

# Dark Model Adaptation: Semantic Image Segmentation from Daytime to Nighttime

Dengxin Dai<sup>1</sup> and Luc Van Gool<sup>1,2</sup>

**Abstract**—This work addresses the problem of semantic image segmentation of nighttime scenes. Although considerable progress has been made in semantic image segmentation, it is mainly related to daytime scenarios. This paper proposes a novel method to *progressive* adapt the semantic models trained on daytime scenes, along with large-scale annotations therein, to nighttime scenes via the bridge of twilight time — the time between dawn and sunrise, or between sunset and dusk. The goal of the method is to alleviate the cost of human annotation for nighttime images by transferring knowledge from standard daytime conditions. In addition to the method, a new dataset of road scenes is compiled; it consists of 35,000 images ranging from daytime to twilight time and to nighttime. Also, a subset of the nighttime images are densely annotated for method evaluation. Our experiments show that our method is effective for knowledge transfer from daytime scenes to nighttime scenes, without using extra human annotation.

## I. INTRODUCTION

Autonomous vehicles will have a substantial impact on people’s daily life, both personally and professionally. For instance, automated vehicles can largely increase human productivity by turning driving time into working time, provide personalized mobility to non-drivers, reduce traffic accidents, or free up parking space and generalize valet service [3]. As such, developing automated vehicles is becoming the core interest of many, diverse industrial players. Recent years have witnessed great progress in autonomous driving [14], resulting in announcements that autonomous vehicles have driven over many thousands of miles and that companies aspire to sell such vehicles in a few years. All this has fueled expectations that fully automated vehicles are coming soon. Yet, significant technical obstacles must be overcome before assisted driving can be turned into full-fledged automated driving, a prerequisite for the above visions to materialize.

While perception algorithms based on visible light cameras are constantly getting better, they are mainly designed to operate on images taken at daytime under good illumination [23], [27]. Yet, outdoor applications can hardly escape from challenging weather and illumination conditions. One of the big reasons that automated cars have not gone mainstream yet is because it cannot deal well with nighttime and adverse weather conditions. Camera sensors can become untrustworthy at nighttime, in foggy weather, and in wet weather. Thus, computer vision systems have to function well

also under these adverse conditions. In this work, we focus on semantic object recognition for nighttime driving scenes.

Robust object recognition using visible light cameras remains a difficult problem. This is because the structural, textural and/or color features needed for object recognition sometimes do not exist or highly disburbed by artificial lights, to the point where it is difficult to recognize the objects even for human. The problem is further compounded by camera noise [32] and motion blur. Due to this reason, there are systems using far-infrared (FIR) cameras instead of the widely used visible light cameras for nighttime scene understanding [31], [11]. Far-infrared (FIR) cameras can be another choice [31], [11]. They, however, are expensive and only provide images of relatively low-resolution. Thus, this work adopts visible light cameras for semantic segmentation of nighttime road scenes. Another reason of this choice is that large-scale datasets are available for daytime images by visible light cameras [8]. This makes model adaptation from daytime to nighttime feasible.

High-level semantic tasks is usually tackled by learning from many annotations of real images. This scheme has achieved a great success for good weather conditions at daytime. Yet, the difficulty of collecting and annotating images for all other weather and illumination conditions renders this standard protocol problematic. To overcome this problem, we depart from this traditional paradigm and propose another route. Instead, we choose to *progressively* adapt the semantic models trained for daytime scenes to nighttime scenes, by using images taken at the twilight time as intermediate stages. The method is based on progressively self-learning scheme, and its pipeline is shown in Figure 1.

The main contributions of the paper are: 1) a novel model adaptation method is developed to transfer semantic knowledge from daytime scenes to nighttime scenes; 2) a new dataset, named *Nighttime Driving*, consisting of images of real driving scenes at nighttime and twilight time, with 35,000 unlabeled images and 50 densely annotated images. These contributions will facilitate the learning and evaluation of semantic segmentation methods for nighttime driving scenes. *Nighttime Driving* is available at <http://people.ee.ethz.ch/~daid/NightDriving/>.

