# SPG: Unsupervised Domain Adaptation for 3D Object Detection via Semantic Point Generation

Qiangeng Xu[1][†]      Yin Zhou[2]      Weiyue Wang[2]      Charles R. Qi[2]      Dragomir Anguelov[2]

[1]University of Southern California          [2]Waymo, LLC

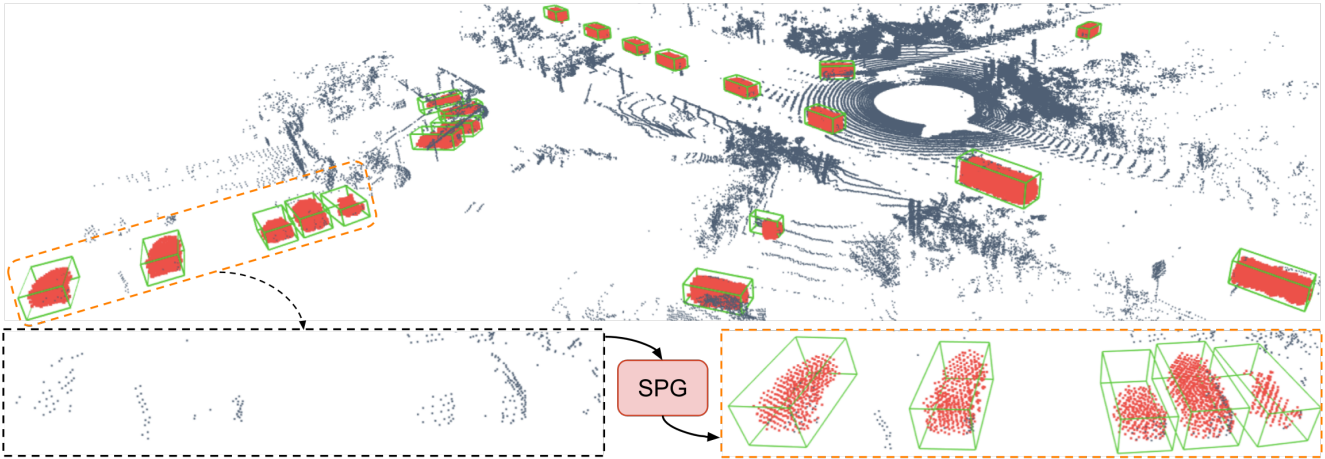{qiangenx}@usc.edu    {yinzhou,weiyuewang,rqi,dragomir}@waymo.com

Figure 1: Our Semantic Point Generation (SPG) recovers the foreground regions by generating semantic points (red). Combined with the original cloud, these semantic points can be directly used by modern LiDAR-based detectors and help improve the detection results (green boxes).

## Abstract

*In autonomous driving, a LiDAR-based object detector should perform reliably at different geographic locations and under various weather conditions. While recent 3D detection research focuses on improving performance within a single domain, our study reveals that the performance of modern detectors can drop drastically cross-domain. In this paper, we investigate unsupervised domain adaptation (UDA) for LiDAR-based 3D object detection. On the Waymo Domain Adaptation [54] dataset, we identify the deteriorating point cloud quality as the root cause of the performance drop. To address this issue, we present Semantic Point Generation (SPG), a general approach to enhance the reliability of LiDAR detectors against domain shifts. Specifically, SPG generates semantic points at the predicted foreground regions and faithfully recovers missing parts of the foreground objects, which are caused by phenomena such as occlusions, low reflectance or weather interference. By merging the semantic points with the original points, we obtain an augmented point cloud, which can be directly con-*

*sumed by modern LiDAR-based detectors. To validate the wide applicability of SPG, we experiment with two representative detectors, PointPillars [25] and PV-RCNN [49]. On the UDA task, SPG significantly improves both detectors across all object categories of interest and at all difficulty levels. SPG can also benefit object detection in the original domain. On the Waymo Open Dataset [54] and KITTI [18], SPG improves 3D detection results of these two methods across all categories. Combined with PV-RCNN [49], SPG achieves state-of-the-art 3D detection results on KITTI.*

## 1. Introduction

A robust autonomous driving system requires its LiDAR-based detector to reliably handle different environmental conditions, *e.g.*, geographic locations and weather conditions. While 3D detection has received increasing interest in recent years, most existing works [79, 7, 10, 11, 16, 23, 25, 26, 27, 29, 30, 36, 41, 42, 49, 50, 51, 53, 63, 64, 65, 67, 68, 62, 78] have focused on the performance in a single domain, where training and test data are captured in similar conditions. It is still an open question how to generalize a 3D detector to different domains, where the environment

---

[†]Work done during internship at Waymo LLC.

| Dataset | Rainy frames | Avg. number of missing points per frame | Avg. number of points per vehicle | 3D L1 AP |
|---------|-------------|------------------------------------------|------------------------------------|----------|
| OD Val | 0.5 % | 23.0K | 306.2 | 56.54 |
| Kirk Dry | 0.0 % | 25.1K | 303.6 | 55.98 |
| Kirk Val | 100.0% | 42.8K | 222.3 | 34.74 |

Table 1: The statistics of OD and Kirk. Each frame contains at most 163.8K points. Kirk Dry is formed by frames with dry weather in Kirk training set.



(a) OD RGB Image      (b) Kirk RGB Image
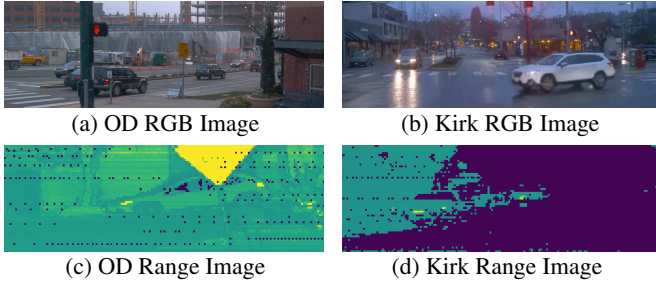
(c) OD Range Image      (d) Kirk Range Image

Figure 2: Examples of RGB and range image (intensity channel) in OD validation set and Kirk validation set. The dark regions in the range images indicate missed LiDAR returns. The regions of "missing points" are irregular in shape.

varies significantly. In this paper, we address the domain gap caused by the deteriorating point cloud quality and aim to improve 3D object detection in the setting of unsupervised domain adaptation (UDA). We use the Waymo Domain Adaptation dataset [54] to analyze the domain gap and introduce semantic point generation (SPG), a general approach to *enhance the reliability of LiDAR detectors against domain shift*. SPG is able to improve detection quality in both the target domain and the source domain and can be naturally combined with modern LiDAR-based detectors.

## 1.1. Understanding the Domain Gap

*Waymo Open Dataset* (OD) is mainly collected in California and Arizona, and *Waymo Kirkland Dataset* (Kirk) [54] is collected in Kirkland. We consider OD as the source domain and Kirk as the target domain. To understand the possible domain gap, we take a PointPillars [25] model trained on the OD training set and compare its 3D vehicle detection performance on OD validation set and those on Kirk validation set. We observe a drastic performance drop of 21.8 points in 3D average precision (AP) (see Table 1).

We first confirm that there is no significant difference in object size between two domains. Then by investigating the meta data in the datasets, we find that only 0.5% of LiDAR frames in OD are collected under rainy weather, but almost all frames in Kirk share the rainy weather attribute. To rule out other factors, we extract all dry weather frames in Kirk training set and form a "Kirk Dry" dataset. Because the the rain drop changes the surface property of objects,

there are twice amount of missing LiDAR points per frame in Kirk validation set than in OD or Kirk Dry (see Table 1). As a result, vehicles in Kirk receive around 27% fewer LiDAR point observations than those in OD (see statistics and more details in the supplemental). In Figure 2, we visualize two range images from OD and Kirk, respectively. We can observe that in the rainy weather, a significant number of points are missing and the distribution of missing points is more irregular compared to the dry weather.

To conclude, the major domain gap between OD and Kirk is the deteriorating point cloud quality, which is caused by the rainy weather condition. In the target domain, we name this phenomenon as the "**missing point**" problem.

## 1.2. Previous Methods to Address the Domain Gap

Multiple studies propose to align the features across domains. Most of them focus on 2D tasks [37, 17, 56, 14] or object-level 3D tasks [77, 45]. Applying feature alignment [9, 20, 35] requires a redesign of the model or loss of a detector. Our goal is to seek a general solution to benefit recently reported LiDAR-based detectors[25, 49, 79, 50, 19].

Another direction is to apply transformations to the data from one domain to match the data from another domain. A naive approach is to randomly down-sample the point cloud but this not only fails to satisfactorily simulate the pattern of missing points (Figure 2d) but also hurts the performance on the source domain. Another approach is to up-sample the point cloud [73, 71, 28] in the target domain, which can increase point density around observed regions. However, those methods have a limited capability in recovering the 3D shape of very partially observed objects. Moreover, upsampling the entire point cloud will lead to a significantly higher latency. A third approach is to leverage style transfer techniques: [80, 40, 12, 20, 48, 21, 47] render point clouds as 2D pseudo images and enforce the renderings from different domains to be resemblant in style. However, these methods introduce an information bottleneck during rasterization [79] and they are not applicable to modern point-based 3D detectors [49].

## 1.3. SPG for Closing the Domain Gap

The "missing point" problem deteriorates the point cloud quality and reduces the number of point observations, thus undermining the detection performance. To address this issue, we propose Semantic Point Generation (SPG). Our approach aims to learn the semantic information of the point cloud and performs foreground region prediction to identify voxels that are inside foreground objects. Based on the predicted foreground voxels, SPG generates points to recover the foreground regions. Since these points are discriminatively generated at foreground objects, we denote them by **semantic points**. These semantic points are merged with the original points into an augmented point cloud, which is

then fed to a 3D detector.

The contributions of this paper are two-fold:

1. We present an in-depth analysis of unsupervised domain adaptation (UDA) for LiDAR 3D detectors across different geographic locations and weather conditions. Our study reveals that the rainy weather can severely deteriorate the quality of LiDAR point clouds and lead to drastic performance drop for modern detectors.

2. We propose semantic point generation (SPG). To our best knowledge, it is the first learning-based model that targets UDA for point cloud 3D detection. Specifically, SPG has the following merits:

- SPG can generate semantic points that faithfully recover the foreground regions suffering from the "missing point" problem. SPG can significantly improve performance over poor-quality point clouds in the target domain while also benefiting source domain, for representative 3D detectors, including PointPillars [25] and PV-RCNN [49].
- SPG also improves the performance for the general 3D object detection task. We verify its effectiveness on KITTI [18] for the aforementioned 3D detectors.
- SPG is a general approach and can be easily combined with modern off-the-shelf LiDAR-based detectors.
- Our approach is light-weight and efficient. Introducing less than 6% additional points, SPG only adds a marginal complexity to a 3D detector.

## 2. Related Work

### 2.1. Unsupervised Domain Adaptation

Unsupervised domain adaptation (UDA) aims to generalize a model to a novel (target) domain by using label information only from the source domain. The two domains are generally related, but there exists a distribution shift (domain gap). Most methods focus on learning aligned feature representations across domains. To reach this goal, [2] proposes Maximum Mean Discrepancy (MMD) while [38] proposes Transfer Component Analysis (TCA). [33] designs a Joint Distribution Adaptation to close the distribution shift while [32, 34] utilize a shared Hilbert space. Without using explicit distance measures, deep learning models [17, 56, 14, 44, 46] use adversarial training to get indistinguishable features between domains.

**Unsupervised Domain Adaptation for 2D Detection** The object detection task is sensitive to local geometric features. [9, 20] hierarchically align the features between domains. Most of these works focus on UDA for 2D detection. With the current advances of unpaired style transfer methods [40, 80], studies such as [48, 21] translate the image from source domain to target domain or vice versa.

**Unsupervised Domain Adaptation for 3D Tasks** Most of the UDA methods focus on 2D tasks, only a few studies explore the UDA in 3D. [77, 45] align the global and local features for object-level tasks. To reduce the sparsity, [59] projects the point cloud to 2D view, while [47] projects the point cloud to birds-eye view (BEV). [15] creates a car model set and adapts their features to the detection object features. However, this study targets general car 3D detection on a single point cloud domain. [57] is the first published study targeting UDA for 3D LiDAR detection. They identify the vehicle size as the domain gap between KITTI[18] and other datasets. So they resize the vehicles in the data. In contrast, we identify the point cloud quality as the major domain gap between Waymo's two datasets[54]. We use a learning-based approach to close the domain gap.

