

Cycle and Semantic Consistent Adversarial Domain Adaptation for Reducing Simulation-to-Real Domain Shift in LiDAR Bird’s Eye View

Alejandro Barrera¹, Jorge Beltrán¹, Carlos Guindel¹, Jose Antonio Iglesias² and Fernando García¹

Abstract—The performance of object detection methods based on LiDAR information is heavily impacted by the availability of training data, usually limited to certain laser devices. As a result, the use of synthetic data is becoming popular when training neural network models, as both sensor specifications and driving scenarios can be generated ad-hoc. However, bridging the gap between virtual and real environments is still an open challenge, as current simulators cannot completely mimic real LiDAR operation. To tackle this issue, domain adaptation strategies are usually applied, obtaining remarkable results on vehicle detection when applied to range view (RV) and bird’s eye view (BEV) projections while failing for smaller road agents. In this paper, we present a BEV domain adaptation method based on CycleGAN that uses prior semantic classification in order to preserve the information of small objects of interest during the domain adaptation process. The quality of the generated BEVs has been evaluated using a state-of-the-art 3D object detection framework at KITTI 3D Object Detection Benchmark. The obtained results show the advantages of the proposed method over the existing alternatives.

I. INTRODUCTION

Nowadays, perception is a crucial task for autonomous vehicles. Research in this field demands accurate sensors and algorithms to perform safe and precise navigation. LiDAR stands as an ideal candidate to directly describe the scene geometry by a dense point cloud representation.

Despite the recent increment in the amount of labeled data, public datasets may not be sufficient to train models able to grasp a complete understanding of the situations that they may encounter in operation due to the domain shift problem. Variations in sensor positions, device specifications (e.g., number and distribution of planes), or even the geographic region [1] can lead to a significant performance drop in supervised learning approaches. Moreover, the annotation type (point-wise, 3D boxes, etc.) could also differ between source and target samples, and collecting well-annotated data for custom applications is prohibitively expensive.

Hence, synthetic data stands as an enticing option to provide on-demand and accurate data which can be modified and extended almost infinitely. Despite the realism of current simulators’ sensor and world models available today, algorithms trained with these data usually fail to generalize in a real environment.

Domain adaptation (DA) techniques have been explored to bridge the aforementioned gaps between domains. Therefore, some approaches have attempted to directly adapt raw

LiDAR information to other data distributions [2]. Nevertheless, due to the sparsity, irregularity, and unstructured distribution of LiDAR data and the high number of points contained in each cloud, on-board perception applications often use efficient projections such as the range view (RV) and the bird’s eye view (BEV), for which DA alternatives have also been proposed [3], [4].

In many of these works, CycleGAN [5] and its cycle consistency mechanism have reported an excellent performance on image-level domain adaptation and content preservation for these 2D projected-based representations. Whilst this method can produce realistic adaptations of big and medium-sized vehicles, we argue that further refinement [6], [7] can help preserve scarcely represented objects, which are normally vulnerable road users such as pedestrians and cyclists.

This work proposes an approach to enhance the style transfer of BEV representations from a synthetic scalable source domain, generated using a simulator, to a real target domain. The conversion, shown in Fig. 1, makes use of an adversarial generative network adapted to BEV representations and endowed with semantic segmentation consistency to help preserve object instances in the scene. To the best of our knowledge, this is the first method addressing unsupervised domain adaptation between unpaired BEV images using a CycleGAN with multi-class semantic regularization.

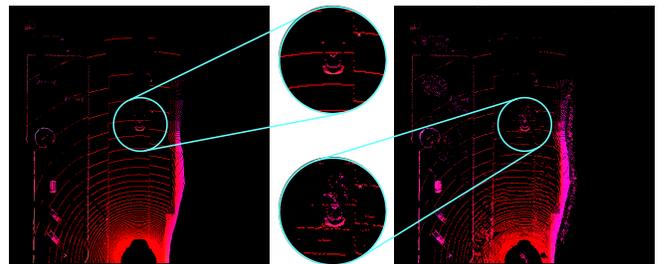


Fig. 1. On the left, raw synthetic BEV recorded in the CARLA simulator. On the right, the result of the proposed domain adaptation. Zoomed regions are provided to better observe the differences.