## II. RELATED WORK

### A. Semantic Understanding of Nighttime Scenes

A lot of work for nighttime object detection/recognition has focused on human detection, by using FIR cameras [31], [10] or visible light cameras [18], or a combination of both [5]. There are also notable examples for detecting other

<sup>1</sup>Dengxin Dai and Luc Van Gool are with the Toyota TRACE-Zurich team at the Computer Vision Lab, ETH Zurich, 8092 Zurich, Switzerland [firstname.lastname@vision.ee.ethz.ch](mailto:firstname.lastname@vision.ee.ethz.ch)

<sup>2</sup>Luc Van Gool is also with the Toyota TRACE-Leuven team at the Dept of Electrical Engineering ESAT, KU Leuven 3001 Leuven, Belgium

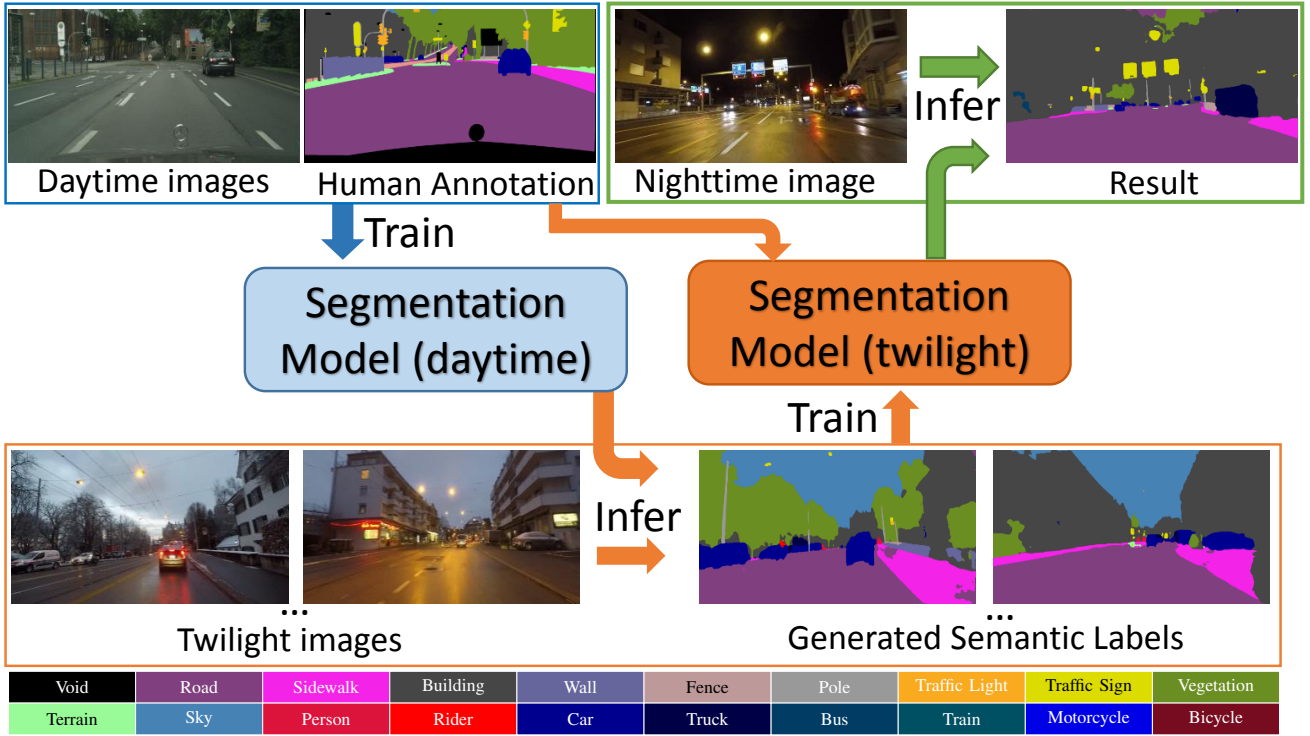


Fig. 1. The pipeline of our approach for semantic segmentation of nighttime scenes, by transferring knowledge from daytime scenes via the bridge of twilight time scenes.

road traffic objects such as cars [19] and their rear lights [29]. Another group of work is to develop methods robust to illumination changes for robust road area detection [2] and semantic labeling [25]. Most of the research in this vein had been conducted before deep learning was widely used.

Semantic understanding of visual scenes have recently undergone rapid growth, making accurate object detection feasible in images and videos in daytime scenes [6], [21]. It is natural to raise the question of how to extend those sophisticated methods to other weather conditions and illumination conditions, and examine and improve the performance therein. A recent effort has been made for foggy weather [27]. This work would like to initiate the same research effort for nighttime.

### B. Model Adaptation

Our work bears resemblance to works from the broad field of transfer learning. Model adaptation across weather conditions to semantically segment simple road scenes is studied in [20]. More recently, domain adaptation based approach was proposed to adapt semantic segmentation models from synthetic images to real environments [33], [16], [28], [27], [7]. The supervision transfer from daytime scenes to nighttime scenes in this paper is inspired by the stream of work on model distillation/imitation [15], [12], [9]. Our approach is similar in that knowledge is transferred from one domain to the next by distilled from the previous domain. The concurrent work in [26] on adaptation of semantic models from clear weather condition to light fog then to dense fog is closely related to ours.

