

Real-time 3D object proposal generation and classification under limited processing resources

Xuesong Li, Jose Guivant, Subhan Khan

Abstract—The task of detecting 3D objects is important to various robotic applications. The existing deep learning-based detection techniques have achieved impressive performance. However, these techniques are limited to run with a graphics processing unit (GPU) in a real-time environment. To achieve real-time 3D object detection with limited computational resources for robots, we propose an efficient detection method consisting of 3D proposal generation and classification. The proposal generation is mainly based on point segmentation, while the proposal classification is performed by a lightweight convolution neural network (CNN) model. To validate our method, KITTI datasets are utilized. The experimental results demonstrate the capability of proposed real-time 3D object detection method from the point cloud with a competitive performance of object recall and classification.

Index Terms—Point cloud segmentation, object detection, Optimization

I. INTRODUCTION

Robots are gradually undertaking more and more tasks in our daily lives, such as driving [1] and providing service for hospitals [2]. Introducing fully autonomous robots into our daily lives implies that robots share the spatial world with people. Due to this situation, it is imperative for robots to recognize all objects in the surrounding environments and understand the scenes they visualize [3]. Over the past ten years, approaches based on deep learning have achieved remarkable progress in vision-based object detection [4], [5], and 3D object recognition from the point cloud [6], which makes the perception systems in robots more reliable and accurate than ever before. However, the deep learning model requires that a robot be equipped with energy-intensive processing capabilities, such as a GPU, for detecting 3D objects in real-time. Furthermore, a large number of small mobile robots are not able to dedicate those energy resources due to limitations in power supply and energy storage. Fortunately, certain relevant parts of the processing can be saved if proper hybrid approaches are exploited, which implies that power requirements can be reduced. In this work, we propose an efficient 3D object proposal generation and classification approach, which is capable of detecting objects in real-time with limited processing resources.

Three-dimensional object detection can provide accurate spatial location and geometrical shapes of targets, benefiting a robot’s perception system. The traditional approach has been to segment out objects of interest (OOIs) and then apply some

simple classifiers for further classification, such as a naive Bayes classifier [7] or support vector machine (SVM) classifier [8]. With the prevalence of deep learning, many methods based on CNN have been proposed for 3D object detection [9], [10], [11]. Despite their high detection accuracy, the deep learning-based methods require a large amount of annotated data for training purposes, and powerful processors for both training and real-time detection. However, many robots do not usually have powerful processing capabilities, especially small and medium-sized robots. When using these methods, a relevant part of the processing effort is dedicated to extracting point features, and generating 3D proposals. However, unlike 2D proposals in the image, 3D proposals are usually well scattered in 3D space. Taking advantage of this characteristic, we propose an efficient segmentation method for searching for 3D proposals, which reduces the computation effort significantly without sacrificing the quality of the estimated proposals. In addition, an efficient and lightweight deep learning model is developed for learning pointwise features for the classification task. Hence, our detection framework consists of two main parts: proposal generation and proposal classification.

The proposal generation firstly involves a pruning step for the removal of points estimated to be part of the ground. Because it is assumed that the OOIs are not part of the ground, the segmentation process can ignore all those points assumed to belong to that surface. The majority of the approaches used for estimating a road or ground surface assume that the surface is flat or possesses a similarly simple shape; consequently, these approaches attempting to approximate the surface with a plane. However, this assumption is not consistent with the real-world scenario, in which the ground/road contains several slopes within an area. Consequently, to achieve better estimation, the ground/road surface is assumed to be piece-wise multi-linear (PWML) [12], or at least piece-wise constant (PWC). In these cases, the area is divided into smaller but sufficiently large regions, in which dominant almost-horizontal, piece-wise approximating patches are assumed to be part of the ground/road surface in that region of the PWML or PWC partition. A dominant patch takes on the role of a local ground/road reference for classifying points as being part of saliences, depressions in the terrain, or the ground/road.

The segmentation process is based on the continuity of the surfaces. Each inferred segment is considered as a potential proposal, except for some improper segments, which are too large or too small. There are several parameters involved in the proposal generation step; these parameters are tuned once through a process based on particle swarm optimization (PSO) [13]. The efficient PointNet [24], [25] used to perform the

X. Li (e-mail: xuesong.li@unsw.edu.au); J. Guivant (e-mail: j.guivant@unsw.edu.au); S. Khan (e-mail: subhan.khan@unsw.edu.au). All above authors are with the School of Mechanical Engineering, University of New South Wales, Sydney, NSW 2052, AU

classification task consists mainly of multi-layer perceptron (MLP) on each point, which is equivalent to a $1 \times 1 \times 1$ convolutional kernel in 3D, and it can perform the real-time classification on a standard CPU. Moreover, every proposal defines its own size; consequently, no bounding box regression task is required. Hence, the model is simpler and faster. In addition, it requires smaller training datasets than those needed for end-to-end CNN-based detection methods [10], [14]. It is even possible to train the model with limited processing resources, such as those available in a common CPU.

The performance of the proposed method is evaluated using well-known KITTI datasets [29] via the consideration of three metrics: recall of proposals, the accuracy of the classification, and the processing time. Experimental results show that the method can attain high recall and accuracy in a real-time fashion and under the limited processing resources, typically found in small and medium-sized robots. This work makes two main contributions. The first contribution is the approach for inferring the ground, which is approximated via a PWC representation. Compared with other ground extraction method [17], [18], [19] in the traffic context, this representation does not require the assumption that ground in the point cloud is one even plane, and is able to deal efficiently with uneven terrains. The second contribution is the proposed framework for detection, which can achieve real-time performance, e.g. by offering processing times of less than 0.1 seconds per massive point cloud, without utilizing a GPU. Consequently, this approach can be exploited by small robots equipped with simple CPUs with 1 or 2 cores, or by low-scale GPUs.

The rest of the paper is organized as follows. Section II introduces the related work, and Section III illustrates how to build an efficient 3D object detection method. Additional procedures for model optimization are presented in Section IV, while experiments involving the proposed method are described in Section V. Section VI concludes our work and summarizes its contributions.

II. RELATED WORK

The proposed 3D object detection framework can be split roughly into two parts: point cloud segmentation, and proposal classification. Related work concerning the point cloud segmentation will be firstly reviewed. Then a literature review of 3D object detection will be undertaken, to compare the proposed 3D object detection framework with those approaches.

A. Point cloud segmentation

Point cloud segmentation is used to infer objects of interest (OOIs) from a point cloud and is based on the assumption that 3D objects are sufficiently separated and there are no relevant intersections between adjacent objects. There are usually two main procedures involved: ground removal, and the clustering of the remaining set of points. Himmelsbach et al. [15] represented the terrain with a circle of infinite radius, centered around the robot itself, and split this circle into a number of circular sectors with fixed interval angles; Then the ground is searched by extracting horizontal lines at a low altitude,

from the point sets of every segment. The remaining non-ground points are mapped into a 2D occupancy grid image, and clustering is done efficiently by finding connected components of occupied grid cells. Moosmann et al. [16] constructed an undirected graph in which every point is treated as a node. The normal of the local surface is estimated for every node against its neighbouring points, and the similar local surfaces, which either form local convexities or are approximately vertical normal vectors, are merged together to grow the ground plane or obstacle area.

