

PSA-Det3D: Pillar Set Abstraction for 3D object Detection

Zhicong Huang^a, Jingwen Zhao^a, Zhijie Zheng^a, Dihu Chen^{a,*}, Haifeng Hu^a

^a*School of Electronics and Information Technology, Sun Yat-sen University, Guangzhou, China*

Abstract

Small object detection for 3D point cloud is a challenging problem because of two limitations: (1) Perceiving small objects is much more difficult than normal objects due to the lack of valid points. (2) Small objects are easily blocked which breaks the shape of their meshes in 3D point cloud. To alleviate these problems, we design a point-based detection network PSA-Det3D which mainly consists of a pillar set abstraction (PSA) and a foreground point compensation (FPC). The PSA improves the query approach of set abstraction to expand its receptive field of the network, which benefits the point-wise feature aggregation for small objects. The FPC fuses the foreground points and the estimated centers to optimize the candidate points selection. This module improves the detection performance for occluded objects effectively. The experiments on the KITTI 3D detection benchmark show that our proposed PSA-Det3D outperforms existing point-based algorithms for small object detection.

1. Introduction

LiDAR-based 3D object detection is widely applied in autonomous driving and robotics [1, 6, 10]. Existing detectors have achieved high accuracy for car category, but the detection performance for small objects such as pedestrian and cyclist is still unsatisfactory. Thus, how to improve the detection of small objects in 3D point cloud is an important yet challenging problem [7]. Some methods apply the multimodal fusion-based approaches in the 3D object detection network [15, 27, 30]. Due to the detailed information provided by the cameras and LiDAR sensors, they achieve high accuracy for both normal and small objects. However, the complex fusion networks increase the computational cost of these methods and thus limit their application.

Although the researches on point-based methods [3, 23, 32, 33] have achieved remarkable progress, the limitations of small objects are not sufficiently considered yet. There are two obvious limitations for LiDAR-based small object detection: (1) Perceiving small objects is much more difficult than normal objects because the sparse LiDAR-based point clouds usually do not provide sufficient information of small objects. However, the existing point-based methods [3, 23, 33] usually apply the native

*Corresponding author;
Email address: stscdh@mail.sysu.edu.cn (Dihu Chen)

set abstraction (SA) [20] to encode the features, either for normal objects or small objects. (2) Small objects are easily blocked in 3D point cloud. Since the scanning beam of LiDAR sensor is not penetrating, the occlusion not only reduces the number of the valid points for small objects but also breaks the shape of their meshes. As the semantic information of the occluded objects are affected, the normal approaches [3, 23] based on semantic features can hardly detect such objects normally.

Considering the above limitations, we first propose a pillar set abstraction (PSA) to learn informative point-wise features for small objects. Motivated by the pillar-based methods [12, 13, 29], our PSA adopts a pillar query operation before feature aggregation to expand the receptive field of the encoder network. Compared with the traditional SA, PSA can help to improve the performance of detection network, especially for small objects, with only a small increase on computational cost.

Furthermore, we present a foreground point compensation (FPC) to locate the objects under occlusion. In our FPC, we not only segment the foreground points but also estimate the centers of all possible objects. The centroid information can effectively compensate for incorrect segmentation results caused by occlusion. We fuse the centers and the foreground points to generate the detection proposals. Both the PSA and FPC are utilized in our point-based object detection network, PSA-Det3D.

In summary, our main contributions include:

- (1) We propose a pillar set abstraction to learn informative features for sparse point clouds of small objects, which adopts a pillar query operation before feature aggregation to expand the receptive field of the network.
- (2) We propose a foreground point compensation to locate the occluded objects. This module optimizes the selection of candidate points by the fusion of foreground points and estimated center points.
- (3) We propose a novel 3D object detection network, PSA-Det3D, which is evaluated on KITTI 3D object detection benchmark. Extensive experimental results confirm the outstanding accuracy of PSA-Det3D for small object detection.

