

Rethinking Pseudo-LiDAR Representation

Xinzhu Ma¹, Shinan Liu², Zhiyi Xia³, Hongwen Zhang⁴,
Xingyu Zeng², and Wanli Ouyang¹

¹ The University of Sydney, SenseTime Computer Vision Research Group, Australia
{xima0693@uni., wanli.ouyang@}sydney.edu.au

² SenseTime Research, China

{liushinan, zengxingyu}@sensetime.com

³ Dalian University of Technology, China

xiazhiyi99@mail.dlut.edu.cn

⁴ Institute of Automation, Chinese Academy of Sciences, China
hongwen.zhang@cripac.ia.ac.cn

Abstract. The recently proposed pseudo-LiDAR based 3D detectors greatly improve the benchmark of monocular/stereo 3D detection task. However, the underlying mechanism remains obscure to the research community. In this paper, we perform an in-depth investigation and observe that the efficacy of pseudo-LiDAR representation comes from the coordinate transformation, instead of data representation itself. Based on this observation, we design an image based CNN detector named PatchNet, which is more generalized and can be instantiated as pseudo-LiDAR based 3D detectors. Moreover, the pseudo-LiDAR data in our PatchNet is organized as the image representation, which means existing 2D CNN designs can be easily utilized for extracting deep features from input data and boosting 3D detection performance. We conduct extensive experiments on the challenging KITTI dataset, where the proposed PatchNet outperforms all existing pseudo-LiDAR based counterparts. Code has been made available at: <https://github.com/xinzhuma/patchnet>.

Keywords: image-based 3D detection, data representation, image, pseudo-LiDAR, coordinate transformation

1 Introduction

3D object detection has received increasing attention from both industry and academia because of its wide applications in various fields such as autonomous driving and robotics. Existing algorithms largely rely on LiDAR sensors, which provide accurate 3D point clouds of the surrounding environment. Although these approaches achieve impressive performance, the excessive dependence on expensive equipment restricts their application prospects.

Compared with fast developing LiDAR-based algorithms, 3D detection [7,6,20] results produced from only RGB images lag considerably behind. This can be attributed to the ill-posed nature of the problem, where a lack of explicit knowledge about the unobserved depth dimension significantly increases the task complexity. An intuitive solution is that we can use a Convolutional Neural Network

(CNN) to predict the depth map [1,10,13] and then use it to augment the input data if we do not have the available depth information. Although the estimated depth map is helpful to 3D scene understanding, the performance improvement brought by it is still limited.

Several recently proposed algorithms [35,24,36] transform the estimated depth map into pseudo-LiDAR representation, and then apply LiDAR based methods to the transformed data. Surprisingly, this simple yet effective method achieves significant improvement in detection accuracy on the challenging KITTI dataset. However, it is unclear why such a representation can bring so much performance improvement. According to the empirical explanation of proponents, the choice of representations is the critical success factor of 3D detection systems. Compared with image representation, they believe that pseudo-LiDAR is more suitable for describing the 3D structure of objects, which is the main reason for performance improvement. However, in the absence of direct evidence, the correctness of this statement is still open to doubt.

In this paper, we aim to explore the essential reasons of this phenomenon. Specifically, on the basis of prior works, we carefully construct an image representation based detector named PatchNet-vanilla, which is an equivalent implementation of pseudo-LiDAR [35] except for the representation of input data. With this detector, we can compare the influence of these two kinds of representations on 3D detection task in depth. Different from the arguments of other works [35,24,36], we observe that the performances of PatchNet-vanilla and pseudo-LiDAR [35] are completely matched, which means that data representation has no effect on 3D detection performance. Moreover, we perform ablation studies on the input data, and observe that the real thing matters is coordinate transformation from image coordinate system to the LiDAR coordinate system, which implicitly encodes the camera calibration information into input data.

PatchNet-vanilla also hints us that pseudo-LiDAR representation is not necessary to improve the accuracy of image based 3D detection. By integrating the generated 3D coordinates as additional channels of input data, our 3D detector gets promising performance. More importantly, this approach can be easily generalized to other image based detectors. Also notice that, as a kind of non-grid structured data, pseudo-LiDAR signals commonly need point-wise CNNs [29,30] to process. However, the development of these technologies still lags behind the standard CNNs. From this point of view, the image-based detectors should outperform their counterparts based on pseudo-LiDAR. To confirm this hypothesis, PatchNet was proposed by extending our original model (e.g., using more powerful backbone network [15,16]), and outperforms other pseudo-LiDAR based detectors on KITTI dataset. In addition, there are other benefits from using images directly as the network’s inputs, such as allowing us to train an end-to-end 3D detector. Based on above reasons, we argue that image representation based 3D detectors have greater development potential.

To summarize, the contributions of this paper are as follows: First, through sufficient experimental demonstration, we confirm that the reason why the pseudo-

LiDAR representation is effective is not the data representation itself, but the coordinate system transformation. Second, we find that pseudo-LiDAR representation is not necessary to improve detection performance. After integrating spatial coordinates, image representation based algorithms can also achieve the competitive if not superior the same performance. Third, thanks to more powerful image-based deep learning technologies, we achieve the state-of-the-art performance and show the potential of image representation based 3D detectors.

2 Related Work

2.1 3D detectors based on image representation

Most of the early works in this scope share the same paradigm with 2D detectors [12,32,40,9,22,21,14]. However, estimating the 3D coordinates (x, y, z) of the object center is much more complicated since there is ambiguity to locate the absolute physical position from only image appearances. Mono3D [6] focus on 3D object proposals generation using prior knowledge (e.g., object size, ground plane). Deep3DBox [26] introduces geometric constraints based on the fact that the 3D bounding box should fit tightly into 2D detection bounding box. DeepMANTA [4] encodes 3D vehicle information using key points, since vehicles are rigid objects with well known geometry. Then the vehicle recognition in DeepMANTA can be considered as key points detection. An expansion stage of ROI-10D [25] takes the advantage of depth information provided by an additional depth estimator [10,5], which itself is learned in a self-supervised manner. In Multi-Fusion [38], a multi-level fusion approach is proposed to exploit disparity estimation results from a pre-trained module for both the 2D box proposal generation and the 3D prediction part of their network. MonoGRNet [31] consists of four subnetworks for progressive 3D localization and directly learning 3D information based on solely semantic cues. MonoDIS [34] disentangles the loss for 2D and 3D detection and jointly trains these two tasks in an end-to-end manner. M3D-RPN [2] is the current state-of-the-art with image representation as input, using multiple 2D convolutions of non-shared weights to learn location-specific features for joint prediction of 2D and 3D boxes. The above approaches utilize various prior knowledge, pre-train models or more powerful CNN designs, but they do not try to use pseudo-LiDAR data to improve their performance. Our work aims to improve the detection accuracy of image-based methods by extracting useful information from pseudo-LiDAR data, which is complementary to these approaches.