### 2.2. Point Cloud Transformation

One way to improve point cloud quality is to suitably transform the point cloud. Studies of point cloud upsampling [73, 71, 28] can transfer a low density point cloud to a high density one. However, they need high density point cloud ground truth during training. These networks can densify the point cloud in the observed regions. But in our case, we also need to recover regions with no point observation, caused by "missing points".

Point cloud completion networks [74, 6, 66, 61] aim to complete the point cloud. Specialized in object-level completion, these models assume a single object has been manually located and the input only consists of the points on this object. Therefore, these models do not fit our purpose of object detection. Point cloud style transfer models [4, 3] can transfer the color theme and the object-level geometric style for the point cloud. However, these models do not focus on preserving local details with high-fidelity. Therefore, their transformation cannot directly help 3D detection.

## 3. Semantic Point Generation

In the input point cloud $PC_{raw} = \{p_1, p_2, ..., p_N\} \in \mathbb{R}^{3+F}$, each point has three channels of xyz and $F$ properties (*e.g.*, intensity, elongation). Figure 3 illustrates the SPG-aided 3D detection pipeline. SPG takes raw point cloud $PC_{raw}$ as input and generates a set of semantic points in the predicted foreground regions. Then, these semantic points are combined with the original point cloud into an augmented point cloud $PC_{aug}$, which is fed into a point cloud detector to obtain object detection results.

As shown in Figure 4, SPG voxelizes $PC_{raw}$ into an evenly spaced 3D voxel grid, and learns the point cloud semantics for these voxels. For each voxel, the network predicts the probability confidence $\tilde{P}^f$ of it being a foreground voxel (contained in a foreground object bounding box). In each foreground voxel, the network generates a **semantic point** $\tilde{sp}$ with point features $\tilde{\psi} = [\tilde{\chi}, \tilde{f}]$. $\tilde{\chi} \in \mathbb{R}^3$ is the xyz coordinate of $\tilde{sp}$ and $\tilde{f} \in \mathbb{R}^F$ is the point properties.

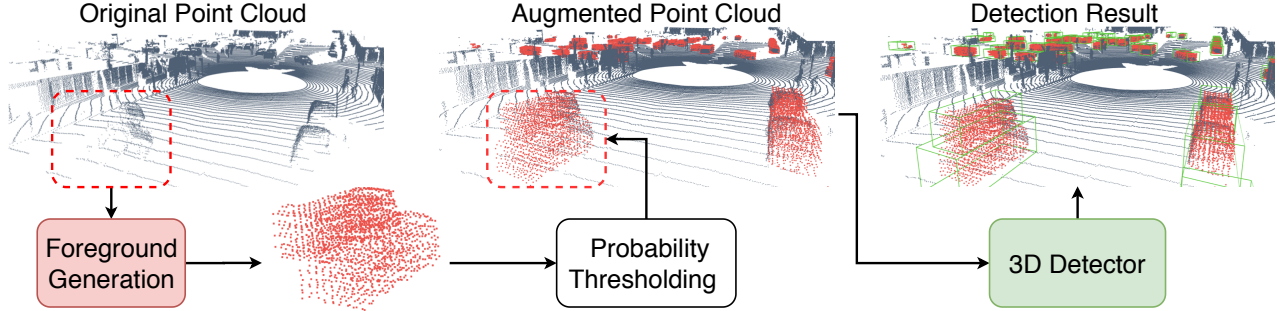To faithfully recover the foreground regions of the ob-

Figure 3: Illustration of SPG-aided 3D detection pipeline. SPG voxelizes the entire point cloud and generates prediction for each voxel (both occupied and empty) within the generation areas. After applying a probability thresholding, we take the top voxels with highest foreground probability and add a semantic point (red) at the predicted location in each of these voxels. These points are merged with the original point cloud and fed into the selected 3D point cloud detector.

served objects, we define a **generation area**. Only voxels occupied or neighbored by the observed points are considered within the generation area. We also filter out semantic points with $\tilde{P}^f$ less than $P_{thresh}$, then take $K$ semantic points $\{\tilde{s}p_1, \tilde{s}p_2, ..., \tilde{s}p_K\}$ with the highest $\tilde{P}^f$ and merge them with the original point cloud $PC_{raw}$ to get $PC_{aug}$. In practice, we use $P_{thresh} = 0.5$.

To enable SPG to be directly used by modern LiDAR-based detectors, we encode the augmented point cloud $PC_{aug}$ as $\{\hat{p}_1, \hat{p}_2, ..., \hat{p}_N, \tilde{s}p_1, \tilde{s}p_2, ..., \tilde{s}p_K\} \in \mathbb{R}^{3+F+1}$. Here we add another property channel to each point, indicating the confidence in foreground prediction: $\tilde{P}^f$ is used for the semantic points, and 1.0 for the original raw points.

### 3.1. Training Targets

To train SPG, we need to create two supervisions: 1) $y^f$, the class label if a voxel (either occupied or empty) is a foreground voxel, which supervises $\tilde{P}^f$; 2) $\psi \in \mathbb{R}^{3+F}$, the regression target for semantic point features $\tilde{\psi}$.

As visualized in Figure 4, we mark a point as a foreground point if it is inside an object bounding box. Voxels contained in a foreground bounding box are marked as foreground voxels $V^f$. For voxel $v_i$, we assign $y_i^f = 1$ if $v_i \in V^f$ and $y_i^f = 0$ otherwise. If $v_i$ is an occupied foreground voxel, we set $\psi_i = [\bar{\chi}_i, \bar{f}_i]$ as the regression target, where $\bar{\chi}_i \in \mathbb{R}^3$ is the centroid (xyz) of all foreground points in $v_i$ while $\bar{f}_i \in \mathbb{R}^F$ is the mean of their point properties (e.g. intensity, elongation).

### 3.2. Model Structure

The lower part of Figure 4 illustrates the network architecture. SPG uses a light-weight encoder-decoder network [79, 25], which is composed of three modules:
1) The Voxel Feature Encoding module [79] aggregates points inside each voxel by using several MLPs. Similar to [25, 49], these voxel features are later stacked into pillars and projected onto a birds-eye view feature space;
2) The Information Propagation module applies 2D convolutions on the pillar features. As shown in Figure 4, the

semantic information in the occupied pillars (dark green) is populated into the neighboring empty pillars (light green), which enables SPG to recover the foreground regions in the empty space.
3. The Point Generation module maps the pillar features to the corresponding voxels. For each voxel $v_i$ in the generation area, the module creates a semantic point $\tilde{s}p_i$ with encoding $[\tilde{\chi}_i, \tilde{f}_i, \tilde{P}_i^f]$, in which $\tilde{\chi}_i$ is the point location, $\tilde{f}_i$ is the point properties, and $\tilde{P}_i^f$ is the foreground probability.

### 3.3. Foreground Region Recovery

The above pipeline supervises SPG to generate semantic points in the occupied voxels. However, it is also crucial to recover the empty voxels caused by the "missing points" problem. To generate semantic points in the empty areas, SPG employs two strategies:
- "Hide and Predict", which produces the "missing points" on the source domain during training and guides SPG to recover the foreground object shape in the empty space.
- "Semantic Area Expansion", which leverages the foreground/background voxel labels derived from the bounding boxes and encourages SPG to recover more unobserved foreground regions in each bounding box.

#### 3.3.1 Hide and Predict

SPG voxelizes $PC_{raw} \in \mathbb{R}^{3+F}$ into a voxel set $V = \{v_1, v_2, ..., v_M\}$. Before passing $V$ to the network, we randomly select $\gamma\%$ of the occupied voxels $V_{hide} \subset V$ and hide all their points. During training, SPG is required to predict the foreground/background label $y^f$ for all voxels in $V$, even though it only observes points in $|V - V_{hide}|$. The predicted point features $\tilde{\psi}$ in $V_{hide}^f$ should match the corresponding ground-truth $\psi$ calculated by these hidden points.

This strategy brings two benefits: 1. Hiding points region by region mimics the missing point pattern in the target domain; 2. The strategy naturally creates the training targets for semantic points in the empty space. Section 4.4 shows the effectiveness of this strategy. Here we set $\gamma = 25$.
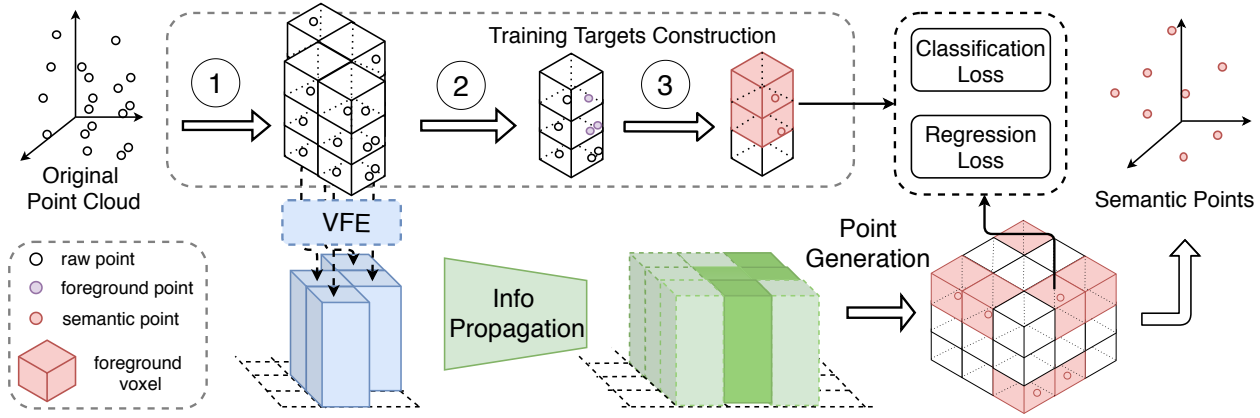
Figure 4: Training targets construction and SPG model architecture. Three steps to create the semantic point training targets: 1.Voxelization; 2. Foreground points searching 3. Label assignment and ground-truth point feature calculation. SPG includes: the Voxel Feature Encoding module (VFE), the Information Propagation module, and the Point Generation module.
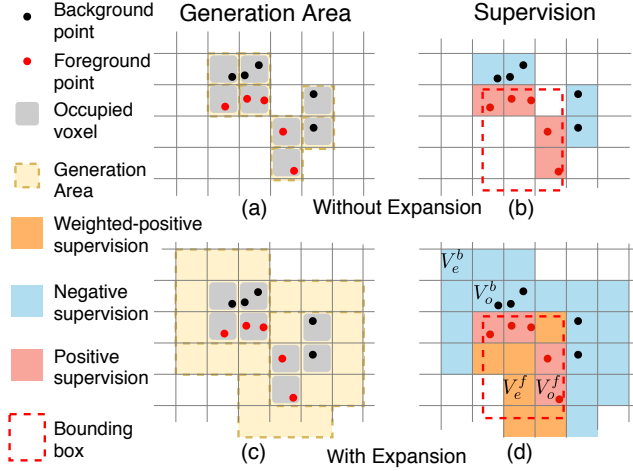


Figure 5: Visualization of "Semantic Area Expansion". (a) and (c) show the occupied voxels and the generation area, respectively. (b) and (d) show the supervision strategies.

### 3.3.2 Semantic Area Expansion

In section 1.1, we find the poor point cloud quality leads to insufficient points on each object and substantially degrades the detection performance. To remedy this problem, we allow SPG to expand the generation area to the empty space. Figure 5 a and c show the examples of the generation area with and without the expansion, respectively.

Without the expansion, we can use the ground-truth knowledge of foreground points to supervise SPG only on the occupied voxels (Figure 5 b). However, with the expansion, there is no foreground point inside these empty voxels. Therefore, as shown in Figure 5 d, we design a supervision scheme as follows:
1. For both occupied and empty background voxels $V_o^b$ and $V_e^b$, we impose negative supervision and set label $y^f = 0$.
2. For the occupied foreground voxels $V_o^f$, we set $y^f = 1$.
3. For the empty voxels inside a bounding box $V_e^f$, we set their foreground label $y^f = 1$ and assign a weighting factor



(a) Without expansion      (b) With expansion

Figure 6: Comparisons between generated semantic points (red) with and without "Semantic Area Expansion".