The goal is to avoid, or at least reduce, the need for labeled samples from the target domain, thus enabling the deployment of high-complexity models on custom setups. In particular, to validate the effectiveness of the approach, we use adapted synthetic data to train a state-of-the-art BEV-based 3D object detection method, BirdNet+ [8], which is later deployed and tested on the KITTI object detection benchmark [9].

¹Authors are with the Dept. of Systems Engineering and Automation, Universidad Carlos III de Madrid, Spain alebarre@pa.uc3m.es, [cguindel](mailto:cguindel@ing.uc3m.es), [jbeltran](mailto:jbeltran@ing.uc3m.es), [fegarcia](mailto:fegarcia@ing.uc3m.es)

²J.A. Iglesias is with the Dept. of Computer Science and Engineering, Universidad Carlos III de Madrid, Spain jiglesia@inf.uc3m.es

II. RELATED WORK

The advent of deep learning has led to a significant research interest in domain adaptation (DA). Among deep DA methods, adversarial-based ones generally use a domain classifier to extract domain-invariant features through the confusion of the source and target domain boundaries [10]. The optimum result from these networks is to minimize the domain distance to maximize the domain discriminator error, producing data that the discriminator cannot distinguish from real [11], [12]. On the other hand, in the reconstruction-based DA category, [13] combines different losses to recombine style and content from two separate images.

CycleGAN [5] combines both solutions using an adversarial loss and a reconstruction loss (cycle consistency loss) to address the image-to-image translation problem when paired training data is not available (unsupervised DA or UDA). Although the CycleGAN reconstruction task shows promising results in a wide variety of scenarios, CYCADA [6] extends its capabilities using both image space alignment and latent representation space alignment. Besides, it incorporates a task to encourage content consistency enforcing relevant semantics to match before and after adaptation. This semantic consistency has proven vital in multimodal scenarios because the invertibility provided by the cycle does not necessarily preserve the arrangement of the classes from the original source domain, as shown in [7]. However, unlike our proposal, this method requires labels of both the source and target domains, as it operates in a supervised fashion.

All these methods are designed for 2D vision tasks where RGB images are the protagonists. However, when it comes to LiDAR point cloud representations, some adaptations are required. In order to work with point clouds, the most straightforward alternative to preserve all the LiDAR information is to use raw clouds to perform point-wise DA and set-level DA [2]. By the same token, PointDAN [14] studies local-level and global-level point cloud alignments by the use of self-adaptive attention nodes.

Although such methods are able to preserve all the LiDAR information, their execution time and memory requirements make them inefficient when it comes to a full point cloud. In this context, projection-based methods gain popularity due to their adaptability to the well-studied 2D approaches. Unfortunately, this also entails the inevitable loss of spatial information. In this field, ePointDA [3] uses range view representations from simulation and real domains to bridge the domain gap at pixel-level and feature-level. LiDARNet [15] combines multiple tasks such as boundary extraction, cycle consistency, and domain invariance to address a full-scene semantic segmentation task on real range view images. BEV-Seg [16] uses multiple camera angles, with RGB and depth images from a simulator to create a semantically enriched point cloud to find BEV semantic segmentation.

Regarding BEV projection, [4] generates from simulation data realistic scenarios to transfer annotations from each other, and [17] shows the capabilities of a similar method on a BEV-based detector.

As can be seen, many of the previous works focus their efforts on simulation-to-real domain adaptation (SRDA). The main reason is to avoid the very challenging annotation task, which is a time and money-consuming task. Considering this issue, simulators such as CARLA [18], which counts with multiple modeled sensors, or LiDAR-based datasets such as GTA-V [19] and SynthCity [20] have been developed. Furthermore, some works improve existing synthetic data adding well-modeled obstacles where needed [21].

III. PROPOSED METHOD

This section provides a detailed explanation of the proposed approach, which is depicted in Fig. 2. Two different sets of bird’s eye view (BEV) images, encoding LiDAR information, are used as input to carry out a transformation between unpaired synthetic BEV point cloud data and real data. This fact will make it possible to use annotations from synthetic data in place of real data, and therefore, expand and improve the diversity of the objects for the detection task.