2.2 3D detectors based on pseudo-LiDAR representation

Recently, several approaches [24,35,36,39] greatly boost the performance of monocular 3D detection task. What they have in common is that they first estimate the depth map from the input RGB image and transform it into pseudo-LiDAR (point cloud) by leveraging the camera calibration information. Specifically, [35]

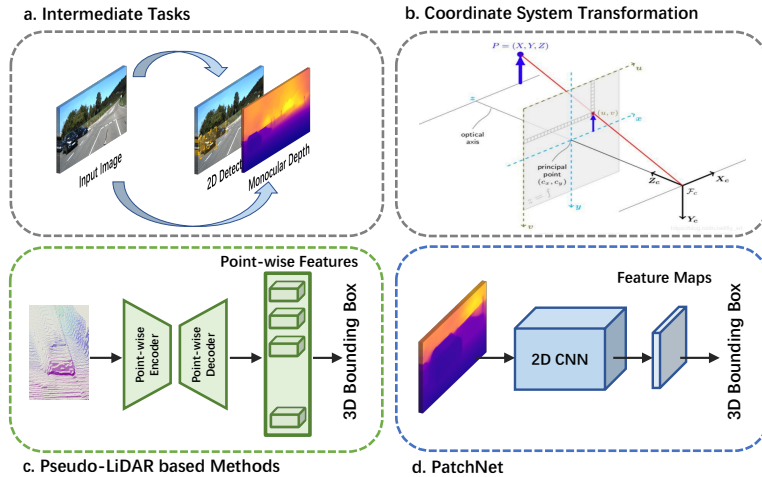


Fig. 1. Comparison of pseudo-LiDAR based methods [24,35,36] and PatchNet. They both generate intermediate tasks using off-the-shelf models (a), and project the image coordinates to the world coordinates (b). Pseudo-LiDAR based methods treat these data as LiDAR signals, and use point-wise network to predict results from them (c). However, PatchNet organizes them as image representation for subsequent processing (d).

adopt off-the-shelf LiDAR-based 3D detectors [28,18] to process the generated pseudo-LiDAR signals directly. AM3D [24] proposes a multi-modal features fusion module to embed the complementary RGB cues into the generated pseudo-LiDAR representation. Besides, [24] also proposes a depth prior based background points segmentation module to avoid the problems caused by the inaccuracy of point cloud annotation. [36] proposes a 2D-3D bounding box consistency loss which can alleviate the local misalignment issue. However, such methods rely heavily on the accuracy of depth map. Overall, pseudo-LiDAR based detectors achieve impressive accuracy in 3D detection task, however, the underlying mechanism is still obscure to the research community. In this paper, we perform an in-depth investigation on this issue. Besides, pseudo-LiDAR based detectors treat generated 3D data as point cloud and use PointNet for processing the point cloud, while our PatchNet organizes them as image and facilitates the use of 2D CNN for processing the data.

3 Delving into pseudo-LiDAR representation

In this section, we investigate the influence of pseudo-LiDAR representation on 3D detection accuracy. In particular, we first give a brief review of pseudo-LiDAR based detectors and introduce the technical details of its image based equivalent

detector. Then, we analyse whether data representation is the internal reason of performance improvement by comparing the performance of these two detectors.

3.1 Review of pseudo-LiDAR based detectors

Here we take pseudo-LiDAR [35] as example for analysis, and the paradigm of [35] can be summarized as follows:

Step 1: Depth estimation. Given a single monocular image (or stereo pairs) as input, [35] predict the depth d for each image pixel (u, v) using a stand alone CNN (Fig 1(a)).

Step 2: 2D detection. Another CNN is adopted to generate 2D object region proposals (Fig 1(a)).

Step 3: 3D data generation. First, regions of interests (RoIs) are cropped from the depth map generated from Step 1, according to the region proposals generated from Step 2. Then, the 3D coordinates of pixels of each RoI can be recovered by:

$$\begin{cases} z = d, \\ x = (u - C_x) \times z/f, \\ y = (v - C_y) \times z/f, \end{cases} \quad (1)$$

where f is the focal length of the camera, (C_x, C_y) is the principal point (Fig 1(b)).

Step 4: 3D object detection. Pseudo-LiDAR based approaches treat the 3D data generated from Step 3 as LiDAR signals, and use point-wise CNN to predict result from them (Fig 1(c)). In particular, they are treated as an unordered point set $\{x_1, x_2, \dots, x_n\}$ with $x_i \in \mathbb{R}^d$, and processed by PointNet, which defines a set function f that maps a set of points to a output vector:

$$f(x_1, x_2, \dots, x_n) = \gamma \left(\underset{i=1, \dots, n}{\mathbf{MAX}} \{h(x_i)\} \right) \quad (2)$$

where γ and h are implemented by multi-layer perceptron (MLP) layers.

3.2 PatchNet-vanilla: equivalent implementation of pseudo-LiDAR

Analysis. The most significant difference between the pseudo-LiDAR based approaches [24,35] and other approaches lies in the representation of depth map. The authors of [24,35] argue that pseudo-LiDAR representation is more suitable for describing the 3D structure of objects, which is the main reason behind the high accuracy of their models. To verify this, we conduct an image representation based detector, i.e., PatchNet-vanilla, which is identical to pseudo-LiDAR [35] except for the input representation.

