

LiDARNet: A Boundary-Aware Domain Adaptation Model for Lidar Point Cloud Semantic Segmentation

Peng Jiang¹ and Srikanth Saripalli¹

Abstract— We present a boundary-aware domain adaptation model for Lidar point cloud semantic segmentation. Our model is designed to extract both the domain private features and the domain shared features using shared weight. We embedded Gated-SCNN into the shared features extractors to help it learn boundary information while learning other shared features. Besides, the CycleGAN mechanism is imposed for further adaptation. We conducted experiments on real-world datasets. The source domain data is from the Semantic KITTI dataset, and the target domain data is collected from our own platform (a warthog) in off-road as well as urban scenarios. The two datasets have differences in channel distributions, reflectivity distributions, and sensors setup. Using our approach, we are able to get a single model that can work on both domains. The model is capable of achieving the state of art performance on the source domain (Semantic KITTI dataset) and get 44.0% mIoU on the target domain dataset.

I. INTRODUCTION

Semantic segmentation of point cloud has wide applications in robotics and autonomous vehicles. In recent years, the prevalence of deep learning has lead to a significant development of semantic segmentation in both 2D and 3D data. However, training a supervised deep learning model requires a large amount of annotated data, and the performance of well-trained models can be hurt because of a slight departure from training data. Cameras produce 2-Dimensional data with color information, which is widely used for semantic segmentation [1], [2], [3], [4]. However, cameras are easily affected by light sources and weather. Lidar can work during day and night. And, the 3D point cloud also provides spatial information for complex scenes [5].

With the development of deep learning, semantic segmentation based on two-dimensional images has made great progress. Compared with the 2-dimensional image, point cloud has characteristics of sparsity, randomness, and irregularity, which makes directly applying neural network on point cloud difficult. However, by pre-processing and carefully design representation, we are able to leverage the neural network to learn from point cloud [6], [7], [8], [5], [9], [10]. Irrespective of the neural network model, the training process always requires a large volume of data. Multiple large datasets for image-based semantic segmentation exist [11], [12], [13]. There is only one point-wise labeled semantic Lidar dataset [14]. Other than requiring a large amount of labeled data, neural networks are also poor at

generalizing learned knowledge to new datasets or environments. Moreover, a slight change in the dataset can cause failure. Furthermore, data annotation is a high workforce cost and time-consuming task. These mentioned problems motivate researchers to approach the semantic segmentation problem of the unlabeled dataset by domain adaptation. Most of the existing work on domain adaptation for semantic segmentation task has focused on the adaptation of images of urban scenes and recently made efforts to adapt models from simulated images to real-world images [15], [16].



Fig. 1: Our platform for collecting data

In this paper, we propose an end-to-end trainable boundary-aware domain adaptation model for the Lidar point cloud semantic segmentation task. Inspired by [17], the model has two branches: domain private branch and domain shared branch. The domain private branch is able to extract features that can be used to distinguish the two domains, while the domain shared branch can extract features that contain shared information from both domains. Intuitively, the semantic segmentation information should be included in the domain shared features. Besides, boundary information as a part of the semantic segmentation information is low-level information but is much easier to learn. Furthermore, if we were able to get the boundary information, we can use it to penalize the segmentation. Therefore, we add Gated-SCNN [4] to the domain shared branch. Gated-SCNN is a two-stream CNN architecture that has a shape stream of learning boundary-related information effectively. Because the target data is unlabeled, we use generative adversarial networks(GAN) to help the domain shared branch to learn the semantic information [18]. GAN is an architecture con-

¹Peng Jiang (maskjp@tamu.edu) and Srikanth Saripalli (ssaripalli@tamu.edu) are with the J. Mike Walker '66 Department of Mechanical Engineering, Texas A&M University, College Station, TX 77840, USA

sisting of a generator and a discriminator. The generator tries to generate data that can cheat the discriminator, while the discriminator keeps learning to discriminate against the fake data generated by the generator [19]. In our case, the segmentation model is the generator that can generate labels based on target data. The discriminator is a classifier that can discriminate against the real labels (source labels) and the generated labels (predicted target labels). However, the GAN loss is not enough to eliminate the gap between the two domains. To further reduce the effect of domain shift, we introduce the CycGAN mechanism [20], [15] in the models. CycGAN mechanism learns two mappings between the source domain and target domain $\{G : S \rightarrow T\}$ and $\{G : T \rightarrow S\}$. We verified our adaptation model using the Semantic KITTI dataset [14] and our Lidar dataset. Semantic KITTI dataset was collected on a Volkswagen Passat B6 with Velodyne Lidar. Our dataset was collected on Warthog, a mobile robot with an Ouster Lidar (see Fig.1). The two Lidars have different channel configurations and setting up on vehicles. From our experiments, our method is capable of achieving the state of art performance on the source domain (Semantic KITTI dataset) and recover 86% performance on the target domain dataset based on IoU metrics.

The remainder of this paper is organized as follows. Section II discusses relevant literature and methods on semantic segmentation and domain adaptation. Section III introduce our approach and new model. Next, section IV performs a experimental evaluation on our model. Finally, section V concludes the paper and discuss future works.

II. RELATED WORK

A. Semantic segmentation

Semantic segmentation is the task of assigning an object label to each basic unit of data like a pixel or a point in a point cloud. Before the prevalence of deep learning, traditional segmentation methods of images and point cloud mainly relied on handcrafted features from geometric constraints and assumed prior knowledge [5], [21]. With the development of deep learning, image-based semantic segmentation has made significant advances [22], [23]. However, due to the irregularity and lack of structure in point clouds, it is not very easy to apply deep learning on point clouds directly. There are mainly two categories of neural network models that can learn 3D-dimensional information and geometry information. One category of the models directly learn from raw points cloud [6], [7]. Another category of models converts point clouds into structured formats (i.e., images and voxel grids), that allow 2D/3D convolutional operations to be used on the converted data. Instead of using 3D representation, using a spherical projection of the Lidar point cloud allows us to get a denser representation and use 2D Convolutional Network. For semantic segmentation of Lidar point cloud, RangeNet++ achieves the state of results (mean IoU: 52.2%) [8], this work represents Lidar point cloud as range image and uses GPU-base nearest neighbor search methods to refine the results.