A. Input Representation

Our BEV representation follows the one proposed in [8]; however, we dispense with the LiDAR intensity data, for which realistic values are difficult to obtain. Thus, three distance-invariant channels are used: maximum height, normalized density within each cell, and binary occupancy. In our experiments, we use a voxel size of 10 cm to obtain handleable representations and a data range of 50 m forward and 22.5 m to each side in order to represent the area where the majority of annotations are available in the KITTI dataset.

For the generation of the synthetic BEV images, we rely on the CARLA simulator [18]. This simulator provides multiple realistic scenarios, agent models, and sensors generated by the graphics engine Unreal. We use a semantically enhanced LiDAR modeled after the KITTI dataset device and let the domain adaptation model the noise.

B. Architecture Description

Our adversarial-based approach, which provides a translation of both representations guided by cycle, identity, and strong pixel-level semantic-aware consistencies, is built upon the CycleGAN architecture [5].

Adversarial network. Given the source domain X (synthetic BEV images) and the target domain Y (real data), an adversarial network aims to map the data distributions of each domain, $x \sim p_d(x)$ and $y \sim p_d(y)$, to the other. First, the architecture is composed of two classifiers named domain discriminators, D_X and D_Y , that learn the characteristic features of each domain separately. Then a set of generators G and F is designed to translate $G : X \rightarrow Y$ and $F : Y \rightarrow X$. Finally, the discriminators D_X and D_Y will provide the necessary feedback of the undergoing mapping $G(x)$ and $F(y)$ until they can not distinguish the domains.

Generators G and F , following the architecture in [13], are organized as follows: first, the encoder module reduces the initial resolution smoothly by a downsampling factor of 4. Afterward, the transformer, which leads the conversion

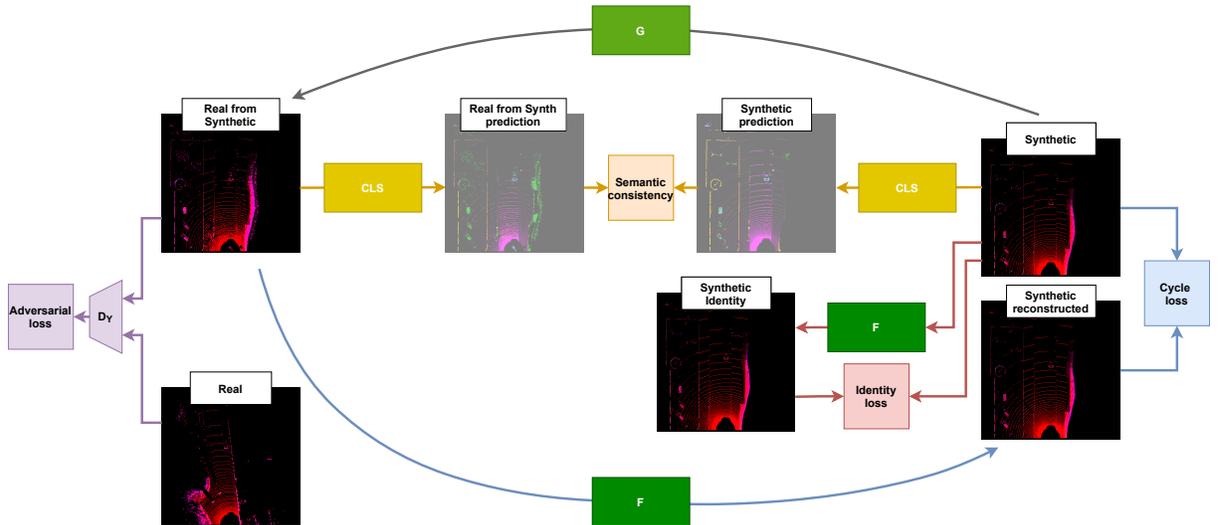


Fig. 2. Overview of the proposed approach. We depict the adversarial loss, to transform the domains, in purple, the cycle consistency mechanism in blue, the identity consistency in red, and the semantic constraint in orange. Blocks G and F stand for the generators $G : X \rightarrow Y$ and $F : Y \rightarrow X$, D_Y for the discriminator of the real domain, and CLS for the semantic segmentation network. For clarity reasons, only losses where the synthetic representation is involved are shown. The real cycle is designed to mirror the synthetic one.

between domains, is built with nine residual blocks that follow a structure similar to that of the encoder, but to which strong dropout is applied during training to provide noise [22]. Finally, the decoder mimics the encoder structure, but it contains fractionally strided convolutions and one convolution to map features to RGB values, followed by a tanh activation function.

