# Pixel and feature level based domain adaptation for object detection in autonomous driving

Yuhu Shan, Wen Feng Lu, Chee Meng Chew

## Abstract

Annotating large-scale datasets to train modern convolutional neural networks is prohibitively expensive and time-consuming for many real tasks. One alternative is to train the model on labeled synthetic datasets and apply it in the real scenes. However, this straightforward method often fails to generalize well mainly due to the domain bias between the synthetic and real datasets. Many unsupervised domain adaptation (UDA) methods were introduced to address this problem but most of them only focused on the simple classification task. This paper presents a novel UDA model which integrates both image and feature level based adaptations to solve the cross-domain object detection problem. We employ objectives of the generative adversarial network and the cycle consistency loss for image translation. Furthermore, region proposal based feature adversarial training and classification are proposed to further minimize the domain shifts and preserve the semantics of the target objects. Extensive experiments are conducted on several different adaptation scenarios, and the results demonstrate the robustness and superiority of the proposed method.

*Keywords:* Autonomous driving, Convolutional neural network, Generative adversarial network, Object detection, Unsupervised domain adaptation

## 1. Introduction

Object detection aims to assign each object a bounding box along with class label, e.g., "pedestrian", "bicycle", "motorcycle" or "car" in an image. It plays an important role in modern autonomous driving systems since it is crucial to detect other traffic participants as shown in Fig. 1. Despite the performance of object detection algorithms has been greatly improved since the introduction of AlexNet [1] in 2012, it is still far from satisfactory when it comes to the practical applications, mainly due to the limited data and expensive labeling cost. Supervised learning algorithms based on deep neural networks require large number of fine labeled images, which are extremely difficult to acquire in real cases. For example, it takes almost ninety minutes to annotate one image from the Cityscapes dataset [2] for driving scene understanding. Even it is already one of the largest driving scene datasets, there are only 2975 training images with fine labels. One promising method to address this problem is to train models on synthetic datasets. Fortunately, with the great progress achieved in graphics and simulation infrastructure, many large-scale datasets with high quality annotations have been produced to simulate different real scenes. However, models trained purely on the rendered images cannot be generalized to real ones, because of the domain shift [3, 4] problem.

During the past several years, many researchers have proposed various unsupervised domain adaptation (UDA) methods [5, 6, 7, 8] to solve this problem. However, most of them only focused on the simple classification task, which is not suitable for more complex vision tasks such as object detection or semantic segmentation. In this paper, we present a new UDA model with adaptations both in pixel and feature
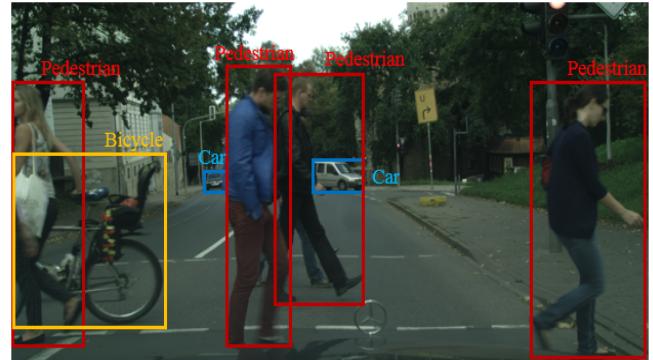


Figure 1: Detection task in autonomous driving system.

spaces to deal with the complex object detection problem in autonomous driving. In the setting of UDA, we generalize the model trained on source dataset with ground truth labels to target dataset without any annotations. Perhaps the most similar works to the proposed work are [9, 10]. In [9], its image generation model is trained based on generative adversarial network (GAN), which is simultaneously combined with the task model for object classification and pose estimation. In [10], the authors attempted to solve the cross-domain semantic segmentation problem based on CycleGAN [11] and traditional segmentation networks. However, this method needs to train an extra segmentation network to preserve the semantics of translated images, which slows down the whole training process. Actually, it is not necessary to pay the same attention to all the image pixels for some specific tasks like object detection. Target objects, such as "car" or "pedestrian" are more important than

other objects or stuffs like "building" or "sky". Therefore, region proposal based feature adversarial training and classification are proposed in this paper to further minimize the domain shifts and preserve the semantics of the target objects. The developed pixel and feature level based domain adaptation modules can be integrated together and trained end-to-end to pursue better detection performance. Qualitative and quantitative results conducted on several datasets show the robustness of the proposed method.

## 2. Related work

### 2.1. Object detection

Object detection [12, 13] as a fundamental problem in computer vision has achieved great progress since 2012 with the development of deep neural networks. Based on Alexnet, many different convolutional neural networks (CNN), such as VGGnet [14], GoogLeNet [15], ResNet [16], DenseNet [17], etc., were proposed to learn more powerful deep features from data. Object detection algorithms also benefit from these architectures since better features are also helpful for other vision tasks. Apart from the different network architectures, recent CNN-based object detectors can be mainly divided into two categories: single-stage detectors (YOLO [18, 19] and SSD [20, 21]) and two-stage detectors (Faster R-CNN [22] and R-FCN [23], etc.). Single stage detectors directly predict object labels and bounding box coordinates within one image based on the default anchor boxes, which is fast but not accurate enough. In contrast, two stage detectors firstly generate large number of region proposals based on the CNN features, and then recognize the proposals with heavy head. Therefore, it has better performance than single stage detectors, but the detection speed is slower. Recently, many researchers also attempt to train the detector with only weak labels [24, 25] or a few ground-truth labels [26] to deal with more practical application scenarios.