Douillard et al. [17] formulated the ground surface detection as non-ground outlier rejection and proposed a Gaussian Process Incremental Sample Consensus, based on the variance of the height of an outlier from the mean of a Gaussian distribution, to propagate the ground surface labels in an iterative manner to eventually include all ground points, while the remaining non-ground points are subsequently segmented into different clusters. Point segmentation approaches, after the removal of all ground points, often bring about issues of under-segmentation or over-segmentation. To alleviate segmentation errors, Held et al. [18] built a probabilistic model to combine the spatial, temporal, and semantic cues for every segment. The probabilistic model iterates through every individual segment when conducting splitting and through all pairs of segments when merging. Zermas et al. [19] started with an approximate coarse ground plane based on a deterministically extracted set of seed points with low height values, and iteratively refined the ground plane fitting by selecting seed points belonging to previously estimated ground surface. After removing the ground points, line segments are found efficiently for every scan and subsequently merged across neighbouring scans. One of the drawbacks in all these methods is that they include the assumption that one even plane should be able to fit the ground in the point cloud, but this assumption may fail in various real-life situations. In contrast, our ground removal approach adopts multiple small planes to fit the ground surface, resulting in a more robust approach.

B. 3D object detection

The conventional pipeline in 3D object detection involves segmenting a point cloud, then classifying the clusters into objects. After point cloud segmentation, Teichman et al. [20] adopted two classifiers, i.e., the segment and holistic classifiers, to classify feature vectors. The results were then combined with a discrete Bayes filter in a boosting framework. Wang et al. [21] proposed a feature-centric voting algorithm which convolutes the 3D point space in a sliding window with a voting scheme in a manner similar to feature transformation, enabling it to find new features describing object locations and orientations. An SVM classifier was subsequently applied to classify 3D candidates into different categories. As CNN-based approaches have become more and more dominant in computer vision tasks, a number of related 3D object detection approaches have been proposed. Chen et al. [22] proposed multi-view 3D object detection and converted the point cloud into bird-view representations, consisting of multiple height maps, one density map, and one intensity map. The mature

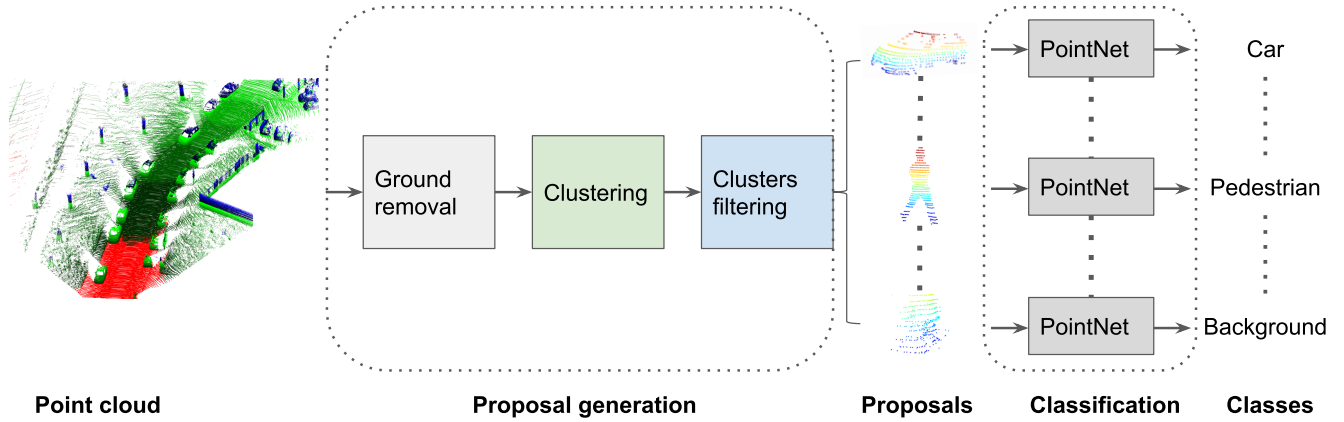


Fig. 1: The framework of our efficient 3D detection approach. Proposals are extracted from the raw point cloud with three procedures, ground removal, clustering, and cluster filtering. The PointNet is then employed to perform a classification for every cluster.

2D CNN framework used for object detection in images was applied to these bird-view images to generate 3D object candidates.

Ku et al. [14] designed AVOD-FPN, which generates 3D proposals on the bird-view images of the point cloud with a 2D CNN, then crops the corresponding features from RGB and bird-view images for each proposal and predicts the 3D bounding boxes based on fused features. Zhou et al. [23] adopted PointNet to extract features from the low resolution raw point cloud for each voxel, then employed a shallow 3D CNN network to convert 3D feature maps into 2D feature maps for 3D object detection. Other methods [26], [10] use sparse CNN [27], [28] to learn 3D feature maps from sparse 3D data, then compress them into 2D feature maps for object detection. Even though a large part of the computation for the sparse CNN is shifted to a CPU, these methods still require a GPU to achieve real-time performance. Although they can achieve desirable detection performance, these deep-learning-based approaches require a large training dataset and powerful processor, and cannot run in real-time without a GPU. To maintain both efficiency and accuracy, our framework employs point cloud segmentation to generate proposals and design an efficient classification model based on a CNN.

III. OBJECT DETECTION METHODOLOGY

The framework of the proposed detection approach can be seen in Fig. 1. It consists mainly of two parts: proposal generation and proposal classification. The proposal generation consists of ground removal, clustering, and cluster filtering. All of these steps are described in the following sections.

A. Ground removal

Pre-processing the data by removing points that belong to the ground has been shown empirically to improve segmentation performance significantly. The typical method employed in urban ground detection is to estimate an approximate 3D plane via model-fitting techniques, such as Random Sample Consensus (RANSAC). However, a simple plane is usually

not sufficient for approximating the ground surface. A more powerful representation involves using a PWML approximation, such as in [12], or even more constrained cases, such as PWC with fixed partitions. The PWC assumes that in each region of the partition, the 'altitude' of the nominal ground is constant (i.e. the nominal surface of the terrain is flat and horizontal). In order to infer the altitude of the nominal surface in a subregion, a histogram of the altitudes of the points whose XY coordinates belong to this subregion can be generated. Additional basic rules, which consider the continuity of the properties of adjacent subregions, enable this process to avoid being misled by plateaus. The PWC requires very low processing effort, and performs better than the more sophisticated RANSAC.

Fig. 2 summarizes ground modelling based on a fixed-size PWC approximation. The coverage in XY ($W \times L$) of the entire point cloud, is partitioned into multiple sub-regions, $W_{sub} \times L_{sub}$. For each subregion, a histogram of points' heights (given simply by the points Z coordinates) is built. It is assumed that the ratio between the number of ground points and the number of non-ground points in every sub-region is larger than a certain user-defined threshold. According to this rule, the minimum bin value, whose frequency in the histogram is over this threshold, is assumed to be ground height. This operation can be implemented efficiently in parallel. The pseudo-code of the algorithm can be found in A. Since the assumption may fail when a sub-region contains the flat plane belonging to an object, such as the roof of a car, we design the post-processing step to amend such false detections. The post-processing will compare the height of every sub-region with those of its adjacent regions, and its ground height value will be updated with the lowest height found in the neighbourhood. Once the ground planes are inferred, points can be classified as either being part of the ground ('close enough') or not. Based on this classification, only the points of interest (those above the ground) are kept for the posterior proposal generation of OOIs.