B. Unsupervised Domain adaptation

Unsupervised Domain adaptation is a subfield of transfer learning, to learn a discriminative model in the presence of domain shift between domains. Most of the existing work on domain adaptation dealing with semantic segmentation tasks focus on the image of urban scenes and transfer models from simulation rendered image domain to real-world image domain [15], [16]. One class of methods develop adaptation layers to reduce the distribution disparity between two domains in order to improve the performance of the target model [24], [25]. Another class of methods utilizes GAN to train a feature extractor to extract features that cannot be discriminated against by discriminator [26], [27]. Zhang et al. [16] learn global label distributions over images and local distributions over landmark superpixels based on the assumption that urban scenes have strong idiosyncrasies. This method will fail if the data do not have strong idiosyncrasies. The most relevant works to ours are SqueezeSegV2 [10], which uses geodesic correlation alignment and progressive domain calibration to perform domain adaptation between real data and simulation data. However, the adaptation is only for road objects (vehicle and pedestrians) detection, not for full semantic segmentation. Tzeng et al. [26] proposes adversarial discriminative domain adaptation (ADDA), which combines discriminative modeling, untied weight sharing, and a GAN loss. Many of the current states of art domain adaptation models are based on this framework. Hoffman et al. [15] proposed a discriminatively-trained cycle consistent adversarial domain adaptation model (CyCADA), which can adapt representation at both the pixel-level and feature-level.

III. OUR APPROACH

A. Input Representation

In this paper, we focus on the domain adaptation from one Lidar dataset to another Lidar dataset. The ring structure of the Lidar point cloud allows us to make the spherical projection (Eq.1). In Eq.1, r is the range, (x, y, z) are the coordinates, (w, h) are width and height of the image, f is the angle of the field of view of Lidar, and f_{up} is the up angle of the field of view. After making the spherical projection of the Lidar point cloud, we can get a range image and a point index image. Based on these two images, we get a 3D coordinate map. Because these are 2D images, we are then able to use a 2D convolutional network to perform semantic segmentation on the Lidar point cloud. Similar to [8], we uses the range image, reflectivity map, and co-ordinate maps (see Fig.5(a-c)) as inputs to our model. Furthermore, we also use a normal map (see 5(d)) as an input. According to [28], the normal map can help to the networks to perform semantic segmentation.

$$\begin{pmatrix} u \\ v \end{pmatrix} = \begin{pmatrix} \frac{1}{2}[1 - \arctan(y, x)\pi^{-1}]w \\ [1 - (\arcsin(zr^{-1}) + f_{up})f^{-1}]h \end{pmatrix} \quad (1)$$

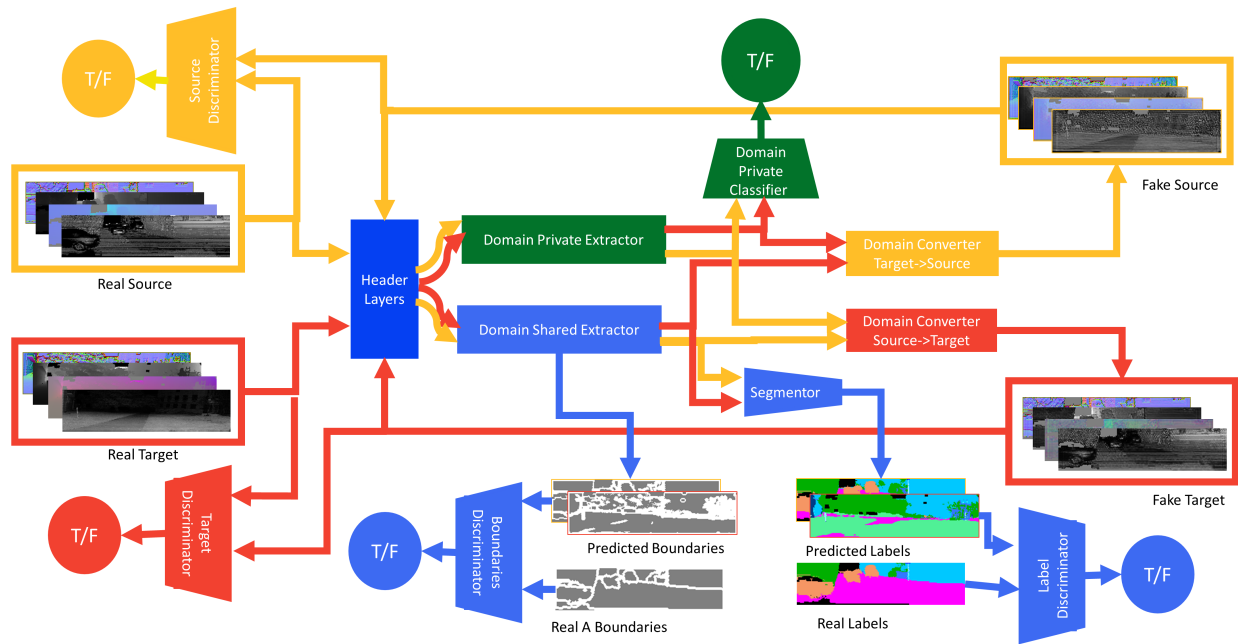


Fig. 2: Information Flow Graph: The data (source and target) are first fed into the header model f_H . After passing the header model, the processed data is fed into two branches: one branch is composed of a domain private extractor f_P and a domain private classifier f_D which can differentiate the input from the two domains, another branch is domain shared extractor f_C which extracts the common feature between the two domains including semantic information Y and boundaries information B . A segmentor f_{Seg} predicts labels \hat{Y} based on the features from domain shared extractor. The predicted boundaries is sent to a boundaries discriminator D_B , and the predicted is sent to a labels discriminator D_Y . Next, the domain private features and domain shared features are fed into domain converters ($f_{S \rightarrow T}$ converts source data into target domain, $f_{T \rightarrow S}$ converts target data into source domain). The conversion are learning through an adversarial learning procedure. Therefore, the converted data are separately fed into domain discriminators (D_T and D_S). Along with this, the converted data is also fed back to the header model to repeat the above procedures.