2. Related work

2.1. Voxel-based Methods

Voxel-based methods convert the irregular point clouds into dense voxels or pillars before sending them to the network. VoxelNet [34] is the pioneer architecture that uses a 3D convolution network to detect 3D object from point clouds. SECOND [31] implements a fast backbone with 3D submanifold convolution network which significantly improved the efficiency of voxel-based network. In PointPillars [12], a pillar-based aggregation operation is proposed to obtain the pseudo-image feature map from bird’s eye view (BEV). It simplifies both the voxelization process and the architecture of detection network to improve the computational efficiency. TANet [16] uses triple attention mechanism in the feature encoding network to improve the detection performance for small objects and the robustness. HotSpotNet [4] converts the voxel-wise features into a new representation and predicts boxes with an anchor-free pipeline. Voxelization reduces the complexity and irregularity of the input data, but the generated voxels lose detailed information of the raw point clouds simultaneously, which is detrimental for small object detection.

2.2. Point-based Methods

Point-based methods learn 3D features from raw point clouds directly without pre-processing. Existing point-based detectors [2, 3, 18, 20, 23, 28, 32, 33] usually learn the point-wise features using PointNet++ and predict 3D proposals based on these features. PointNet [2] successfully applies the neural network to aggregate the point-wise features through the multi-layer perceptron (MLP) networks. They also conduct experiments on different 3D computer vision tasks to prove the generalizability of PointNet. PointNet++[20] enhances the performance of feature extraction by embedding the farthest point sampling (FPS) algorithm and the grouping operation based on PointNet. PointRCNN[23] is a two-stage object detection algorithm proposed by Shi et al. in 2019. It adopts PointNet++ for the point-wise feature extraction and predicts 3D boxes based on the foreground point segmentation. 3D-CenterNet[28] predicts the position of the object centers and estimates the bounding boxes around them. Pointformer [18] use a local and global transformer network for feature learning in a object detection model. SASA [3] present a semantics-augmented set abstraction module to maintain more foreground points in the point-based encoder. Zhang et al. [33] implement a fast single-stage detector IA-SSD by using an instance-aware downsampling strategy to select foreground points of interest and predict the centroids of objects.

2.3. Point-Voxel Fusion Methods

In addition to the above methods, there is another stream that uses the voxel-based and the point-based networks in a single detector [5, 11, 14, 21, 22]. Fast PointRCNN [5] adopts a point-based network to fuse the voxel-wise features from a voxel-based backbone for the optimization of 3D proposals. In PV-RCNN [22], the voxel-wise features are mapped to each sampled key point through the voxel set extraction module where the obtained point-wise features are further used to optimize the proposals. SA-SSD [11] projects the voxel-wise features back to the point-wise features to exploit point-wise supervisions in an auxiliary network. The auxiliary network guides the voxel-wise features in the backbone network to be aware of the object structure, thus improving the detection performance while introducing no extra computation in the inference stage. M3DETR [9] learns multi-representation features of point cloud and fuse the features using a multi-scale transformer network. The point-voxel fusion methods contribute to a better detection performance, but how to effectively fuse the features of the two modalities is still an open and challenging problem.

3. Our Method

As shown in Fig. 1, PSA-Det3D mainly comprises three components: (a) a point-based backbone with PSA layers; (b) a foreground point compensation (FPC) block; and (c) an RCNN which generates the final detection results. The input point set P first passes through the backbone to produce the point-wise features F . In this block, several stacked PSA layers are utilized to build the feature encoder. Then the point-wise features are sent to the FPC block. We predict 3D bounding boxes BB and calculate the foreground confidence by a segmentation following PointRCNN [23] in the FPC first. Then we employ a center estimation as well as a fusion algorithm to optimize the

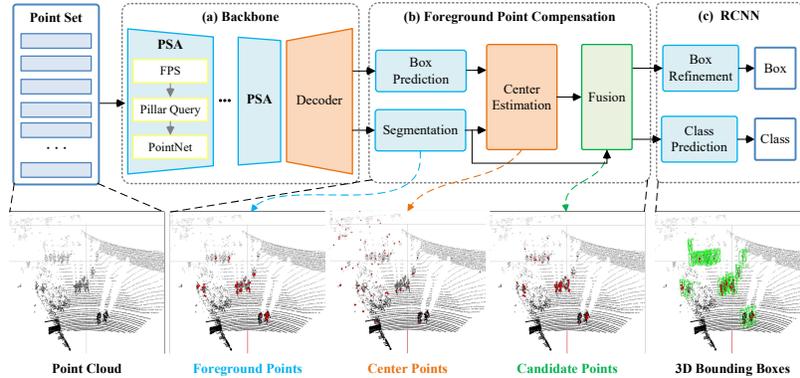


Figure 1: Outline of PSA-Det3D. The input points pass through (a) backbone network with PSA layer to produce the point-wise features. The (b) foreground point compensation is used to predict 3D proposals based on the features. The final detection results are generated in the (c) RCNN.

candidate points and generate 3D proposals. At last, the proposals are refined by an RCNN. The details of the PSA and FPC will be described below.