On the other hand, discriminator networks are inspired in the PatchGAN architecture [22], where four sets of convolution + instance normalization + LeakyReLU layers downsample the initial resolution and then, two non-strided convolutions predict on 70×70 -pixel overlapping image patches the domain to which they belong.

Cycle consistency. Although adversarial training produces some resemblance between each domain results, it still lacks a mechanism to preserve the structure from the original domain after the conversion. Therefore, in [5], a cyclical training is proposed to prevent G and F from contradicting each other; thus, the mapping $G(F(y))$ will attempt to reproduce the content in y and $F(G(x))$ the content in x .

Identity consistency. One more constrain is imposed on the generators to minimize their distance to an identity mapping when real samples of the opposite domain are provided: $y \approx G(y)$ and $x \approx F(x)$. This idea generally preserves better the information contained in each channel, which may otherwise suffer from undesirable blending under adversarial training [5]. Working with a BEV representation requires retaining the per-pixel structure of the channels so that crucial information, such as object height or cell density, is not modified in the process.

Semantic consistency. Unlike [17], Sem-GAN [7] demonstrates that the previous methods (i.e., cycle and identity consistency) do not necessarily maintain object identities, but

they focus on the whole image adaptation instead. It implies that small details in a BEV representation such as poles and pedestrians could disappear during the transformation. In order to encourage source voxels to keep their structure while being translated to the target domain, we follow the approach in CYCADA [6]. The idea consists of training a semantic segmentation network on the synthetic domain beforehand and using its predictions to ensure high semantic consistency after conversion, thus preserving the fine-grained content and the style of the input.

As our method is unsupervised, we only use synthetic labels in the process. In addition, for the semantic segmentation task, we chose a network that has previously provided state-of-the-art results in LiDAR projection-based representations such as SalsaNext [23]. Although not dedicated to operating in BEV representations, its predecessor was able to use both representations indistinctly. This network, arranged in an encoder-decoder fashion, is composed of a contextual module that stacks three residual blocks to fuse features of two different receptive fields. Afterward, the network follows a U-Net-based structure concatenating residual blocks i with the $n - i$ blocks. Similarly, the decoder utilizes a sequence that mimics the encoder; however, it is preceded by a pixel-shuffle layer.

C. Loss Functions

To train the proposed adversarial network, two different cost functions are minimized. On the one hand, each discriminator D tries to reduce, independently of the generator, an adversarial loss \mathcal{L}_{adv_D} ; hence, for D_Y :

$$\mathcal{L}_{adv_{D_Y}} = \mathbb{E}_{y \sim p_d(y)} [(D_Y(y) - 1)^2] + \mathbb{E}_{x \sim p_d(x)} [D_Y(G(x))^2] \quad (1)$$

In this way, the classifier is trained to distinguish between its domain ($D_Y(y) \approx 1$) and the domain representation created

by the opposite generator ($D_Y(G(x)) \approx 0$).

On the other hand, the generators’ multi-task training loss is given by the following equation, where each component has a weight λ and is computed in both directions; i.e., $X \rightarrow Y$ and $Y \rightarrow X$:

$$\mathcal{L} = \mathcal{L}_{\text{adv}_{\text{gen}}} + \lambda_{\text{cyc}}\mathcal{L}_{\text{cyc}} + \lambda_{\text{idt}}\mathcal{L}_{\text{idt}} + \lambda_{\text{sem}}\mathcal{L}_{\text{sem}} \quad (2)$$