Implementation. The steps 1, 2 and 3 in PatchNet-vanilla are the same as that in the pseudo-LiDAR based detectors. Therefore, they have the same estimated depth, 2D detection results and generated 3D data. The main difference is the Step 4, which will be analyzed in details. Specifically, in PatchNet-vanilla,

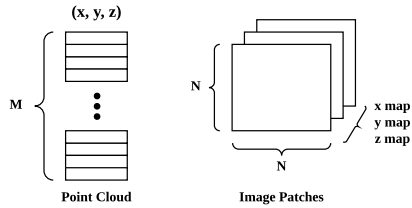


Fig. 2. Illustration of input data. Pseudo-LiDAR based approaches use point cloud (*left*) as input, while PatchNet use image patches (*right*) as input. We set $M = N \times N$ so that these two kinds of input data contain the same amount of information.

the generated 3D data is organized as image representation (see Fig 2), where each pixel location with 3 channels, i.e. (x, y, z) in Eq. 1. Different from point-wise CNN used in pseudo-LiDAR counterparts, 2D CNN is used for processing the input data in PatchNet-vanilla (Fig 1(d)). Note that we can define a same function as Eq. 2 using 2D convolution with 1×1 receptive field and global max pooling. This scheme is also adopted in the official implementation¹ of PointNet.

3.3 Preliminary conclusion

Table 1. Comparison of different input representation. Experiments are conducted on KITTI *validation* set. * indicates the method is reproduced by ourself. Metric is $AP|_{R_{11}}$ of the **Car** category.

Method	Modality	3D detection			BEV detection		
		Easy	Moderate	Hard	Easy	Moderate	Hard
pseudo-LiDAR [35]	pseudo-LiDAR	28.2	18.5	16.4	40.6	26.3	22.9
pseudo-LiDAR*	pseudo-LiDAR	28.9	18.4	16.2	41.0	26.2	22.8
PatchNet-vanilla	image	28.7	18.4	16.4	40.8	26.1	22.8

The performances of PatchNet-vanilla and pseudo-LiDAR are reported in Tab. 1, where we reproduce pseudo-LiDAR to eliminate the impact of implementation details. As can be seen, PatchNet-vanilla achieves almost the same accuracy as pseudo-LiDAR, which means the choice of data representation has no substantial impact on 3D detection tasks. Moreover, we perform ablation studies on data content, and observe that coordinate transform is the key factor for performance improvement (experimental results and analysis can be found in Sec. 5.2).

Above observations reveal that pseudo-LiDAR representation is not necessary, and after integrating the generated 3D information, image representation has the same potential. More importantly, compared with point-wise CNNs [29,30], image based representation can utilize the well-studied 2D CNNs

¹ <https://github.com/charlesq34/pointnet>.

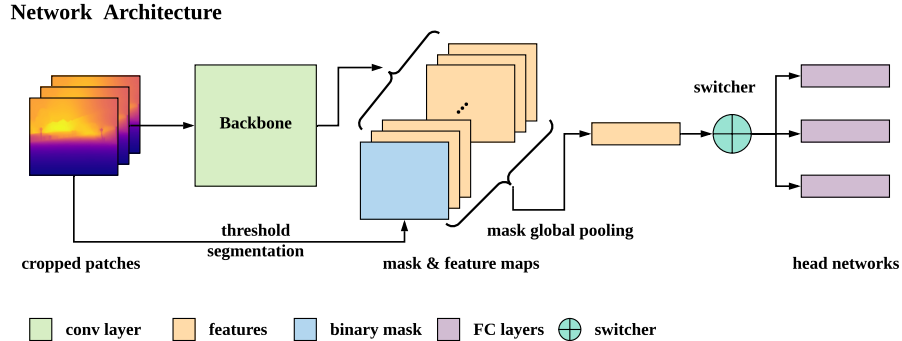


Fig. 3. Illustration of the network architecture. Given an input patch with $\{x, y, z\}$ channels, we first generate a binary mask according to mean depth, and use it to guide pooling layer to extract the features corresponding to foreground object. Then, we assign examples to different head networks according to the prediction difficulty of them.

for developing high-performance 3D detectors. Along this direction, we show how the proposed PatchNet framework is used to further improve the detection performance in Sec. 4.

4 PatchNet

In PatchNet, we first train two deep CNNs on two intermediate prediction tasks (i.e., 2D detection and depth estimation) to obtain position and depth information, which are the same as PatchNet-vanilla and pseudo-LiDAR based detectors (Fig 1(a)). Then, as shown in Fig 3, for each detected 2D object proposal, we crop the corresponding region from the depth map, and recover its spatial information using Eq 1. Next, deep features of RoIs are extracted by backbone network, and filtered by the mask global pooling and foreground mask. Finally, we use a detection head with difficulty assignment mechanism to predict the 3D bounding box parameterized by $(x, y, z, h, w, l, \theta)$.

Backbone Most of existing backbone networks can be used in our method to extract image features. In our implementation, we use the ResNet-18 [15] with Squeeze-and-Excitation (SE) block [16] as the 3D detection backbone. Moreover, we remove all pooling layers in the original SE-ResNet-18 so that its output features have the same size as input image patches. Then we use mask global pooling operation and generated mask to extract features from foreground object.

Mask global pooling The feature maps \mathbf{X} output from the backbone network will be converted to a feature vector by global pooling. Conventional global

pooling takes features of all positions into account and output the global feature. To obtain more robust features, we perform global pooling only on those features within foreground regions so that the final feature is corresponding to those pixels of interest. Specifically, we additionally generate a binary mask \mathbf{M} which indicates the foreground region. This masks will be applied to the feature maps \mathbf{X} to select foreground features before global pooling. Such a mask global pooling encourages the final feature to focus on the regions of interest.

Mask generation Following the prior work [24], the fore/background binary mask \mathbf{M} is obtained by setting a threshold to the depth map. Specifically, we empirically add an offset on the mean depth of each patch and set it as the threshold. The regions with the depth values smaller than this threshold will be regarded as the foreground regions. The binary mask \mathbf{M} has the same resolution as the input image with its values corresponding to foreground regions set as 1 and otherwise 0.

Head Inspired by difficulty-wise evaluation adopted by KITTI dataset, we use three branches to deal with samples of different difficulty levels separately. To select the branch, we need a specific module. Specifically, before sending the feature maps to the three parallel box estimators, we add another branch to predict the difficulty level of each instance.