### C. Road Scene Understanding

Road scene understanding is a crucial enabler for applications such as assisted or autonomous driving. Typical examples include the detection of roads [4], traffic lights [17], cars and pedestrians [8], [27], and tracking of such objects [30], [22]. We refer the reader to the excellent surveys [24]. The aim of this work is to extend/adapt the advanced models developed recently for road scene understanding at daytime to nighttime, without manually annotating nighttime images.

## III. APPROACH

Training a segmentation model with large amount of human annotations should work for nighttime images, similar to what has been achieved for daytime scene understanding [13], [21]. However, applying this protocol to other weather conditions and illumination conditions is problematic as it is hardly affordable to annotate the same amount of data for all different conditions and their combinations. We depart from this protocol and investigate an automated approach to transfer the knowledge from existing annotations of daytime scenes to nighttime scenes. The approach leverages the fact that illumination changes continuously between daytime and nighttime, through the twilight time. Twilight is the time between dawn and sunrise, or between sunset and dusk. Twilight is defined according to the solar elevation angle, which is the position of the geometric center of the sun relative to the horizon [1]. See Figure 2 for an illustration.

During a large portion of twilight time, solar illumination suffices enough for cameras to capture the terrestrial objects and suffices enough to alleviate the interference of artificial

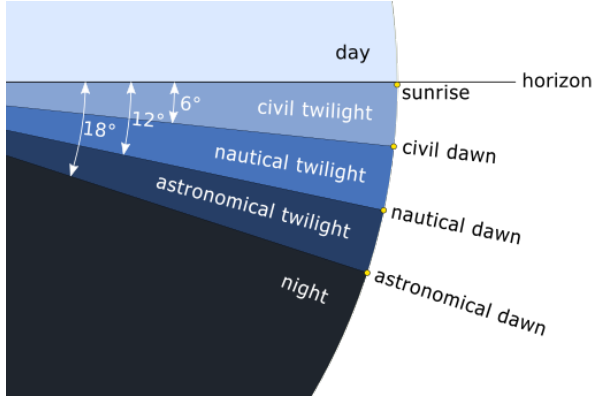


Fig. 2. Twilight is defined according to the solar elevation angle and is categorized into three subcategories: civil twilight, nautical twilight, and astronomical twilight. (picture is from wikipedia).

lights to a limited amount. See Figure 1 for examples of road scenes at twilight time. These observations lead to our conjecture that the domain discrepancy between daytime scenes and twilight scenes, and the domain discrepancy between twilight scenes and nighttime scenes are both smaller than the domain discrepancy between daytime scenes and nighttime scenes. Thus, images captured during twilight time can serve our purpose well — transfer knowledge from daytime to nighttime. That is, twilight time constructs a bridge for knowledge transfer from our source domain daytime to our target domain nighttime.

In particular, we train a semantic segmentation model on daytime images using the standard supervised learning paradigm, and apply the model to a large dataset recorded at civil twilight time to generate the class responses. The three subgroups of twilight are used: civil twilight, nautical twilight, and astronomical twilight [1]. Since the domain gap between daytime condition and civil twilight condition is relatively small, these class responses, along with the images, can then be used to fine-tune the semantic segmentation model so that it can adapt to civil twilight time. The same procedure is continued through nautical twilight and astronomical twilight. We then apply the final fine-tuned model to nighttime images.

This learning approach is inspired by the stream of work on model distillation [15], [9], [12]. Those methods either transfer supervision from sophisticated models to simpler models for efficiency [15], [9], or transfer supervision from the domain of images to other domains such as depth maps [12]. We here transfer the semantic knowledge of annotations of daytime scenes to nighttime scenes via the unlabeled images recorded at twilight time.

Let us denote an image by  $\mathbf{x}$ , and indicate the image taken at *daytime*, *civil twilight time*, *nautical twilight time*, *astronomical twilight time* and *nighttime* by  $\mathbf{x}^0$ ,  $\mathbf{x}^1$ ,  $\mathbf{x}^2$ ,  $\mathbf{x}^3$ , and  $\mathbf{x}^4$ , respectively. The corresponding human annotation for  $\mathbf{x}^0$  is provided and denoted by  $\mathbf{y}^0$ , where  $\mathbf{y}^0(m, n) \in \{1, \dots, C\}$  is the label of pixel  $(m, n)$ , and  $C$  is the total number of classes. Then, the training data consist of labeled

data at daytime  $\mathcal{D}^0 = \{(\mathbf{x}_i^0, \mathbf{y}_i^0)\}_{i=1}^{l^0}$ , and three unlabeled datasets for the three twilight categories:  $\mathcal{D}^1 = \{\mathbf{x}_j^1\}_{j=1}^{l^1}$ ,  $\mathcal{D}^2 = \{\mathbf{x}_k^2\}_{k=1}^{l^2}$ , and  $\mathcal{D}^3 = \{\mathbf{x}_q^3\}_{q=1}^{l^3}$ , where  $l^0$ ,  $l^1$ ,  $l^2$ , and  $l^3$  are the total number of images in the corresponding datasets.