The generators attempt to minimize the adversarial loss $\mathcal{L}_{\text{adv}_{\text{gen}}}$, where the usual negative log-likelihood or binary cross-entropy loss has been modified by a least-squares loss, as in [5]:

$$\mathcal{L}_{\text{adv}_{\text{gen}}} = \mathbb{E}_{x \sim p_d(x)} [(D_Y(G(x)) - 1)^2] + \mathbb{E}_{x \sim p_d(y)} [(D_X(F(y)) - 1)^2] \quad (3)$$

The first reconstruction component, namely cycle consistency \mathcal{L}_{cyc} , is an L1 penalty between the initial sample from one domain to the final representation, after both translations $x \rightarrow G(x) \rightarrow F(G(x)) \approx x$ (and the analogous for domain Y):

$$\mathcal{L}_{\text{cyc}} = \mathbb{E}_{x \sim p_d(x)} [\|F(G(x)) - x\|_1] + \mathbb{E}_{y \sim p_d(y)} [\|G(F(y)) - y\|_1] \quad (4)$$

Secondly, the identity loss is defined as the mean absolute error, L1 loss, to ensure when a generator is fed by a sample from the opposite domain, it can produce a representation of its own domain:

$$\mathcal{L}_{\text{idt}} = \mathbb{E}_{x \sim p_d(x)} [\|G(y) - y\|_1] + \mathbb{E}_{y \sim p_d(y)} [\|F(x) - x\|_1] \quad (5)$$

Finally, regarding the semantic loss, we use the SalsaNext [23] semantic segmentation outputs as a noisy labeler to keep as much context as possible after translation:

$$\mathcal{L}_{\text{sem}} = \mathcal{L}_{\text{wCE}}(\text{CLS}(G(x)), \text{argmax}(\text{CLS}(x))) + \mathcal{L}_{\text{wCE}}(\text{CLS}(F(y)), \text{argmax}(\text{CLS}(y))) \quad (6)$$

where \mathcal{L}_{wCE} represents a weighted cross-entropy loss, which is computed over the pixels of the semantic prediction (CLS) for $G(x)$ and $F(y)$ with the predictions from the source BEVs (x and y , respectively) as weak labels. Only cells with non-zero values in both inputs contribute to the loss in order to preserve both the class and location of the points. Weights are used to increase the importance, by a $2\times$ factor, of the categories of interest (i.e., cars, pedestrians, and cyclists) over the rest of classes (e.g., roads or buildings), which are still included to maintain the geometric consistency of the complete scene.

As stated above, we use a SalsaNext model as semantic predictor (CLS). This model is trained beforehand with synthetic data through the usual weighted multi-category cross-entropy and Lovász-Softmax losses.

IV. EXPERIMENTAL RESULTS

In this section, we present a set of experiments to validate the performance of our domain adaptation approach from synthetic BEV data to real BEV representations. Our implementation will be evaluated on the well-known KITTI object dataset using the 3D detector specified below.

A. 3D Obstacle Detection for Evaluation

Assessing the feasibility of a domain adaptation method on a 3D detector is a common practice nowadays. The main advantage lies in the fact that the results will offer a good estimate of the resemblance of the generated data to the real dataset to be emulated. With this task in mind, we have chosen a 3D detector that uses enriched BEV inputs to provide the object’s location, shape, and category in a two-stage end-to-end fashion. The first stage of the BirdNet+ architecture [8] is built upon a residual-based encoder with per-level skips to preserve global and local content (ResNet-50 and a feature pyramid network). These features are fed into a region proposal network that classifies and refines a default anchor estimation. These non-axis-aligned proposals are dimensionally normalized by an ROI Align layer and forwarded to a second stage composed of a sequence of multiple fully-connected layers, which finally estimate the 3D object parameters.

B. Experimental Setup

As in the original CycleGAN approach [5], the adversarial network is trained from scratch with random weights following a normal distribution $\mathcal{N}(0, 0.02)$ to initialize the weights of every layer in our model. Training data is randomly augmented using different techniques: horizontal flip, point dropout, and additive Gaussian noise with similar distribution to our input representation. Following [12], real labels in (1) and (3) are softened randomly and kept between 0.7 and 1. Additionally, the last 50 samples are used to compute the discriminators’ losses and provide better stability.