Note that all three branches are the same in network architecture, and only different in learned parameters for handling different difficulty levels. Besides, in our implementation, all three branches predict results simultaneously, and two of them are blocked according to the output of difficulty predictor. Theoretically, this does not affect the accuracy of the algorithm, and allows all branches to run in parallel with the cost of extra GPU memory.

Loss function The ground truth box is parameterized by center (x, y, z) , size (w, h, l) and heading angle θ . We adopted the loss function proposed by [28] to our baseline model:

$$\mathcal{L} = \mathcal{L}_{center} + \mathcal{L}_{size} + \mathcal{L}_{heading} + \lambda \cdot \mathcal{L}_{corner} \quad (3)$$

where \mathcal{L}_{center} , \mathcal{L}_{size} , and $\mathcal{L}_{heading}$ respectively denote the loss function for the center, size, and heading angle. λ is an empirical weight, and \mathcal{L}_{corner} is used to alleviate the potential sub-optimal problem. Please refer to [28] for details.

5 Experiments

5.1 Setup

Dataset We evaluate our approach on the challenging KITTI dataset [11], which provides 7,481 images for training and 7,518 images for testing. Detection and

localization (i.e., bird’s-eye-view detection) tasks are evaluated in three different subsets: *easy*, *moderate* and *hard*, according to the occlusion and truncation levels of objects. Since the ground truth for the test set is not available and the access to the test server is limited, we follow the protocol of prior works [6,7,8] to divide the training data into a training set (3,712 images) and a validation set (3,769 images). We will conduct ablation studies based on this split and also report final results on the testing set provided by KITTI server. Due to space limitations, we only report the **Car** detection results of **monocular images** in the main paper. More results about **stereo pairs** and **Pedestrian/Cyclist** can be found in Appendix.

Metric Most of previous works use 11-point interpolated average precision (IAP) metric [11] as follows:

$$AP|_{R_{11}} = \frac{1}{11} \sum_{r \in R_{11}} \max_{\tilde{r} \geq r} \rho(\tilde{r}). \quad (4)$$

Recently, to avoid ostensible boost in performance, KITTI and [34] call for a new 40-point IAP ($AP|_{R_{40}}$) with the exclusion of “0” and four-times denser interpolated prediction for better approximation of the area under the Precision/Recall curve. For fair and comprehensive comparisons with previous and future works, we show both $AP|_{R_{11}}$ and $AP|_{R_{40}}$ in the following experiments.

5.2 Investigation of pseudo-LiDAR representation

Table 2. 3D object detection results on KITTI *validation* set. Metrics are AP_{3D} and AP_{BEV} of the **Car** category with 11 recall positions. * indicates method is reproduced by ourselves.

Method	Modality	3D detection			BEV detection		
		Easy	Moderate	Hard	Easy	Moderate	Hard
pseudo-LiDAR [35]	pseudo-LiDAR	28.2	18.5	16.4	40.6	26.3	22.9
pseudo-LiDAR*	pseudo-LiDAR	28.9	18.4	16.2	41.0	26.2	22.8
AM3D [24]	pseudo-LiDAR	32.2	21.1	17.3	43.8	28.4	23.9
PatchNet-vanilla	image	28.7	18.4	16.4	40.8	26.1	22.8
PatchNet-AM3D	image	32.8	20.9	17.3	43.5	28.2	23.6
PatchNet	image	35.1	22.0	19.6	44.4	29.1	24.1
Improvement	-	+2.9	+0.9	+2.3	+0.6	+0.7	+0.2

Analysis of data representation As shown in Tab. 2, PatchNet-vanilla shows a comparable results with pseudo-LiDAR, which indicates that *data representation is not the key factor to improve the performance of 3D detectors*. To further validate this claim, we also adjust our image representation based detector based on AM3D, where we achieve a matched performance again.

Table 3. Comparison between different input data on KITTI *validation* set. Metrics are AP_{3D} and AP_{BEV} of the **Car** category with 11 recall positions.

input	AP_{3D}			AP_{BEV}		
	Easy	Moderate	Hard	Easy	Moderate	Hard
{z}	4.51	3.48	3.03	6.31	4.50	3.98
{x, z}	27.1	18.3	15.8	35.9	23.4	18.3
{x, y, z}	35.1	22.0	19.6	44.4	29.1	24.1
{u, v, z}	24.6	15.7	14.6	33.2	21.3	16.7

Analysis of data content We conduct an ablation study on the effect of input channels and report the results in Tab. 5.2. We can see from the results that, using only depth as an input, it is almost impossible to obtain accurate 3D bounding boxes. If other coordinates are used, the accuracy of predicted boxes improves greatly, which validates the importance of generated spatial features. It should be noted that in the absence of y-axis data, this detection accuracy is much worse than our full model. This shows that all coordinates are useful for the 3D detection.

In pseudo-LiDAR, the coordinate (u, v) for images is projected to the world coordinate (x, y) using the camera information. Experimental results in Tab. 5.2 also compares the effectiveness of different coordinate systems. According to experimental results, world coordinate (x, y) , which utilizes the camera information, performs much better than image coordinate (u, v) . Through the above experiments, we can observe that *that real thing matters is coordinate system transformation, instead of data representation itself*.

5.3 Boosting the performance of PatchNet

Backbone Compared with point-wise backbone nets commonly used in (pseudo) LiDAR based methods, standard 2D backbones such as [15,16,37] can extract more discriminative features, which is a natural advantage of image based detectors. We investigate the impact of different backbones on proposed PatchNet, and the experimental results are summarized in Tab. 4 (*left*). The original PointNet has only 8 layers. For fair comparison, we construct a PointNet with 18 layers, which is denoted by PointNet-18 in Tab. 4. Compared with PointNet-18, using

2D convolution backbones can improve the accuracy of 3D boxes, especially for *hard* setting. This is because these cases are usually occluded/truncated or far away from the camera, and estimating the pose of them is more dependent on context information. However, it is evident that the point-wise CNNs are hard to extract local features of data efficiently. From this perspective, image representation based detectors have greater development potentials. Besides, we can see from Tab. 4 (*right*) that the accuracy does not improve much when the CNN has more layers from ResNeXt-18 to ResNeXt-50. Compared with ResNeXt-50, ResNeXt-101 performs worse, which can be attributed to over-fitting. All the CNNs are trained from scratch.