$\alpha$, where $\alpha < 1$.
4. We only impose point features supervision $\psi$ at occupied foreground voxels $V_o^f$.

To investigate the effectiveness of the expansion, we train a model on the OD training set and evaluate it on the Kirk validation set. The expansion results in **510%** more semantic points on foreground objects, which mitigates the "missing points" problem caused by environmental interference and occlusions. Figure 6 shows the generation results with and without the expansion. The supervision scheme encourages SPG to learn the extended shape of vehicle parts and enables SPG to fill in more foreground space with semantic points. We also conduct ablation studies (Section 4.4) to show the effectiveness of the proposed strategy.

### 3.4. Objectives

We use two loss functions, *i.e.,* foreground area classification loss $L_{cls}$ and feature regression loss $L_{reg}$.

To supervise $\tilde{P}^f$ with label $y^f$, we use Focal loss [31] to mitigate the background-foreground class imbalance. $L_{cls}$ can be decomposed as focal losses on four categories of voxels: the occupied voxels $V_o$, the empty background voxels $V_e^b$, the empty foreground voxels $V_e^f$ and the hidden voxels $V_{hide}$. The labeling strategy for these categories is described in Section 3.3.2.

$$L_{cls} = \frac{1}{|V_o \cup V_e^b|} \sum_{V_o \cup V_e^b} L_{focal}$$
$$+ \frac{\alpha}{|V_e^f|} \sum_{V_e^f} L_{focal} + \frac{\beta}{|V_{hide}|} \sum_{V_{hide}} L_{focal} \quad (1)$$

We use Smooth-L1 loss [20] for point feature $\tilde{\psi}$ regression, and supervise on the semantic points in occupied foreground voxels $V_o^f$ and the hidden foreground voxels $V_{hide}^f$.

$$L_{reg} = \frac{1}{|V_o^f|} \sum_{V_o^f} L_{smooth-L1}(\tilde{\psi}, \psi)$$
$$+ \frac{\beta}{|V_{hide}^f|} \sum_{V_{hide}^f} L_{smooth-L1}(\tilde{\psi}, \psi) \quad (2)$$

Please note that we are only interested in the $L_{cls}$ and $L_{reg}$ on voxels inside the generation area. We find $\alpha = 0.5$ and $\beta = 2.0$ achieves the best result.

## 4. Experiments

In this section, we first evaluate the effectiveness of SPG as a general UDA approach for 3D detection, based on the Waymo Domain Adaptation Dataset [54]. In addition, we show that SPG can also improve results for top-performing 3D detectors on the source domain[54, 18]. To demonstrate the wide applicability of SPG, we choose two representative detectors: 1) PointPillars [25], popular among industrial-grade autonomous driving systems; 2) PV-RCNN [49], a high performance LiDAR-based 3D detector [18, 54]. We perform two groups of model comparisons under the setting of unsupervised domain adaptation (UDA) and general 3D object detection: group 1, PointPillars vs. SPG + Point-Pillars; group 2, PV-RCNN vs. SPG + PV-RCNN. SPG can also be combined with range image-based detectors [36, 78, 1] by applying ray casting to the generated points. However, we leave this as future work.

**Datasets** The Waymo Domain Adaptation dataset 1.0 [54] consists of two sub datasets, the *Waymo Open Dataset* (OD) and the *Waymo Kirkland Dataset* (Kirk). OD provides 798 training segments of 158,361 LiDAR frames and 202 validation segments of 40,077 frames. Captured across California and Arizona, 99.40% of its frames have dry weather. Kirk is a smaller dataset including 80 training segments of 15,797 frames and 20 validation segments of 3,933 frames. Captured in Kirkland, 97.99% its LiDAR frames have rainy weather. To examine a detector's reliability when entering a new environment, we conduct UDA experiments without using the data in Kirk during training.

*KITTI* [18] contains 7481 training samples and 7518 testing samples. Following [8], we divide the training data into a *train* split and a *val* split containing 3721 and 3769 LiDAR frames, respectively.

**Implementation and Training Details** We use a single lightweight network architecture on all experiments. As

shown in Figure 4, our Voxel Feature Encoding[79] module includes a single layer point-wise MLP and a voxel-wise max-pooling [43, 79]. The Information Propagation module includes two levels of CNN layers. The first level includes three CNN layers with stride 1. The second level includes one CNN layer with stride 2 and four subsequent CNN layers with stride 1, then up-sampled back to the original resolution. Each layer has an output dimension of 128. From the BEV feature map, the Point Generation module uses one FC layer to produce $\tilde{P}^f$ and another FC layer to generate the features $\tilde{\psi}$ for the voxels in each pillar. SPG and each detector are trained separately.

We implement PointPillars following [25] and use the PV-RCNN code provided by [49] (the training settings on OD 1.0 are obtained via direct communication with the author). On the Waymo Domain Adaptation Dataset [54], we set the voxel dimensions to (0.32m, 0.32m, 0.4m) for Point-Pillars and (0.2m, 0.2m, 0.3m) for PV-RCNN. On KITTI, we set the voxel dimensions to (0.16m, 0.16m, 0.2m) and (0.2m, 0.2m, 0.3m) for PointPillars and PV-RCNN, respectively. By default, the generation area includes voxels within 6 steps of any occupied voxel. After probability thresholding, we preserve up to 8000 semantic points for the Waymo Domain Adaptation Dataset and 6000 for KITTI.

### 4.1. Evaluation on the Waymo Open Dataset

We perform two groups of model comparisons by training them on the OD training set and evaluating them on both the OD validation set and the Kirk validation set.

**Evaluation Metrics** The Kirk 1.0 validation set only provides the evaluation labels for the vehicle and the pedestrian classes. We use the official evaluation tool released by [54]. The IoU thresholds for vehicles and pedestrians are 0.7 and 0.5. In Table 2 we report both 3D and BEV AP on two difficulty levels. More results with distance breakdown are shown in the supplemental material.

**Target Domain** On Kirk, we observe that SPG brings remarkable improvements over both detectors across all object types. Averaged over two difficulty levels, SPG improves PointPillars on Kirk vehicle 3D AP by 6.7% and BEV AP by 8.8%. For PV-RCNN, SPG improves Kirk pedestrian 3D AP by 5.6% and BEV AP by 5.7%.

**Source Domain** Unlike most UDA methods [9, 21, 48] that only optimize the performance on the target domain, SPG also consistently improves the results on the source domain. Averaged across both difficulty levels, SPG improves OD vehicle 3D AP for PointPillars by 5.4% and improves OD pedestrian 3D AP for PV-RCNN by 1.6%.

**Comparison with Alternative Strategies** We compare SPG with alternative strategies that also target the deteriorating point cloud quality. We employ PointPillars as the baseline and choose LEVEL_1 vehicle 3D AP as the main

| Difficulty | Method | Target Domain - Kirk | | | | Source Domain - OD | | | |
| | | Vehicle | | Pedestrian | | Vehicle | | Pedestrian | |
| | | 3D AP | BEV AP | 3D AP | BEV AP | 3D AP | BEV AP | 3D AP | BEV AP |
|---|---|---|---|---|---|---|---|---|---|
| LEVEL_1 | PointPillars | 34.65 | 51.88 | 20.65 | 22.33 | 57.27 | 72.26 | 55.20 | 63.82 |
| | SPG + PointPillars | **41.56** | **60.44** | **23.72** | **24.83** | **62.44** | **77.63** | **56.06** | **64.66** |
| | *Improvement* | *+6.91* | *+8.56* | *+3.07* | *+2.50* | *+5.17* | *+5.37* | *+0.86* | *+0.84* |
| LEVEL_2 | PointPillars | 31.67 | 47.93 | 17.66 | 18.40 | 52.96 | 69.09 | 51.33 | 60.13 |
| | SPG + PointPillars | **38.15** | **56.94** | **19.57** | **20.67** | **58.54** | **74.90** | **52.33** | **60.93** |
| | *Improvement* | *+6.48* | *+9.01* | *+1.91* | *+2.27* | *+5.58* | *+5.81* | *+1.00* | *+0.80* |
| LEVEL_1 | PV-RCNN | 55.16 | 70.38 | 24.47 | 25.39 | 74.01 | 85.13 | 65.34 | 70.35 |
| | SPG + PV-RCNN | **58.31** | **72.56** | **30.82** | **31.92** | **75.27** | **87.38** | **66.93** | **70.37** |
| | *Improvement* | *+3.15* | *+2.18* | *+6.35* | *+6.53* | *+1.26* | *+2.25* | *+1.59* | *+0.02* |
| LEVEL_2 | PV-RCNN | 45.81 | 60.13 | 17.16 | 17.88 | 64.69 | 76.84 | 56.03 | 60.81 |
| | SPG + PV-RCNN | **48.70** | **62.03** | **22.05** | **22.65** | **65.98** | **78.05** | **57.68** | **60.88** |
| | *Improvement* | *+2.89* | *+1.90* | *+4.89* | *+4.77* | *+1.29* | *+1.21* | *+1.65* | *+0.07* |

Table 2: Results on the Waymo Open Dataset 1.0 and the Kirkland Dataset. Results for PointPillars are based on our own implementation following [25]. We use the PV-RCNN source code and obtain training settings for the Waymo Open Dataset [54] via direct communication with the author.

metric on the Kirk validation set, during UDA. Three strategies are implemented: 1. RndDrop, where we randomly drop 17% of the points in the source domain during training. This dropout ratio is chosen for the number of points in the source and target domain to match (see Table 1). 2. K-frames, where we use $K$ consecutive historical frames in both the source domain and the target domain. The points in the first $K - 1$ are transformed into the last frame according to the ground-truth ego-motion, so that the last frame has $K$ times the number of points. 3. Adversarial Domain Adaptation (ADA), where we follow [17] and add a domain classification loss on the pillar features of PointPillars.

As shown in Table 3, although "RndDrop" enforces the quantity of missing points in the source domain to match with that in the target domain, the pattern of missing points still differs from the reality (see Figure 2), which limits the improvement to only 0.80% in 3D AP. To remedy the "missing points" problem, "3-frames" contains real points from 3 frames and "5-frames" contains points from 5 frames. With around 800K points per scene, "5-frames" significantly improves the single-frame baseline. However, aggregating multiple frames inevitably increases the memory usage and the processing time. ADA improves 3D AP to 36.34 on the target domain, but we observe an AP drop of 1.52 in the source domain. Remarkably, SPG can outperform "5-frames", by adding only 8000 semantic points, which is less than 6% of the points in a single frame.

| Method | Baseline | RndDrop | 3-frames | 5-frames | ADA | SPG |
|---|---|---|---|---|---|---|
| 3D AP | 34.65 | 35.45 | 38.00 | 38.51 | 36.34 | **41.56** |

Table 3: Comparisons of different strategies targeting at the deteriorating point cloud quality. The models are trained on OD and evaluated on Kirk. The metric is LEVEL_1 Vehicle 3D AP. We use PointPillars[25] as the baseline.

| Method | Reference | Car - 3D AP | | | |
| | | Easy | Mod. | Hard | Avg. |
|---|---|---|---|---|---|
| SA-SSD[19] | CVPR 2020 | 88.75 | 79.79 | 74.16 | 80.90 |
| 3D-CVF[72] | ECCV 2020 | 89.20 | 80.05 | 73.11 | 80.79 |
| CIA-SSD[60] | AAAI 2021 | 89.59 | 80.28 | 72.87 | 80.91 |
| Asso-3Ddet[15] | CVPR 2020 | 85.99 | 77.40 | 70.53 | 77.97 |
| Voxel R-CNN[13] | AAAI 2021 | **90.90** | 81.62 | 77.06 | 83.19 |
| PV-RCNN[49] | CVPR 2020 | 90.25 | 81.43 | 76.82 | 82.83 |
| **SPG**+PV-RCNN | - | 90.50 | **82.13** | **78.90** | **83.84** |

Table 4: Car detection Results on the KITTI test set. See the full list of comparisons in the supplemental.