For our experiments, we fix $\lambda_{\text{cyc}} = \lambda_{\text{idt}} = 10$ and $\lambda_{\text{sem}} = 0.5$ in (2) to weigh all the losses. We use the Adam solver for the optimization with momentum [0.5, 0.99] and train up to 50 epochs. The number of epochs, batch size, learning rate (LR), and learning rate decay of all networks involved in this work are indicated in Table I.

TABLE I
TRAINING PARAMETERS FOR THE DIFFERENT MODELS.

Network	Epochs	Batch	LR	LR decay
DA network	50	1	0.0001	None
SalsaNext	100	1	0.01	0.01 every epoch
BirdNet+	~12	4	0.01	0.1 at the 10 th epoch

Our synthetic data have been extracted from five different towns delivered by the CARLA simulator. Between 1000 to 1500 samples per town were extracted with a delay of 0.5 simulated seconds to provide more variability to our results. It is worth noting that in one of the maps, we only use cyclists and pedestrians to deal with the unbalanced data.

The scenes contain a limited and approximately constant number of agents, which are spawned and destroyed outside the field of view of our BEV representation to avoid inconsistencies. Besides, we divide the CARLA vehicle class into more fine-grained categories (i.e., car, motorcycle,

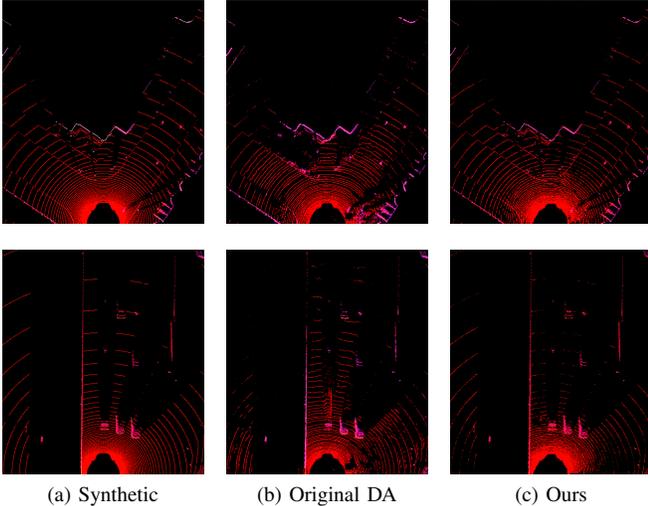


Fig. 3. Domain adaptation results with the original DA (CycleGAN) and the proposed approach.

and bicycle), matching the KITTI labeling criteria for two-wheelers (so that they include both the vehicle and the rider). Additionally, 3D labels for parked vehicles, which are not reported by CARLA, are estimated using the dense BEV semantic representation provided by a top-view camera.

In total, 6878 synthetic BEV images, endowed with semantic and 3D box labels, are extracted. They are used to train both SalsaNext and the proposed domain adaptation framework (together with the KITTI training set as target). Finally, the 3D detection model is trained with the adapted images (considering only the car, pedestrian, and cyclist categories) and tested on the KITTI validation set, defined as in [8].

C. Overall assessment

To shed light on the importance of the proposed consistency in the domain adaptation, Fig. 3 shows a sample from the CARLA simulator, the output after the adversarial training proposed in [5], and the output with our model, including all the losses described in Sec. III-C.

The noise introduced in the original DA approach (Fig. 3b) adds unrealistic LiDAR measurements and changes the pixel value of some areas. On the other hand, our method (Fig. 3c) seems to eliminate some points from the ground that may impact the segmentation task, but it preserves significantly better the semantic identity of each individual pixel while modeling the noise of the real LiDAR.

For the evaluation of the performance of the BirdNet+ detector trained with the adapted data, we follow the 3D and BEV detection tasks from the KITTI object detection benchmark [9]. The strong IoU requirements between detections and labels imposed by the official evaluation do not fit well the localization uncertainty shown by models trained only with synthetic data; therefore, we employ less strict thresholds: 50% (cars) and 30% (pedestrians and cyclists). Table II shows the considerable performance gap in the domain shift between the synthetic and the real in our

representations. It is worth noting that the original DA method focuses either on medium and big-sized obstacles (i.e., cars) and structures rather than small objects, which, in the end, are partially occluded by the noise generated. In our approach, the disappearance or modification of these elements penalizes the generator, that becomes more aware of the semantic context of each point. Thus, our method preserves better the details in the scene, easing the task of 3D object detection in BEV representations.