The method consists of eight steps and it is summarized below.

- 1: train a segmentation model with daytime images and the human annotations:

$$\min_{\phi^0} \frac{1}{l^0} \sum_{i=1}^{l^0} L(\phi^0(\mathbf{x}_i^0), \mathbf{y}_i^0), \quad (1)$$

where  $L(\cdot, \cdot)$  is the cross entropy loss function;

- 2: apply segmentation model  $\phi^0$  to the images recorded at civil twilight time to obtain “noisy” semantic labels:  $\hat{\mathbf{y}}^1 = \phi^0(\mathbf{x}^1)$ , and augment dataset  $\mathcal{D}^1$  to  $\hat{\mathcal{D}}^1$ :  $\hat{\mathcal{D}}^1 = \{(\mathbf{x}_j^1, \hat{\mathbf{y}}_j^1)\}_{j=1}^{l^1}$ ;
- 3: instantiate a new model  $\phi^1$  by duplicating  $\phi^0$ , and then fine-tune (retrain) the semantic model on  $\mathcal{D}^0$  and  $\hat{\mathcal{D}}^1$ :

$$\phi^1 \leftarrow \phi^0, \quad (2)$$

and

$$\min_{\phi^1} \left( \frac{1}{l^0} \sum_{i=1}^{l^0} L(\phi^1(\mathbf{x}_i^0), \mathbf{y}_i^0) + \frac{\lambda^1}{l^1} \sum_{j=1}^{l^1} L(\phi^1(\mathbf{x}_j^1), \hat{\mathbf{y}}_j^1) \right), \quad (3)$$

where  $\lambda^1$  is a hyper-parameter balancing the weights of the two data sources;

- 4: apply segmentation model  $\phi^1$  to the images recorded at nautical twilight time to obtain “noisy” semantic labels:  $\hat{\mathbf{y}}^2 = \phi^1(\mathbf{x}^2)$ , and augment dataset  $\mathcal{D}^2$  to  $\hat{\mathcal{D}}^2$ :  $\hat{\mathcal{D}}^2 = \{(\mathbf{x}_k^2, \hat{\mathbf{y}}_k^2)\}_{k=1}^{l^2}$ ;
- 5: instantiate a new model  $\phi^2$  by duplicating  $\phi^1$ , and fine-tune (train) semantic model on  $\mathcal{D}^0$ ,  $\hat{\mathcal{D}}^1$  and  $\hat{\mathcal{D}}^2$ :

$$\phi^2 \leftarrow \phi^1, \quad (4)$$

and then

$$\min_{\phi^2} \left( \frac{1}{l^0} \sum_{i=1}^{l^0} L(\phi^2(\mathbf{x}_i^0), \mathbf{y}_i^0) + \frac{\lambda^1}{l^1} \sum_{j=1}^{l^1} L(\phi^2(\mathbf{x}_j^1), \hat{\mathbf{y}}_j^1) + \frac{\lambda^2}{l^2} \sum_{k=1}^{l^2} L(\phi^2(\mathbf{x}_k^2), \hat{\mathbf{y}}_k^2) \right), \quad (5)$$

where  $\lambda^1$  and  $\lambda^2$  are hyper-parameters regulating the weights of the datasets;

- 6: apply segmentation model  $\phi^2$  to the images recorded at astronomical twilight data to obtain “noisy” semantic labels:  $\hat{\mathbf{y}}^3 = \phi^2(\mathbf{x}^3)$ , and augment dataset  $\mathcal{D}^3$  to  $\hat{\mathcal{D}}^3$ :  $\hat{\mathcal{D}}^3 = \{(\mathbf{x}_q^3, \hat{\mathbf{y}}_q^3)\}_{q=1}^{l^3}$ ;
- 7: instantiate a new model  $\phi^3$  by duplicating  $\phi^2$ , and fine-tune (train) the semantic model on all four datasets  $\mathcal{D}^0$ ,  $\hat{\mathcal{D}}^1$ ,  $\hat{\mathcal{D}}^2$  and  $\hat{\mathcal{D}}^3$ :

$$\phi^3 \leftarrow \phi^2, \quad (6)$$

and then

$$\min_{\phi^3} \left( \frac{1}{l^0} \sum_{i=1}^{l^0} L(\phi^3(\mathbf{x}_i^0), \mathbf{y}_i^0) + \frac{\lambda^1}{l^1} \sum_{j=1}^{l^1} L(\phi^3(\mathbf{x}_j^1), \hat{\mathbf{y}}_j^1) \right. \\ \left. + \frac{\lambda^2}{l^2} \sum_{k=1}^{l^2} L(\phi^3(\mathbf{x}_k^2), \hat{\mathbf{y}}_k^2) + \frac{\lambda^3}{l^3} \sum_{q=1}^{l^3} L(\phi^3(\mathbf{x}_q^3), \hat{\mathbf{y}}_q^3) \right), \quad (7)$$

where  $\lambda^1$ ,  $\lambda^2$  and  $\lambda^3$  are hyper-parameters regulating the weights of the datasets;