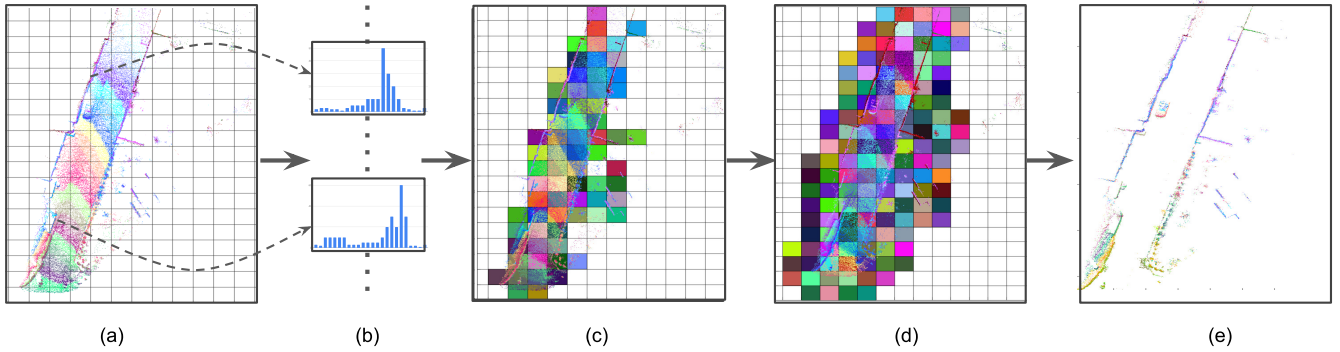


Fig. 2: The pipeline for the ground detection methodology: (a) a bird-view image of the entire original point cloud covering a full local area of size $80m$ by $70m$, every element of the partition is a fixed-size sub-region with size $4m$ by $3.5m$; (b) the histogram of points' heights in every sub-region; (c) and (d) depict the detected ground planes before and after post-processing; (e) depicts the point cloud shown in (a) after the ground points are removed.

B. Point cloud clustering

The remaining points, i.e. the points 'over the ground', belong mostly to the OOIs, such as vehicles, cyclists, pedestrians, and background objects (buildings and plants). These objects are usually well separated when the ground points are removed. Therefore, we can perform clustering on these remaining points to find the potential 3D proposals.

Two methods are used to do point cloud clustering: distance-based and scan-based clustering. For the distance-based clustering method, the only criterion used to decide whether two points belong to the same cluster or not is their Euclidean distance. Hence, the distances between all point pairs are calculated and compared with a pre-defined distance threshold T_d . A random point P_0 is first selected and put into a separate list C . Then, the distances of other points from P_0 are calculated; Those points with distances from P_0 less than T_d , are removed and added to the same list C . Similar operations are performed repeatedly between points in C and the rest of the points until no more points are added to list C , and C is treated as a cluster. After this first cluster is complete, a new list will be created to search for the next cluster and so on until no more points are left. The algorithm details are shown in A. Distance-based clustering, which does not require the ordering information among points, can be generalized to any point cloud dataset. However, processing a large number of point clouds takes a considerable amount of time, since every point is sequentially required to do a comparison against the rest of the points.

The 3D point cloud, generated by the LiDAR, has a multi-layer structure and every layer, also called a scan, produced from the same LiDAR ring includes contiguous points. Based on this characteristic, the scan-based clustering method is proposed to speed up processing. The points in every scan, or single layer, are first segmented into line segments, then these line segments are compared with their vertical neighbouring segments. Finally, the line segments which are close together are merged into one segment.

It is assumed that points traverse in a raster in a counter-clockwise fashion, starting from the top scan-line. H_d is used to decide whether two points in the same scan belongs to the same segment or not, and V_d is set to distinguish two points in

neighbouring scans. We save all clusters, called *Global_segs*, in the form of dictionary variables, in which the key and value are label number and corresponding points, respectively. All line segments begin with the first scan and are saved into a dictionary with the label as the key and the point group as the value, called *Above_segs*.

The first *Above_segs* is appended to *Global_segs*. Then line segments are extracted from the second scan and saved in another dictionary, *Current_segs*. Next, every line segment seg_i in *Current_segs* is compared with all segments in *Above_segs*. If *Above_segs* has no segment whose distance from seg_i is smaller than V_d . Then, the key for seg_i in *Current_segs* is assigned a new label number, and this new key and its points are added to *Global_segs*. If *Above_segs* includes one connection with seg_i , no new key is generated, and the key for seg_i is the same as the one in *Above_segs*; all points in seg_i are appended to the point group with the same key in the *Above_segs* and *Global_segs*. If *Above_segs* includes more than two connections with seg_i , the minimum key for the connected segments in *Above_segs* is selected as the unique key for these connected clusters.

Merging operations are implemented in the *Above_segs* and *Global_segs*. After processing a current scan, the *Current_segs* becomes the *Above_segs* and the algorithm starts to process the next scan. The pseudo-code of scan-based clustering can be found in A. The author in [19] designed a complicated label management system to handle label conflict and merge issues, while this paper employs the dictionary to save the points in all clusters directly, making our algorithm easier to implement.

In contrast to the distance-based clustering method, the scan-based approach takes the advantage of the ordering of the points along the LiDAR scan-lines and convert clustering problem in 3D space into the line segmentation, and only the Euclidean distance between one point and its neighbouring points along the scan line is required, resulting in much lower algorithm complexity. Therefore, it can achieve a real-time running performance. But the scan-based approach is only applicable to point cloud which is generated by the Lidar or other sensors with fixed multi-layer scanning structure, and is

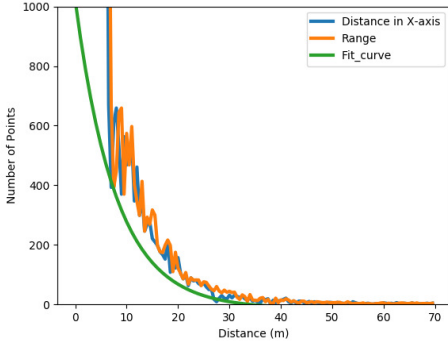


Fig. 3: The number of points within a proposal with respect to distance. The blue curve represents the number of points versus the distance along the X-axis to LiDAR, while the red curve represents the number of points against the range. The green curve is the fitted curve.

not as generalized as the distance-based method.

C. Proposals filtering

Three-dimensional proposals, generated by the clustering method, could also be for some background objects, such as super-large and tiny objects, which could be the buildings or leaves, respectively. To remove some improper background proposals and reduce the number of proposals, we set up two criteria to filter background proposals: 1) the size of the desired proposal should be within a certain range, and 2) the number of points in the non-occluded proposals should have a minimum. For the first criterion, we choose the maximum values for length and width to filter out very large objects, and the minimum value for height is defined to remove objects which are too flat, since objects of interest, such as pedestrians, cars, and cyclists, have certain height value.

As for the second criterion, the occlusion should be inferred first. The horizontal angle span of every 3D proposal is calculated against the location of the LiDAR and is compared with the angle spans of other proposals. A proposal is labelled as a non-occlusion proposal if it satisfies one of the following cases: 1) its angle scan does not overlap those of other proposals, and 2) its angle scan has overlapped with those of other proposals, but it is closer to the LiDAR than the other proposals. The details of this occlusion labelling algorithm are presented in A. To estimate the minimum number of points in non-occluded proposals, we investigate the relationship between the number of points and the distance from proposals to LiDAR, as shown in Fig. 3. The number of points and the distance to LiDAR are collected from all the ground truth bounding boxes in the training dataset provided by the KITTI benchmark [29]. The distance, which can be represented as either the range or the distance along X-axis to LiDAR, is sampled with an interval of 0.5 m , and its corresponding number of points is the minimum number of points in all the ground truth bounding boxes, whose centre points are located in the distance interval. The exponential function is used to

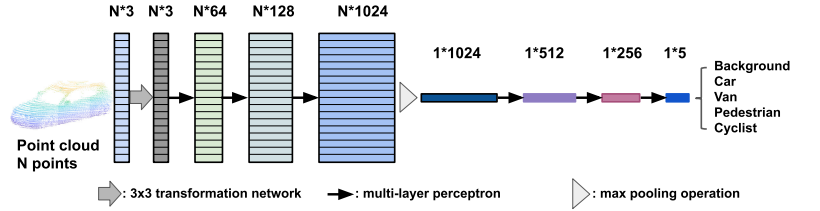


Fig. 4: Classification network. A MLP is applied on each point to increase the feature channel, and the parameters in the same layer is shared. Then max pooling is used to extract the global feature vector and output the final class scores

fit these sample points, see the green line in Fig. 3. Three-dimensional proposals with fewer than the minimum number of points are discarded.