## 4.2. Evaluation on the KITTI Dataset

In this section, we show besides the usefulness in UDA (Sec. 4.1) the proposed SPG can also boost performance in another popular 3D detection benchmark (i.e. KITTI [18]). We follow the training and evaluation protocols in [25, 49].

**KITTI Test Set** As shown in Table 4, SPG significantly improves PV-RCNN on Car 3D detection. As of Mar. 3rd, 2021, our method ranks the **1st** on KITTI car 3D detection among all published methods (4th among all submitted approaches). Moreover, SPG demonstrates strong robustness in detecting hard objects (truncation up to 50%). Specifically, SPG surpasses all submitted methods on the hard category by a big margin and achieves the **highest** overall 3D AP of 83.84% (averaged over Easy, Mod. and Hard).

**KITTI Validation Set** We summarize the results in Table 5. We train each group of models using the recommended settings of baseline detectors [25, 49].

SPG remarkably improves both PointPillars and PV-RCNN on all object types and difficulty levels. Specifically, for PointPillars, SPG improves the 3D AP of car detection by 2.02%, 2.97%, 3.67% on easy, moderate, and hard levels, respectively. For PV-RCNN, SPG improves the 3D AP

| Method | Car - 3D AP | | | Car - BEV AP | | | Pedestrian - 3D AP | | | Pedestrian - BEV AP | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Easy | Mod. | Hard | Easy | Mod. | Hard | Easy | Mod. | Hard | Easy | Mod. | Hard |
| PointPillars | 87.75 | 78.39 | 75.18 | 92.03 | 88.05 | 86.66 | 57.30 | 51.41 | 46.87 | 61.59 | 56.01 | 52.04 |
| SPG + PointPillars | **89.77** | **81.36** | **78.85** | **94.38** | **89.92** | **87.97** | **59.65** | **53.55** | **49.24** | **65.38** | **59.48** | **55.32** |
| *Improvement* | *+2.02* | *+2.97* | *+3.67* | *+2.35* | *+1.87* | *+1.31* | *+2.35* | *+2.14* | *+2.47* | *+3.79* | *+3.47* | *+3.28* |
| PV-RCNN | 92.10 | 84.36 | 82.48 | 93.02 | 90.33 | 88.53 | 64.26 | 56.67 | 51.91 | 67.97 | 60.52 | 55.80 |
| SPG + PV-RCNN | **92.53** | **85.31** | **82.82** | **94.99** | **91.11** | **88.86** | **69.66** | **61.80** | **56.39** | **71.79** | **64.50** | **59.51** |
| *Improvement* | *+0.43* | *+0.95* | *+0.34* | *+1.97* | *+0.78* | *+0.33* | *+5.40* | *+5.13* | *+4.48* | *+3.82* | *+3.98* | *+3.71* |

Table 5: Comparisons on the KITTI validation set. Average Precision (AP) is computed over 40 recall positions. The baseline results[49, 55] are obtained based on publically released models. See more results (including Cyclist) in the supplemental.

of pedestrian detection by 5.40%, 5.13%, 4.48% on easy, moderate and hard levels, respectively.

### 4.3. Model Efficiency

We evaluate the efficiency of SPG on the KITTI *val* split (Table 6). SPG contains 0.39 million parameters while adding less than 17 milliseconds latency to the detectors. This indicates that SPG is highly efficient for industrial-grade deployment on a stringent computation budget.

| Detectors | PointPillars | | PV-RCNN | | - |
|---|---|---|---|---|---|
| With SPG | No | Yes | No | Yes | Yes |
| Latency (ms) | 23.56 | 36.67 | 139.96 | 156.85 | 16.82 |
| Parameters | 4.83M | 5.22M | 13.12M | 13.51M | 0.39M |

Table 6: Latency and model parameters. "M" stands for million. The last column shows the results of standalone SPG. The evaluation is based on a 1080Ti GPU with a batch size of 1. The latency is averaged over the KITTI *val* split.

### 4.4. Ablation Studies

| Model | Expansion | Hide & Predict | Foreground Confidence | 3D AP | *Improve* |
|---|---|---|---|---|---|
| Baseline | — | — | — | 34.65 | — |
| SPG | — | — | ✓ | 35.89 | *+1.24* |
| SPG | — | 25% | ✓ | 38.09 | *+3.44* |
| SPG | ✓($\alpha$=0.0) | 25% | ✓ | 38.96 | *+4.31* |
| SPG | ✓($\alpha$=1.0) | 25% | ✓ | 38.42 | *+3.77* |
| SPG | ✓($\alpha$=0.5) | — | ✓ | 39.22 | *+4.57* |
| SPG | ✓($\alpha$=0.5) | 25% | — | 37.96 | *+3.31* |
| SPG(ours) | ✓($\alpha$=0.5) | 25% | ✓ | **41.56** | ***+6.91*** |

Table 7: Ablation studies of SPG. The models are trained on OD and evaluated on Kirk. The metric is LEVEL_1 Vehicle 3D AP. We use PointPillars[25] as our baseline.

We conduct ablation studies on "Semantic Area Expansion", "Hide and Predict" and whether to add foreground confidence ($\tilde{P}^f$) as a point property and show all of them can benefit detection quality (see Table 7). We also change the weighting factor $\alpha$ on the empty foreground voxels $V_e^f$. A larger $\alpha$ encourages more point generation in the empty foreground space. However, in reality, an object typically does not occupy the entire space within a bounding

| $P_{thresh}$ | 0.3 | 0.4 | 0.5 | 0.6 | 0.7 |
|---|---|---|---|---|---|
| 3D AP | 39.39 | 40.09 | **41.56** | 41.18 | 40.89 |

Table 8: Ablation studies on the probability threshold $P_{thresh}$ (only keep the semantic point if $\tilde{P}^f > P_{thresh}$). Our best SPG model uses $P_{thresh} = 0.5$. The metric is LEVEL_1 Vehicle 3D AP on the Kirk validation set.

box. Therefore, over-aggressively generating points does not help improve the performance (see $\alpha = 1.0$).

**Probability Thresholding** In Table 8, we show the effect of choosing different thresholds during probability thresholding. While a higher $P_{thresh}$ only keeps semantic points with high foreground probability, a lower $P_{thresh}$ admits more points, but may introduce points to the background. We find the threshold of 0.5 achieves the best results.

## 5. Conclusions

In this paper, we investigate unsupervised domain adaptation for LiDAR-based 3D detectors across different geographic locations and weather conditions. We observe that rainy weather can severely deteriorate the point cloud quality and cause drastic performance drop for modern 3D detectors, based on the Waymo Domain Adaptation dataset. The proposed SPG method addresses this issue as a novel unsupervised domain adaptation (UDA) task without using any training data from the new domain. This setting allows us to rigorously test 3D detectors against real-world challenges autonomous vehicles may experience due to diverse conditions (e.g., different levels of fog/rain/snow beyond what one may effectively train for) during the trip.

Utilizing two strategies "Hide and Predict" and "Semantic Area Generation", SPG generates semantic points to recover the shape of foreground objects with a negligible overhead (only adding 6% extra points) and can be conveniently integrated with modern LiDAR-based detectors. We test SPG with two detectors: PointPillars and PV-RCNN. For unsupervised domain adaptation, SPG achieves significant performance gains on the challenging target domain. On Waymo Open dataset and KITTI, SPG also consistently benefits detection quality on the source domain.

## 6. Acknowledgement

We would like to thank Boqing Gong for the helpful discussions. We also thank Jingwei Ji for the careful proofreading.

## References

[1] Alex Bewley, Pei Sun, Thomas Mensink, Drago Anguelov, and Cristian Sminchisescu. Range conditioned dilated convolutions for scale invariant 3d object detection. In *Conference on Robot Learning*, 2020. 6

[2] Karsten M Borgwardt, Arthur Gretton, Malte J Rasch, Hans-Peter Kriegel, Bernhard Schölkopf, and Alex J Smola. Integrating structured biological data by kernel maximum mean discrepancy. *Bioinformatics*, 22(14):e49–e57, 2006. 3

[3] Xu Cao, Weimin Wang, and Katashi Nagao. Neural style transfer for point clouds. *arXiv preprint arXiv:1903.05807*, 2019. 3

[4] Xu Cao, Weimin Wang, Katashi Nagao, and Ryosuke Nakamura. Psnet: A style transfer network for point cloud stylization on geometry and color. In *The IEEE Winter Conference on Applications of Computer Vision*, pages 3337–3345, 2020. 3

[5] Qi Chen, Lin Sun, Zhixin Wang, Kui Jia, and Alan Yuille. Object as hotspots: An anchor-free 3d object detection approach via firing of hotspots. In *European Conference on Computer Vision*, pages 68–84. Springer, 2020. 18

[6] Xuelin Chen, Baoquan Chen, and Niloy J Mitra. Unpaired point cloud completion on real scans using adversarial training. In *ICLR*, 2020. 3

[7] Xiaozhi Chen, Huimin Ma, Ji Wan, Bo Li, and Tian Xia. Multi-view 3d object detection network for autonomous driving. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1907–1915, 2017. 1

[8] Xiaozhi Chen, Huimin Ma, Ji Wan, Bo Li, and Tian Xia. Multi-view 3d object detection network for autonomous driving. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6526–6534, 2017. 6

[9] Yuhua Chen, Wen Li, Christos Sakaridis, Dengxin Dai, and Luc Van Gool. Domain adaptive faster r-cnn for object detection in the wild. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3339–3348, 2018. 2, 3, 6

[10] Yilun Chen, Shu Liu, Xiaoyong Shen, and Jiaya Jia. Fast point r-cnn. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 9775–9784, 2019. 1, 18

[11] Yilun Chen, Shu Liu, Xiaoyong Shen, and Jiaya Jia. Dsgn: Deep stereo geometry network for 3d object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12536–12545, 2020. 1

[12] Jaehoon Choi, Taekyung Kim, and Changick Kim. Self-ensembling with gan-based data augmentation for domain adaptation in semantic segmentation. In *Proceedings of the IEEE international conference on computer vision*, pages 6830–6840, 2019. 2

[13] Jiajun Deng, Shaoshuai Shi, Peiwei Li, Wengang Zhou, Yanyong Zhang, and Houqiang Li. Voxel r-cnn: Towards high performance voxel-based 3d object detection. *arXiv:2012.15712*, 2020. 7, 18

[14] Jiahua Dong, Yang Cong, Gan Sun, and Dongdong Hou. Semantic-transferable weakly-supervised endoscopic lesions segmentation. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 10712–10721, 2019. 2, 3

[15] Liang Du, Xiaoqing Ye, Xiao Tan, Jianfeng Feng, Zhenbo Xu, Errui Ding, and Shilei Wen. Associate-3ddet: Perceptual-to-conceptual association for 3d point cloud object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13329–13338, 2020. 3, 7, 18

[16] Xinxin Du, Marcelo H Ang, Sertac Karaman, and Daniela Rus. A general pipeline for 3d detection of vehicles. In *2018 IEEE International Conference on Robotics and Automation (ICRA)*, pages 3194–3200. IEEE, 2018. 1

[17] Yaroslav Ganin and Victor Lempitsky. Unsupervised domain adaptation by backpropagation. In *International conference on machine learning*, pages 1180–1189. PMLR, 2015. 2, 3, 7

[18] Andreas Geiger, Philip Lenz, Christoph Stiller, and Raquel Urtasun. Vision meets robotics: The kitti dataset. *The International Journal of Robotics Research*, 32(11):1231–1237, 2013. 1, 3, 6, 7

[19] Chenhang He, Hui Zeng, Jianqiang Huang, Xian-Sheng Hua, and Lei Zhang. Structure aware single-stage 3d object detection from point cloud. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2020. 2, 7, 18

[20] Zhenwei He and Lei Zhang. Multi-adversarial faster-rcnn for unrestricted object detection. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 6668–6677, 2019. 2, 3, 6

[21] Han-Kai Hsu, Chun-Han Yao, Yi-Hsuan Tsai, Wei-Chih Hung, Hung-Yu Tseng, Maneesh Singh, and Ming-Hsuan Yang. Progressive domain adaptation for object detection. In *The IEEE Winter Conference on Applications of Computer Vision*, pages 749–757, 2020. 2, 3, 6

[22] Tengteng Huang, Zhe Liu, Xiwu Chen, and Xiang Bai. Epnet: Enhancing point features with image semantics for 3d object detection. In *European Conference on Computer Vision*, pages 35–52. Springer, 2020. 18