### 2.2. Unsupervised domain adaptation

Unsupervised domain adaptation aims to solve the learning problem in a target domain without labels by leveraging training data in a different but related source domain with ground truth labels. Pan et al. [27] proposed Transfer Component Analysis (TCA), a kernel method based on Maximum Mean Discrepancy (MMD) [28], to learn better feature representation across domains. Based on TCA, [29] provided a new method of Joint Distribution Adaptation to jointly adapt both the marginal distribution and the conditional distribution, which was robust for substantial distribution difference. Jiang et al. [30] proposed to close the second moments of the source and target domain distributions to obtain better adaptation performance. Recently, with the advent of deep learning, many works were proposed to learn deep domain invariant features within the neural networks. For example, Long et al. [31, 32] proposed to embed hidden network features in a reproducing kernel Hilbert space and explicitly measure the difference between the two domains with MMD and its variants. Sun et al. [33] tried to minimize domain shift by matching the second order statistics of feature

distributions between the source and the target domains. Rather than explicitly modeling the term to measure the domain discrepancy, another stream of works utilized adversarial training to implicitly find the domain invariant feature representations. Ganin et al. [34, 35] added one more domain classifier to the deep neural network model to classify the domains of the inputs. Adversarial training is conducted through reversing the gradients from the domain classification loss. Rather than using shared feature layers, Tzeng et al. [36] proposed to learn indistinguishable features for the target domain data by training a separate network with objectives similar to the traditional GAN model [37]. Then it is combined with the classifier trained on the source domain for recognition tasks. In [38], the authors argued that each domain should have its own specific features and only part of the features were shared between the different domains. Therefore, they explicitly modeled the private and shared domain representations and only conducted adversarial training on these share ones.

Despite many methods have been proposed to solve the UDA problem, most of them only focus on the simple classification task. Very limited works were conducted related to more complex tasks such as object detection or semantic segmentation. As far as we know, [39] is the first work to deal with the domain adaptation problem for object detection. They conducted adversarial training on features from convolutional and fully connected layers, along with the regular training for detection task. A domain label consistency check was used as a regularizer to help the network to learn better domain invariant features. Although good detection results were achieved in [39], we argue that conducting adaptations on both feature and image pixel spaces would be a better alternative since adapting high-level features can fail to model the low-level image details. Therefore, the proposed work is also closely related to image translation.

### 2.3. Image-to-image translation

Currently, many works have been done to convert images into another style. [40, 41, 42, 43, 44] conducted image translation based on the assumption of available paired training images, which is not fit to the problem of unsupervised domain adaptation. Several other recent works tried to solve the problem with unpaired images. [45, 46] shared part of the network weights to learn the joint distribution of multi-modal images. PixelDA [9] proposed a novel GAN based architecture to convert images across domains. However, this method needs prior knowledge about which parts of the image contributing to the content similarity loss. Neural style transfer [47, 48, 49] is another kind of method to convert one image into another style while preserving its own contents through optimizing the pixel values during back-propagation process. One shortcoming of style transfer is that it only targets on translation between two specific images while not at the dataset level. Recently proposed CycleGAN model [11] is a promising method for unpaired image translation. The authors utilized cycle consistency loss to regularize the generative model and hence to preserve structural information of the transferred image. However, this method can only guarantee that, one area if occupied by one object before the
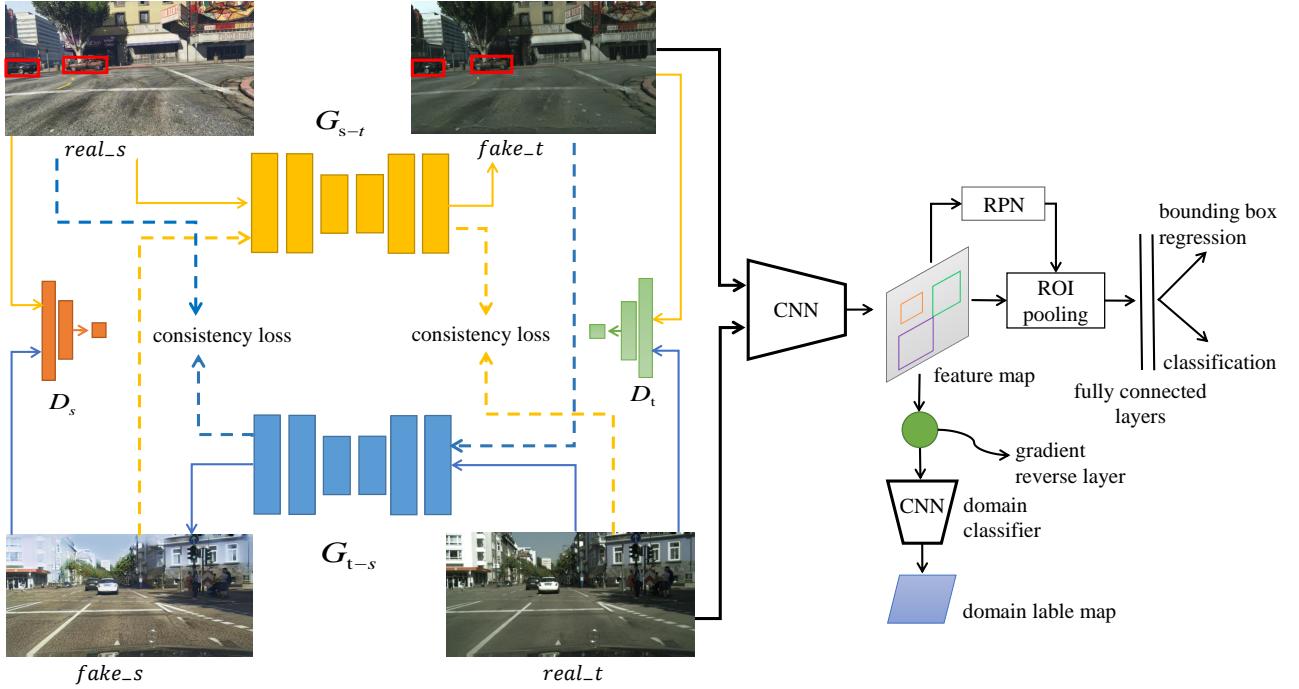
Figure 2: Overall architecture of the proposed learning method. On the left, we show the pixel-level based image translation module, in which source image is firstly converted to the target domain before it is used to train the detection network. On the right, detection network is trained together with the feature-level based adversarial training.

translation, will also get occupied after the generation process. Semantics of the pixels are not guaranteed to be consistent with this only cycle consistence loss.