Table 4. Comparisons of different backbone nets on KITTI *validation* set. Metrics are $AP_{3D|R_{11}}$ for 3D detection task of the **Car** category with IoU threshold = 0.7. Other settings are same as PatchNet-vanilla.

Backbone	Easy	Moderate	Hard
PointNet-18	31.1	20.5	17.0
ResNet-18	33.2	21.3	19.1
ResNeXt-18	33.4	21.2	19.2
SE-ResNet-18	33.7	21.5	19.2

Backbone	Easy	Moderate	Hard
ResNeXt-18	32.7	21.2	19.2
ResNeXt-50	32.9	21.4	17.3
ResNeXt-101	31.1	20.9	17.0

Mask global pooling. In the PatchNet, we design the mask global pooling operation to force the feature maps must be extracted from a set of pixels of interests, which can be regarded as a hard attention mechanism. Tab. ?? shows the effectiveness of this operation, e.g., mask global pooling (max) can improve $AP_{3D|R_{11}}$ by 1.4% for moderate setting and by 2.7% for easy setting, and max pooling is slightly better than avg pooling. Besides, the visualization result shown in Fig. ?? intuitively explains the reason for the performance improvement. Specifically, most activation units filtered by mask global pooling correspond to foreground goals, while the ones from standard global max pooling will have many activation units on the background.

It should be noted that the background points provide contextual information in our model, but they are not involved in [28,35] as input for PointNet.

Instance assignment. We use a stand alone module to predict the ‘difficulty’ of each instance, and assign it to its corresponding head network. Tab. 5.3 shows the ablation study of this mechanism. First, we can find that the accuracy of outputs increases with instance assignment. Interestingly, considering that not all cases we can get the annotations of ‘difficulty’, we use a simple alternative: using the distance from object to camera to represent the ‘difficulty’ of objects (our default setting), and the threshold used in this experiment is (30, 50). Experiment shows that this scheme get a similar performance as predicted difficulty levels.

Table 5. Ablation study of mask global pooling on KITTI *validation* set. Metrics are AP_{3D} and AP_{BEV} of the **Car** category with 11 recall positions. Other settings are same as PatchNet (full model).

pooling type	AP_{3D}			AP_{BEV}		
	Easy	Moderate	Hard	Easy	Moderate	Hard
standard max	32.4	20.6	17.7	41.3	27.0	21.6
mask avg	34.6	21.6	19.3	43.5	28.7	23.3
mask max	35.1	22.0	19.6	44.4	29.1	24.1

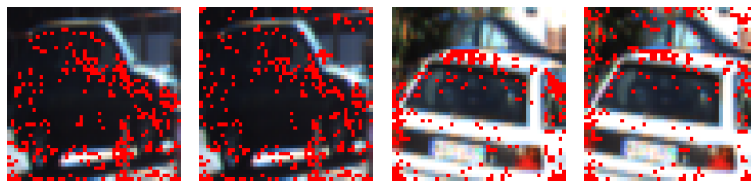


Fig. 4. Qualitative comparison of max global pooling on KITTI *validation* set. The left/right image in each image pair marks the units activated by mask/standard global pooling.

Table 6. Ablation study of instance assignment on KITTI *validation* set. Metrics are AP_{3D} and AP_{BEV} of the **Car** category with 11 recall positions.

assignment	switcher	AP_{3D}			AP_{BEV}		
		Easy	Moderate	Hard	Easy	Moderate	Hard
-	-	33.7	21.5	19.2	42.5	28.2	23.5
✓	difficulty	34.7	22.1	19.5	44.1	29.0	24.2
✓	distance	35.1	22.0	19.6	44.4	29.1	24.1

5.4 Comparing with state-of-the-art methods

As shown in Tab. 7, we report our 3D detection results on the car category on KITTI dataset, where the proposed PatchNet ranks 1st among all published methods (ranked by *moderate* setting). Overall, our method achieves superior result over other state-of-the-art methods across all settings except for *easy* level of *testing* set. For instance, we outperform the current state-of-the-art AM3D [24] by **0.65/1.56/2.34** under *hard* setting on the listed three metrics, which is the most challenging cases in the KITTI dataset. Besides, the proposed method outperforms existing pseudo-LiDAR based approaches. Note we use the same depth estimator (DORN) as [24,35,39,3] and the pipeline of proposed method is much simpler than pseudo-LiDAR based counterparts [3,39]. This shows the

Table 7. 3D detection performance of the **Car** category on KITTI dataset. For *testing* set, only $AP|_{R_{40}}$ is provided by the official leaderboard. For *validation* set, we report both $AP|_{R_{40}}$ and $AP|_{R_{11}}$ for better comparisons. IoU threshold is set to 0.7. * indicates method is based on pseudo-LiDAR data. Methods are ranked by *moderate* setting (same as KITTI leaderboard). We highlight the best results in **bold**.

Method	testing ($AP _{40}$)			validation($AP _{40}$)			validation ($AP _{11}$)		
	Easy	Mod.	Hard	Easy	Mod.	Hard	Easy	Mod.	Hard
OFTNet [33]	1.61	1.32	1.00	-	-	-	4.07	3.27	3.29
FQNet [23]	2.77	1.51	1.01	-	-	-	5.98	5.50	4.75
ROI-10D [25]	4.32	2.02	1.46	-	-	-	10.25	6.39	6.18
GS3D [19]	4.47	2.90	2.47	-	-	-	13.46	10.97	10.38
Shift R-CNN [27]	6.88	3.87	2.83	-	-	-	13.84	11.29	11.08
Multi-Fusion [38]	7.08	5.18	4.68	-	-	-	22.03	13.63	11.60
MonoGRNet [31]	9.61	5.74	4.25	-	-	-	13.88	10.19	7.62
Decoupled-3D* [3]	11.08	7.02	5.63	-	-	-	26.95	18.68	15.82
MonoPSR [18]	10.76	7.25	5.85	-	-	-	12.75	11.48	8.59
MonoPL* [36]	10.76	7.50	6.10	-	-	-	31.5	21.00	17.50
SS3D [17]	10.78	7.68	6.51	-	-	-	14.52	13.15	11.85
MonoDIS [34]	10.37	7.94	6.40	11.06	7.60	6.37	18.05	14.98	13.42
M3D-RPN [2]	14.76	9.71	7.42	-	-	-	20.27	17.06	15.21
PL-AVOD* [35]	-	-	-	-	-	-	19.5	17.2	16.2
PL-FPointNet* [35]	-	-	-	-	-	-	28.2	18.5	16.4
AM3D* [24]	16.50	10.74	9.52	28.31	15.76	12.24	32.23	21.09	17.26
PatchNet	15.68	11.12	10.17	31.6	16.8	13.8	35.1	22.0	19.6