[23] Hendrik Königshof, Niels Ole Salscheider, and Christoph Stiller. Realtime 3d object detection for automated driving using stereo vision and semantic information. In *2019 IEEE Intelligent Transportation Systems Conference (ITSC)*, pages 1405–1410. IEEE, 2019. 1

[24] Jason Ku, Melissa Mozifian, Jungwook Lee, Ali Harakeh, and Steven L Waslander. Joint 3d proposal generation and object detection from view aggregation. In *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 1–8. IEEE, 2018. 18

[25] Alex H Lang, Sourabh Vora, Holger Caesar, Lubing Zhou, Jiong Yang, and Oscar Beijbom. Pointpillars: Fast encoders

for object detection from point clouds. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 12697–12705, 2019. 1, 2, 3, 4, 6, 7, 8, 16, 17, 18

[26] Buyu Li, Wanli Ouyang, Lu Sheng, Xingyu Zeng, and Xiaogang Wang. Gs3d: An efficient 3d object detection framework for autonomous driving. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1019–1028, 2019. 1

[27] Peiliang Li, Xiaozhi Chen, and Shaojie Shen. Stereo r-cnn based 3d object detection for autonomous driving. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7644–7652, 2019. 1

[28] Ruihui Li, Xianzhi Li, Chi-Wing Fu, Daniel Cohen-Or, and Pheng-Ann Heng. Pu-gan: a point cloud upsampling adversarial network. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 7203–7212, 2019. 2, 3

[29] Ming Liang, Bin Yang, Yun Chen, Rui Hu, and Raquel Urtasun. Multi-task multi-sensor fusion for 3d object detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7345–7353, 2019. 1, 18

[30] Ming Liang, Bin Yang, Shenlong Wang, and Raquel Urtasun. Deep continuous fusion for multi-sensor 3d object detection. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 641–656, 2018. 1

[31] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision*, pages 2980–2988, 2017. 5

[32] Mingsheng Long, Yue Cao, Jianmin Wang, and Michael Jordan. Learning transferable features with deep adaptation networks. In *International conference on machine learning*, pages 97–105. PMLR, 2015. 3

[33] Mingsheng Long, Jianmin Wang, Guiguang Ding, Jiaguang Sun, and Philip S Yu. Transfer feature learning with joint distribution adaptation. In *Proceedings of the IEEE international conference on computer vision*, pages 2200–2207, 2013. 3

[34] Mingsheng Long, Han Zhu, Jianmin Wang, and Michael I Jordan. Unsupervised domain adaptation with residual transfer networks. In *Advances in neural information processing systems*, pages 136–144, 2016. 3

[35] Haifeng Luo, Kourosh Khoshelham, Lina Fang, and Chongcheng Chen. Unsupervised scene adaptation for semantic segmentation of urban mobile laser scanning point clouds. *ISPRS Journal of Photogrammetry and Remote Sensing*, 169:253–267, 2020. 2

[36] Gregory P Meyer, Ankit Laddha, Eric Kee, Carlos Vallespi-Gonzalez, and Carl K Wellington. Lasernet: An efficient probabilistic 3d object detector for autonomous driving. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 12677–12686, 2019. 1, 6

[37] Pietro Morerio, Jacopo Cavazza, and Vittorio Murino. Minimal-entropy correlation alignment for unsupervised deep domain adaptation. *arXiv preprint arXiv:1711.10288*, 2017. 2

[38] Sinno Jialin Pan, Ivor W Tsang, James T Kwok, and Qiang Yang. Domain adaptation via transfer component analysis. *IEEE Transactions on Neural Networks*, 22(2):199–210, 2010. 3

[39] Su Pang, Daniel Morris, and Hayder Radha. Clocs: Camera-lidar object candidates fusion for 3d object detection. *arXiv preprint arXiv:2009.00784*, 2020. 18

[40] Taesung Park, Alexei A Efros, Richard Zhang, and Jun-Yan Zhu. Contrastive learning for unpaired image-to-image translation. *arXiv preprint arXiv:2007.15651*, 2020. 2, 3

[41] Alex D Pon, Jason Ku, Chengyao Li, and Steven L Waslander. Object-centric stereo matching for 3d object detection. In *2020 IEEE International Conference on Robotics and Automation (ICRA)*, pages 8383–8389. IEEE, 2020. 1

[42] Charles R Qi, Wei Liu, Chenxia Wu, Hao Su, and Leonidas J Guibas. Frustum pointnets for 3d object detection from rgb-d data. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 918–927, 2018. 1, 18

[43] Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 652–660, 2017. 6

[44] Can Qin, Lichen Wang, Yulun Zhang, and Yun Fu. Generatively inferential co-training for unsupervised domain adaptation. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*, pages 0–0, 2019. 3

[45] Can Qin, Haoxuan You, Lichen Wang, C-C Jay Kuo, and Yun Fu. Pointdan: A multi-scale 3d domain adaption network for point cloud representation. In *Advances in Neural Information Processing Systems*, pages 7192–7203, 2019. 2, 3

[46] Kuniaki Saito, Kohei Watanabe, Yoshitaka Ushiku, and Tatsuya Harada. Maximum classifier discrepancy for unsupervised domain adaptation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3723–3732, 2018. 3

[47] Khaled Saleh, Ahmed Abobakr, Mohammed Attia, Julie Iskander, Darius Nahavandi, Mohammed Hossny, and Saeid Nahvandi. Domain adaptation for vehicle detection from bird's eye view lidar point cloud data. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*, pages 0–0, 2019. 2, 3

[48] Yuhu Shan, Wen Feng Lu, and Chee Meng Chew. Pixel and feature level based domain adaptation for object detection in autonomous driving. *Neurocomputing*, 367:31–38, 2019. 2, 3, 6

[49] Shaoshuai Shi, Chaoxu Guo, Li Jiang, Zhe Wang, Jianping Shi, Xiaogang Wang, and Hongsheng Li. Pv-rcnn: Point-voxel feature set abstraction for 3d object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10529–10538, 2020. 1, 2, 3, 4, 6, 7, 8, 16, 17, 18

[50] Shaoshuai Shi, Xiaogang Wang, and Hongsheng Li. Pointr-cnn: 3d object proposal generation and detection from point cloud. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 770–779, 2019. 1, 2, 18

[51] Shaoshuai Shi, Zhe Wang, Jianping Shi, Xiaogang Wang, and Hongsheng Li. From points to parts: 3d object detection from point cloud with part-aware and part-aggregation network. *arXiv preprint arXiv:1907.03670*, 2019. 1

[52] S. Shi, Z. Wang, J. Shi, X. Wang, and H. Li. From points to parts: 3d object detection from point cloud with part-aware and part-aggregation network. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pages 1–1, 2020. 18

[53] Weijing Shi and Raj Rajkumar. Point-gnn: Graph neural network for 3d object detection in a point cloud. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1711–1719, 2020. 1, 18

[54] Pei Sun, Henrik Kretzschmar, Xerxes Dotiwalla, Aurelien Chouard, Vijaysai Patnaik, Paul Tsui, James Guo, Yin Zhou, Yuning Chai, Benjamin Caine, Vijay Vasudevan, Wei Han, Jiquan Ngiam, Hang Zhao, Aleksei Timofeev, Scott Ettinger, Maxim Krivokon, Amy Gao, Aditya Joshi, Yu Zhang, Jonathon Shlens, Zhifeng Chen, and Dragomir Anguelov. Scalability in perception for autonomous driving: Waymo open dataset, 2019. 1, 2, 3, 6, 7, 14

[55] OpenPCDet Development Team. Openpcdet: An opensource toolbox for 3d object detection from point clouds. https://github.com/open-mmlab/OpenPCDet, 2020. 8, 16, 17

[56] Eric Tzeng, Judy Hoffman, Kate Saenko, and Trevor Darrell. Adversarial discriminative domain adaptation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7167–7176, 2017. 2, 3

[57] Yan Wang, Xiangyu Chen, Yurong You, Li Erran Li, Bharath Hariharan, Mark Campbell, Kilian Q Weinberger, and Wei-Lun 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*, pages 11713–11723, 2020. 3

[58] Zhixin Wang and Kui Jia. Frustum convnet: Sliding frustums to aggregate local point-wise features for amodal 3d object detection. In *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 1742–1749. IEEE, 2019. 18

[59] Bichen Wu, Xuanyu Zhou, Sicheng Zhao, Xiangyu Yue, and Kurt 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)*, pages 4376–4382. IEEE, 2019. 3

[60] Sijin Chen Li Jiang Chi-Wing Fu Wu Zheng, Weiliang Tang. Cia-ssd: Confident iou-aware single-stage object detector from point cloud. In *AAAI*, 2021. 7, 18

[61] Haozhe Xie, Hongxun Yao, Shangchen Zhou, Jiageng Mao, Shengping Zhang, and Wenxiu Sun. Grnet: Gridding residual network for dense point cloud completion. *arXiv preprint arXiv:2006.03761*, 2020. 3

[62] Qiangeng Xu, Xudong Sun, Cho-Ying Wu, Panqu Wang, and Ulrich Neumann. Grid-gcn for fast and scalable point cloud learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5661–5670, 2020. 1

[63] Zhenbo Xu, Wei Zhang, Xiaoqing Ye, Xiao Tan, Wei Yang, Shilei Wen, Errui Ding, Ajin Meng, and Liusheng Huang. Zoomnet: Part-aware adaptive zooming neural network for 3d object detection. In *AAAI*, pages 12557–12564, 2020. 1

[64] Yan Yan, Yuxing Mao, and Bo Li. Second: Sparsely embedded convolutional detection. *Sensors*, 18(10):3337, 2018. 1, 18

[65] Bin Yang, Wenjie Luo, and Raquel Urtasun. Pixor: Realtime 3d object detection from point clouds. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 7652–7660, 2018. 1

[66] Yaoqing Yang, Chen Feng, Yiru Shen, and Dong Tian. Foldingnet: Point cloud auto-encoder via deep grid deformation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 206–215, 2018. 3

[67] Zetong Yang, Yanan Sun, Shu Liu, and Jiaya Jia. 3dssd: Point-based 3d single stage object detector. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11040–11048, 2020. 1, 18

[68] Zetong Yang, Yanan Sun, Shu Liu, Xiaoyong Shen, and Jiaya Jia. Std: Sparse-to-dense 3d object detector for point cloud. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1951–1960, 2019. 1, 18

[69] Yangyang Ye, Houjin Chen, Chi Zhang, Xiaoli Hao, and Zhaoxiang Zhang. Sarpnet: Shape attention regional proposal network for lidar-based 3d object detection. *Neurocomputing*, 379:53–63, 2020. 18

[70] Hongwei Yi, Shaoshuai Shi, Mingyu Ding, Jiankai Sun, Kui Xu, Hui Zhou, Zhe Wang, Sheng Li, and Guoping Wang. Segvoxelnet: Exploring semantic context and depth-aware features for 3d vehicle detection from point cloud. In *2020 IEEE International Conference on Robotics and Automation (ICRA)*, pages 2274–2280. IEEE, 2020. 18

[71] Wang Yifan, Shihao Wu, Hui Huang, Daniel Cohen-Or, and Olga Sorkine-Hornung. Patch-based progressive 3d point set upsampling. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5958–5967, 2019. 2, 3

[72] Jin Hyeok Yoo, Yecheol Kim, Ji Song Kim, and Jun Won Choi. 3d-cvf: Generating joint camera and lidar features using cross-view spatial feature fusion for 3d object detection. *arXiv preprint arXiv:2004.12636*, 3, 2020. 7, 18

[73] Lequan Yu, Xianzhi Li, Chi-Wing Fu, Daniel Cohen-Or, and Pheng-Ann Heng. Pu-net: Point cloud upsampling network. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2790–2799, 2018. 2, 3

[74] Wentao Yuan, Tejas Khot, David Held, Christoph Mertz, and Martial Hebert. Pcn: Point completion network. In *2018 International Conference on 3D Vision (3DV)*, pages 728–737. IEEE, 2018. 3

[75] Dingfu Zhou, Jin Fang, Xibin Song, Chenye Guan, Junbo Yin, Yuchao Dai, and Ruigang Yang. Iou loss for 2d/3d object detection. In *2019 International Conference on 3D Vision (3DV)*, pages 85–94. IEEE, 2019. 18

[76] Dingfu Zhou, Jin Fang, Xibin Song, Liu Liu, Junbo Yin, Yuchao Dai, Hongdong Li, and Ruigang Yang. Joint 3d instance segmentation and object detection for autonomous

driving. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020. 18

[77] Xingyi Zhou, Arjun Karpur, Chuang Gan, Linjie Luo, and Qixing Huang. Unsupervised domain adaptation for 3d keypoint estimation via view consistency. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 137–153, 2018. 2, 3

[78] Yin Zhou, Pei Sun, Yu Zhang, Dragomir Anguelov, Jiyang Gao, Tom Ouyang, James Guo, Jiquan Ngiam, and Vijay Vasudevan. End-to-end multi-view fusion for 3d object detection in lidar point clouds. In *Conference on Robot Learning*, pages 923–932, 2020. 1, 6

[79] Yin Zhou and Oncel Tuzel. Voxelnet: End-to-end learning for point cloud based 3d object detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4490–4499, 2018. 1, 2, 4, 6

[80] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *Proceedings of the IEEE international conference on computer vision*, pages 2223–2232, 2017. 2, 3

In this supplementary material, we provide detailed analysis about the statistics of the Waymo Domain Adaptation Dataset in Section A; the robustness analysis of the foreground voxel classifier in Section B; the derivation of the dropout rate used in the RndDrop method in Section C; more results on the Waymo Domain Adaptation Dataset in Section D; more results on KITTI in Section E; and more visualization of the semantic point generation in Section F.