3.1. Pillar Set Abstraction

The SA of PointNet++ is widely applied to learn the point-wise features in existing point-based methods [3, 23, 28, 32, 33]. In SA, the input points are divided into several balls around the sampled key points in 3D space through a ball query operation [20]. We use the point cloud of a pedestrian as an example to simulate this operation. As shown in Fig. 2, the ball query of SA generates balls around the key points to aggregate their local features. However, this operation is restricted by the number of sampled points and the ball radius adopted [17]. If the number of the key points reduces to two or even one, the ball-based grouping operation cannot keep a high coverage anymore. This is probably one of the main reasons for the decreased accuracy on small objects because of insufficient context information. Simply increasing the radius of the balls cannot solve the problem because it would introduce severe noise to the aggregated features and the fine-grained 3D information may lose.

The pillar-based methods [12, 29] provide a special query operation that partitions points into pillars by increasing the boundaries of voxels in the z-direction. Motivated by this, we propose a novel pillar set abstraction (PSA) which implements the pillar-based grouping operation in the SA of PointNet++ [20]. Specifically, we use the FPS algorithm to sample key points and PointNet to learn the features in the PSA layer. After the FPS, we employ a pillar-based grouping operation to partition points into several groups. This operation is represented by the following formula:

$$\begin{aligned}
 GA(k) &= \{p \in P \mid HD(p, k) < r_0\}, k \in K \\
 HD(p, k) &= \sqrt{(p(x) - k(x))^2 + (p(y) - k(y))^2} \\
 P &= \{p_n\}_{n=0}^{N_i-1}, K = \{k_n\}_{n=0}^{N_o-1}, N_o \leq N_i
 \end{aligned} \tag{1}$$

where GA is the array of pillar-based grouping result, P is the set of input points, N_i is the number of input points, K is the set of key points, N_o is the number of sampled points, $HD(p, k)$ is the horizontal distance between points p and k .

The pillar-based grouping of our PSA is also simulated for comparison. As shown in Fig. 2, the generated pillar captures most of the points even if there is only one key point used. In other words, the PSA provides more valid points to the MLPs and thus helps to learn informative features.

3.2. Foreground Point Compensation

After feature learning, the point-based detector needs to select candidate points for 3D proposal generation. Existing methods [3, 18, 23] usually adopt the foreground point segmentation to select the most valuable points. However, the foreground scores cannot fully reflect the quality of the predicted proposals [24] because the relationship between 3D proposals and points is not considered. Additionally, the occlusion of small objects will inevitably affects the accuracy of the segmentation due to the lack of semantic information. An increasing number of approaches produce the candidate points by locating the object centers [19, 28, 33], which is an effective improvement compared with the foreground point segmentation. Therefore, our proposed FPC fuses both the foreground points and the object centers to optimize the selection of candidate points. As shown in Fig. 1, we adopt a box prediction and segmentation to predict preliminary 3D proposals and semantic scores. Next, we present a center estimation algorithm to predict object centers, which are then used for candidate points selection in a fusion module. We detail the center estimation in algorithm 1.

The inputs include the size of input points N , the foreground confidences C_{FG} and the set of generated 3D bounding boxes BB . The mask of centers M_{center} is calculated based on the intrinsic relationship between the predicted 3D boxes and the coordinates of the point cloud. First, the point with the highest foreground confidence is selected as the first center t_0 . The corresponding 3D box $BB[t_0]$ is used to identify the other points in the same box which are marked as ignored points in the mask $M_{ignored}$. Then, we continue to select the next unmarked point with the highest confidence, and perform the same operation until all points have been marked as centers or ignored points.

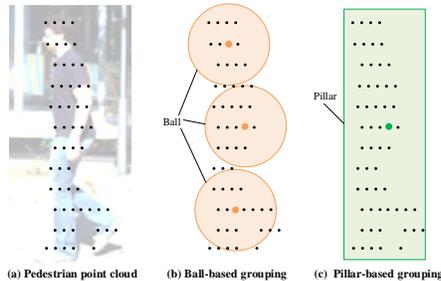


Figure 2: Both (b) ball query operation and (c) pillar query operation are simulated on the (a) point cloud of a pedestrian.