effectiveness of our design. We also observe that proposed model lags behind AM3D [24] under the *easy* setting on *testing* set. This may be attributed to the differences of the 2D detectors. We emphasize that *easy* split contains the least number of examples, so the performance of this setting is prone to fluctuations. Also note that these three splits are containment relationships (e.g., *hard* split contains all instances belong to *easy* and *moderate* setting).

5.5 Qualitative results

We visualize some representative outputs of our PatchNet model in Fig. 5. We can observe that for simple cases in reasonable distance, our model outputs remarkably accurate 3D bounding boxes. Relatively, for distant objects, our estimates of their size and heading angle are still accurate, although it is difficult to determine its center.

On the other hand, we do observe several failure patterns, which indicate possible directions for future efforts. First, our method often makes mistakes with truncated/occluded objects, and often manifests itself as inaccurate heading estimates. Second, sometimes our 2D detector misses object due to strong occlusion, which will cause these samples to be ignored in subsequent steps.



Fig. 5. Qualitative results on KITTI *validation* set. Red boxes represent our predictions, and green boxes come from ground truth. LiDAR signals are only used for visualization. Best viewed in color with zoom in.

6 Conclusions

In this paper, a novel network architecture, namely PatchNet, is proposed to explore the fundamental cause why pseudo-LiDAR representation based 3D detectors achieve promising performance. Different from other works, we argue that the key factor is projecting the image coordinates to the world coordinates by the camera parameters, rather than the point cloud representation itself. More importantly, the world coordinate representation can be easily integrated into image representation, which means we can further boost the performance of 3D detector using more flexible and mature 2D CNN technologies. Experimental results on KITTI dataset demonstrate our argument and show potential of image representation based 3D detector. We hope these novel viewpoints provide insights to monocular/stereo 3D object detection community, and promote the development of new 2D CNN designs for image based 3D detection.

7 Acknowledgement

This work was supported by SenseTime, the Australian Research Council Grant DP200103223, and Australian Medical Research Future Fund MRFAI000085.

References

1. Alhashim, I., Wonka, P.: High quality monocular depth estimation via transfer learning. arXiv e-prints **abs/1812.11941**, arXiv:1812.11941 (2018), <https://arxiv.org/abs/1812.11941>
2. Brazil, G., Liu, X.: M3d-rpn: Monocular 3d region proposal network for object detection. In: The IEEE International Conference on Computer Vision (ICCV) (October 2019)
3. Cai, Y., Li, B., Jiao, Z., Li, H., Zeng, X., Wang, X.: Monocular 3d object detection with decoupled structured polygon estimation and height-guided depth estimation. arXiv preprint arXiv:2002.01619 (2020)
4. Chabot, F., Chaouch, M., Rabarisoa, J., Teuliere, C., Chateau, T.: Deep manta: A coarse-to-fine many-task network for joint 2d and 3d vehicle analysis from monocular image. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 2040–2049 (2017)
5. Chang, J.R., Chen, Y.S.: Pyramid stereo matching network. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 5410–5418 (2018)
6. Chen, X., Kundu, K., Zhang, Z., Ma, H., Fidler, S., Urtasun, R.: Monocular 3d object detection for autonomous driving. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 2147–2156 (2016)
7. Chen, X., Kundu, K., Zhu, Y., Berneshawi, A.G., Ma, H., Fidler, S., Urtasun, R.: 3d object proposals for accurate object class detection. In: Advances in Neural Information Processing Systems. pp. 424–432 (2015)
8. Chen, X., Ma, H., Wan, J., Li, B., Xia, T.: Multi-view 3d object detection network for autonomous driving. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (July 2017)
9. Dai, J., Li, Y., He, K., Sun, J.: R-fcn: Object detection via region-based fully convolutional networks. In: Advances in neural information processing systems. pp. 379–387 (2016)
10. Fu, H., Gong, M., Wang, C., Batmanghelich, K., Tao, D.: Deep ordinal regression network for monocular depth estimation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 2002–2011 (2018)
11. Geiger, A., Lenz, P., Urtasun, R.: Are we ready for autonomous driving? the kitti vision benchmark suite. In: 2012 IEEE Conference on Computer Vision and Pattern Recognition. pp. 3354–3361. IEEE (2012)
12. Girshick, R.: Fast r-cnn. In: Proceedings of the IEEE international conference on computer vision. pp. 1440–1448 (2015)
13. Godard, C., Mac Aodha, O., Brostow, G.J.: Unsupervised monocular depth estimation with left-right consistency. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 270–279 (2017)
14. He, K., Gkioxari, G., Dollár, P., Girshick, R.: Mask r-cnn. In: Proceedings of the IEEE international conference on computer vision. pp. 2961–2969 (2017)
15. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2016)
16. Hu, J., Shen, L., Sun, G.: Squeeze-and-excitation networks. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2018)
17. Jrgensen, E., Zach, C., Kahl, F.: Monocular 3d object detection and box fitting trained end-to-end using intersection-over-union loss. CoRR **abs/1906.08070** (2019), <http://arxiv.org/abs/1906.08070>