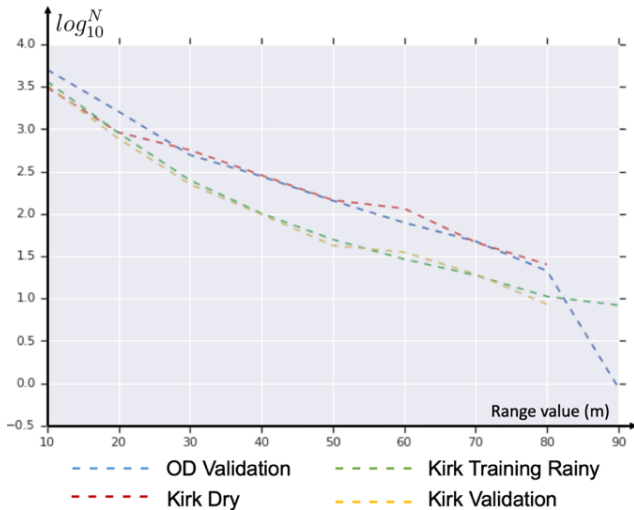## A. Statistics of the Waymo Domain Adaptation Dataset



Figure 7: The average number of raw points per vehicle across different ranges. On the x axis, the range value stands for the distance between the center of a bounding box and the LiDAR sensor. The y axis shows the value after applying $log_{10}$ on the number of points $N$. "Kirk Dry" is extracted from the Kirk Training set and contains frames captured in the dry weather.

We collect the statistics about the average number of points in a vehicle bounding box across different ranges. The range value is calculated as the euclidean distance between the LiDAR sensor and the center of a bounding box. We investigate four sets of point clouds:
- The OD Validation set, in which $99.5\%$ of the frames are collected in the dry weather.
- The Kirk Dry set, which consists of all the frames with the dry weather condition from the Kirk training set.
- The Kirk Training Rainy set, which consists of all the frames with the rainy weather condition from the Kirk training set.
- The Kirk Validation set, in which all the frames are collected in the rainy weather.

As shown in Figure 7, the point clouds with similar weather conditions share similar numbers of points per ob-

ject, even though they are collected at different locations. Specifically, the vehicle objects of the two "dry datasets", *i.e.,* the Kirk Dry set and the OD Validation set, have similar numbers of points across all ranges. The vehicle objects of the two "rainy datasets" *i.e.,* the Kirk Training Rainy set and the Kirk Validation set, share similar statistics.

In addition, the point clouds captured in **the dry weather** (the OD Validation set and the Kirk Dry set) have more points on each object than those collected in **the rainy weather** (the Kirk Training Rainy set and the Kirk Validation set). Please note that we have applied $log_{10}$ to the number of points for better visualization. The difference in the number of points is substantial between two weather conditions across all ranges.

## B. The Robustness of the Foreground Voxel Classifier

In order to generalize detectors to different domains, it is crucial to correctly classify foreground voxels so that semantic points can be reliably generated. Table 9 lists the evaluation results of the foreground voxel classifier. The

| Train | Eval | Accuracy | Precision | Recall | AP |
|-------|------|----------|-----------|--------|-----|
| OD Train | OD Val | 99.3 % | 90.9 % | 92.9 % | 86.7 % |
| OD Train | Kirk Val | 98.9 % | 88.4 % | 88.2 % | 78.3 % |

Table 9: Foreground voxel classification results of our SPG. The model is trained on the OD training set and then it is evaluated on the OD validation set and Kirk validation set, respectively. The accuracy, precision and recall are evaluated by setting $\tilde{P}^f > 0.5$.

results in Table 9 are averaged among all voxels in the foreground regions. Our SPG is trained on the OD training set. Then it is evaluated on the OD validation set and the Kirk validation set, respectively. The classification of a voxel is correct if its prediction score $\tilde{P}^f > 0.5$ when $y^f = 1.0$ or $\tilde{P}^f < 0.5$ when $y^f = 0.0$. The accuracy, precision and recall are all calculated under this setting. The AP is calculated using 40 recall thresholds. The results show that SPG achieves high performance in both domains.

## C. Dropout Rate of the RndDrop Method

In the experiment section, we implement a baseline method RndDrop, where we randomly drop out $17\%$ of points for point clouds from the source domain during training. This dropout ratio is chosen to match the ratio of missing points in the target domain. We calculate $(\overline{N}_{src} - \overline{N}_{tgt})/\overline{N}_{src} = 17\%$, where $\overline{N}_{src} = 121.2K$ is the average number of points per scene in the source domain and $\overline{N}_{tgt} = 100.4K$ is the average number of points per scene in the target domain.

## D. More Results on the Waymo Domain Adaptation Dataset

The evaluation tool [54] provides the average precision for three distance-based breakdowns: 0 to 30 meters, 30 to 50 meters, and beyond 50 meters. The AP is calculated using 100 recall thresholds.

We perform two groups of model comparisons in the setting of UDA: Group 1. PointPillars vs. SPG + PointPillars; Group 2. PV-RCNN vs. SPG + PV-RCNN. We train all models on the OD training set and evaluate them on both the OD validation set and the Kirk validation set. Table 10 and 11 show the comparisons on vehicle 3D AP and vehicle BEV AP, respectively. Table 12 and Table 13 show the comparisons in pedestrian 3D AP and pedestrian BEV AP, respectively. In most cases, SPG improves the detection performance across all ranges for both vehicles and pedestrians.

## E. More Results on KITTI

We provide more 3D object detection results on KITTI. There are two commonly used metric standards for evaluating the detection performance: 1) R11, where the AP is evaluated with 11 recall positions; 2) R40, where the AP is evaluated with 40 recall positions. In addition to the improvement on car and pedestrian detection, SPG also significantly boosts the performance in cyclist detection. Based on R11, Table 14 and Table 15 show the results in 3D AP and BEV AP for three object types, respectively. Based on R40, Table 16 and Table 17 show the results in 3D AP and BEV AP for three object types, respectively.

We show more comparisons on the KITTI test set in Table 18.

## F. More Visualization of Semantic Point Generation

In Figure 9, we illustrate more augmented point clouds, where the raw points are rendered in the grey color and the generated semantic points are highlighted in red.

| | | Target Domain - Kirk | | | | Source Domain - OD | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | Vehicle 3D AP (IoU = 0.7) | | | | Vehicle 3D AP (IoU = 0.7) | | | |
| Difficulty | Method | Overall | 0-30m | 30-50m | 50-Inf | Overall | 0-30m | 30-50m | 50-Inf |
| | PointPillars | 34.65 | 63.13 | 24.56 | 7.65 | 57.27 | 84.39 | 52.97 | 28.22 |
| LEVEL_1 | SPG + PointPillars | **41.56** | **68.26** | **31.91** | **13.08** | **62.44** | **86.18** | **58.13** | **35.40** |
| | *Improvement* | *+6.91* | *+5.13* | *+7.35* | *+5.43* | *+5.17* | *+1.79* | *+5.16* | *+7.18* |
| | PointPillars | 31.67 | 59.26 | 22.09 | 7.08 | 52.96 | 82.30 | 50.74 | 24.6 |
| LEVEL_2 | SP + PointPillar | **38.15** | **64.57** | **28.66** | **11.96** | **58.54** | **85.75** | **56.02** | **31.02** |
| | *Improvement* | *+6.48* | *+5.31* | *+6.57* | *+4.88* | *+5.58* | *+3.45* | *+5.28* | *+6.42* |
| | PV-RCNN | 55.16 | 76.68 | 47.96 | 27.59 | 74.01 | 91.39 | 70.94 | 49.51 |
| LEVEL_1 | SPG+PV-RCNN | **58.31** | **77.81** | **51.65** | **31.29** | **75.27** | **92.36** | **73.47** | **51.03** |
| | *Improvement* | *+3.15* | *+1.13* | *+3.69* | *+3.70* | *+1.26* | *+0.97* | *+2.53* | *+1.52* |
| | PV-RCNN | 45.81 | 71.31 | 38.83 | 20.52 | 64.69 | 88.95 | 64.80 | 37.37 |
| LEVEL_2 | SPG + PV-RCNN | **48.70** | **72.41** | **42.16** | **23.52** | **65.98** | **91.62** | **65.61** | **39.87** |
| | *Improvement* | *+2.89* | *+1.10* | *+3.33* | *+3.00* | *+1.29* | *+2.67* | *+0.81* | *+2.50* |

Table 10: The unsupervised domain adaptation vehicle detection results on both Waymo Open Dataset (OD) and Kirkland Dataset (Kirk). We show the vehicle 3D AP results in this table. The AP distance breakdowns are provided by the official evaluation tool.

| | | Target Domain - Kirk | | | | Source Domain - OD | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | Vehicle BEV AP (IoU = 0.7) | | | | Vehicle BEV AP (IoU = 0.7) | | | |
| Difficulty | Method | Overall | 0-30m | 30-50m | 50-Inf | Overall | 0-30m | 30-50m | 50-Inf |
| | PointPillars | 51.88 | 75.56 | 46.04 | 25.55 | 72.26 | 92.23 | 71.35 | 51.11 |
| LEVEL_1 | SPG + PointPillars | **60.44** | **80.89** | **53.73** | **38.24** | **77.63** | **93.39** | **75.96** | **61.16** |
| | *Improvement* | *+8.56* | *+5.33* | *+7.69* | *+12.69* | *+5.37* | *+1.16* | *+4.61* | *+10.05* |
| | PointPillars | 47.93 | 71.18 | 42.41 | 23.47 | 69.09 | 91.83 | 68.87 | 45.53 |
| LEVEL_2 | SPG + PointPillars | **56.94** | **77.13** | **49.99** | **35.04** | **74.90** | **93.06** | **73.96** | **54.51** |
| | *Improvement* | *+9.01* | *+5.95* | *+7.58* | *+11.57* | *+5.81* | *+1.23* | *+5.09* | *+8.98* |
| | PV-RCNN | 70.38 | 84.27 | 65.31 | 52.98 | 85.13 | 95.99 | 84.02 | 72.19 |
| LEVEL_1 | SPG + PV-RCNN | **72.56** | **84.43** | **68.79** | **58.49** | **87.38** | **97.54** | **86.63** | **74.59** |
| | *Improvement* | *+2.18* | *+0.16* | *+3.48* | *+5.51* | *+2.25* | *+1.55* | *+2.61* | *+2.40* |
| | PV-RCNN | 60.13 | 78.10 | 54.36 | 40.67 | 76.84 | 93.29 | 76.64 | 58.29 |
| LEVEL_2 | SPG + PV-RCNN | **62.03** | **78.86** | **56.47** | **44.94** | **78.05** | **94.45** | **80.25** | **59.56** |
| | *Improvement* | *+1.90* | *+0.76* | *+2.11* | *+4.27* | *+1.21* | *+1.16* | *+3.61* | *+1.27* |