- 8: apply model  $\phi^3$  to nighttime images to perform the segmentation:  $\hat{\mathbf{y}}^4 = \phi^3(\mathbf{x}^4)$ .

We term our method Gradual Model Adaptation. During training, in order to balance the weights of different data sources (in Equation 3, Equation 5 and Equation 7), we empirically give equal weight to all training datasets. An optimal value can be obtained via cross-validation. The optimization of Equation 3, Equation 5 and Equation 7 are implemented by feeding to the training algorithm a stream of hybrid data, for which images in the considered datasets are sampled proportionally according to the parameters  $\lambda^1$ ,  $\lambda^2$ , and  $\lambda^3$ . In this work, they all set to 1, which means all datasets are sampled at the same rate.

Rather than applying the model trained on daytime images directly to nighttime images, Gradual Model Adaptation breaks down the problem to three progressive steps to adapt the semantic model. In each of the step, the domain gap is much smaller than the domain gap between daytime domain and nighttime domain. Due to the unsupervised nature of this domain adaptation, the algorithm will also be affected by the noise in the labels. The daytime dataset  $\mathcal{D}^1$  is always used for the training, to balance between noisy data of similar domains and clean data of a distinct domain.

#### IV. EXPERIMENTS

##### A. Data Collection

*Nighttime Driving* was collected during 5 rides with a car inside multiple Swiss cities and their suburbs using a GoPro Hero 5 camera. We recorded 5 large video sequence with length of about 2 hours. The video recording starts from daytime, goes through twilight time and ends at full nighttime. The video frames are extracted at a rate of one frame per second, leading to 35,000 images in total. According to [1] and the sunset time of each recording day, we partition the dataset into five parts: daytime, civil twilight time, nautical twilight time, astronomical twilight time, and nighttime. They consist of 8000, 8750, 8750, 8750, and 9500 images, respectively.

We manually select 50 nighttime images of diverse visual scenes, and construct the test set of *Nighttime Driving* therefrom, which we term *Nighttime Driving-test*. The aforementioned selection is performed manually in order to guarantee that the test set has high diversity, which compensates for its relatively small size in terms of statistical significance of evaluation results. We annotate these images with fine pixel-level semantic annotations using the 19 evaluation classes

TABLE I  
PERFORMANCE COMPARISON BETWEEN THE VARIANTS OF OUR  
METHOD TO THE ORIGINAL SEGMENTATION MODEL.

Model	Fine-tuning on twilight data	Mean IoU
Refinenet [21]	—	35.2
Refinenet	$\phi^1$ ( $\rightarrow$ civil)	38.6
Refinenet	$\phi^2$ ( $\rightarrow$ civil $\rightarrow$ nautical)	39.9
Refinenet	$\phi^3$ ( $\rightarrow$ civil $\rightarrow$ nautical $\rightarrow$ astronomical)	<b>41.6</b>
Refinenet	$\rightarrow$ all twilight (1-step adaptation)	39.1

of the Cityscapes dataset [8]: *road, sidewalk, building, wall, fence, pole, traffic light, traffic sign, vegetation, terrain, sky, person, rider, car, truck, bus, train, motorcycle and bicycle*. In addition, we assign the *void* label to pixels which do not belong to any of the above 19 classes, or the class of which is uncertain due to insufficient illumination. Every such pixel is ignored for semantic segmentation evaluation.

##### B. Experimental Evaluation

Our model of choice for experiments on semantic segmentation is the RefineNet [21]. We use the publicly available *RefineNet-res101-Cityscapes* model, which has been trained on the daytime training set of Cityscapes. In all experiments of this section, we use a constant base learning rate of  $5 \times 10^{-5}$  and mini-batches of size 1.

Our segmentation experiment showcases the effectiveness of our model adaptation pipeline, using twilight time as a bridge. The models which are obtained after the initial adaptation step are further fine-tuned on the union of the daytime Cityscapes dataset and the previously segmented twilight datasets, where the latter sets are labeled by the adapted models one step ahead.

We evaluate four variants of our method and compare them to the original segmentation model trained on daytime images directly. Using the pipeline described in Section III, three models can be obtained, in particular  $\phi^1$ ,  $\phi^2$ , and  $\phi^3$ .