TABLE II
BIRDNET+ DETECTION PERFORMANCE (AP BEV % AND AP 3D %) ON THE KITTI VAL SET FOR THE DIFFERENT TRAINING DATA SOURCES.

	Car		Pedestrian		Cyclist	
	BEV	3D	BEV	3D	BEV	3D
Oracle (KITTI)	81.94	67.04	50.17	43.90	42.74	39.89
Synthetic	52.91	46.82	18.11	17.91	22.37	21.85
Original DA	61.53	48.08	10.58	07.45	18.82	16.56
Ours	53.79	48.61	26.21	25.81	29.88	29.75

The validity of the method is further confirmed in the qualitative results depicted in Fig. 4. As can be seen, our method adjusts the elevation of cars (first row) and pedestrians (third row) in a better way due to the fact that it preserves the pixel values better than the others whilst generating fewer artifacts. In addition, our method provides more recall in small classes (second and third rows) such as pedestrians and cyclists; however, it occasionally fails to distinguish between them. It is clear that, although detection capabilities are naturally limited by the domain gap, our method demonstrates a significant improvement over its predecessor and the synthetic-only approach in the 3D detection task.

V. CONCLUSIONS

In this work, it has been shown that the enforcement of semantic consistency in GAN-based domain adaptation of BEV projections benefits the preservation of the original layout of the elements in the synthetic scene during style transfer. The performance of the presented framework has been assessed using a state-of-the-art object detection network in the challenging KITTI Benchmark. Contrary to the baseline method, our results improve by a wide margin those obtained when training with raw synthetic data, being especially significant the difference in the detection of smaller road participants.

In future works, a lossless LiDAR style transfer will be studied so that any kind of object detection network can be used regardless of its input representation. To this aim, two different approaches will be explored: first, by means of several simultaneous generators per domain dedicated to different projections, which will ultimately allow reconstructing the LiDAR point cloud; and second, by designing a method able to perform domain adaptation over the raw 3D information.



Fig. 4. Qualitative results in KITTI validation set produced by BirdNet+ using different training data. From left to right: raw synthetic data, cycle consistent DA, cycle and semantic consistent DA.

ACKNOWLEDGMENT

Research conducted within the project PEVAUTO-CM-UC3M. The research project PEVAUTO-CM-UC3M has been funded by the call “Programa de apoyo a la realización de proyectos interdisciplinarios de I+D para jóvenes investigadores de la Universidad Carlos III de Madrid 2019-2020 under the frame of the Convenio Plurianual Comunidad de Madrid-Universidad Carlos III de Madrid. We gratefully acknowledge the support of NVIDIA Corporation with the donation of the GPUs used for this research.