## 3. Method

We tackle the problem of unsupervised cross-domain object detection with the assumption of available source images $X_s$ with ground truth labels $Y_s$ and target images $X_t$ without any labels. Our target is to train the detection network with source dataset, which also need to perform well on the target dataset. The whole framework is shown in Fig. 2. Our model consists of two modules: 1) pixel-level domain adaptation (PDA) mainly based on CycleGAN, and 2) feature-level domain adaptation (FDA) based on Faster R-CNN. The two modules can be integrated together and trained in an end-to-end way to pursue better performance. Source images are firstly converted into the target image style. Then the transformed images are used to train the object detector and domain classifier, along with the sampled images from the target dataset.

### 3.1. Pixel-level based domain adaptation

We firstly introduce the pixel-level based domain adaptation module. As shown in Fig. 2, two symmetric generative networks $G_{s-t}$ and $G_{t-s}$ are employed to generate images $fake\_t$ and $fake\_s$ separately in two domains. Another two discriminators $D_s$ and $D_t$ are trained to distinguish the real sampled and fake generated images. The whole training process runs

in a min-max manner, in which the generators always try to generate images which cannot be distinguished from the real ones. The discriminators are trained simultaneously to be good at classifying real and fake images. The whole objectives of generator $G_{s-t}$ and discriminator $D_t$ can be formulated as Eq. (1).

$$
\begin{aligned}
L_{GAN}(D_t, G_{s-t}) = \; & \mathbb{E}_{t \sim X(t)}[log D_t(x)] \\
& + \mathbb{E}_{s \sim X(s)}[log(1 - D_t(G_{s-t}(s)))].
\end{aligned}
\tag{1}
$$

Similar equation can also be formulated for $G_{t-s}$ and $D_s$ as $L_{GAN}(D_s, G_{t-s})$, which is ignored here. This kind of GAN objectives can, in theory, learn the mapping functions $G_{s-t}$ and $G_{t-s}$ to produce images identically sampled from the data distribution of $X_t$ and $X_s$. However, it also faces the problem of mode collapse [37] and losing the structural information of the source images. To address these problems, cycle consistency loss is adopted here to force the image $fake\_t$ generated by $G_{s-t}$ to have identical result as $real\_s$ after it is sent into generator $G_{t-s}$ and vice versa. The whole cycle consistency loss is formulated in Eq. (2),

$$
\begin{aligned}
L_{cyc}(G_{t-s}, G_{s-t}) = \; & \mathbb{E}_{t \sim X(t)}[\|G_{s-t}(G_{t-s}(x_t)) - x_t\|_1] \\
& + \mathbb{E}_{s \sim X(s)}[\|G_{t-s}(G_{s-t}(x_s)) - x_s\|_1].
\end{aligned}
\tag{2}
$$

Therefore, the full objective for CycleGAN training is

$$
\begin{aligned}
L_{cyc-gan}(G_{t-s}, G_{s-t}, D_t, D_s) = \; & L_{GAN}(D_t, G_{s-t}) + \\
& L_{GAN}(D_s, G_{t-s}) + \lambda_{cyc} L_{cyc}(G_{t-s}, G_{s-t}),
\end{aligned}
\tag{3}
$$

3

with the target of solving

$$G_{s-t}^*, G_{t-s}^* = \arg \min_{G_{s-t}, G_{t-s}} \max_{D_s, D_t} L_{cyc-gan}(G_{t-s}, G_{s-t}, D_t, D_s). \tag{4}$$

### 3.2. Feature-level based domain adaptation

Our detection network is based on the famous framework of Faster R-CNN. Specifically, region proposal network (RPN) is trained to generate region proposals and Fast R-CNN [13] trained for bounding box classification and regression. A small fully convolutional network is newly added to the framework for further domain adversarial training. Specifically, the inputs are features extracted from the final convolutional features by the corresponding region proposals. Gradients generated by the detection and reversed domain classification losses will flow into the shared convolution layers to learn domain-invariant features for object detection.

*Losses for Faster R-CNN.* Assuming there are $m$ categories in the detection task, the region classification layer will output $m + 1$ dimension probability distribution for each region proposal, $p_{obj} = (p_{obj}^0, p_{obj}^1, ..., p_{obj}^m)$, with one more category for the background. $p_{obj}^*$ is used to represent the ground truth label for the region proposal. Box coordinate $t_{obj} = (t_x, t_y, t_w, t_h)$ is predicted for each possible class by the bounding box regression layer to approach the ground truth regression target $t_{obj}^*$. Here $t_{obj}$ and $t_{obj}^*$ are normalized as [22] for better training of the network. Similarly, we use $p_{rp}$ and $p_{rp}^*$ for the training of RPN. Since RPN is only trained to discriminate object/non-object, labels $p_{rp}^*$ can only be 1 or 0. The full objectives for Faster R-CNN training can be formulated as Eq. (5),

$$\begin{aligned} L_{det}(p, p^*, t, t^*) = L_{cls-det}(p_{obj}, p_{obj}^*) + [p_{obj}^* \geq 1] \\ L_{loc-det}(t_{obj}, t_{obj}^*) + L_{cls-rpn}(p_{rp}, p_{rp}^*) \\ + [p_{rp}^* \geq 1]L_{loc-rpn}(t_{rp}, t_{rp}^*), \end{aligned} \tag{5}$$

in which $L_{cls-det}$ and $L_{cls-rpn}$ indicate the cross-entropy loss for classification. $L_{loc-det}$ and $L_{loc-rpn}$ represent the smooth $L_1$ loss [13] for bounding box regression. The Iverson bracket indicator function $[p^* \geq 1]$ evaluates to 1 when the true object category $p^* \geq 1$ and 0 otherwise.