18. Ku, J., Mozifian, M., Lee, J., Harakeh, A., Waslander, S.L.: Joint 3d proposal generation and object detection from view aggregation. In: 2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS). pp. 1–8. IEEE (2018)
19. Li, B., Ouyang, W., Sheng, L., Zeng, X., Wang, X.: Gs3d: An efficient 3d object detection framework for autonomous driving. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 1019–1028 (2019)
20. Li, P., Chen, X., Shen, S.: Stereo r-cnn based 3d object detection for autonomous driving. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (June 2019)
21. Lin, T.Y., Dollar, P., Girshick, R., He, K., Hariharan, B., Belongie, S.: Feature pyramid networks for object detection. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (July 2017)
22. Lin, T.Y., Goyal, P., Girshick, R., He, K., Dollár, P.: Focal loss for dense object detection. In: Proceedings of the IEEE international conference on computer vision. pp. 2980–2988 (2017)
23. Liu, L., Lu, J., Xu, C., Tian, Q., Zhou, J.: Deep fitting degree scoring network for monocular 3d object detection. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 1057–1066 (2019)
24. Ma, X., Wang, Z., Li, H., Zhang, P., Ouyang, W., Fan, X.: Accurate monocular 3d object detection via color-embedded 3d reconstruction for autonomous driving. In: The IEEE International Conference on Computer Vision (ICCV) (October 2019)
25. Manhardt, F., Kehl, W., Gaidon, A.: Roi-10d: Monocular lifting of 2d detection to 6d pose and metric shape. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2019)
26. Mousavian, A., Anguelov, D., Flynn, J., Kosecka, J.: 3d bounding box estimation using deep learning and geometry. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 7074–7082 (2017)
27. Naiden, A., Paunescu, V., Kim, G., Jeon, B., Leordeanu, M.: Shift r-cnn: Deep monocular 3d object detection with closed-form geometric constraints. In: 2019 IEEE International Conference on Image Processing (ICIP). pp. 61–65. IEEE (2019)
28. Qi, C.R., Liu, W., Wu, C., Su, H., Guibas, L.J.: Frustum pointnets for 3d object detection from rgb-d data. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2018)
29. Qi, C.R., Su, H., Mo, K., Guibas, L.J.: Pointnet: Deep learning on point sets for 3d classification and segmentation. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 652–660 (2017)
30. Qi, C.R., Yi, L., Su, H., Guibas, L.J.: Pointnet++: Deep hierarchical feature learning on point sets in a metric space. In: Advances in neural information processing systems. pp. 5099–5108 (2017)
31. Qin, Z., Wang, J., Lu, Y.: Monogrnet: A geometric reasoning network for monocular 3d object localization. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 33, pp. 8851–8858 (2019)
32. Ren, S., He, K., Girshick, R., Sun, J.: Faster r-cnn: Towards real-time object detection with region proposal networks. In: Advances in neural information processing systems. pp. 91–99 (2015)
33. Roddick, T., Kendall, A., Cipolla, R.: Orthographic feature transform for monocular 3d object detection. arXiv preprint arXiv:1811.08188 (2018)

34. Simonelli, A., Bulo, S.R., Porzi, L., Lopez-Antequera, M., Kotschieder, P.: Disentangling monocular 3d object detection. In: The IEEE International Conference on Computer Vision (ICCV) (October 2019)
35. Wang, Y., Chao, W.L., Garg, D., Hariharan, B., Campbell, M., Weinberger, K.Q.: Pseudo-lidar from visual depth estimation: Bridging the gap in 3d object detection for autonomous driving. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2019)
36. Weng, X., Kitani, K.: Monocular 3d object detection with pseudo-lidar point cloud. In: IEEE International Conference on Computer Vision (ICCV) Workshops (Oct 2019)
37. Xie, S., Girshick, R., Dollár, P., Tu, Z., He, K.: Aggregated residual transformations for deep neural networks. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 1492–1500 (2017)
38. Xu, B., Chen, Z.: Multi-level fusion based 3d object detection from monocular images. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2018)
39. You, Y., Wang, Y., Chao, W.L., Garg, D., Pleiss, G., Hariharan, B., Campbell, M., Weinberger, K.Q.: Pseudo-lidar++: Accurate depth for 3d object detection in autonomous driving. arXiv preprint arXiv:1906.06310 (2019)
40. Zhou, D., Zhou, X., Zhang, H., Yi, S., Ouyang, W.: Cheaper pre-training lunch: An efficient paradigm for object detection. arXiv preprint arXiv:2004.12178 (2020)

A Overview

This document provides additional analysis and extra experiments to the main paper. Specifically, in Sec. B, we analyse the latency of proposed method and compare it with some pseudo-LiDAR based methods. Sec. C shows the performance of stereo images while Sec. D gives the results of Pedestrian and Cyclist detection. Finally, Sec. E presents more visualization examples.

B Runtime Analysis

In this section, we will analyze the latency of our PatchNet and compare it with some existing methods [24,35,36] based on pseudo-LiDAR representation. In general, all the four methods can be divided into three main stages. In our designs, the processing flows of PatchNet-vanilla and pseudo-LiDAR are the same, but the representations of inputs are different. So the runtime of these two methods are almost the same, which is shown as follow (tested on a single 1080 GPU):

Table 8. Runtime of PachNet-vanilla and pseudo-LiDAR.

2D detection	Depth estimation	3D detection
60ms	400ms	28ms

PatchNet shares the same 2D detector and depth estimator (note the runtime of different depth estimators varies greatly, see **KITTI Benchmark** for details), and we show its runtime of 3D detection stage for different backbone models as follows:

Table 9. Runtime of PatchNet in 3D detection stage.

Backbone	PointNet-18	ResNet-18	ResNeXt-18	SE-ResNet-18
runtime	12ms	23ms	18ms	26ms

Although we add some extra operations in PatchNet, the runtime of the baseline model (PointNet-18) is 12ms while the runtime of PatchNet-vanilla is 28ms. This is mainly because we remove the foreground segmentation net and use a dynamic threshold to segment the foreground, which can save about 18ms. For the best backbone, the runtime is only 26ms, which has similar runtime of pseudo-LiDAR for 3d detection.

Besides, although PatchNet and [24] use the same segmentation method, [24] add another ResNet-34 to extract image features. For [36], it adds a 2D instance segmentation net, which will bring lots of computing overhead (e.g., about 200ms for Mask RCNN [14]).

In summary, PatchNet is more efficient than [24,36] and has the similar run time as [35].