Algorithm 1 Range Image Preprocessing

input $I_{mask}, I_{range}, I_{label}, I_{reflectivity}$.

$$\hat{I}_{mask} = Closing(I_{mask})$$

$$I_{stripe} = \hat{I}_{mask} - I_{mask}$$

$$\hat{I}_{range} = NSInpainting(I_{range}, I_{stripe})$$

$$\hat{I}_{reflectivity} = NSInpainting(I_{reflectivity}, I_{stripe})$$

$$\hat{I}_{label} = NNInpainting(I_{label}, I_{stripe})$$



Fig. 3: (a) the original label; (b) the inpainted label;

In our experiment, Semantic KITTI dataset is the source dataset, which the model mainly learn from. However, the quality of the KITTI data is low. A lot of stripe pattern appears on projected images and labels (see Fig.3(a)), which affects the domain adaptation process. To reduce the negative effect, we pre-process the data (see Fig.3(b)) as shown in algorithm 1. We first perform Closing morphological operator on mask image to close stripe pattern. We then get the positions of stripes on the image by using the closed

mask image to subtract the original image. Next, based on the stripe positions, we inpaint the reflectivity image and range image using the Navier-Stokes inpainting algorithm [29], and inpaint the label based on the nearest neighbors of the stripes. The algorithm is based on the Navier-Stokes equations for fluid dynamics. In this algorithm, the image intensity is considered as a "stream function" for a 2-dimensional incompressible flow, and the Laplacian of the image intensity is treated as the vorticity of the fluid. The intensity flows from the exterior region into the region to be inpainted by following the vector field of the stream function.

B. Network Structure

Generally, we expect that there exists a model that can complete a task for data from close domains. However, feature difference between two close domains causes a model, which learns from one domain (called source domain), to not perform well on another domain (called target domain). Therefore, we expect a method that can adapt a model from one domain to another domain. If the target domain does not provide ground truth, the problem is called unsupervised domain adaptation. In this paper, the task is semantic segmentation. In this problem, we use X_S denotes source data, Y_S denotes source labels, and X_T denotes target data, but

target labels are not accessible. Intuitively, two close domains should contain shared information, which is necessary for completing the task. In our case, the two Lidar datasets must have the same objects, such as buildings, vegetables, cars. Nevertheless, the data from two datasets also have different features, for example, different noise distribution, point cloud distribution, and reflectivity distribution. Therefore, inspired by [17], we made an assumption that data has the domain private features and domain shared features. We also assume that semantic information is contained in the domain shared features. Therefore, we want to get a feature extractor, which can extract the domain shared features and a segmentator that can use the domain shared features to perform semantic segmentation. In this paper, we designed an end-to-end trainable model that can extract the domain shared features and then utilize the extracted common features to perform the semantic segmentation. After training, the model can perform semantic segmentation on both the source domain and the target domain. Fig.2 shows the information flow in the model. The model has a header model f_H , which is composed of several layers and performs pre-processing on the data. The data is first fed into the header model. After passing the header model, the processed data is fed into two branches: one branch is composed of a domain private extractor f_P and a domain private classifier f_D which can differentiate the input from the two domains, another branch is a domain shared component extractor f_C which extracts the shared feature between the two domains including semantic information Y . A segmentor f_{Seg} can utilize the domain shared features from domain shared extractor to predicts labels \hat{Y} . Meanwhile, the domain shared extractor is forced to learn boundary maps B through Gated-SCNN [4]. The design is based on the intuition that the shared domain features should include boundary information. The predicted boundaries are sent to a boundaries discriminator D_B , and the predicted labels are sent to a labels discriminator D_Y . To further adapt the features of two domains, we adopt the CycGAN mechanism [20], [15] in the model. The domain private features and domain shared features are fed into domain converters to convert the data from one domain to another domain ($f_{S \rightarrow T}$ converts source data into target domain, $f_{T \rightarrow S}$ converts target data into source domain). The conversion is learning through an adversarial learning procedure. Therefore, the converted data are separately fed into domain discriminators (D_T and D_S). To close the cycle of CycGAN, we feed the converted data back to the header model to repeat the above procedures.

C. Multi-task Learning

The domain adaptation procedure is essentially a multi-task learning procedure. The tasks include domain private feature classification, boundaries extraction, semantic segmentation, and domain mutual conversion. The complete loss function can be written as follows:

$$L = \lambda_P L_P + \lambda_B L_B + \lambda_{Seg} L_{Seg} + \lambda_M L_M \quad (2)$$

where L_P , L_B , L_{Seg} and L_M correspond to the loss of domain private feature classification, boundaries extraction, semantic segmentation, and domain mutual conversion, and λ_P , λ_B , λ_{Seg} and λ_M are hyperparameters that control the weighting between losses.