*Losses for FDA training.* As shown in Fig. 2, domain classifier learns to classify $fake\_t$ as label $d = 0$ and $real\_t$ as label $d = 1$. Then the gradients are reversed before they flow into the shared convolutional layers. Loss of the domain adversarial training is shown in Eq. (6),

$$L_{domain} = -\sum_{i,j}[d \log p_{i,j} + (1 - d)\log(1 - p_{i,j})], \tag{6}$$

in which $p_{i,j}$ indicates the classification output on the *i-th* region proposal of the *j-th* image. Based on the above equations, we can formulate the full training objectives in Eq. (7), where $\lambda_1, \lambda_2$ are the weights to balance the different losses,

$$L_{full} = L_{det} + \lambda_1 L_{domain} + \lambda_2 L_{cyc-gan}. \tag{7}$$

## 4. Implementation

### 4.1. Datasets

To test the validity of the proposed method and also compare with current state-of-the-art (SOTA) work [39], we choose Cityscapes, KITTI [50], Foggy-Cityscapes [51], VKITTI-Rainy [52] and Sim10k [53] datasets for the experiments. Cityscapes dataset has 2975 training images and 500 images for validation. Eight classes of common traffic participants are annotated with instance labels. KITTI is another famous dataset for benchmarking different vision tasks in autonomous driving. There are 7481 labeled training images with bounding boxes for categories "car", "pedestrian" and "cyclist". Foggy-Cityscapes, VKITTI-Rainy and Sim10k are synthetic datasets which simulate different driving scenes. Particularly, Foggy-Cityscapes and VKITTI-Rainy are rendered based on the real Cityscapes and KITTI datasets to simulate the foggy and rainy weathers. Sim10k has 10,000 training images which are collected from the computer game of GTA5 and annotated automatically by access to the original game engine. In Sim10k dataset, only objects "car" are annotated with bounding boxes for the detection task. Therefore, we only calculate and compare the "car" detection result for the experiments based on Sim10k in this paper. Five exampled images randomly sampled from the above five datasets are shown in Fig. 3 to show their domain differences.

### 4.2. Implementation details

We use U-Net structure [54] with skip connections between layers for the two generators in the pixel adaptation module and PatchGAN [43] for the other two discriminators. Instance normalization is adopted since it is more effective as stated in original CycleGAN paper. For the detection network, we use VGG16 [14] as the backbone and a small fully convolutional network for domain adversarial training. Inputs for the feature-level based adversarial training are the cropped Conv5 features of VGG16 based on the 64 region proposals generated by the RPN module. In our practical training process, we choose to firstly pre-train the detection and image translation networks separately and then conduct the end-to-end training based on these two pre-trained models. This mainly considers the fact that most of the generated images are quite noisy in the start training stage of the pixel-level adaptation module. We train the PDA module with Adam optimizer and an initial learning rate of 0.0002. After 30 epochs, the learning rate linearly decays to be zero in the following training process for another 30 epochs. FDA module is trained together with the object detection network with initial learning rate of 0.001 based on the standard SGD algorithm. After 6 epochs, we reduce the learning rate to 0.0001 and train the network for another 3 epochs. Gradients from the domain classifier are reversed before they flow into the shared CNN layers during the back-propagation. For the end-to-end training, all the above initial learning rates are scaled down by ten times. We then finetune the whole network for another 10 epochs with $\lambda_1$ and $\lambda_2$ set as 0.5.

(a) Cityscapes     (b) Foggy-Cityscapes     (c) Sim10k

(d) VKITTI-Rainy     (e) KITTI

Figure 3: Sampled images from the five datasets to show the domain bias.

## 5. Results and discussion

We show our experimental results under three adaptation scenarios: "synthetic to real", "cross different weathers" and "cross different cameras". Specifically, experiments are conducted on "Sim10k → Cityscapes" and "Sim10k → KITTI" for the scenario of "synthetic to real", "Cityscapes → Foggy-Cityscapes" and "KITTI → VKITTI-Rainy" for "cross different weathers", "Cityscapes → KITTI" and "KITTI → Cityscapes" for "cross different cameras". For each dataset pair, we conduct three experiments with considering the PDA, FDA modules and the final end-to-end training. To test the validity of the PDA module, we only train the network for image translation, and use the translated images to train a pure detection network. For the FDA training, we directly use the source images as inputs for the detection training. Randomly sampled images from the target domain are combined with the source images for the feature-level based adversarial training. Finally, we integrate the two modules together and train the whole network in an end-to-end way. All the results are evaluated with the commonly used metrics of Average Precision (AP) and mean Average Precision (mAP). IoU threshold is set as 0.5.

Table 1: Detection results (AP (%) of "car") are evaluated on the Cityscapes validation and KITTI training datasets by using Sim10k as the source dataset.

|  | Sim10k → Cityscapes | Sim10k → KITTI |
| --- | --- | --- |
| Faster R-CNN | 30.1 | 52.7 |
| SOTA[39] | 39.0 | – |
| Faster R-CNN w/ PDA | 37.8 | 58.4 |
| Faster R-CNN w/ FDA | 33.8 | 55.3 |
| Faster R-CNN w/ (PDA+FDA) | **39.6** | **59.3** |

### 5.1. Synthetic to real

In this scenario, object detector is trained on the synthetic dataset and evaluated on the real datasets.

*Sim10k → Cityscapes.* To study the efficacy of the proposed method, we firstly consider to remove the domain bias between the computer synthetic dataset of Sim10k and the real dataset of Cityscapes. Table 1 shows our results with different adaptation modules. Our baseline result with Faster R-CNN is trained purely on the source Sim10k and evaluated on Cityscapes directly. The calculated mAP is 30.1% with the VGG16 backbone. Compared with the baseline and current SOTA, the proposed method obtains 9.5% and 0.6% performance gains, respectively. Specifically, the proposed feature-level based adversarial training can improve the performance of baseline to 33.8%. The PDA module can bring 7.78% gains. When an end-to-end training is conducted, we can obtain better results than current SOTA.

*Sim10k → KITTI.* To further check the proposed method's robustness to the scenario of "synthetic to real", experiment is conducted by using the KITTI dataset as the target dataset. Since there are no other works reporting the detection results under this specific setting, we only compare our results with the baseline ones. All the results are evaluated on the KITTI training split just like [39]. As shown in Table 1, the baseline Faster R-CNN network has an AP of 52.7%. It can then be improved to 55.3% and 58.4% by using the proposed FDA and PDA modules separately. Through conducting the end-to-end training, we can get the highest result of 59.3%.