Algorithm 1: Center Estimation algorithm

Input: the size of input points: N ; the input points: $P = \{p_n\}_{n=0}^{N-1}$; the foreground confidences: $C_{FG} = \{C_i\}_{n=0}^{N-1}$; the set of bounding boxes: $BB = \{bb_n\}_{n=0}^{N-1}$;
Data: Center point mask: $M_{center}[0 \dots N - 1]$; Ignored point mask: $M_{ignored}[0 \dots N - 1]$

```
1 Initialize  $M_{center}$  with zeros;  
2 Initialize  $M_{ignored}$  with zeros;  
3 for  $i = 0$  to  $N - 1$  do  
4    $idx = \text{argmax}(C_{FG})$ ;  
5    $C_{FG}[idx] = 0$ ;  
6   if  $M_{center}[idx] == 0$  and  $M_{ignored}[idx] == 0$  then  
7     for  $j = 0$  to  $N - 1$  do  
8       if  $P[j] \in BB[idx]$  then  
9          $M_{ignored}[j] = 1$ ;  
10      end  
11     end  
12      $M_{center}[idx] = 1$ ;  
13   end  
14 end  
15 return  $M_{center}$ 
```

To generate the refined scores of 3D proposals, we design an algorithm to fuse the center mask M_{center} with the foreground confidences C_{FG} . The final score generated through the FPC C_{FPC} is formulated as follows:

$$C_{FPC} = C_{FG} + r \cdot C_{FG} \cdot M_{center} \quad (2)$$

where r is a hyperparameter that determines the fusion ratio.

As shown in Fig. 3, we demonstrate the selection of candidate points during the FPC, where the groundtruth boxes are represented by red boxes. As the red points represent the selection of the segmented foreground points, we can see that some occluded pedestrians are not detected by the segmentation due to the occlusion in 3D point cloud. These pedestrians are then recognized by the estimated centers which are drawn as blue points. Finally, our FPC generates the proper candidate points by fusing the foreground points as well as the estimated center points, which are illustrated as gold points.

3.3. Refinement and loss function

We adopt the box refinement layer and the loss function from [23]. The 3D bounding box regression loss L_{reg} can be formulated as

$$L_{reg} = \frac{1}{N_{pos}} \sum_{p \in pos} (L_{bin}^{(p)} + L_{res}^{(p)}) \quad (3)$$

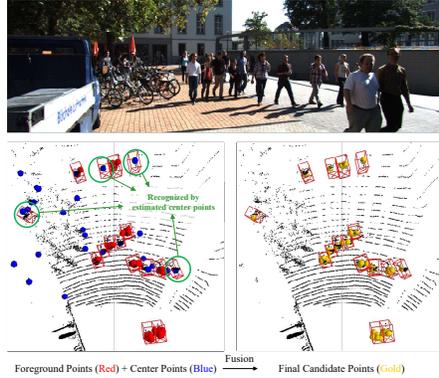


Figure 3: Illustration of the FPC. We fuse the foreground points and the estimated center points to generate the optimized candidate points.

Table 1: Performance comparison of 3D object detection with existing methods on KITTI *test* split. All results are evaluated by mean Average Precision with 40 recall positions via the official KITTI evaluation server.

Methods	Car 3D AP (%)			Ped. 3D AP (%)			Cyc. 3D AP (%)		
	Easy	Moderate	Hard	Easy	Moderate	Hard	Easy	Moderate	Hard
VoxelNet [34]	77.47	65.11	57.73	39.48	33.69	31.50	61.22	48.36	44.37
SECOND [31]	84.65	75.96	68.71	45.31	35.52	33.14	75.83	60.82	53.67
PointPillars [12]	82.58	74.31	68.99	51.45	41.92	38.89	77.10	58.65	51.92
PointRCNN [23]	86.96	75.64	70.70	47.98	39.37	36.01	74.96	58.82	52.53
Pointformer [18]	87.13	77.06	69.25	50.67	42.43	39.60	75.01	59.80	53.99
IA-SSD [33]	88.34	80.13	75.04	46.51	39.03	35.61	78.35	61.94	55.70
SASA [3]	88.76	82.16	77.16	-	-	-	-	-	-
PSA-Det3D	87.46	78.80	74.47	49.72	42.81	39.58	75.82	61.79	55.12

Table 2: Quantitative comparison of different approaches on the *val* split of the KITTI dataset. Our PSA-Det3D achieves the best performance for small objects including “pedestrian” and “cyclist”.