D. Proposals classification

The PointNet-based, efficient, and lightweight classification model is designed to identify the class of 3D region proposals [24]. There are numerous techniques available in the existing literature for processing raw point clouds in every proposal for classification, such as 3D CNN on a 3D voxelization grid [26], or a 2D CNN for bird-view images [22]. In contrast to these methods, the PointNet can consume the raw points directly without altering the data representation, and it also employs an efficient and shared MLP to process the features on each point independently (i.e. pointwise $1 \times 1 \times 1$ CNN) and max-pooling operation, which can achieve real-time performances under constrained processing resources. Therefore, we choose the PointNet to build our proposal classification network, as shown in Fig. 4.

In terms of the classification model, the input consists of 3D coordinates (x, y, z) along with N points for one proposal. The raw point-cloud can be encapsulated into an $N \times 3$ matrix. The affine transformation on the point set (such as rotation) should not change the classification results for a geometric object. Therefore, a 3×3 transformation matrix generated by a spatial transformer network is first learned to transform the input points in order to make the model invariant to certain transformations. This type of transformation matrix is known as input 'TNet' [24]. The PointNet in [24] also includes a 64×64 feature transformation matrix to align the features, but this matrix and corresponding feature 'TNet' is omitted in our classification model, as it increases a lot of computation with limited accuracy gain. To increase the channel features from 3 to 1024, three layers of MLP are applied with weights shared among all the N point features. A max-pooling network is used to aggregate point features in order to generate 1×1024 global feature vector, and three layers of MLP are applied in order to decrease the channel features from 1024 to 256, then outputs the classification scores for 5 classes: background, car, pedestrian, van, and cyclist. In addition, batch normalization [30] is applied to all layers via ReLU [31] and dropout layers [32] are used for the fully-connected neural network.

IV. PARAMETERS LEARNING

The proposed method consists of two sets of parameters: the parameters (P) in the proposal segmentation model, and the weights (W) in the classification model. To learn P , the PSO is applied, while gradient descent is employed to optimize W . This section will cover the optimization of these parameters.

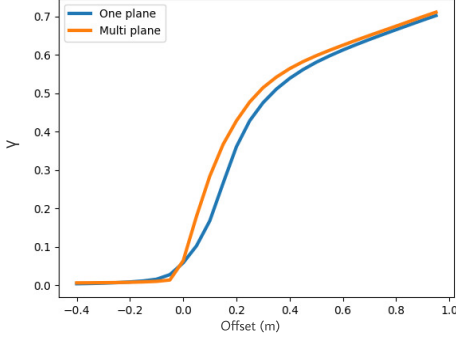


Fig. 5: The ratio between the number of removed points and all points versus the offset distance.

A. Proposal segmentation optimization

The parameters in the proposal segmentation model include the distance threshold H_d for the line segmentation done in one scan, merging threshold V_d between adjacent scans, and filtering distance D_o to the ground plane. The objective of the proposed segmentation method is to achieve the highest possible recall of objects via using PSO to optimize parameters. The PSO searches the space of an objective function by adjusting the trajectories of individual agents, called particles. Every particle, x_i , represents a solution, namely, $x_i = [H_d, V_d, D_o]$ in our case, and the objective function is the recall of objects, as shown in Equation (1).

$$R = \frac{TP}{TP + FN}, \quad (1)$$

where TP is the number of true positives, i.e. the correctly segmented objects, and FN is the number of false negatives, i.e. the missed objects.

The movement of a swarming particle consists of two components: the stochastic and deterministic movements. Each particle is attracted towards the position of a current global best g^* and its own best location x_i^* in history, while, at the same time, it tends to move randomly. For each iteration, the particle x_i^t is updated by Equations (2) and (3). The pseudo-code of PSO can be found in A.

$$x_i^{t+1} = x_i^t + v_i^{t+1} \quad (2)$$

$$v_i^{t+1} = \alpha \cdot v_i^t + \lambda \cdot r_1 \cdot (g^* - x_i^t) + \theta \cdot r_2 \cdot (x_i^* - x_i^t), \quad (3)$$

where α , λ and θ are acceleration constants for different terms, and r_1 and r_2 are the random number used to simulate diversity in the quality solutions.

B. Classification model learning

Considering that the number of labels in the different categories is balanced, we define a naive classification loss function, i.e. negative log-likelihood loss, with the following equation:

$$L_{cls} = -\frac{1}{N} \sum_{i=1}^N \log(\hat{y}_i), \quad (4)$$

where N is the number of training samples in one batch, and \hat{y}_i is the predicted probability of i belonging to the labelled class.

Since there exists a differentiable equation between the loss function (4) and parameters W , the gradient-based optimization method can be used to optimize the parameters W . The gradient descent method, i.e. AdamOptimizer [33], is applied to do the optimization in our classification model with an initial learning rate of 0.0002 and an exponential decay factor of 0.8 for every 18570 steps.

V. EXPERIMENTAL RESULTS

A. Implementation details

1) *Dataset*: The KITTI benchmark dataset is adopted to validate our proposed method. It includes 7481 training pairs for the 3D LIDAR point cloud, RGB image, ground truth annotation, and calibration parameters. The optimization of the PSO and classification model training are conducted on this training dataset using different split ratios between the training and validation datasets.

2) *Detection method details*: The offset distance D_o to the detected ground plane for removing ground points is set to 0.26 m. For the task of clustering the point cloud, the T_d in Euclidean-based clustering is assigned the value 0.5 m, and the H_d and V_d in the scan-based method are also specified as 0.49 and 0.58 m, respectively. These parameters are optimized by the PSO. When calculating the objective fitness, i.e. recall value, the intersection over union (IoU) threshold for true positives is set to 0.25. For the ground removal, the user-defined ground ratio threshold for finding the ground height in every subregion is 0.05, and the width of bin in histogram is 0.15 m.

For the classification task, the number of points in one proposal is randomly sampled to 100 points, and all points are normalized by subtracting the central point from each other and then dividing this difference by the furthest distance from the center. All convolutional kernel sizes are $1 \times 1 \times 1$ in classification model, 3 layers of CNNs are used to increase the channel number of features, and the corresponding feature maps sizes are 300×64 , 300×128 and 300×1024 . After the max-pooling operation, 3 layers of fully connected networks compress the feature (1×1024 to 1×5) through two middle features (1×512 and 1×256).

3) *Hardware platform*: The processors in our hardware platform are 8 Intel(R) Xeon(R) CPU E5-1620 v3 @ 3.50GHz, and every core has 2 threads. The memory capacity is 16 Gigabyte. The GPU card is a NVIDIA TITAN Xp with a 12 Gigabyte memory.

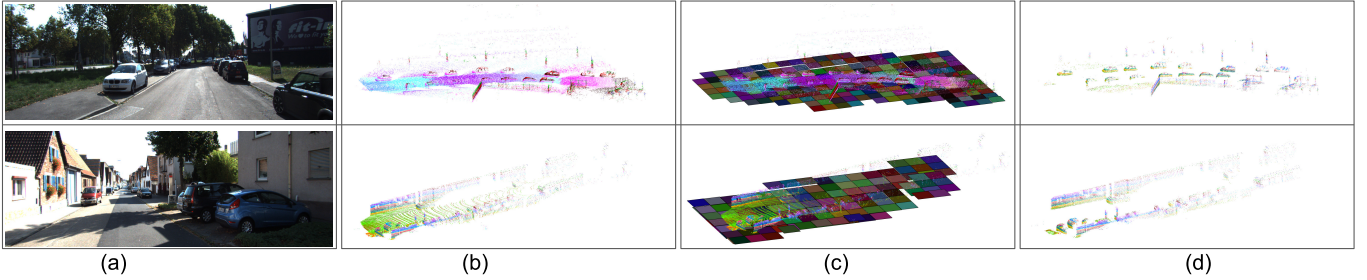
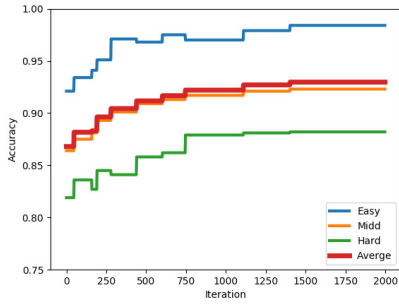
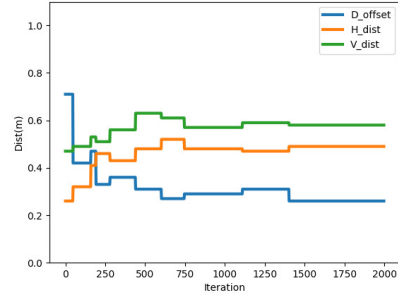


Fig. 6: Visualization of the ground removal: (a) RGB image, (b) raw point cloud, (c) raw point cloud with ground planes, and (d) point cloud with ground points removed.