C Stereo Images

Pseudo-LiDAR representation is also widely used in the field of stereo 3D detection task. In order to verify that the proposed method is still work with binocular images, we replace the monocular depth maps with the stereo ones (we use PSMNet [5] as our stereo depth estimator and get the pre-trained model from [35]) and test the performance on KITTI *validation* set using $AP|_{R_{11}}$ for better comparison with previous works. As shown in the Tab. 10, PatchNet-vanilla has almost the same accuracy as pseudo-LiDAR, while PatchNet achieves better performances. We also report the $AP|_{R_{40}}$ for reference.

Table 10. Stereo 3D detection performance of the **Car** category on KITTI *validation* dataset. IoU threshold is set to 0.7. We highlight the best results in **bold**.

Method	3D Detection			BEV Detection		
	Easy	Moderate	Hard	Easy	Moderate	Hard
3DOP [7]	6.55	5.07	4.10	12.63	9.49	7.59
Multi-Fusion [38]	-	9.80	-	-	19.54	-
Stereo-RCNN [20]	54.1	36.7	31.1	68.5	48.3	41.5
Pseudo-LiDAR [35]	59.4	39.8	33.5	72.8	51.8	44.0
PatchNet-vanilla	60.8	40.1	33.6	72.7	51.2	43.8
PatchNet	65.9	42.5	38.5	74.5	52.9	44.8
PatchNet-vanilla@ $AP _{R_{40}}$	61.4	37.6	31.6	73.5	49.8	41.7
PatchNet@ $AP _{R_{40}}$	66.0	41.1	34.6	76.8	52.8	44.3

D Pedestrian and Cyclist

For better comparison, we also report **Pedestrian/Cyclist** detection performance for 3D detection task on KITTI *validation* set in this part. Specifically, we conduct these experiments using both monocular and stereo images with $AP|_{R_{11}}$ as metric. It can be seen from Tab. 11 that the proposed model also get better performance than [35] with each setting. Note that results of pseudo-LiDAR are evaluated by ourselves using its official code, since pseudo-LiDAR did not provide Pedestrian/Cyclist detection results for monocular images.

Besides, the accuracy of **Pedestrian/Cyclist** detection fluctuate greatly compared with **Car** detection. This fluctuation of performance is mainly caused by insufficient training samples (there are only 2,207/734 training samples for **Pedestrian/Cyclist** in KITTI *training* set, while it provides 14,357 **Car** instances). This problem can be reduced by introducing more training data or more effective data augmentation strategies.

Table 11. 3D detection performance of the **Pedestrian/Cyclist** category on KITTI *validation* dataset. Metric is $AP|_{R_{11}}$ and IoU threshold is set to 0.5. We highlight the best results in **bold**.

Method	Category	Monocular			Stereo		
		Easy	Moderate	Hard	Easy	Moderate	Hard
Pseudo-LiDAR [35]	Pedestrian	7.32	6.19	5.64	33.8	27.4	24.0
PatchNet	Pedestrian	9.82	7.86	6.84	38.8	30.1	26.5
Pseudo-LiDAR [35]	Cyclist	5.49	3.85	3.82	41.3	25.2	24.9
PatchNet	Cyclist	8.14	4.84	4.62	46.8	29.0	26.8

E More Qualitative Examples

In this part, we compare the monocular images and stereo pairs by some representative qualitative results in Fig. 6. First, we can find that stereo images can detect objects more accurately, which is generally reflected to the better depth estimation, instead of size or heading estimation. Then, for most of close range objects, in terms of visual experience, monocular images are not inferior to stereo images (although there are still some failure case among those instances).

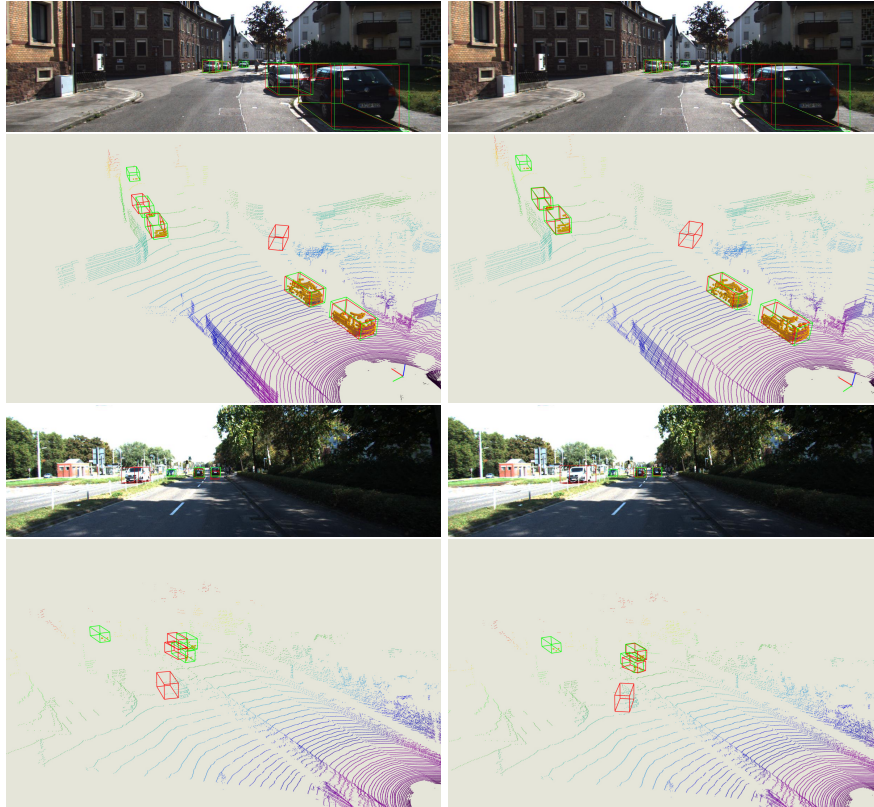


Fig. 6. Qualitative results on KITTI *validation* set. *Left*: monocular detection results. *Right*: stereo detection results. Red boxes represent our predictions, and green boxes come from ground truth. LiDAR signals are only used for visualization. Best viewed in color with zoom in.