Methods	Car 3D AP (%)			Ped. 3D AP (%)			Cyc. 3D AP (%)		
	Easy	Moderate	Hard	Easy	Moderate	Hard	Easy	Moderate	Hard
VoxelNet [34]	81.97	65.46	62.85	57.86	53.42	48.87	67.17	47.65	45.11
SECOND [31]	90.55	81.61	78.61	55.94	51.14	46.17	82.96	66.74	62.78
PointPillars [12]	87.75	78.40	75.18	57.30	51.41	46.87	81.58	62.94	58.98
PointRCNN [23]	92.07	80.61	78.35	64.01	57.34	51.10	92.10	73.03	68.52
Pointformer [18]	90.05	79.65	78.89	-	-	-	-	-	-
IA-SSD [33]	91.88	83.41	80.44	61.22	56.77	51.15	88.42	70.14	65.99
SASA [3]	92.19	85.76	83.11	67.46	61.52	55.75	91.57	73.89	69.63
PSA-Det3D	91.34	80.53	78.09	68.35	63.09	57.09	92.23	74.11	69.97

where N_{pos} is the number of foreground points, $L_{bin}^{(p)}$ is the bin-based regression losses and $L_{res}^{(p)}$ is smooth L1 loss of regression.

4. Experiments

Our 3D detector is evaluated on the KITTI object detection benchmark [8] which provides 7481 training samples and 7518 test samples. The objects are labeled as three levels of difficulty (“easy”, “moderate” and “hard”), according to their occlusion level

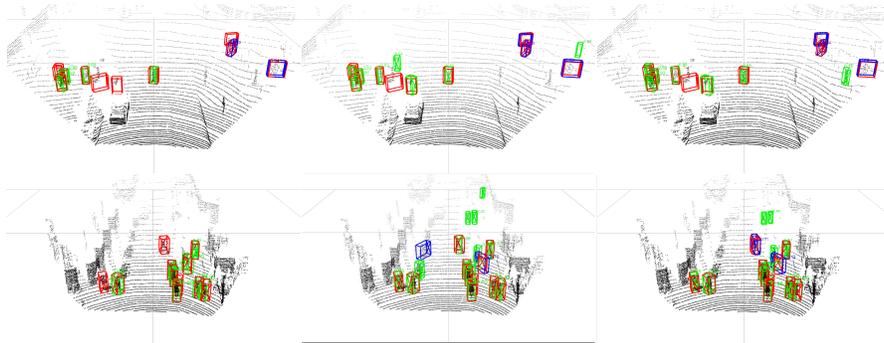


Figure 4: The detection result PointRCNN (*left*), SASA (*middle*) and PSA-Det3D (*right*). The green and blue boxes indicate the predicted pedestrians and cyclists respectively while the groundtruth boxes are represented by the red boxes.

and truncation ratio. Based on the common practice in related research, we divide training data into “train” split (3712 samples) and “val” split (3769 samples) for training and evaluating detection networks.

4.1. Implementation Detail

Our PSA-Det3D is implemented based on the OpenPCDet framework [26]. As only objects appearing on the image plane are labeled, the point clouds of KITTI dataset are first filtered through a projection from 3D to 2D. Next, we randomly sample 16384 points from the remaining points as the input of our detector. In terms of network parameters, the sizes of key points in four PSA layers are set to 8192, 1024, 256 and 64 respectively. The receptive radius of the PSA layers is set to 0.1, 0.5, 1.0 and 2.0 in meters. We implement the FPC module in the proposal generation layer after the foreground point segmentation and the 3D box prediction. The fusion ratio r of formula 2 is set to 0.5.

During training, we adopt three commonly used data augmentation methods, including random flipping about the X axis or Y axis, random scaling, and random rotation about the Z axis. The network is trained on 2 TITAN Xp GPUs for total 80 epochs with Adam optimizer [25] and one cycle learning rate schedule.