REFERENCES

- [1] Y. Wang, X. Chen, Y. You, L. E. Li, B. Hariharan, M. Campbell, K. Q. Weinberger, and W.-L. Chao, “Train in Germany, test in the USA: Making 3D object detectors generalize,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 11 713–11 723.
- [2] S. Peng, X. Xi, C. Wang, R. Xie, P. Wang, and H. Tan, “Point-based multilevel domain adaptation for point cloud segmentation,” *IEEE Geoscience and Remote Sensing Letters*, 2020.
- [3] S. Zhao, Y. Wang, B. Li, B. Wu, Y. Gao, P. Xu, T. Darrell, and K. Keutzer, “ePointDA: An end-to-end simulation-to-real domain adaptation framework for LiDAR point cloud segmentation,” *arXiv preprint arXiv:2009.03456*, 2020.
- [4] A. E. Sallab, I. Sobh, M. Zahran, and N. Essam, “LiDAR sensor modeling and data augmentation with GANs for autonomous driving,” *arXiv preprint arXiv:1905.07290*, 2019.
- [5] J.-Y. Zhu, T. Park, P. Isola, and 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.
- [6] J. Hoffman, E. Tzeng, T. Park, J.-Y. Zhu, P. Isola, K. Saenko, A. Efros, and T. Darrell, “Cycada: Cycle-consistent adversarial domain adaptation,” in *International Conference on Machine Learning*, 2018, pp. 1989–1998.
- [7] A. Cherian and A. Sullivan, “Sem-GAN: semantically-consistent image-to-image translation,” in *2019 IEEE Winter Conference on Applications of Computer Vision (WACV)*. IEEE, 2019, pp. 1797–1806.
- [8] A. Barrera, C. Guindel, J. Beltrán, and F. García, “BirdNet+: end-to-end 3D object detection in LiDAR bird’s eye view,” in *2020 IEEE 23rd International Conference on Intelligent Transportation Systems (ITSC)*, 2020, pp. 2985–2990.
- [9] A. Geiger, P. Lenz, and R. Urtasun, “Are we ready for autonomous driving? the KITTI Vision Benchmark Suite,” in *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2012, pp. 3354–3361.
- [10] Y. Ganin, E. Ustinova, H. Ajakan, P. Germain, H. Larochelle, F. Laviolette, M. Marchand, and V. Lempitsky, “Domain-adversarial training of neural networks,” *The journal of machine learning research*, vol. 17, no. 1, pp. 2096–2030, 2016.
- [11] I. J. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, “Generative adversarial networks,” *arXiv preprint arXiv:1406.2661*, 2014.
- [12] T. Salimans, I. Goodfellow, W. Zaremba, V. Cheung, A. Radford, and X. Chen, “Improved techniques for training GANs,” *arXiv preprint arXiv:1606.03498*, 2016.
- [13] J. Johnson, A. Alahi, and L. Fei-Fei, “Perceptual losses for real-time style transfer and super-resolution,” in *European conference on computer vision*. Springer, 2016, pp. 694–711.
- [14] C. Qin, H. You, L. Wang, C.-C. J. Kuo, and Y. Fu, “PointDAN: A multi-scale 3D domain adaption network for point cloud representation,” *arXiv preprint arXiv:1911.02744*, 2019.
- [15] P. Jiang and S. Saripalli, “LiDARNet: A boundary-aware domain adaptation model for lidar point cloud semantic segmentation,” *arXiv preprint arXiv:2003.01174*, 2020.
- [16] M. H. Ng, K. Radia, J. Chen, D. Wang, I. Gog, and J. E. Gonzalez, “BEV-Seg: Bird’s eye view semantic segmentation using geometry and semantic point cloud,” *arXiv preprint arXiv:2006.11436*, 2020.
- [17] K. Saleh, A. Abobakr, M. Attia, J. Iskander, D. Nahavandi, M. Hossny, and S. Nahvandi, “Domain adaptation for vehicle detection from bird’s eye view LiDAR point cloud data,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops*, 2019.
- [18] A. Dosovitskiy, G. Ros, F. Codevilla, A. Lopez, and V. Koltun, “CARLA: An open urban driving simulator,” in *Conference on Robot Learning*, 2017, pp. 1–16.
- [19] B. Wu, X. Zhou, S. Zhao, X. Yue, and K. Keutzer, “SqueezeSegV2: Improved model structure and unsupervised domain adaptation for road-object segmentation from a LiDAR point cloud,” in *2019 International Conference on Robotics and Automation (ICRA)*, 2019, pp. 4376–4382.
- [20] D. Griffiths and J. Boehm, “SynthCity: A large scale synthetic point cloud,” *arXiv preprint arXiv:1907.04758*, 2019.
- [21] J. Beltrán, I. Cortés, A. Barrera, J. Urdiales, C. Guindel, F. García, and A. de la Escalera, “A method for synthetic LiDAR generation to create annotated datasets for autonomous vehicles perception,” in *2019 IEEE Intelligent Transportation Systems Conference (ITSC)*. IEEE, 2019, pp. 1091–1096.
- [22] P. Isola, J.-Y. Zhu, T. Zhou, and 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.
- [23] T. Cortinhal, G. Tzelepis, and E. E. Aksoy, “SalsaNext: Fast, uncertainty-aware semantic segmentation of LiDAR point clouds,” in *International Symposium on Visual Computing*, 2020, pp. 207–222.