The domain private feature classification task is a binary classification problem, which uses standard binary cross-entropy loss 3.

$$L_{BCE}(Y, \hat{Y}) = E_{y \sim Y} [y \log(\hat{y}) + (1 - y) \log(1 - \hat{y})] \quad (3)$$

Therefore, the classification loss is 4, where $\delta(X) = \{1 : x \in X_S; 0 : x \in X_T\}$.

$$L_P = L_{BCE}(\delta(X), f_D(f_P(f_H(X)))) \quad (4)$$

For the boundaries extraction task, we have access to labels of source data, which allows us to get the boundaries of source data B_S . We then use standard binary cross-entropy (BCE) loss on predicted boundary maps \hat{B}_s of source data. From experiments, the network inclines to generate blank results if there was no penalty on target data. Therefore, we added a GAN loss to encourage the network to predict boundaries for target data too. We express GAN loss as 5

$$L_{GAN}(G, D_Y, X, Y) = E_{y \sim Y} [\log D_Y(y)] + E_{x \sim X} [\log(1 - D_Y(G(x)))] \quad (5)$$

Then, the complete loss function of boundary extraction task can be written as 6, where $G_B(x)$ equals $f_C(f_H(x))$, and $\lambda_{B_{GAN}}$, $\lambda_{B_{BCE}}$ are hyper-parameters for balancing the effect between the GAN loss and BCE loss.

$$L_{bd} = \lambda_{B_{BCE}} L_{BCE}(B_S, \hat{B}_S) + \lambda_{B_{GAN}} L_{GAN}(G_B, D_B, X_T, B_T) \quad (6)$$

For the semantic segmentation task, L_{Seg} consist of two parts: L_{Seg}^S of source data and L_{Seg}^T of target data. We employ standard cross-entropy (CE) loss with dual boundary regularizer [4] on predicted labels of source data. The dual boundary regularizer penalizes the boundaries error between the ground truth label and predicted label and also encourages the boundary predictions to be consistent with the predicted boundaries.

$$L_{Seg}^S = \lambda_{SS1} L_{CE}(Y_S, \hat{Y}_S) + \lambda_{SS2} L_{dual}(Y_S, \hat{Y}_S, \hat{B}_S) \quad (7)$$

Where $G_{Seg}(X) = f_{seg}(f_C(f_H(X)))$.

For target data, we employ GAN loss to learn segmentation [18]. Besides, from our observation, learning the boundary map is easier than learning semantic segmentation. Therefore, we consider the boundary prediction \hat{B}_T from the domain shared extractor to be the true boundary of the semantic labels. Then, the predicted labels \hat{Y}_T can be penalized by boundary prediction. We add a Laplacian layer in the model to extract the boundary of the predicted labels and use the L1 loss to measure the difference. In the end, the segmentation loss of source data is 8

$$L_{Seg}^T = \lambda_{ST1} L_{GAN}(G_{Seg}, D_{Seg}, X_T, Y_S) + \lambda_{ST2} E_{x_t \sim X_T}^{b_t \sim \hat{B}_T} [\|Laplacian(G_{Seg}(x_t)) - \hat{b}_t\|] \quad (8)$$

In order to further eliminate the effect of domain difference, we introduce the CycleGAN mechanism into our model, which leads to the fourth task: domain mutual conversion task. We expect that through learning the domain mutual conversion, the model can find the interior relationship between two domains. The mutual conversion task requires two mapping functions: $G_{S \rightarrow T}$ maps data from the source domain to target domain, $G_{T \rightarrow S}$ maps data from target domain to source domain. The two mappings function can be expressed as 9.

$$\begin{aligned} G_{T \rightarrow S}(X) &= f_{T \rightarrow S}(f_P(f_H(x_t)), f_C(f_H(X))) \\ G_{S \rightarrow T}(X) &= f_{S \rightarrow T}(f_P(f_H(x_s)), f_C(f_H(X))) \end{aligned} \quad (9)$$

Based on the two mapping function, we can define the domain mutual conversion loss as follows:

$$L = \lambda_{Minv} L_{inv} + \lambda_{Mcy} L_{cy} \quad (10)$$

Where L_{inv} and L_{cy} represent domain invariance loss and cycle consistency loss.

The domain invariance means that the data domain will not be changed if it passes through its domain convertor. For example, we will get data in source domain $X_{S(S)}$ after data from source domain X_S pass the mapping function $G_{T \rightarrow S}$. This invariance character of the mapping function can be learned through the following function:

$$\begin{aligned} L_{inv}(G_{S \rightarrow T}, G_{T \rightarrow S}, X_S, X_T) &= \\ E_{x_s \sim X_S} [\| \hat{x}_{s(s)} - x_s \|_1] & \\ + E_{x_t \sim X_T} [\| \hat{x}_{t(t)} - x_t \|_1] & \\ + L_{GAN}(G_{S \rightarrow T}, D_T, \hat{X}_{S(S)}, X_T) & \\ + L_{GAN}(G_{T \rightarrow S}, D_S, \hat{X}_{T(T)}, X_S) & \end{aligned} \quad (11)$$

On the other end, cycle consistency means that after data passes two different mapping functions, its domain should be in its original domain. For example, source domain data X_S first passes the mapping function $G_{S \rightarrow T}$. The converted results $X_{T(S)}$ passes the mapping function $G_{T \rightarrow S}$. We will finally get data $X_{S(T(S))}$, which should be in the source domain. The cycle consistency loss can be defined as:

$$\begin{aligned} L_{cy}(G_{S \rightarrow T}, G_{T \rightarrow S}, X_S, X_T) &= \\ E_{x_s \sim X_S} [\| \hat{x}_{s(t(s))} - x_s \|_1] & \\ + E_{x_t \sim X_T} [\| \hat{x}_{t(s(t))} - x_t \|_1] & \\ + E_{x_s \sim X_S} [\| G_{Seg}(\hat{x}_{s(t(s))}) - G_{Seg}(x_s) \|_1] & \\ + E_{x_t \sim X_T} [\| G_{Seg}(\hat{x}_{t(s(t))}) - G_{Seg}(x_t) \|_1] & \end{aligned} \quad (12)$$

