xu, S. Ramos, D. Vazquez, and A. López. Hierarchical adaptive structural svm for domain adaptation. *CoRR*, arXiv:1408.5400, 2014. 3

[66] M. Yamada, L. Sigal, and Y. Chang. Domain adaptation for structured regression. *International journal of computer vision*, 109(1-2):126–145, 2014. 2

[67] M. D. Zeiler. ADADELTA: an adaptive learning rate method. *arXiv preprint arXiv:1212.5701*, 2012. 5

[68] C. Zhang, L. Wang, and R. Yang. Semantic segmentation of urban scenes using dense depth maps. In *European Conference on Computer Vision*, pages 708–721. Springer, 2010. 3

[69] H. Zhao, J. Shi, X. Qi, X. Wang, and J. Jia. Pyramid Scene Parsing Network. *arXiv preprint arXiv:1612.01105*, 2016. 3

## GTA→Cityscapes

The main text above has been accepted to IEEE International Conference on Computer Vision (ICCV) 2017. After the paper submission, we have been continuously working on the project and have got more results. We include them below to complement the experiments in the main text.

The new experiment is basically the same as the one in the main text except that we replace SYNTHIA with the GTA dataset [47]. GTA is a synthetic, vehicle-egocentric image dataset collected from the open world in the realistically rendered computer game Grand Theft Auto V (GTA, or GTA5). It contains 24,996 images, whose semantic segmentation annotations are fully compatible with the classes used in Cityscapes. Hence we use all the 19 official training classes in our experiment. The results are shown in Table 3.

As in the main text, the same observations about our approach apply here. Additionally, we note that the results are overall better than those adapting from SYNTHIA to Cityscapes. This is not surprising, because the GTA images are more photo-realistic than SYNTHIA's.

Table 3: Comparison results for the semantic segmentation of the Cityscapes images [12] by adapting from GTA [47].

| Method    % | IoU | bike | fence | wall | t-sign | pole | mbike | t-light | sky | bus | rider | veg | terrain | train | bldg | car | person | truck | sidewalk | road |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| NoAdapt [27] | 21.1 | 0.0 | 3.1 | 7.4 | 1.0 | 16.0 | 0.0 | 10.4 | 58.9 | 3.7 | 1.0 | 76.5 | 13 | 0.0 | 47.7 | 67.1 | 36 | 9.5 | 18.9 | 31.9 |
| FCN Wld [27] | 27.1 | 0.0 | 5.4 | **14.9** | 2.7 | 10.9 | 3.5 | 14.2 | 64.6 | 7.3 | 4.2 | **79.2** | 21.3 | 0.0 | 62.1 | **70.4** | **44.1** | 8.0 | **32.4** | 70.4 |
| NoAdapt | 22.3 | 13.8 | 8.7 | 7.3 | **16.8** | **21.0** | 4.3 | 14.9 | 64.4 | 5.0 | **17.5** | 45.9 | 2.4 | 6.9 | 64.1 | 55.3 | 41.6 | 8.4 | 6.8 | 18.1 |
| **Ours (I)** | 23.1 | 9.5 | 9.4 | 10.2 | 14.0 | 20.2 | 3.8 | 13.6 | 63.8 | 3.4 | 10.6 | 56.9 | 2.8 | **10.9** | 69.7 | 60.5 | 31.8 | 10.9 | 10.8 | 26.4 |
| **Ours (SP)** | 27.8 | **15.6** | 11.7 | 5.7 | 12.0 | 9.2 | 12.9 | 15.5 | 64.9 | 15.5 | 9.1 | 74.6 | 11.1 | 0.0 | 70.5 | 56.1 | 34.8 | 15.9 | 21.8 | 72.1 |
| SP Lndmk | 21.4 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | **82.9** | 10.0 | 0.0 | 74.5 | 22.5 | 0.0 | 69.9 | 52.7 | 13.1 | 11.2 | 8.0 | 61.8 |
| SP | 26.8 | 0.3 | 4.1 | 7.6 | 0.0 | 0.2 | 0.9 | 0.0 | 81.6 | **25.3** | 3.5 | 73.0 | **32.1** | 0.0 | 71.0 | 61.9 | 26.2 | **30.4** | 19.2 | 71.8 |
| **Ours (I+SP)** | **28.9** | 14.6 | **11.9** | 6.0 | 11.1 | 8.4 | **16.8** | 16.3 | 66.5 | 18.9 | 9.3 | 75.7 | 13.3 | 0.0 | **71.7** | 55.2 | 38.0 | 18.8 | 22.0 | **74.9** |