(a) The convergence graph of recall versus iteration.



(b) The convergence graph of the global particle versus iterations.

Fig. 7: The convergence graphs for the PSO algorithm.

Because this efficient detection method is proposed for robots with limited processing capabilities, we set an affinity mask to restrict our code to one core when generating the experimental data to make sure that our method can run well on small or medium robot platforms.

B. Ground removal

In order to compare our ground removal approach with RANSAC qualitatively, we set up a criterion γ , i.e. the ratio between the number of removed points N_g and the number of all points N_a , see Equation (5). The majority of points can be removed if the plane fits the real landscape of the ground plane. Therefore, given the same offset distance D_o to the detected ground plane, the higher γ is, the more effective the method is at removing the ground points. The qualitative experimental results and visualization of the ground removal performance can be found in Fig. 5 and 6.

$$\gamma = \frac{N_g}{N_a} \quad (5)$$

The γ is calculated by averaging all the ratio values of 1000 point clouds randomly selected from the training dataset. From Fig. 5, we can see that, when the offset distance is small, such as $0 < D_o < 0.3$ m, the γ ratio for the proposed method is higher than that for RANSAC. This finding validates the effectiveness of the proposed method, as it can effectively remove more points on the ground surface and match the detected ground planes of the real landscape. In the case of

a larger offset, it does not matter whether the detected plane fits the landscape closely or not. Thus, the ratios of the two methods converge towards each other.

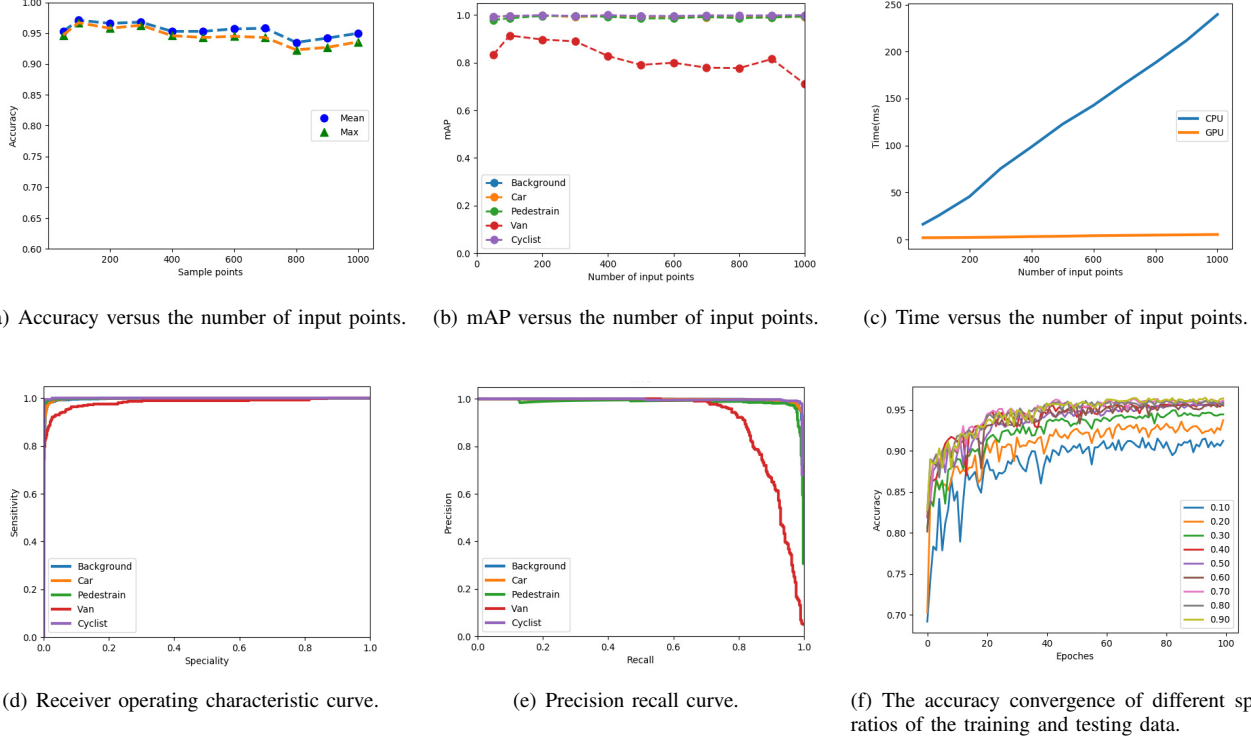
C. Proposal generation

TABLE I: The performance of the proposal method for different combinations.

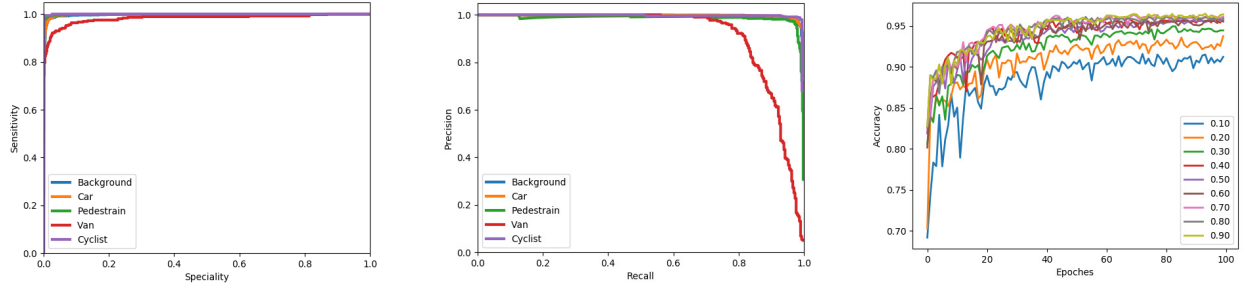
Method modality	Recall(%)	Number	Time (s)
Clustering_1+Filtering	92.9	55	0.038
Clustering_1	94.0	146	0.032
Clustering_2+Filtering	93.3	78	2.351
Clustering_2	94.5	212	2.344

Clustering_1 uses the scan-based clustering method, and clustering_2 uses the distance-based approach (the time spent removing the ground is also included in clustering), the second column indicates the average recall value for each combination, and the third column lists the average number of proposals produced by each combination.

In order to make experimental data, which is meaningful for robots with limited processing capabilities, we force our code into one core when running the experiments. Proposal generation consists of three steps: ground removal, clustering, and proposal filtering. For the scan-based clustering method, the processing times for these steps are 0.002, 0.03, and 0.006 s, respectively. More experimental data is presented in Table I, which shows that the scan-based method is much faster than the distance-based one. The distance-based method is slower because it calculates the distances from one point to all available points and searches a large number of adjacent points



(a) Accuracy versus the number of input points. (b) mAP versus the number of input points. (c) Time versus the number of input points.