IV. EXPERIMENTAL EVALUATION

A. Dataset.

We evaluated our algorithm on two datasets: the source dataset is semantic KITTI dataset [14], which is labeled from the KITTI dataset collected around the mid-size city of Karlsruhe, in rural areas and on highways. The data was collected using a Volkswagen Passat B6 with a Velodyne HDL-64E. The target dataset was collected by us on Texas A& M University, College Station campus, and RELLIS

campus. The data was collected using a Warthog with an Ouster OS1-64 Lidar. The robot was driven on the road and sidewalks (see Fig.4(b)). Besides the campus's scene, we also collect data in the off-road environment (see Fig.4(a)). The Semantic KITTI dataset has 23201 labeled scans. Our dataset has 18987 unlabeled scans and 1200 labeled scans, which were labeled according to [14]. Along with the different platforms and environments, the lidar sensors are also different. KITTI dataset is collected using Velodyne HDL-64, while our dataset is collected from Ouster OS1-64. The two sensors have different vertical fields of view and vertical angular resolutions. Velodyne HDL-64 has a field of view $26.33(+2/-24.33)$ degrees, and its vertical angular resolution is $1/3$ degree from $+2$ to -8.33 and $1/2$ degree from -8.83 to -24.33 . The field of view of Ouster is $45.0(+22.5/-22.5)$ degrees, and the vertical angular has a normal distribution (0.703 degree). Moreover, the reflectivity is different. The reflectivity of Velodyne is much noisier than Ouster.



Fig. 4: Dataset Collection Environment: (a) is off-road environment of RELLIS campus; (b) is part of College Station Campus;

B. Evaluation Metrics.

To evaluate performance of our model, we use the widely used intersection-over-union (IoU) metric, mIoU, over all classes [30], given by

$$mIoU = \frac{1}{C} \sum_{c=1}^c \frac{TP_c}{TP_c + FP_c + FN_c} \quad (13)$$

where TP_c, FP_c and FN_c represent the number of true positive, false positive and false negative predictions for class c and C is the number of classes.

C. Implementation Details

In our experiments, all our networks are implemented using PyTorch. Training is done on an NVIDIA GTX 1080-Ti. The lidar point cloud was projected as 64×1028 resolution range images. While training, the projected data was randomly cropped into a size of 64×256 , with a batch size of 2. We use group normalization to replace batch normalization because of our small batch size [9]. Moreover, we also replaced the ReLU activation function as SeLU function and Dropout as AlphaDropout to avoid "dying ReLU". Furthermore, we initialize the model uses the method described in [31].

TABLE I: IoU of semantic segmentation results. RangeNet and Source were tested on Semantic KITTI. The results of RangeNet++ is cited from [8]. The other four models were tested on our dataset. "Baseline" means that the model was trained on Semantic KITTI dataset. "Adapt" refers to our domain adaptation model. "+normal" represents the input includes normal map.

Model	car	bicycle	person	road	building	fence	vegetation	trunk	terrain	pole	traffic-sign	mIoU
RangeNet++	0.914	0.257	0.388	0.918	0.874	0.586	0.805	0.555	0.646	0.479	0.559	0.522
Source	0.913	0.130	0.047	0.925	0.830	0.130	0.769	0.555	0.676	0.426	0.216	0.511
Baseline	0.143	0.017	0.019	0.365	0.282	0.070	0.265	0.075	0.182	0.193	0.277	0.172
Baseline+normal	0.331	0.024	0.140	0.432	0.474	0.161	0.571	0.145	0.214	0.294	0.243	0.276
Adapt	0.402	0.161	0.096	0.782	0.437	0.371	0.482	0.164	0.281	0.213	0.182	0.325
Adapt+normal	0.534	0.134	0.183	0.773	0.607	0.768	0.726	0.342	0.381	0.328	0.371	0.440

D. Quantitive Evaluation

In order to evaluate our approach, we need to create a baseline. Table.I row 1 shows the state of art segmentation results of semantic segmentation on the Semantic KITTI dataset from RangeNet++[8], but the evaluation of RangeNet++ used 20 classes. Compared with the Sematic KITTI dataset, our dataset is more unstructured. Hence, we reduce the types of classes from 20 to 11, especially the types of grounds (sidewalks, road). Therefore, we only cite the corresponding results from [8] on Table.I.

To create a baseline for the adaptation performance, we first removed the domain adaptation part of our model and kept the semantic segmentation part of our model G_{Seg} . We trained this segmentor on the Semantic KITTI dataset. The results are shown on Table.I row 2. The mean IoU of the model is 51.1%, which is close to the performance of RangeNet++. Then, we evaluate the model on our dataset. From Table.I row 3 column 13, we notice that the performance of the model drops to 17.2% in terms of the mIoU. It is worth noting that if the input representation includes the normal map, the model trained on the source domain can perform better on the target domain. See Table.I row 4, the performance of all classes except traffic-sign increase in terms of IoU. However, the performance is only 54% of the original performance.

To compare the effect of the normal map in adaptation, we trained two adaptation models: one uses range image, co-ordinate maps, reflectivity image, and normal map as input; another only uses range image, co-ordinate maps, and reflectivity image as input. The results are shown on Table.I row 5 and 6. After adaptation, the performance of both models increased compared with models without adaptation. But the model with the normal map as input achieved better perform on all classes except bicycle and road. The model achieved a 44.0% mIoU, which is 86% of the performance on source data. The improvement is especially apparent on large objects. When comparing the final adaptation result (see Table.I row 6) with with baseline (see Table.I row 3), IoU of the road increases by 40.8%, IoU of the building increases by 32.5%, IoU of the vegetation increases by 46.1%.