We also compare to an alternative adaptation approach which generates labels (by using the original model trained on daytime data) for all twilight images at once and fine-tunes the original daytime segmentation model once. To put in another word, the three-step progressive model adaptation is reduced to a one-step progressive model adaptation.

**Quantitative Results.** The overall intersection over union (IoU) over all classes of the semantic segmentation by all methods are reported in Tables I. The table shows that all variants of our adaptation method improve the performance of the original semantic model trained with daytime data. This is mainly due to the fact that twilight time fall into the middle ground of daytime and nighttime, so the domain gaps from twilight to the other two domains are smaller than the direct domain gap of the two.

Also, it can be seen from the table that our method benefits from the progressive adaptation in three steps, i.e. from daytime to civil twilight, from civil twilight to nautical twilight, and from nautical twilight to astronomical twilight. The complete pipeline outperforms the two incomplete alternatives. This means that the gradual adaptation closes

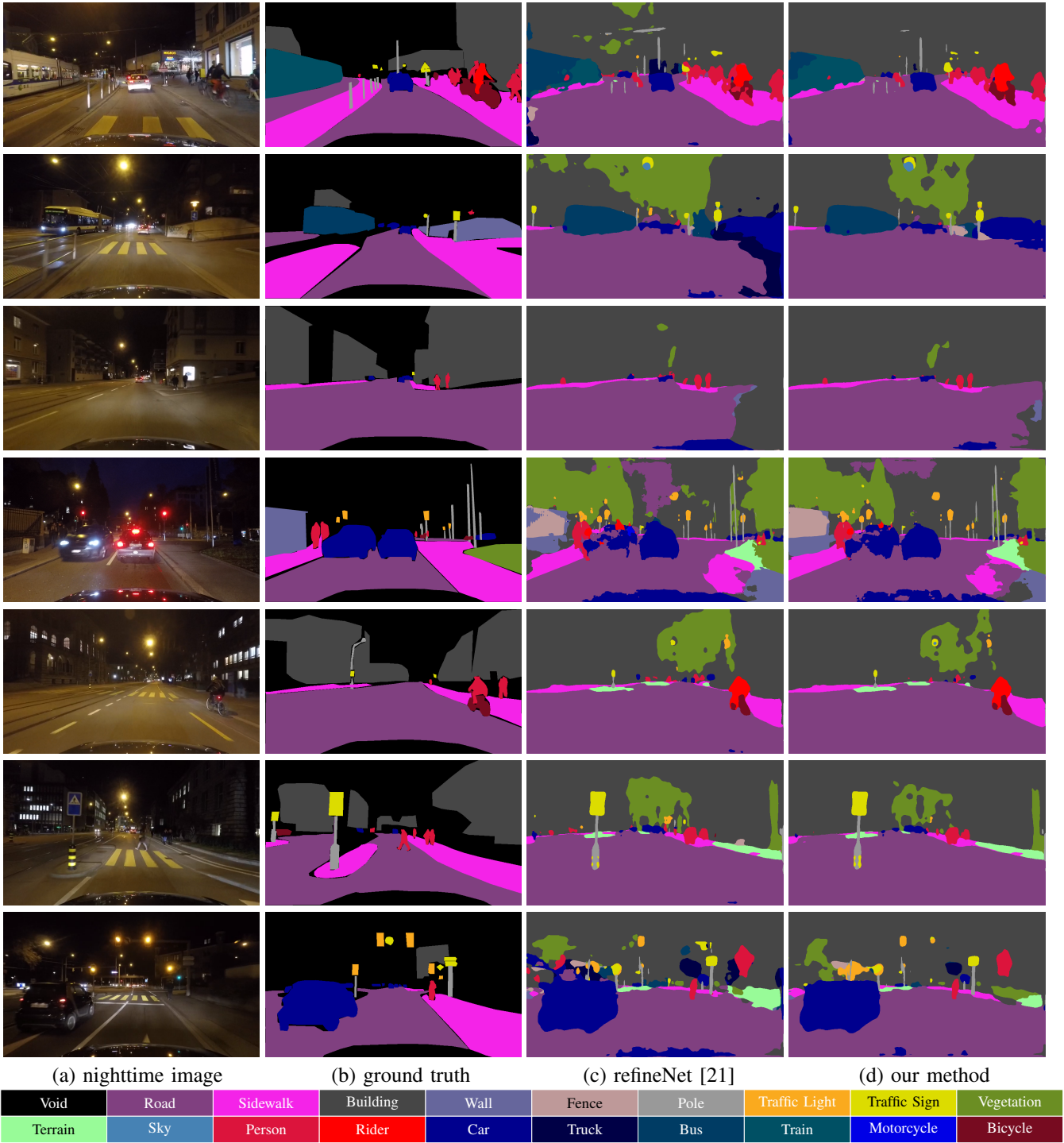


Fig. 3. Qualitative results for semantic segmentation on *Nighttime Driving-test*.

the domain gap progressively. As the model is adapted one more step forward, the gap to the target domain is further narrowed. Data recorded through twilight time constructs a trajectory between the source domain (daytime) and the target domain (nighttime) and makes daytime-to-nighttime knowledge transfer feasible.