(d) Receiver operating characteristic curve. (e) Precision recall curve. (f) The accuracy convergence of different split ratios of the training and testing data.

Fig. 8: The analysis of the performance of the classification model. (a) Accuracy distribution ('Max' represents the best accuracy selected between 100 epochs, while the 'Mean' is calculated by averaging the accuracy of the last 50 epochs) versus the number of sample points. (b) The mAP of every class changes with the number of input points. (c) Running time performance on the CPU and GPU devices, given a batch size of 32. (d) The ROC plot for every category of detected objects. (e) The PRC plot for every detected class. (f) The split ratio of the training and testing data affects the speed of convergence and classification accuracy.

in each iteration. The performance of the proposal generation is very similar to that of the distance-based method, since both of them employ distance as the discriminative criterion regardless of the clustering mechanism used.

The proposal generation method with clustering has a higher recall (%) and three times the number of proposals compared to the combination of clustering and filtering, since the filtering procedure can filter out a number of false positives and reduce the number of proposals. However, it also removes some true positives and leads to a decrease in the recall value. In terms of the processing time, the clustering takes up to 0.032 s, while the filtering only consumes 0.006 s and can decrease the number of proposals greatly. A large number of proposals will require a long processing time for the next classification task. Therefore, for a task involving limited computational resources, the filtering procedure is very important for achieving a real-time performance.

To enhance the performance of the proposal generation method, D_o , H_d , and V_d , are optimized using the PSO algorithm. For every generation in the PSO, the recall value is calculated by running the proposal generation algorithm through 500 training samples randomly selected from the entire training dataset. A total of 1000 generation and 50 particles are set for optimization. The convergences of the recall value and

all parameters are illustrated in Fig. 7. All parameters are randomly initialized between 0 and 1.2 m. From Fig. 7 (b), we can see that the parameters fluctuate at the beginning, then slowly converge to their optimal values. Since these parameters are constrained by each other, the combination of parameter values will decide the final fitness value. All parameter convergences fluctuate instead of moving in one direction. The convergence graph shows that the best values for D_o , H_d , and V_d are 0.26, 0.49, and 0.58 m, respectively.

D. Classification model

To reduce the overfitting and generalize the model, two basic techniques for data augmentation are introduced:

a) *Rotation*: Every point cloud in the entire block is rotated along the Z-axis and about the origin [0,0,0] by θ , where θ is drawn from the uniform distribution $[-\pi/4, \pi/4]$.

b) *Scaling*: Similar to the global scaling operation used in image-based data augmentation, the whole point cloud block is zoomed in or out. The XYZ coordinates of the points in the whole block are multiplied by a random variable drawn from the uniform distribution [0.95, 1.05].

We have investigated how the number of sample points N in a proposal affects the classification accuracy and demonstrate their relationship in Fig. 8 (a). It can be observed that the



Fig. 9: Some samples of our detection approach. The bottom row shows the detection results in the point cloud, and the top row depicts the corresponding bounding boxes. The 2D bounding boxes in the images are generated by projecting 3D bounding boxes in Lidar coordinates into the image coordinates using the calibration parameters provided in the KITTI benchmark [29]. Here, the pedestrians are each denoted with a blue bounding box, the cyclists are each represented by a yellow bounding box, the cars are each in red bounding box.

TABLE II: The performance of the proposal method for different combinations against other methods.

Method	Modality	Recall(%)	Number of proposal	Processing time (s)	
				GPU	CPU
Clustering+Filtering+Classification	Lidar	92.9	55	0.041	0.082*
Clustering+Classification	Lidar	94.0	146	0.035	0.148*
AVOD-FPN [14]	Lidar+Camera	95.2	100	0.1	3.6
Voxelnet [23]	Lidar	–	–	0.23	5.9
SECOND [10]	Lidar+Camera	–	–	0.05	4.8
MV3D [22]	Lidar+Camera	94.4	50	0.36	11.6

The third column lists the average recall values at IoU threshold of 0.25 and the fourth column contains the average number of proposals produced by the method; the times spent on the GPU and CPU are the average times spent processing one point cloud; the time with * is generated with only one core of CPU, otherwise, the 8 cores of CPU are used to report the time. The clustering method used is scan-based clustering. '–' indicates unavailable data.

best accuracies of 0.971 (Maximum) and 0.967 (Mean) are achieved when $N=100$ for one proposal. The classification performance could deteriorate if more points (such as 1000 or more) are added to a proposal. The training samples contains 18000 background objects, i.e. 14357 cars, 1297 vans, 734 cyclists, and 2207 pedestrians, and Fig. 8 (b), (d), and, (e) show the mean Average Precision (mAP) curve, the Receiver Operating Characteristic (ROC) curve, and the Precision-Recall Curve (PRC), respectively, for these samples. These curves demonstrate the detailed classification performance of our model for every class. It can be seen that the classification performance for vans is worse than those for other classes. Furthermore, the scores for the background, car, pedestrian, and cyclist are similar. The shape of the van is similar to that of the car, while the number of vans in the training sample is much smaller than the number of the car. Therefore, the classification model for identifying the van cannot be well-trained, leading to the poor classification performance for the van class. As for the cyclist class, even though the number of cyclists in the training sample is small, the shape of this class is distinct from those of other classes, which makes it easy for the classification model to capture its features. Thus, the model achieves a good classification performance for the cyclist class.

The running time for processing 32 proposals versus different numbers of input points can be found in Fig 8 (c). The implementation time for the CPU increases almost linearly with the number of sample points, while the GPU time is almost constant around 0.003 s due to the parallel computations. We find that it takes 0.0254 s to process 32 proposals with 100

points, making it feasible to run our classification model on the CPU-only hardware platform in real-time.

To capture how many data points are needed to train our model, we split the dataset into training and testing sets using different ratios, and the results can be found in Fig. 8 (f). For this experiment, the number of points per proposal is 100, and we find that the classification performances are similar when the split ratio is larger than 0.3. Other end-to-end training detection methods, such as SECOND [10] and 3DBN [26], require a minimum 0.5 split ratio to guarantee consistent performance. In comparison, our framework for detection required less training data and had better capability in generalization because most of the background is removed by the segmentation algorithm, and most of the noise is peeled off, making it easier for this classification model to learn different classes.

E. Overall performance

To evaluate the performance, we assemble two sub-modules, i.e. proposal generation and classification, together into an object detection method. In Fig. 9, the overall detection performance is presented¹. It can be observed that we can find bounding boxes which fit nearby objects well, but the bounding boxes are much smaller than the real sizes of the objects that are far away from the Lidar. For instance, in the top middle image, the 2D bounding boxes for cars, which are far away, only cover the bottom halves of these cars because

¹More visualization results can be found in videos: <https://youtu.be/poSbDQ1LCR0>, <https://youtu.be/7C4kOnLrtig>, <https://youtu.be/QhXHMC4wsMc>, https://youtu.be/LrXg3J_4FKY

the density of points on the surface of the object is inversely proportional to the distance between the object and Lidar. The proposal of a nearby object includes a number of points and these points can represent its real shape, while, for a far object, there are only a small number of points and only a part of the object can be captured.

The comparison of the computational efficiency of our method with other 3D detection methods is provided in Table II, which shows that the proposed method is much faster than other existing methods and can achieve a real-time performance on either a GPU (0.035 s) or CPU (0.082 s) hardware platform. As for the performance of the object proposal, AVOD-FPN [14] and MV3D [22] can obtain better recall value than ours, but their efficiency is very low, especially on the CPU platform, since they employ the CNN backbone networks to extract features images for generating proposals.