E. Qualitative Discussion

In Fig.5, we provide qualitative results of our method on our dataset and Semantic KITTI. Target (a-d) and Source(a-d) are the inputs of the model. Fake Target (a-d) and Fake

Source(a-d) are the conversion results. Row (b) shows the reflectivity images. In the source domain, the reflectivity of near objects has high quality, but the reflectivity of far objects are almost black. The converted target data (Fake Source) has brighter reflectivity of far objects but becomes noisier. In Fig.6, we provide segmentation results from four samples. We can see that without the normal map, both (b) and (d) misclassified buildings (blue) as roads (purple). All models can recognize the person at the right end of the image. He is the data collector, who exists in all the data of our dataset, and is very close to the Lidar. However, Semantic KITTI does not have this pattern; therefore, the models cannot learn how to classify it, and hence all models misclassify. Fig.7 shows 30 registered segmented Lidar scan. (a) and (c) are ground truth results, while (b) and (d) are the predictions. (a) and (b) shows our off-road data. The environment is almost all vegetation. A misclassification is the red part, which represents two moving people. This misclassification also happens in our urban data Fig.7(d). Moreover, in (d), we can also see that the model has difficulty in distinguishing road (purple) and terrain (light green). For further reseach, We are also in the process of open-source our labeled Lidar data as well as our dataset along with our source code.

V. CONCLUSIONS

In this paper, we propose a boundary-aware domain adaptation approach for semantic segmentation of the lidar point cloud. We design a model that can extract domain shared features and domain private features. We utilize the Gated-SCNN to enable the domain shared feature extractor to keep boundary information in the domain shared features and utilize the learned boundary to refine the segmentation results. Our experiments show improvement of around 26.8% mIoU given the baseline. In future work, we further explore more effective point cloud representation and more efficient architecture to learn the general geometry information.

REFERENCES

- [1] L. C. Chen, Y. Zhu, G. Papandreou, F. Schroff, and H. Adam, "Encoder-decoder with atrous separable convolution for semantic image segmentation," in *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, vol. 11211 LNCS. Springer Verlag, 2018, pp. 833–851.
- [2] K. He, G. Gkioxari, P. Dollar, and R. Girshick, "Mask R-CNN," *Proceedings of the IEEE International Conference on Computer Vision*, vol. 2017-October, pp. 2980–2988, dec 2017.

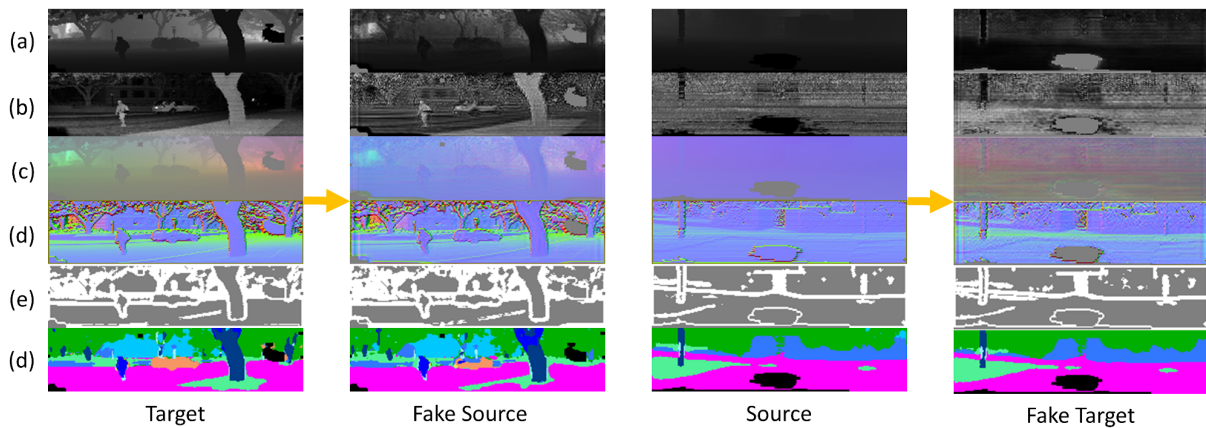


Fig. 5: Input Representation and Corresponding Results: (a) are the range images; (b) are the reflectivity images; (c) are the co-ordinate maps; (d) are the normal maps; (e) are the predicted labels; (d) are the predicted boundaries. Columns 1 and 3 are the original target and source data; Columans 2 and 4 are the conversion results of cyc mechanism.

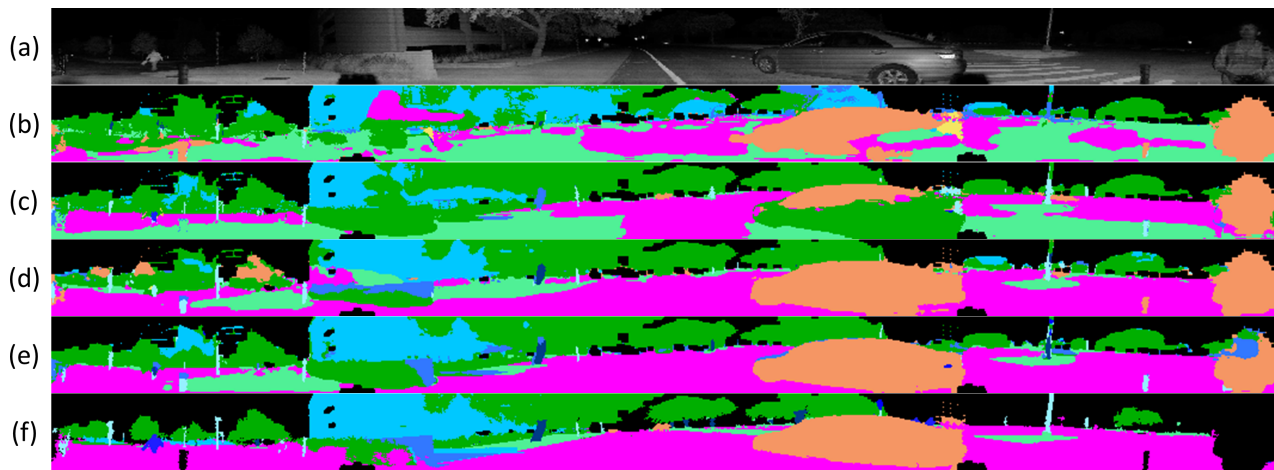


Fig. 6: Segmentation results. (a) is the reflectivity image; (b) is the baseline results; (c) is the baseline+normal results; (d) is the adaptation results; (e) is the adaptation+normal results; (f) is the ground truth