Table 11: The unsupervised domain adaptation vehicle detection results on both Waymo Open Dataset (OD) and Kirkland Dataset (Kirk). We show the vehicle BEV AP results in this table. The AP distance breakdowns are provided by the official evaluation tool.

| Difficulty | Method | Target Domain - Kirk Pedestrian 3D AP (IoU = 0.5) | | | | Source Domain - OD Pedestrian 3D AP (IoU = 0.5) | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | Overall | 0-30m | 30-50m | 50-Inf | Overall | 0-30m | 30-50m | 50-Inf |
| | PointPillars | 20.65 | 43.98 | 9.27 | 3.24 | 55.20 | 69.24 | 52.04 | 32.72 |
| LEVEL_1 | SPG + PointPillars | **23.72** | **50.19** | **9.11** | **5.57** | **56.06** | **69.32** | **53.12** | **34.73** |
| | *Improvement* | *+3.07* | *+6.21* | *-0.16* | *+2.33* | *+0.86* | *+0.08* | *+1.08* | *+2.01* |
| | PointPillars | 17.66 | 40.67 | 7.40 | 2.32 | 51.33 | 65.85 | 49.32 | 29.29 |
| LEVEL_2 | SPG + PointPillars | **19.57** | **46.42** | **7.44** | **3.99** | **52.33** | **65.63** | **50.10** | **31.25** |
| | *Improvement* | *+1.91* | *+5.75* | *+0.04* | *+1.67* | *+1.00* | *-0.22* | *+0.78* | *+1.96* |
| | PV-RCNN | 24.47 | 39.69 | 14.24 | 8.05 | 65.34 | 72.23 | 64.89 | 50.04 |
| LEVEL_1 | SPG + PV-RCNN | **30.82** | **48.04** | **18.80** | **13.39** | **66.93** | **73.55** | **66.60** | **50.82** |
| | *Improvement* | *+6.35* | *+8.35* | *+4.56* | *+5.34* | *+1.59* | *+1.32* | *+1.71* | *+0.78* |
| | PV-RCNN | 17.16 | 36.39 | 9.64 | 3.51 | 56.03 | 66.88 | 56.58 | 35.76 |
| LEVEL_2 | SPG + PV-RCNN | **22.05** | **44.07** | **12.91** | **5.77** | **57.68** | **68.28** | **58.29** | **37.64** |
| | *Improvement* | *+4.89* | *+7.68* | *+3.27* | *+2.26* | *+1.65* | *+1.40* | *+1.71* | *+1.88* |

Table 12: The unsupervised domain adaptation pedestrian detection results on both Waymo Open Dataset (OD) and Kirkland Dataset (Kirk). We show the pedestrian 3D AP results in this table. The AP distance breakdowns are provided by the official evaluation tool.

| Difficulty | Method | Target Domain - Kirk Pedestrian BEV AP (IoU = 0.5) | | | | Source Domain - OD Pedestrian BEV AP (IoU = 0.5) | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | Overall | 0-30m | 30-50m | 50-Inf | Overall | 0-30m | 30-50m | 50-Inf |
| | PointPillars | 22.33 | 45.00 | 10.50 | 3.49 | 63.82 | 76.33 | 61.90 | 42.81 |
| LEVEL_1 | SPG + PointPillars | **24.83** | **51.44** | **10.80** | **5.71** | **64.66** | **76.11** | **62.69** | **44.98** |
| | *Improvement* | *+2.50* | *+6.44* | *+0.30* | *+2.22* | *+0.84* | *-0.22* | *+0.79* | *+2.17* |
| | PointPillars | 18.40 | 41.63 | 8.58 | 2.49 | 60.13 | 73.34 | 58.77 | 38.83 |
| LEVEL_2 | SPG + PointPillars | **20.67** | **47.56** | **8.98** | **4.11** | **60.93** | **72.94** | **59.54** | **41.11** |
| | *Improvement* | *+2.27* | *+5.93* | *+0.40* | *+1.62* | *+0.80* | *-0.40* | *+0.77* | *+2.28* |
| | PV-RCNN | 25.39 | 40.23 | 14.72 | 9.76 | 70.35 | 76.22 | 70.49 | 56.77 |
| LEVEL_1 | SPG + PV-RCNN | **31.92** | **49.06** | **19.87** | **14.87** | **70.37** | **75.86** | **72.29** | **57.47** |
| | *Improvement* | *+6.53* | *+8.83* | *+5.15* | *+5.11* | *+0.02* | *-0.36* | *+1.80* | *+0.70* |
| | PV-RCNN | 17.88 | 36.89 | 9.97 | 4.23 | 60.81 | 69.22 | 61.86 | 41.32 |
| LEVEL_2 | SPG + PV-RCNN | **22.65** | **44.57** | **13.48** | **6.38** | **60.88** | **70.62** | **63.65** | **43.27** |
| | *Improvement* | *+4.77* | *+7.68* | *+3.51* | *+2.15* | *+0.07* | *+1.40* | *+1.79* | *+1.95* |

Table 13: The unsupervised domain adaptation pedestrian detection results on both Waymo Open Dataset (OD) and Kirkland Dataset (Kirk). We show the pedestrian BEV AP results in this table. The AP distance breakdowns are provided by the official evaluation tool.

| Method | Car - 3D AP (R11) | | | Pedestrian - 3D AP (R11) | | | Cyclist - 3D AP (R11) | | |
|---|---|---|---|---|---|---|---|---|---|
| | Easy | Mod. | Hard | Easy | Mod. | Hard | Easy | Mod. | Hard |
| PointPillars[25] | 86.46 | 77.28 | 74.65 | 57.75 | 52.29 | 47.90 | 80.05 | 62.68 | 59.70 |
| SPG + PointPillars | **87.98** | **78.54** | **77.32** | **59.91** | **54.58** | **50.34** | **81.58** | **65.70** | **62.28** |
| *Improvement* | *+1.52* | *+1.26* | *+2.67* | *+2.16* | *+2.29* | *+2.44* | *+1.53* | *+3.02* | *+2.58* |
| PVRCNN[49] | 89.35 | 83.69 | 78.70 | 64.60 | 57.90 | 53.23 | 85.22 | 70.47 | 65.75 |
| SPG + PVRCNN | **89.81** | **84.45** | **79.14** | **69.04** | **62.18** | **56.77** | **86.82** | **73.35** | **69.30** |
| *Improvement* | *+0.46* | *+0.76* | *+0.44* | *+4.44* | *+4.28* | *+3.54* | *+1.60* | *+2.88* | *+3.55* |

Table 14: Result comparisons on the KITTI validation set. The results are evaluated by the Average Precision with 11 recall positions. The baseline detectors, PointPillars and PV-RCNN, are directly evaluated by using the checkpoints released by [49, 55].

| Method | Car - BEV AP (R11) | | | Pedestrian - BEV AP (R11) | | | Cyclist - BEV AP (R11) | | |
|---|---|---|---|---|---|---|---|---|---|
| | Easy | Mod. | Hard | Easy | Mod. | hard | Easy | Mod. | Hard |
| PointPillars[25] | 89.65 | 87.17 | 84.37 | 61.63 | 56.27 | 52.60 | 82.27 | 66.25 | 62.64 |
| SPG + PointPillars | **90.07** | **88.00** | **86.63** | **65.16** | **59.86** | **56.07** | **86.02** | **71.93** | **65.69** |
| *Improvement* | *+0.42* | *+0.83* | *+2.26* | *+3.53* | *+3.59* | *+3.47* | *+3.75* | *+5.68* | *+3.05* |
| PVRCNN[49] | 90.09 | 87.90 | 87.41 | 67.01 | 61.38 | 56.10 | 86.79 | 73.55 | 69.69 |
| SPG + PVRCNN | **90.41** | **88.49** | **87.74** | **71.19** | **64.37** | **59.88** | **92.54** | **74.43** | **70.99** |
| *Improvement* | *+0.32* | *+0.59* | *+0.33* | *+4.18* | *+2.99* | *+3.78* | *+5.75* | *+0.88* | *+1.30* |

Table 15: Result comparisons on the KITTI validation set. The results are evaluated by the Average Precision with 11 recall positions. The baseline detectors, PointPillars and PV-RCNN, are directly evaluated by using the checkpoints released by [49, 55].

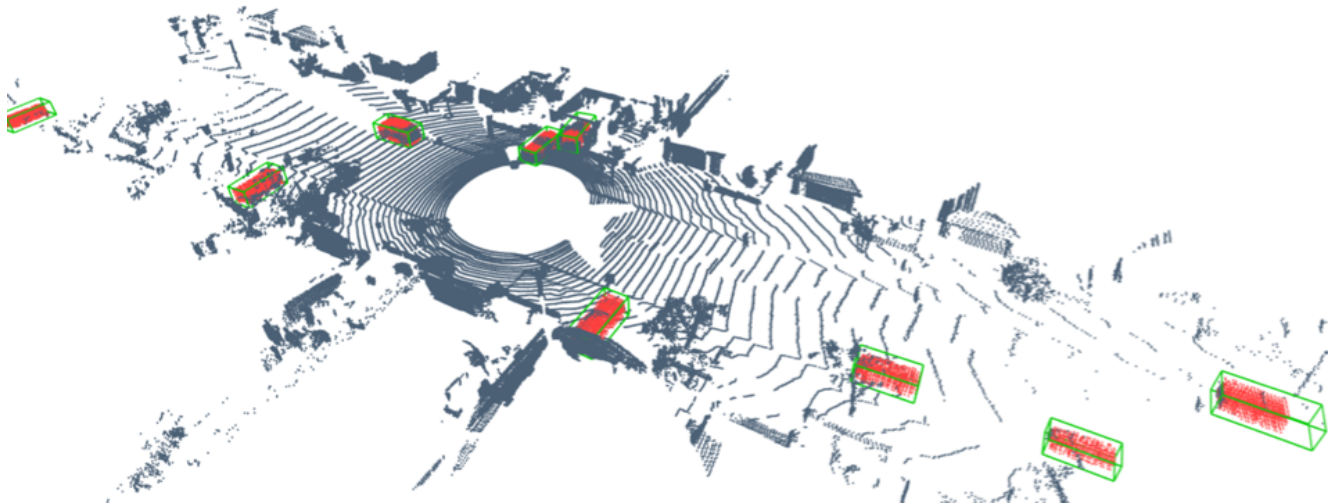| Method | Car - 3D AP (R40) | | | Pedestrian - 3D AP (R40) | | | Cyclist - 3d AP (R40) | | |
|---|---|---|---|---|---|---|---|---|---|
| | Easy | Mod. | Hard | Easy | Mod. | Hard | Easy | Mod. | Hard |
| PointPillars[25] | 87.75 | 78.39 | 75.18 | 57.30 | 51.41 | 46.87 | 81.57 | 62.94 | 58.98 |
| SPG+PointPillars | **89.77** | **81.36** | **78.85** | **59.65** | **53.55** | **49.24** | **83.27** | **66.11** | **61.99** |
| *Improvement* | *+2.02* | *+2.97* | *+3.67* | *+2.35* | *+2.14* | *+2.37* | *+1.70* | *+3.17* | *+3.01* |
| PVRCNN[49] | 92.10 | 84.36 | 82.48 | 64.26 | 56.67 | 51.91 | 88.88 | 71.95 | 66.78 |
| SPG+PVRCNN | **92.53** | **85.31** | **82.82** | **69.66** | **61.80** | **56.39** | **91.75** | **74.35** | **69.49** |
| *Improvement* | *+0.43* | *+0.95* | *+0.34* | *+5.40* | *+5.13* | *+4.48* | *+2.87* | *+2.40* | *+2.71* |

Table 16: Result comparisons on the KITTI validation set. The results are evaluated by the Average Precision with 40 recall positions. The baseline detectors, PointPillars and PV-RCNN, are directly evaluated by using the checkpoints released by [49, 55].