Finally, we find that our three-step progressive pipeline outperforms the one-step progressive alternative. This is mainly due to the unsupervised nature of the model adap-

tation: the method learns from generated labels for model adaptation. This means that the accuracy of the generated labels directly affect the quality of the adaptation. The one-step adaptation alternative proceeds more aggressively and in the end learns from more noisy generated labels than our three-step complete pipeline. The three-step model adaptation method generate labels only on data which falls slightly off the training domain of the previous model. Our three-step model adaptation strikes a good balance between



computational cost and quality control.

**Qualitative Results.** We also show multiple segmentation examples by our method (the three-step complete pipeline) and the original daytime RefineNet model in Figure 3. From the two figures, one can see that our method generally yields better results than the original RefineNet model. For instance, in the second image of Figure 3, the original RefineNet model misclassified some *road* area as *car*.

While improvement has been observed, the performance of for nighttime scenes is still a lot worse than that for daytime scenes. Nighttime scenes are indeed more challenging than daytime scenes for semantic understanding tasks. There are more underlying causal factors of variation that generated night data, which requires either more training data or more intelligent learning approaches to disentangle the increased number of factors. Also, the models are adapted in an unsupervised manner. Introducing a reasonable amount of human annotations of nighttime scenes will for sure improve the results. This constitutes our future work.

**Limitation.** Many regions in nighttime images are uncertain for human annotators. Those areas should be treated as a separate, special class; algorithms need to be trained to predict this special class as well. It is misleading to assign a class label to those areas. This will be implemented in our next work. We also argue that street lamps should be considered as a separate class in addition to the classes considered in Cityscapes’ daytime driving.

## V. CONCLUSIONS

This work has investigated the problem of semantic image segmentation of nighttime scenes from a novel perspective. This paper has proposed a novel method to *progressive* adapts the semantic models trained on daytime scenes to nighttime scenes via the bridge of twilight time — the time between dawn and sunrise, or between sunset and dusk. Data recorded during twilight times are further grouped into three subgroups for a three-step progressive model adaptation, which is able to transfer knowledge from daytime to nighttime in an unsupervised manner. In addition to the method, a new dataset of road driving scenes is compiled. It consists of 35,000 images ranging from daytime to twilight time and to nighttime. Also, 50 diverse nighttime images are densely annotated for method evaluation. The experiments show that our method is effective for knowledge transfer from daytime scenes to nighttime scenes without using human supervision.