- [3] C. Hazirbas, L. Ma, C. Domokos, and D. Cremers, "FuseNet: Incorporating Depth into Semantic Segmentation via Fusion-Based CNN Architecture," in *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, vol. 10111 LNCS. Springer Verlag, 2017, pp. 213–228. [Online]. Available: http://link.springer.com/10.1007/978-3-319-54181-5_{-}_14
- [4] T. Takikawa, D. Acuna, V. Jampani, and S. Fidler, "Gated-SCNN: Gated Shape CNNs for Semantic Segmentation," jul 2019. [Online]. Available: <http://arxiv.org/abs/1907.05740>
- [5] J. Zhang, X. Zhao, Z. Chen, and Z. Lu, "A Review of Deep Learning-Based Semantic Segmentation for Point Cloud," pp. 179 118–179 133, 2019.
- [6] C. R. Qi, L. Yi, H. Su, and L. J. Guibas, "PointNet++: Deep hierarchical feature learning on point sets in a metric space," in *Advances in Neural Information Processing Systems*, vol. 2017-Decem. Neural information processing systems foundation, jun 2017, pp. 5100–5109. [Online]. Available: <http://arxiv.org/abs/1706.02413>
- [7] C. R. Qi, H. Su, K. Mo, and L. J. Guibas, "PointNet: Deep learning on point sets for 3D classification and segmentation," in *Proceedings - 30th IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017*, vol. 2017-Janua. Institute of Electrical and Electronics Engineers Inc., nov 2017, pp. 77–85.
- [8] A. Milioto, I. Vizzo, J. Behley, and C. Stachniss, "RangeNet++: Fast and Accurate LiDAR Semantic Segmentation, Tech. Rep. i, 2019. [Online]. Available: <https://github.com/PRBonn/lidar-bonnetal>.
- [9] Y. Wu and K. He, "Group normalization," in *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, vol. 11217 LNCS, mar 2018, pp. 3–19. [Online]. Available: <http://arxiv.org/abs/1803.08494>
- [10] 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 *Proceedings - IEEE International Conference on Robotics and Automation*, vol. 2019-May. Institute of Electrical and Electronics Engineers Inc., may 2019, pp. 4376–4382.
- [11] A. Geiger, P. Lenz, C. Stiller, and R. Urtasun, "Vision meets robotics: The kitti dataset," *International Journal of Robotics Research (IJRR)*, 2013.
- [12] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, and B. Schiele, "The cityscapes dataset for semantic urban scene understanding," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016.
- [13] G. Ros, L. Sellart, J. Materzynska, D. Vazquez, and A. M. Lopez, "The SYNTHIA Dataset: A Large Collection of Synthetic Images for Semantic Segmentation of Urban Scenes," in *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, vol. 2016-Decem. IEEE Computer Society, dec 2016, pp. 3234–3243.
- [14] J. Behley, M. Garbade, A. Milioto, J. Quenzel, S. Behnke, C. Stachniss, and J. Gall, "SemanticKITTI: A Dataset for Semantic