VI. CONCLUSION

In this paper, we have designed an efficient 3D object detection approach for robots with limited computational resources consisting of proposal generation and a lightweight classification model. The existing deep learning-based 3D object detection methods can achieve the desired detection performance but require the robots to be equipped with a powerful processing unit, which is impossible in many small robots. In comparison to existing deep-learning techniques, the proposed method can achieve competitive recall values and classification accuracies. Most conspicuously, the proposed method has the capability of detection in real-time with either a GPU or CPU. In future work, we intend to fuse our work with a semantic segmentation task and combine the information from different resources for some complicated tasks, such as instance segmentation.

VII. ACKNOWLEDGMENT

The presented research was supported by the China Scholarship Council (Grant No. 201606950019).

REFERENCES

- [1] C. Urmsion, J. Anhalt, D. Bagnell, C. Baker, R. Bittner, M. Clark, J. Dolan, D. Duggins, T. Galatali, C. Geyer *et al.*, "Autonomous driving in urban environments: Boss and the urban challenge," *Journal of Field Robotics*, vol. 25, no. 8, pp. 425–466, 2008.
- [2] A. G. Ozkil, Z. B. Fan, S. Dawids, H. Aanes, J. K. Kristensen, and K. H. Christensen, "Service robots for hospitals: A case study of transportation tasks in a hospital," *IEEE International Conference on Automation and Logistics*, pp. 289–294, 2009.
- [3] A. Geiger, M. Lauer, C. Wojek, C. Stiller, and R. Urtasun, "3D traffic scene understanding from movable platforms," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 36, pp. 1012–1025, 2014.
- [4] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *IEEE conference on computer vision and pattern recognition*. IEEE, 2009, pp. 248–255.
- [5] X. Li., N. Kwok., J. E. Guivant., K. Narula., R. Li., and H. Wu., "Detection of imaged objects with estimated scales," in *Proceedings of the 14th International Joint Conference on Computer Vision, Imaging and Computer Graphics Theory and Applications - Volume 5: VISAPP, INSTICC*. SciTePress, 2019, pp. 39–47.
- [6] A. Dai, A. X. Chang, M. Savva, M. Halber, T. A. Funkhouser, and M. Nießner, "ScanNet: Richly-annotated 3D reconstructions of indoor scenes," *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 2432–2443, 2017.
- [7] I. Rish, "An empirical study of the naive Bayes classifier," *IJCAI 2001 Work Empir Methods Artif Intell*, vol. 3, 01 2001.
- [8] M. A. Hearst, "Support vector machines," *IEEE Intelligent Systems*, vol. 13, no. 4, pp. 18–28, Jul. 1998.
- [9] B. Graham, M. Engelcke, and L. van der Maaten, "3D semantic segmentation with submanifold sparse convolutional networks," *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 9224–9232, 2017.
- [10] Y. Yan, Y. Mao, and B. Li, "SECOND: Sparsely embedded convolutional detection," *Journal of Sensors*, vol. 18, no. 10, 2018.
- [11] B. Li, T. Zhang, and T. Xia, "Vehicle detection from 3D Lidar using fully convolutional network," *arXiv preprint arXiv:1608.07916*, 2017.
- [12] A. Robledo, S. Cossell, and J. Guivant, "Outdoor ride: Data fusion of a 3D kinect camera installed in a bicycle," in *Proceedings of Australasian Conference on Robotics and Automation, Australia*, 05 2011.
- [13] N. M. Kwok, Q. P. Ha, D. Liu, G. Fang, and K. C. Tan, "Efficient particle swarm optimization: a termination condition based on the decision-making approach," *IEEE Congress on Evolutionary Computation*, pp. 3353–3360, 2007.
- [14] J. Ku, M. Mozifian, J. Lee, A. Harakeh, and S. L. Waslander, "Joint 3D proposal generation and object detection from view aggregation," *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 1–8, 2017.
- [15] M. Himmelsbach, F. von Hundelshausen, and H.-J. Wünsche, "Fast segmentation of 3D point clouds for ground vehicles," *IEEE Intelligent Vehicles Symposium*, pp. 560–565, 2010.
- [16] F. Moosmann, O. Pink, and C. Stiller, "Segmentation of 3D Lidar data in non-flat urban environments using a local convexity criterion," *IEEE Intelligent Vehicles Symposium*, pp. 215–220, 2009.
- [17] B. Douillard, J. Underwood, N. Kuntz, V. Vlaskine, A. Quadros, P. Morton, and A. Frenkel, "On the segmentation of 3D Lidar point clouds," in *IEEE International Conference on Robotics and Automation*, May 2011, pp. 2798–2805.
- [18] D. Held, D. Guillory, B. Rebsamen, S. Thrun, and S. Savarese, "A probabilistic framework for real-time 3D segmentation using spatial, temporal, and semantic cues," in *Robotics: Science and Systems (RSS)*, June 2016.
- [19] D. Zermas, I. Izzat, and N. Papanikolopoulos, "Fast segmentation of 3D point clouds: A paradigm on Lidar data for autonomous vehicle applications," in *IEEE International Conference on Robotics and Automation (ICRA)*, May 2017, pp. 5067–5073.
- [20] A. Teichman, J. Levinson, and S. Thrun, "Towards 3D object recognition via classification of arbitrary object tracks," *IEEE International Conference on Robotics and Automation*, pp. 4034–4041, 2011.
- [21] D. Z. Wang and I. Posner, "Voting for voting in online point cloud object detection," in *Robotics: Science and Systems*, 2015.
- [22] X. Chen, H. Ma, J. Wan, B. Li, and T. Xia, "Multi-view 3D object detection network for autonomous driving," *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 6526–6534, 2016.
- [23] Y. Zhou and O. Tuzel, "Voxelnet: End-to-end learning for point cloud based 3D object detection," *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 4490–4499, 2017.
- [24] C. R. Qi, H. Su, K. Mo, and L. J. Guibas, "Pointnet: Deep learning on point sets for 3D classification and segmentation," *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 77–85, 2016.
- [25] C. R. Qi, L. Yi, H. Su, and L. J. Guibas, "Pointnet++: Deep hierarchical feature learning on point sets in a metric space," *arXiv preprint arXiv:1706.02413*, 2017.
- [26] X. Li, J. Guivant, N. Kwok, Y. Xu, R. Li, and H. Wu, "Three-dimensional backbone network for 3D object detection in traffic scenes," *CoRR*, vol. abs/1901.08373, 2019.
- [27] B. Graham and L. van der Maaten, "Submanifold sparse convolutional networks," *arXiv preprint arXiv:1706.01307*, 2017.
- [28] B. Graham, "Sparse 3D convolutional neural networks," in *Proceedings of Conference on the British Machine Vision Association (BMVA)*, 2015, pp. 150.1–150.9.
- [29] A. Geiger, P. Lenz, and R. Urtasun, "Are we ready for autonomous driving? the KITTI vision benchmark suite," in *Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2012, pp. 3354–3361.
- [30] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," in *Proceedings of the 32nd International Conference on International Conference on Machine Learning - Volume 37*, 2015, pp. 448–456.
- [31] V. Nair and G. E. Hinton, "Rectified linear units improve restricted Boltzmann machines," in *Proceedings of the 27th International Conference on International Conference on Machine Learning*, ser. ICML'10. USA: Omnipress, 2010, pp. 807–814.

- [32] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, "Dropout: A simple way to prevent neural networks from overfitting," *J. Mach. Learn. Res.*, vol. 15, no. 1, pp. 1929–1958, Jan. 2014.
- [33] D. Kingma and J. Ba, "Adam: A method for stochastic optimization," *International Conference on Learning Representations*, vol. abs/1412.6980, 12 2014.