**Acknowledgement** This work is supported by Toyota Motor Europe via the research project TRACE-Zurich.

## REFERENCES

- [1] Definitions from the us astronomical applications dept (usno). Retrieved 2011-07-22.
- [2] J. M. . Alvarez and A. M. Lopez. Road detection based on illuminant invariance. *IEEE Transactions on Intelligent Transportation Systems*, 12(1):184–193, 2011.
- [3] J. M. Anderson, K. Nidhi, K. D. Stanley, P. Sorensen, C. Samaras, and O. A. Oluwatola. *Autonomous Vehicle Technology: A Guide for Policymakers*. Santa Monica, CA: RAND Corporation, 2016.
- [4] A. Bar Hillel, R. Lerner, D. Levi, and G. Raz. Recent progress in road and lane detection: A survey. *Mach. Vision Appl.*, 25(3):727–745, Apr. 2014.
- [5] Y. Chen and C. Han. Night-time pedestrian detection by visual-infrared video fusion. In *World Congress on Intelligent Control and Automation*, 2008.
- [6] Y. Chen, W. Li, C. Sakaridis, D. Dai, and L. Van Gool. Domain adaptive faster r-cnn for object detection in the wild. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [7] Y. Chen, W. Li, and L. Van Gool. Road: Reality oriented adaptation for semantic segmentation of urban scenes. In *Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [8] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, and B. Schiele. The cityscapes dataset for semantic urban scene understanding. In *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [9] D. Dai, T. Kroeger, R. Timofte, and L. Van Gool. Metric imitation by manifold transfer for efficient vision applications. In *CVPR*, 2015.
- [10] J. Ge, Y. Luo, and G. Tei. Real-time pedestrian detection and tracking at nighttime for driver-assistance systems. *IEEE Transactions on Intelligent Transportation Systems*, 10(2):283–298, 2009.
- [11] A. Gonzalez, Z. Fang, Y. Socarras, J. Serrat, D. Vazquez, J. Xu, and A. M. Lpez. Pedestrian detection at day/night time with visible and fir cameras: A comparison. *Sensors*, 16(6), 2016.
- [12] S. Gupta, J. Hoffman, and J. Malik. Cross modal distillation for supervision transfer. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [13] K. He, G. Gkioxari, P. Dollr, and R. Girshick. Mask r-cnn. In *The IEEE International Conference on Computer Vision (ICCV)*, 2017.
- [14] S. Hecker, D. Dai, and L. Van Gool. End-to-end learning of driving models with surround-view cameras and route planners. In *European Conference on Computer Vision (ECCV)*, 2018.
- [15] G. Hinton, O. Vinyals, and J. Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015.
- [16] J. Hoffman, E. Tzeng, T. Park, J.-Y. Zhu, P. Isola, K. Saenko, A. A. Efros, and T. Darrell. Cycada: Cycle consistent adversarial domain adaptation. In *International Conference on Machine Learning (ICML)*, 2018.
- [17] M. B. Jensen, M. P. Philipsen, A. Mgelmoose, T. B. Moeslund, and M. M. Trivedi. Vision for looking at traffic lights: Issues, survey, and perspectives. *IEEE Transactions on Intelligent Transportation Systems*, 17(7):1800–1815, July 2016.
- [18] J. H. Kim, H. G. Hong, and K. R. Park. Convolutional neural network-based human detection in nighttime images using visible light camera sensors. *Sensors*, 17(5), 2017.
- [19] H. Kuang, K. Yang, L. Chen, Y. Li, L. L. H. Chan, and H. Yan. Bayes saliency-based object proposal generator for nighttime traffic images. *IEEE Transactions on Intelligent Transportation Systems*, 19(3):814–825, 2018.
- [20] E. Levinkov and M. Fritz. Sequential bayesian model update under structured scene prior for semantic road scenes labeling. In *IEEE International Conference on Computer Vision*, 2013.
- [21] G. Lin, A. Milan, C. Shen, and I. Reid. Refinenet: Multi-path refinement networks with identity mappings for high-resolution semantic segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [22] S. Manen, M. Gygli, D. Dai, and L. Van Gool. Pathtrack: Path supervision for efficient video annotation. In *International Conference on Computer Vision (ICCV)*, 2017.
- [23] S. G. Narasimhan and S. K. Nayar. Vision and the atmosphere. *Int. J. Comput. Vision*, 48(3):233–254, July 2002.
- [24] E. Ohn-Bar and M. M. Trivedi. Looking at humans in the age of self-driving and highly automated vehicles. *IEEE Transactions on Intelligent Vehicles*, 1(1):90–104, 2016.
- [25] G. Ros and J. M. Alvarez. Unsupervised image transformation for outdoor semantic labelling. In *IEEE Intelligent Vehicles Symposium (IV)*, 2015.
- [26] C. Sakaridis, D. Dai, S. Hecker, and L. Van Gool. Model adaptation with synthetic and real data for semantic dense foggy scene understanding. In *European Conference on Computer Vision (ECCV)*, 2018.
- [27] C. Sakaridis, D. Dai, and L. Van Gool. Semantic foggy scene understanding with synthetic data. *International Journal of Computer Vision*, 2018.
- [28] S. Sankaranarayanan, Y. Balaji, A. Jain, S. N. Lim, and R. Chellappa. Learning from Synthetic Data: Addressing Domain Shift for Semantic Segmentation. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [29] R. K. Satzoda and M. M. Trivedi. Looking at vehicles in the night: Detection and dynamics of rear lights. *IEEE Transactions on Intelligent Transportation Systems*, 2016.
- [30] S. Sivaraman and M. M. Trivedi. Looking at vehicles on the road:

A survey of vision-based vehicle detection, tracking, and behavior analysis. *IEEE Transactions on Intelligent Transportation Systems*, 14(4):1773–1795, 2013.

- [31] F. Xu, X. Liu, and K. Fujimura. Pedestrian detection and tracking with night vision. *IEEE Transactions on Intelligent Transportation Systems*, 6(1):63–71, 2005.
- [32] Y. Xu, Q. Long, S. Mita, H. Tehrani, K. Ishimaru, and N. Shirai. Real-time stereo vision system at nighttime with noise reduction using simplified non-local matching cost. In *IEEE Intelligent Vehicles Symposium (IV)*, 2016.
- [33] Y. Zhang, P. David, and B. Gong. Curriculum domain adaptation for semantic segmentation of urban scenes. In *International Conference on Computer Vision (ICCV)*, 2017.