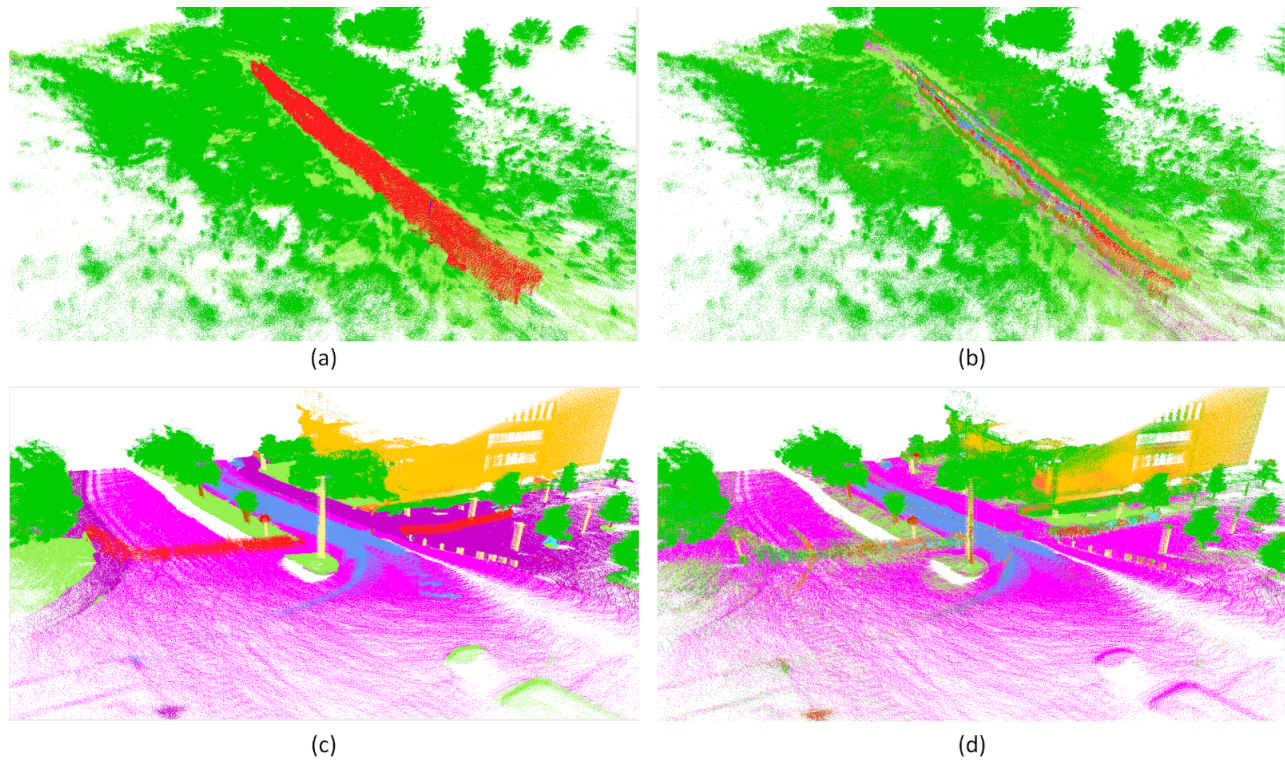


Fig. 7: Point Cloud Segmentation Results. (a) is the ground truth point cloud in off-road environment; (b) is the predicted point cloud in off-road environment; (c) is the ground truth point cloud of Texas A&M University; (d) is the predicted point cloud of Texas A&M University

Scene Understanding of LiDAR Sequences,” apr 2019. [Online]. Available: <http://arxiv.org/abs/1904.01416>

[15] 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 *35th International Conference on Machine Learning, ICML 2018*, vol. 5, nov 2018, pp. 3162–3174. [Online]. Available: <http://arxiv.org/abs/1711.03213>

[16] Y. Zhang, P. David, H. Foroosh, and B. Gong, “A Curriculum Domain Adaptation Approach to the Semantic Segmentation of Urban Scenes,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pp. 1–1, dec 2019. [Online]. Available: <http://arxiv.org/abs/1812.09953>

[17] K. Bousmalis, G. Trigeorgis, N. Silberman, D. Krishnan, and D. Erhan, “Domain separation networks,” in *Advances in Neural Information Processing Systems*, 2016, pp. 343–351.

[18] P. Luc, C. Couprie, S. Chintala, and J. Verbeek, “Semantic Segmentation using Adversarial Networks,” nov 2016. [Online]. Available: <http://arxiv.org/abs/1611.08408>

[19] I. J. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, “Generative adversarial nets,” in *Advances in Neural Information Processing Systems*, vol. 3, no. January. Neural information processing systems foundation, jun 2014, pp. 2672–2680.

[20] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros, “Unpaired Image-to-Image Translation using Cycle-Consistent Adversarial Networks,” mar 2017. [Online]. Available: <http://arxiv.org/abs/1703.10593>

[21] Y. Xie, J. Tian, and X. X. Zhu, “A Review of Point Cloud Semantic Segmentation,” *IEEE Geoscience and Remote Sensing Magazine*, aug 2019. [Online]. Available: <https://arxiv.org/abs/1908.08854><http://arxiv.org/abs/1908.08854>

[22] F. Lateef and Y. Ruichek, “Survey on semantic segmentation using deep learning techniques,” *Neurocomputing*, vol. 338, pp. 321–348, apr 2019. [Online]. Available: <https://doi.org/10.1016/j.neucom.2019.02.003>

[23] A. Garcia-Garcia, S. Orts-Escolano, S. Oprea, V. Villena-Martinez, P. Martinez-Gonzalez, and J. Garcia-Rodriguez, “A survey on deep learning techniques for image and video semantic segmentation,” pp. 41–65, sep 2018.

[24] M. Long, Y. Cao, J. Wang, and M. I. Jordan, “Learning transferable features with deep adaptation networks,” in *32nd International Conference on Machine Learning, ICML 2015*, vol. 1. International Machine Learning Society (IMLS), feb 2015, pp. 97–105. [Online]. Available: <http://arxiv.org/abs/1502.02791>

[25] E. Tzeng, J. Hoffman, N. Zhang, K. Saenko, and T. Darrell, “Deep Domain Confusion: Maximizing for Domain Invariance,” *CoRR*, vol. abs/1412.3, dec 2014. [Online]. Available: <http://arxiv.org/abs/1412.3474>

[26] E. Tzeng, J. Hoffman, K. Saenko, and T. Darrell, “Adversarial discriminative domain adaptation,” in *Proceedings - 30th IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017*, vol. 2017-Janua, 2017, pp. 2962–2971.

[27] C. Yu, J. Wang, Y. Chen, and M. Huang, “Transfer Learning with Dynamic Adversarial Adaptation Network,” in *2019 IEEE International Conference on Data Mining*, sep 2020, pp. 778–786. [Online]. Available: <http://arxiv.org/abs/1909.08184>

[28] S. Gupta, R. Girshick, P. Arbeláez, and J. Malik, “Learning rich features from RGB-D images for object detection and segmentation,” in *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, vol. 8695 LNCS, no. PART 7. Springer Verlag, 2014, pp. 345–360.

[29] M. Bertalmio, A. L. Bertozzi, and G. Sapiro, “Navier-Stokes, fluid dynamics, and image and video inpainting,” in *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, vol. 1, 2001.

[30] M. Everingham, S. M. Eslami, L. Van Gool, C. K. Williams, J. Winn, and A. Zisserman, “The Pascal Visual Object Classes Challenge: A Retrospective,” *International Journal of Computer Vision*, vol. 111, no. 1, pp. 98–136, jun 2014.

[31] K. He, X. Zhang, S. Ren, and J. Sun, “Delving deep into rectifiers: Surpassing human-level performance on imagenet classification,” in *Proceedings of the IEEE International Conference on Computer Vision*, vol. 2015 Inter, feb 2015, pp. 1026–1034. [Online]. Available: <http://arxiv.org/abs/1502.01852>