Table 2: Detection results are evaluated for the adaptations of Cityscapes → Foggy-Cityscapes and KITTI → VKITTI-Rainy.

|  | Cityscapes → Foggy-Cityscapes | KITTI → VKITTI-Rainy |
| --- | --- | --- |
| Faster R-CNN | 18.8 | 40.0 |
| SOTA[39] | 27.6 | – |
| Faster R-CNN w/ PDA | 27.1 | 51.3 |
| Faster R-CNN w/ FDA | 23.6 | 44.5 |
| Faster R-CNN w/ (PDA+FDA) | **28.9** | **52.2** |

5

## 5.2. Cross different weathers

In addition to the "synthetic to real" scenario, domain bias caused by different weathers are also considered in the experiments. Specifically, the proposed method is testified to deal with the foggy and rainy weathers in the driving scenes.

*Cityscapes → Foggy-Cityscapes.* Cityscapes and its foggy version are used as the source and target datasets in this experiment. All the training settings are the same with Section 5.1. The main target here is to adapt the detector to fit to the foggy weather by utilizing the labeled images collected in the clear weather. Experimental results are shown in Table 2. We can finally achieve 10.1% and 1.8% performance gains compared to the baseline result and current SOTA. Both of the two modules can largely improve the detection performance under the foggy circumstance. Specifically, baseline results of Faster R-CNN can be improved from 18.8% to 23.6% with FDA and 27.1% with PDA.

*KITTI → VKITTI-Rainy.* To certify the method's robustness to the rainy weather, VKITTI-Rainy is used as the target dataset, and KITTI as the source dataset. Because previous SOTA work does not report their results under this circumstance, experimental results are only compared with the baseline Faster R-CNN. Consistent improvements are achieved in this rainy condition. Detection results can be finally improved from 40.0% to 52.2% with the end-to-end training. It also can be seen that PDA module performs much better than the FDA module in this case.

Table 3: Detection results (AP (%) of "car") are evaluated for the adaptations of Cityscapes → KITTI and KITTI → Cityscapes.

|  | Cityscapes → KITTI | KITTI → Cityscapes |
|---|---|---|
| Faster R-CNN | 53.5 | 30.2 |
| SOTA[39] | 64.1 | 38.5 |
| Faster R-CNN w/ PDA | 64.4 | 41.1 |
| Faster R-CNN w/ FDA | 58.6 | 34.5 |
| Faster R-CNN w/ (PDA+FDA) | **65.6** | **41.8** |

## 5.3. Cross different cameras

To fully compare with current SOTA work, experiments are further conducted to deal with the domain bias caused by using different cameras. Two real datasets KITTI and Cityscapes are selected as the source and target datasets alternatively. The results are shown in Table 3. With Cityscapes as the source dataset, baseline results of Faster R-CNN can be improved from 53.5% to 64.4% (with PDA) and 58.6% (with FDA). When KITTI is used as the source dataset, the detection results can be improved from 30.2% to 41.1% (with PDA) and 34.5% (with FDA). The proposed method achieves better performance than current SOTA in both adaptation settings, with improvements of 1.5% and 2.3%, respectively.

Table 4: Quantitative ablation mAP results (%) of the different domain adaptation scenarios.

|  | Faster R-CNN | Faster R-CNN w/ PDA | Faster R-CNN w/ FDA | Faster R-CNN w/ PDA+FDA |
|---|---|---|---|---|
| Sim10k → Cityscapes | 30.1 | 37.8 (+7.7) | 33.8 (+3.7) | 39.6 (+9.5) |
| Sim10k → KITTI | 52.7 | 58.4 (+5.7) | 55.3 (+2.6) | 59.3 (+6.6) |
| Cityscapes → Foggy-Cityscapes | 18.8 | 27.1 (+8.3) | 23.6 (+4.8) | 28.9 (+10.1) |
| Cityscapes → KITTI | 53.5 | 64.4 (+10.9) | 58.6 (+5.1) | 65.6 (+12.1) |
| KITTI → Cityscapes | 30.2 | 41.1 (+10.9) | 34.5 (+4.3) | 41.8 (+11.6) |
| KITTI → VKITTI-Rainy | 40.0 | 51.3 (+11.3) | 44.5 (+4.5) | 52.2 (+12.2) |

## 5.4. Analysis and discussion

*PDA versus FDA.* Two modules (PDA and FDA) are proposed to reduce the domain bias between the source and target datasets. A final end-to-end training of the whole network is also conducted to pursue better performance of cross-domain object detection. The quantitative ablation results are further summarized in Table 4 to compare the effectiveness of the different modules.
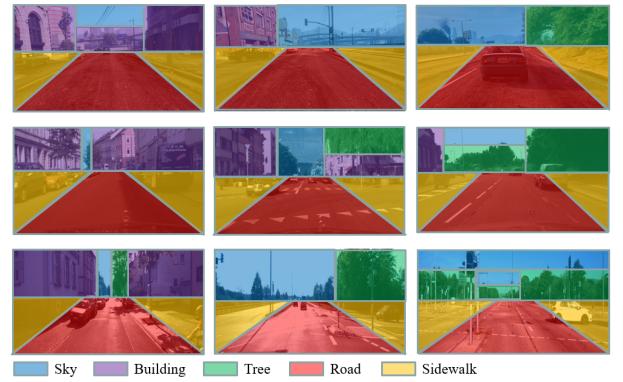


Figure 4: Overall structure layouts of the different driving scenes. From up to bottom, we show the images from Sim10k, Cityscapes and KITTI, respectively. Different driving scenes share the similar layouts along the "road" (*i.e.* the red and yellow areas.)