| Method | Car - BEV AP (R40) | | | Pedestrian - BEV AP (R40) | | | Cyclist - BEV AP (R40) | | |
|---|---|---|---|---|---|---|---|---|---|
| | Easy | Mod. | Hard | Easy | Mod. | Hard | Easy | Mod. | Hard |
| PointPillars[25] | 92.03 | 88.05 | 86.66 | 61.59 | 56.01 | 52.04 | 85.27 | 66.34 | 62.36 |
| SPG + PointPillars | **94.38** | **89.92** | **87.97** | **65.38** | **59.48** | **55.32** | **90.29** | **71.43** | **66.96** |
| *Improvement* | *+2.35* | *+1.87* | *+1.31* | *+3.79* | *+3.47* | *+3.28* | *+5.02* | *+5.09* | *+4.60* |
| PVRCNN[49] | 93.02 | 90.33 | 88.53 | 67.97 | 60.52 | 55.80 | 91.02 | 74.54 | 69.92 |
| SPG + PVRCNN | **94.99** | **91.11** | **88.86** | **71.79** | **64.50** | **59.51** | **93.62** | **76.45** | **71.64** |
| *Improvement* | *+1.97* | *+0.78* | *+0.33* | *+3.82* | *+3.98* | *+3.71* | *+2.60* | *+1.91* | *+1.72* |

Table 17: Result comparisons on the KITTI validation set. The results are evaluated by the Average Precision with 40 recall positions. The baseline detectors, PointPillars and PV-RCNN, are directly evaluated by using the checkpoints released by [49, 55].

| Method | Reference | Modality | Car - 3D AP | | | |
|---|---|---|---|---|---|---|
| | | | Easy | Mod. | Hard | Avg. |
| F-PointNet[42] | CVPR 2018 | LIDAR & RGB | 82.19 | 69.79 | 60.59 | 70.86 |
| AVOD-FPN[24] | IROS 2018 | LIDAR & RGB | 83.07 | 71.76 | 65.73 | 73.52 |
| F-ConvNet[58] | IROS 2019 | LIDAR & RGB | 87.36 | 76.39 | 66.69 | 76.81 |
| UberATG-MMF[29] | CVPR 2019 | LIDAR & RGB | 88.40 | 77.43 | 70.22 | 78.68 |
| EPNet[22] | ECCV 2020 | LiDAR & RGB | 89.81 | 79.28 | 74.59 | 81.23 |
| CLOCs_PVCas[39] | IROS 2020 | LiDAR & RGB | 88.94 | 80.67 | 77.15 | 82.25 |
| 3D-CVF[72] | ECCV 2020 | LiDAR & RGB | 89.20 | 80.05 | 73.11 | 80.79 |
| SECOND[64] | Sensors 2018 | LiDAR | 83.34 | 72.55 | 65.82 | 73.90 |
| PointPillars[25] | CVPR 2019 | LiDAR | 82.58 | 74.31 | 68.99 | 75.30 |
| PointRCNN[50] | CVPR 2019 | LiDAR | 86.96 | 76.50 | 71.39 | 77.77 |
| 3D IoU Loss[75] | 3DV 2019 | LiDAR | 86.16 | 75.64 | 70.70 | 78.28 |
| Fast Point R-CNNs[10] | ICCV 2019 | LiDAR | 85.29 | 77.40 | 70.24 | 77.64 |
| STD[68] | ICCV 2019 | LiDAR | 87.95 | 79.71 | 75.09 | 80.91 |
| SegVoxelNet[70] | ICRA 2020 | LiDAR | 86.04 | 76.13 | 70.76 | 77.64 |
| SARPNET[69] | Neuro Computing 2019 | LiDAR | 85.63 | 76.64 | 71.31 | 77.86 |
| HRI-VoxelFPN[70] | Sensor 2020 | LiDAR | 85.63 | 76.70 | 69.44 | 77.26 |
| HotSpotNet[5] | ECCV 2020 | LiDAR | 87.60 | 78.31 | 73.34 | 79.75 |
| PartA$^2$[52] | TPAMI 2020 | LiDAR | 87.81 | 78.49 | 73.51 | 79.94 |
| SERCNN[76] | CVPR 2020 | LiDAR | 87,74 | 78.96 | 74.14 | 51.03 |
| Point-GNN[53] | CVPR 2020 | LiDAR | 88.33 | 79.47 | 72.29 | 80.03 |
| 3DSSD[67] | CVPR 2020 | LiDAR | 88.36 | 79.57 | 74.55 | 80.83 |
| SA-SSD[19] | CVPR 2020 | LiDAR | 88.75 | 79.79 | 74.16 | 80.90 |
| CIA-SSD[60] | AAAI 2021 | LiDAR | 89.59 | 80.28 | 72.87 | 80.91 |
| Asso-3Ddet[15] | CVPR 2020 | LiDAR | 85.99 | 77.40 | 70.53 | 77.97 |
| Voxel R-CNN[13] | AAAI 2021 | LiDAR | **90.90** | 81.62 | 77.06 | 83.19 |
| PV-RCNN[49] | CVPR 2020 | LiDAR | 90.25 | 81.43 | 76.82 | 82.83 |
| SPG+PV-RCNN (Ours) | - | LiDAR | 90.49 | **82.13** | **78.88** | **83.83** |

Table 18: Car detection result comparisons on the KITTI test set. The results are evaluated by the Average Precision with 40 recall positions on the KITTI benchmark website. We compare with the leader board front runner detectors that are associated with conferences or journals released before our submission. The Avg. AP is calculated by averaging over the APs of Easy, Mod. and Hard. difficulty levels.
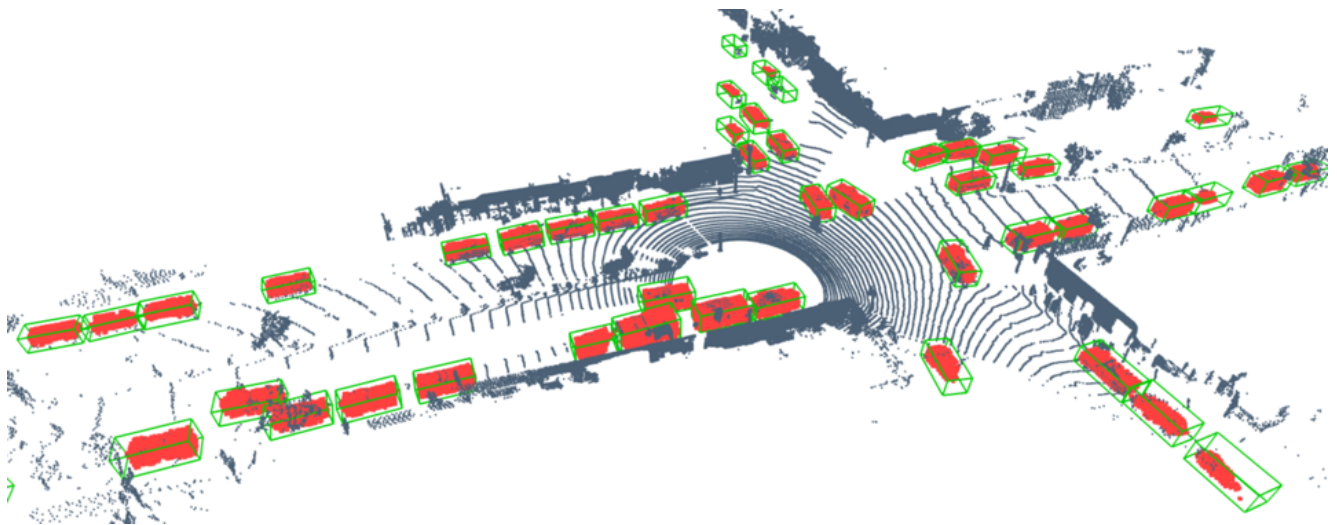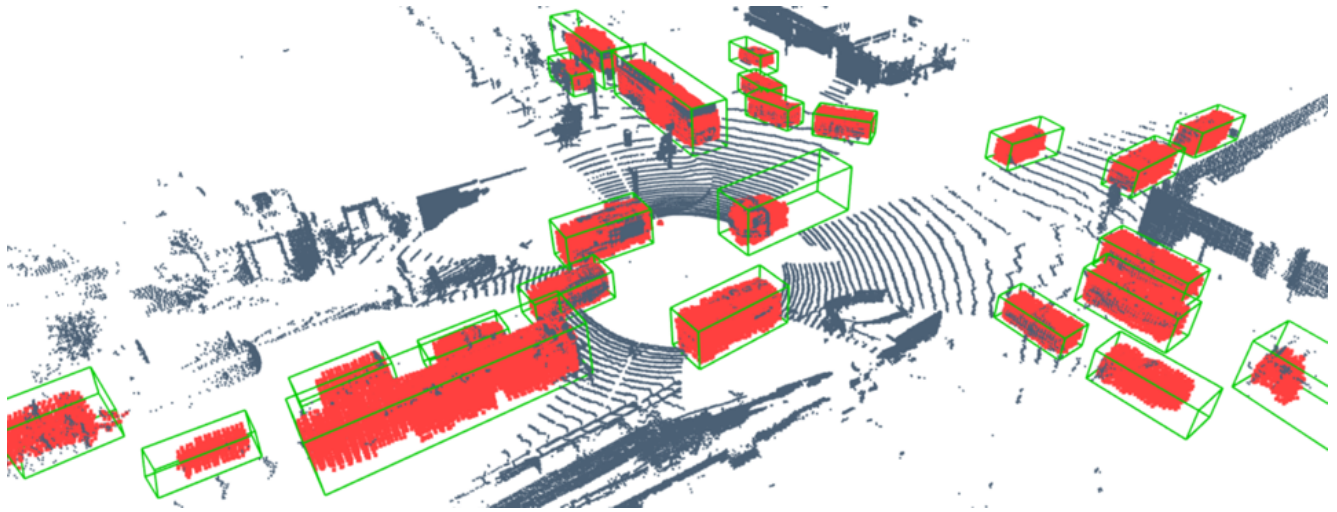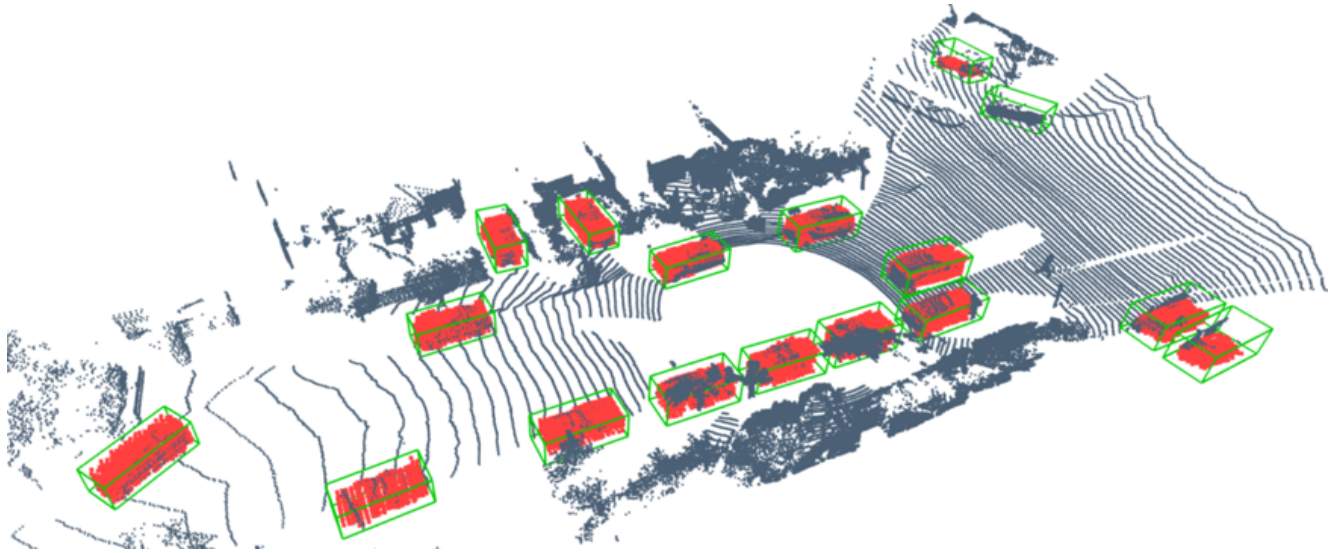
Figure 9: More visualization of generated semantic points. The grey points are original raw points. The red points are the generated semantic points. The green boxes are the predicted bounding boxes.