APPENDIX

Algorithm 1 Output ground planes *Ground_post*

Legend:

W : the width of whole region;
 L : the length of whole region;
 w_{sub} : the width of sub region;
 l_{sub} : the height of sub region;
 $points$: all points in the format [x, y, z] with the size $n \times 3$;
 $ground_num$:
the minimum number of points in ground bin;
 bin_width :
the width of bin in the histogram;
 $subregion_func()$:
the function of getting all points in given sub_region;
 $histogram_build()$:
the function to build histogram;
 $neighbourhood_height_func$:
the function to compare the height in neighbourhood and return the lowest height value;

```

Ground = array( $w_{sub}, l_{sub}$ );
Ground[:, :] = 1000;
for  $i = 1 : ceil(\frac{W}{w_{sub}})$  do
  for  $j = 1 : ceil(\frac{L}{l_{sub}})$  do
     $points_{sub} = subregion\_func(points, i, w_{sub}, j, l_{sub})$ ;
     $hist = histogram\_build(points_{sub}[:, 2], bin\_width)$ ;
     $frequency_z = hist.frequency$ ;
     $bin_z = hist.bin$ ;
    for  $k = 1 : len(frequency_z)$  do
      if  $frequency_z[k] \geq ground\_num$  then
         $Ground[i, j] = height$ ;
        break;
      end if
    end for
  end for
end for
Ground_post = array( $w_{sub}, l_{sub}$ );
Ground_post[:, :] = 1000;
for  $i = 1 : ceil(\frac{W}{w_{sub}})$  do
  for  $j = 1 : ceil(\frac{L}{l_{sub}})$  do
     $low\_height = neighbourhood\_height\_func(i, j, Ground)$ 

     $Ground\_post[i, j] = low\_height$ ;
  end for
end for

```

Algorithm 2 Output clustering results *Cluster*

Legend:

T_d : the minimum distance threshold;
 $points$: all the point cloud, $n \times 3$;
 $dist_func$: distance calculation, $n \times n$;

```

Clusters  $\leftarrow \{\}$ : create empty cluster
 $D_m \leftarrow dist\_func(points)$ 
while  $D_m$  is not empty do
   $stop \leftarrow False$ 
   $C \leftarrow []$ 
  C append 1
   $index\_new \leftarrow 1$ 
  while not stop do
     $index\_t \leftarrow fun\_index\_dist(index\_new, T_d, D_m)$ 
     $index\_t \leftarrow index\_t - C$ 
     $index\_new \leftarrow index\_t$ 
    if  $index\_new$  is empty then
       $stop \leftarrow True$ 
    else
      C append  $index\_t$ 
    end if
  end while
  Cluster  $\leftarrow get\_xyz(C, points)$ 
  delete the  $index\_t$  in points
  delete the  $index\_t$  in  $D_m$ 
end while

```

Algorithm 3 Output clustering results *Cluster*

Legend:

H_d : minimum distance threshold insider one scan;
 V_d : minimum distance threshold between two scans;
 $points$: all the point cloud, $n \times 3$;
 $mini_points$:
 minimum number of connected points between two segments;
 $find_segments()$:
 find all segments in one scan with the distance threshold H_d ;
 $dist_func()$:
 distance calculation between two segments;

 $Clusters \leftarrow []$: create empty cluster $scans \leftarrow find_all_scans(points)$: rearrange the point cloud in the form of Lidar scan. $Global_segs \leftarrow \{\}$: dictionary for segments in global setting $Above_segs \leftarrow \{\}$: dictionary for segments in above scan $Current_segs \leftarrow \{\}$: dictionary for segments in current scan $current_segments \leftarrow []$ $next_segments \leftarrow []$ $global_label \leftarrow 0$ $Current_segs \leftarrow find_segment(scans(1), H_d)$ $Global_segs = Current_segs$ $Above_segs = Current_segs$ **for** $i = 2 : len(scans)$ **do** $Current_segs \leftarrow find_segments(scans(i), H_d)$ $len_current \leftarrow length(Current_segs)$ $temp_dict = \{\}$ **for** $j = 1 : len_current$ **do** $connect_index = find_connection(Current_segs(j), Above_segs, V_d, mini_points)$ $len_connect = length(connect_index)$ **if** $len_connect == 0$ **then** $global_cluster = global_cluster + 1$ $Global_segs[global_label] = Current_segs(j)$ $temp_dict.[global_label] = Current_segs(j)$ **end if****if** $len_connect == 1$ **then** $Global_segs[connect_index(0)].append(Current_segs(j))$ $temp_dict[connect_index(0)] = Current_segs(j)$ **end if****if** $len_connect \geq 2$ **then** $mini_key = fetch_mini_key(connect_index)$ $Global_segs = merge_delete_global(Global_segs, mini_key, connect_index)$ $Above_segs = merge_delete_above(Above_segs, mini_key, connect_index)$ $temp_dict = merge_delete_current(temp_dict, mini_key, connect_index)$ **end if****end for** $Above_segs = temp_dict$ **end for****Algorithm 4** Output occlusion label *Occlusion*

Legend:

θ_t : angle difference threshold to decide whether two adjacent clusters are occluded or not;
 $clusters$: all the clusters which are used to be labeled with occlusion or not;
 $get_all_angle()$:
 function to calculate all the angle with respective to original point;
 $calculate_range()$:
 function to calculate the mean range of one cluster;
 $calculate_angle_overlap$:
 function to build the $n \times n$ symmetric matrix about angle overlap;
 $calculate_diff_range$:
 function to build the $n \times n$ symmetric matrix about range difference;

 $all_angles \leftarrow []$: the list of left and right most angle of every cluster $all_ranges \leftarrow []$: the list of range of every cluster**for** c in $clusters$ **do** $leftmost_ang, rightmost_ang = get_all_angle(c)$ $all_angles.append([leftmost_ang - \theta_t, rightmost_ang + \theta_t])$ $r = calculate_range(c)$ $all_ranges.append(r)$ **end for** $occlusion \leftarrow []$ **for** $i = 1 : len(clusters)$ **do** $ind = find_angle_overlap(i, all_angle)$ $ind.remove(i)$ **if** $len(ind) == 0$ **then** $occlusion.append(0)$ **else****for** j in ind **do****if** $all_ranges[i] < all_ranges[j]$ **then** $occlusion.append(0)$ **else** $occlusion.append(1)$ **end if****end for****end if****end for**

Algorithm 5 Output global best g^* and its particle X^*

Legend:

Num_generation: number of iterations;
Num_particles: number of particles;
Range_particles: domain of particle search
Recall_func(): recall calculation function;
Init_func(): initialization function for all particles;
Prop_generation_func():
 proposal generation function;

// Initialization

$X = [x_1, x_2, \dots, x_n]^T = \text{Init_func}(\text{Range_particles})$: randomly initialize all n particles with the defined domain

$V = [v_1, v_2, \dots, v_n]^T = \text{Init_func}(1)$: randomly initialize velocities for n particles

$g^* = \text{Real_large_number}$: initialize the global best with a large number

$X^* = [x_1^*, x_2^*, \dots, x_n^*] = \text{Real_large_number}$: initialize the global best with a large number *Real_large_value*

for all $t = 0 : \text{Num_generation}$ **do**

for $i = 0 : \text{Num_particles}$ **do**

$\text{proposals} = \text{Prop_generation_func}(x_i^t)$

$\text{recall} = \text{Recall_func}(\text{proposals})$

// Find the personnel and global best particle

if $\text{recall} < x_i^*$ **then**

$x_i^* = \text{recall}$

end if

if $\text{recall} < g^*$ **then**

$g^* = \text{recall}$

end if

end for

// Particle updating step

for $i = 0 : \text{Num_particles}$ **do**

$v_i^{t+1} = \alpha \cdot v_i^t + \lambda \cdot r_1 \cdot (g^* - x_i^t) + \theta \cdot r_2 \cdot (x_i^* - x_i^t)$

$x_i^{t+1} = x_i^t + v_i^{t+1}$

end for

// Particle repair

for $i = 0 : \text{Num_particles}$ **do**

if x_i^{t+1} is out of range **then**

$x_i^{t+1} = \text{Init_func}(\text{Range_particles})$

end if

end for

end for