From Table 4, it can be seen that both of the two modules work effectively to improve the detector's performance under different domain bias. Image pixel-level transformation performs much better than the feature-level based adversarial training, with an average two times of improvements. Since the PDA module is based on CycleGAN, semantics of the image pixels are not guaranteed to be consistent after the translation process. However, the experimental results show that the image translation still works effectively to improve the detector's performance on the target datasets. Two reasons can be explained to this phenomenon. The first one is that the translated images, despite not perfect, can help the detector to learn generic target-oriented CNN features. The second one is that traffic

participants are the main concerns of the object detection task. Semantics of these objects can be mostly kept with PDA since the structure layouts along the "road" are similar among different driving scenes as shown by Fig. 4. Despite the upper parts of the images (*i.e.* the purple, blue and green areas) vary across different scenes, the lower parts (*i.e.* the red and yellow areas) are quite similar. In these similar areas, CycleGAN based PDA would mainly focus on the translation of color and textures.
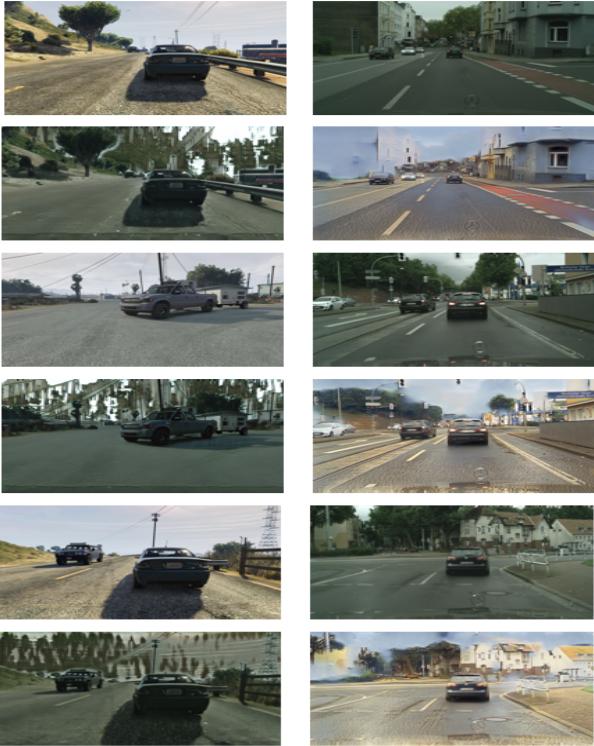


Figure 5: Qualitative image translation results between Sim10k and Cityscapes. On the left, we show the translations from Sim10k to Cityscapes. Translation results from Cityscapes to Sim10k are shown on the right.

Fig. 5 shows more image translation results between Sim10k and Cityscapes to illustrate this phenomenon. Generally, the synthetic dataset of Sim10k owns quite different structure layouts with the real Cityscapes dataset. Most of the images in Sim10k are taken from the scene of high-way, while Cityscapes are mainly taken from the cities. Therefore, it can be seen from Fig. 5 that many image pixels of Sim10k are mapped into "tree" or "building" to fit to the structure layout of Cityscapes. Similarly, pixels of "tree" and "building" are largely mapped into "sky" in the translation from Cityscapes to Sim10k. Nevertheless, objects along the "road" can still be maintained since similar layouts are shared between the two datasets. The translations are mainly focused on color and textures as shown in Fig. 5.

*With/without the end-to-end training.* During the experiments, an end-to-end training of the whole network is further conducted for another 10 epochs. In this training stage, more at-



(a) Original image from Sim10k

(b) Transformed image with PDA

(c) Transformed image with end-to-end training

(d) Cropped images with PDA

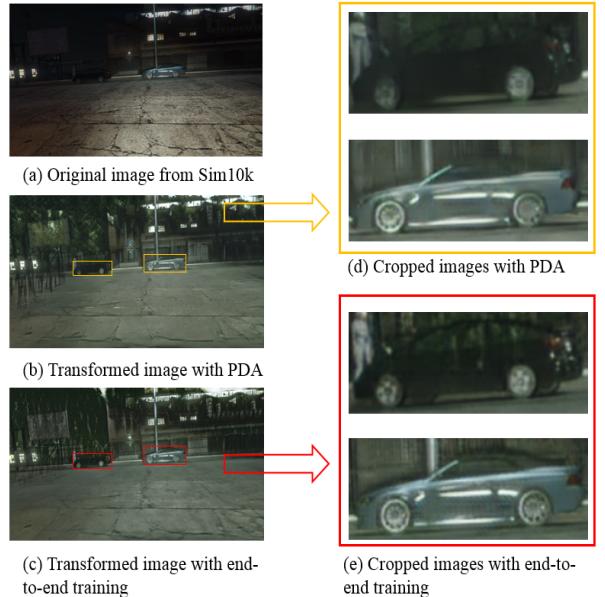(e) Cropped images with end-to-end training

Figure 6: Image translation results are shown with PDA (the yellow rectangular box) and the end-to-end training (the red rectangular box), respectively.

tentions are paid to the target objects. PDA module would also attempt to generate more realistic objects to reduce the adversarial training loss of FDA. Quantitatively, by conducting the end-to-end training, more improvements around 1% ~ 2% can be further achieved for the detection task. Qualitative results are also shown in Fig. 6 for better comparison of the translated images. Results with only PDA are shown in the yellow rectangular boxes and the end-to-end ones in the red rectangular boxes. It can be seen that better results are generated after the end-to-end training. When the target objects are close to the background or out of the share structured layouts, one potential problem is that the objects may get merged with the background or falsely translated into other semantics. With the end-to-end training, more details of the target objects can be maintained, such as the roof of the car.

*Analysis of the qualitative results.* Fig. 7 shows the qualitative image translation results under different scenarios. Row 1 and 2 show the translations between synthetic images of Sim10k and realistic ones of Cityscapes and KITTI. In this "synthetic to real" scenario, structures of the images are quite different. However, semantics of the image contents along the "road" can still be maintained. Row 3 and 4 show the translation from Cityscapes to KITTI and vice versa. Since both of these two datasets are collected from the cities of Germany. Most of the images share the similar structure layouts. Good results can be generated in these two experimental settings. Row 5 and 6 show the results of translating Cityscapes to its foggy version and KITTI to VKITTI-Rainy. In these experiments, robustness of the proposed method is verified to deal with different weathers (*e.g.* foggy or rainy) of the city. Since the structure of a city almost remains the same, domain bias brought by different

7

Figure 7: Qualitative results of the final translated images. From up to bottom, we show the image translations of Sim10k → Cityscapes, Cityscapes → Foggy-Cityscapes, Sim10k → KITTI, Cityscapes → KITTI, KITTI → Cityscapes and KITTI → VKITTI-Rainy. From left to right, we show the source and translated target domain images, alternatively.

weathers mainly comes from the color and texture. In this case, quite good results can be generated by the image translation module.

*Analysis of the inference performance.* The proposed method aims to adapt the source images into the style of the target domain. Then the translated images can be used to train the detector as augmented data. During the inference, the image translation module can be simply removed. Therefore, the proposed method would remain the same inference time with Faster R-CNN. Similar to [39], short side of the input image is scaled to 500 pixels. Using NVIDIA Tesla M40 GPU, the corresponding inference time with VGG16 is 165 ms.

## 6. Conclusions and future work

A new unsupervised domain adaptation method is proposed in this paper to solve the object detection problem in the field of autonomous driving. Extensive experiments are implemented to certify the efficacy of the proposed method. We can achieve better performance than current SOTA work through conducting adaptations both in image pixel and feature spaces. In the future, we will try to solve more complex cross-domain vision tasks such as instance segmentation or depth estimation based on the proposed method.

## References

[1] A. Krizhevsky, I. Sutskever, G. E. Hinton, Imagenet classification with deep convolutional neural networks, in: Advances in neural information processing systems, 2012, pp. 1097–1105.

[2] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, B. Schiele, The cityscapes dataset for semantic urban scene understanding, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2016, pp. 3213–3223.

[3] S. Ben-David, J. Blitzer, K. Crammer, F. Pereira, Analysis of representations for domain adaptation, in: Advances in neural information processing systems, 2007, pp. 137–144.

[4] S. Ben-David, J. Blitzer, K. Crammer, A. Kulesza, F. Pereira, J. W. Vaughan, A theory of learning from different domains, Machine learning 79 (1-2) (2010) 151–175.

[5] R. Gopalan, R. Li, R. Chellappa, Domain adaptation for object recognition: An unsupervised approach, in: Computer Vision (ICCV), 2011 IEEE International Conference on, IEEE, 2011, pp. 999–1006.

[6] W. Jiang, H. Gao, W. Lu, W. Liu, F.-L. Chung, H. Huang, Stacked robust adaptively regularized auto-regressions for domain adaptation, IEEE Transactions on Knowledge & Data Engineering (3) (2019) 561–574.

[7] W. Jiang, W. Liu, F.-l. Chung, Knowledge transfer for spectral clustering, Pattern Recognition 81 (2018) 484–496.

[8] G.-J. Qi, W. Liu, A. C, H. T, Joint intermodal and intramodal label transfers for extremely rare or unseen classes., IEEE transactions on pattern analysis and machine intelligence 39 (7) (2017) 1360–1373.

[9] K. Bousmalis, N. Silberman, D. Dohan, D. Erhan, D. Krishnan, Unsupervised pixel-level domain adaptation with generative adversarial networks, in: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR), Vol. 1, 2017, p. 7.

[10] J. Hoffman, E. Tzeng, T. Park, J.-Y. Zhu, P. Isola, K. Saenko, A. Efros, T. Darrell, Cycada: Cycle-consistent adversarial domain adaptation, in: International Conference on Machine Learning, 2018, pp. 1994–2003.

[11] J.-Y. Zhu, T. Park, P. Isola, A. A. Efros, Unpaired image-to-image translation using cycle-consistent adversarial networks, in: Proceedings of the IEEE international conference on computer vision, 2017, pp. 2223–2232.

[12] R. Girshick, J. Donahue, T. Darrell, J. Malik, Rich feature hierarchies for accurate object detection and semantic segmentation, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2014, pp. 580–587.

[13] R. Girshick, Fast r-cnn, in: Proceedings of the IEEE international conference on computer vision, 2015, pp. 1440–1448.

[14] K. Simonyan, A. Zisserman, Very deep convolutional networks for large-scale image recognition, arXiv preprint arXiv:1409.1556.

[15] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, A. Rabinovich, et al., Going deeper with convolutions, Cvpr, 2015.

[16] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recognition, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2016, pp. 770–778.

[17] G. Huang, Z. Liu, K. Q. Weinberger, L. van der Maaten, Densely connected convolutional networks, in: Proceedings of the IEEE conference on computer vision and pattern recognition, Vol. 1, 2017, p. 3.

[18] J. Redmon, S. Divvala, R. Girshick, A. Farhadi, You only look once: Unified, real-time object detection, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2016, pp. 779–788.

[19] J. Redmon, A. Farhadi, Yolov3: An incremental improvement, arXiv preprint arXiv:1804.02767.

[20] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, A. C. Berg, Ssd: Single shot multibox detector, in: European conference on computer vision, Springer, 2016, pp. 21–37.

[21] C.-Y. Fu, W. Liu, A. Ranga, A. Tyagi, A. C. Berg, Dssd: Deconvolutional single shot detector, arXiv preprint arXiv:1701.06659.

[22] S. Ren, K. He, R. Girshick, J. Sun, Faster r-cnn: Towards real-time object detection with region proposal networks, in: Advances in neural information processing systems, 2015, pp. 91–99.

[23] J. Dai, Y. Li, K. He, J. Sun, R-fcn: Object detection via region-based fully convolutional networks, in: Advances in neural information processing systems, 2016, pp. 379–387.

[24] Z. Jie, Y. Wei, X. Jin, J. Feng, W. Liu, Deep self-taught learning for weakly supervised object localization, in: Proceedings of the IEEE Con-

ference on Computer Vision and Pattern Recognition, 2017, pp. 1377–1385.

[25] P. Tang, X. Wang, A. Wang, Y. Yan, W. Liu, J. Huang, A. Yuille, Weakly supervised region proposal network and object detection, in: Proceedings of the European conference on computer vision (ECCV), 2018, pp. 352–368.

[26] B. Kang, Z. Liu, X. Wang, F. Yu, J. Feng, T. Darrell, Few-shot object detection via feature reweighting, arXiv preprint arXiv:1812.01866.

[27] S. J. Pan, I. W. Tsang, J. T. Kwok, Q. Yang, Domain adaptation via transfer component analysis, IEEE Transactions on Neural Networks 22 (2) (2011) 199–210.

[28] K. M. Borgwardt, A. Gretton, M. J. Rasch, H.-P. Kriegel, B. Schölkopf, A. J. Smola, Integrating structured biological data by kernel maximum mean discrepancy, Bioinformatics 22 (14) (2006) e49–e57.

[29] M. Long, J. Wang, G. Ding, J. Sun, P. S. Yu, Transfer feature learning with joint distribution adaptation, in: Proceedings of the IEEE international conference on computer vision, 2013, pp. 2200–2207.

[30] W. Jiang, C. Deng, W. Liu, F. Nie, F.-l. Chung, H. Huang, Theoretic analysis and extremely easy algorithms for domain adaptive feature learning, in: The 26th International Joint Conference on Artificial Intelligence (IJCAI 2017), 2017.

[31] M. Long, Y. Cao, J. Wang, M. Jordan, Learning transferable features with deep adaptation networks, in: International Conference on Machine Learning, 2015, pp. 97–105.

[32] M. Long, H. Zhu, J. Wang, M. I. Jordan, Unsupervised domain adaptation with residual transfer networks, in: Advances in Neural Information Processing Systems, 2016, pp. 136–144.

[33] B. Sun, J. Feng, K. Saenko, Return of frustratingly easy domain adaptation., in: AAAI, Vol. 6, 2016, p. 8.

[34] Y. Ganin, E. Ustinova, H. Ajakan, P. Germain, H. Larochelle, F. Laviolette, M. Marchand, V. Lempitsky, Domain-adversarial training of neural networks, The Journal of Machine Learning Research 17 (1) (2016) 2096–2030.

[35] Y. Ganin, V. Lempitsky, Unsupervised domain adaptation by backpropagation, in: International Conference on Machine Learning, 2015, pp. 1180–1189.

[36] E. Tzeng, J. Hoffman, K. Saenko, T. Darrell, Adversarial discriminative domain adaptation, in: Computer Vision and Pattern Recognition (CVPR), Vol. 1, 2017, p. 4.

[37] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, Y. Bengio, Generative adversarial nets, in: Advances in neural information processing systems, 2014, pp. 2672–2680.

[38] K. Bousmalis, G. Trigeorgis, N. Silberman, D. Krishnan, D. Erhan, Domain separation networks, in: Advances in Neural Information Processing Systems, 2016, pp. 343–351.

[39] Y. Chen, W. Li, C. Sakaridis, D. Dai, L. Van Gool, Domain adaptive faster r-cnn for object detection in the wild, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2018, pp. 3339–3348.

[40] N. Srivastava, R. R. Salakhutdinov, Multimodal learning with deep boltzmann machines, in: Advances in neural information processing systems, 2012, pp. 2222–2230.

[41] J. Ngiam, A. Khosla, M. Kim, J. Nam, H. Lee, A. Y. Ng, Multimodal deep learning, in: Proceedings of the 28th international conference on machine learning (ICML-11), 2011, pp. 689–696.

[42] J. Yang, J. Wright, T. S. Huang, Y. Ma, Image super-resolution via sparse representation, IEEE transactions on image processing 19 (11) (2010) 2861–2873.

[43] P. Isola, J.-Y. Zhu, T. Zhou, A. A. Efros, Image-to-image translation with conditional adversarial networks, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2017, pp. 1125–1134.

[44] J.-Y. Zhu, R. Zhang, D. Pathak, T. Darrell, A. A. Efros, O. Wang, E. Shechtman, Toward multimodal image-to-image translation, in: Advances in Neural Information Processing Systems, 2017, pp. 465–476.

[45] M.-Y. Liu, O. Tuzel, Coupled generative adversarial networks, in: Advances in neural information processing systems, 2016, pp. 469–477.

[46] Y. Aytar, L. Castrejon, C. Vondrick, H. Pirsiavash, A. Torralba, Crossmodal scene networks, IEEE transactions on pattern analysis and machine intelligence 40 (10) (2017) 2303–2314.

[47] J. Johnson, A. Alahi, L. Fei-Fei, Perceptual losses for real-time style transfer and super-resolution, in: European Conference on Computer Vision, Springer, 2016, pp. 694–711.

[48] H. Zhang, K. Dana, Multi-style generative network for real-time transfer, arXiv preprint arXiv:1703.06953.

[49] L. A. Gatys, A. S. Ecker, M. Bethge, Image style transfer using convolutional neural networks, in: Computer Vision and Pattern Recognition (CVPR), 2016 IEEE Conference on, IEEE, 2016, pp. 2414–2423.

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

[51] C. Sakaridis, D. Dai, L. Van Gool, Semantic foggy scene understanding with synthetic data, International Journal of Computer Vision 126 (9) (2018) 973–992.

[52] A. Gaidon, Q. Wang, Y. Cabon, E. Vig, Virtual worlds as proxy for multiobject tracking analysis, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2016, pp. 4340–4349.

[53] M. Johnson-Roberson, C. Barto, R. Mehta, S. N. Sridhar, K. Rosaen, R. Vasudevan, Driving in the matrix: Can virtual worlds replace humangenerated annotations for real world tasks?, in: Robotics and Automation (ICRA), 2017 IEEE International Conference on, IEEE, 2017, pp. 746–753.

[54] O. Ronneberger, P. Fischer, T. Brox, U-net: Convolutional networks for biomedical image segmentation, in: International Conference on Medical image computing and computer-assisted intervention, Springer, 2015, pp. 234–241.