4.2. Online evaluation on test set

We evaluate the detection performance of our method and several other methods from existing literatures [3, 12, 23, 31, 33, 34] on KITTI *test* split dataset for “car”, “pedestrian” and “cyclist” under three different difficulties. The results on “moderate” are usually adopted as the main indicator for final ranking. Note that, since [3] does not submit the result for “pedestrian” and “cyclist” to KITTI server, we only report its performance for “car”. For the most challenging “pedestrian” detection race track, our method achieves 42.81% “moderate” 3D AP, outperforming Pointformer [18] by 0.38%. Compared with the most recent point-based detector IA-SSD [33], we improve the 3D AP of “pedestrian” category by 3.21%, 3.78% and 3.97% on three difficulties.

4.3. Offline evaluation on val split

Apart from the detection results on the *test* split, we also report the performance comparison on the *val* split of the KITTI dataset in Table 2. It can be seen that our method achieves the best performance for both “pedestrian” and “cyclist”. Comparing with the baseline model PointRCNN, our PSA-Det3D boosts the AP by 3.34%, 5.75%, 5.99% for “pedestrian” and 0.13%, 1.08%, 1.45% for “cyclist” respectively. Since we mainly focus on small object detection in this paper, the performance for “car” is basically the same as that of our baseline model.

We visualize some representative results using different methods in Fig. 4. We can observe that PointRCNN [23] locates most of the objects in the point clouds but there are still some unrecognized objects, which means that the accuracy of PointRCNN on small object detection is unsatisfactory. SASA [3] generates plenty of 3D proposals to figure out more objects, but some incorrect boxes are introduced simultaneously. The performance of PSA-Det3D (*right*) is generally the best since it identifies more objects than PointRCNN while producing fewer false boxes than SASA.

4.4. Ablation Study

To analyze the effectiveness of the proposed PSA and FPC, we conduct comprehensive ablation studies. All models are trained on the *train* set and evaluated on the *val* set for “pedestrian” and “cyclist” categories of KITTI dataset [8].

Table 3: Effectiveness of the PSA layer for small objects with different numbers of key points.

Methods	Query	Number of key points	3D mAP for Ped. and Cyc.(%)	Batch sizes	FPS
PointRCNN [23]	Ball	4096	67.70	21	8.97
		8192	67.72	20	7.85
+PSA	Pillar	4096	69.11	20	8.57
		8192	70.43	19	7.69

4.4.1. Effect of Pillar Set Abstraction

To investigate the effectiveness of the proposed PSA, we replace the encoder network of the PointRCNN with PSA and evaluate their performance and computational cost on KITTI dataset. We report the maximum number of batches that can be parallelized on one TITAN Xp(12GB) as well as the inference speed. As shown in Table

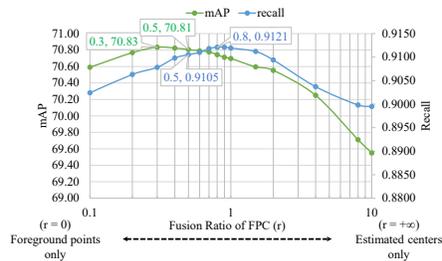


Figure 5: Comparison of different fusion ratio r in FPC.

3, the methods “+PSA” achieves higher accuracy than the baseline model with only a small increase on computational cost. It is worth noting that when the number of key points in the first abstraction layer increases to 8192, the performance of PSA-based method is further boosted by 1.32%.

The results in Table 4 show that the PSA improves the detection performance for small object by 2.73% on 3D mAP. This is mainly because the native point-wise encoder of PointRCNN is difficult to perceive small objects, especially the pedestrian. Note that our PSA significantly improves the 3D AP for “pedestrian” by 4.71%, 4.79% and 4.17% on three difficulties, which indicates that our PSA aggregates more informative features for small objects.

Table 4: Ablation study on the PSA and the FPC.

+PSA	+FPC	3D mAP (%)	Ped. 3D AP (%)			Cyc. 3D AP (%)		
			E.	M.	H.	E.	M.	H.
		67.70	64.01	57.34	51.12	92.10	73.03	68.53
✓		70.43	68.72	62.13	55.29	92.06	74.18	70.19
	✓	70.19	66.98	61.06	55.93	93.64	74.55	69.97
✓	✓	70.81	68.35	63.09	57.09	92.23	74.11	69.97

4.4.2. Fusion Ratio of FPC

According to formula 2, the parameter r determines the fusion ratio between the estimated centers and the segmented foreground points. To investigate the optimal fusion ratio of FPC, we conduct an experiment that compares the performance of PSA-Det3D using different r to figure out the best fusion setting. The comparison results in Fig. 5 indicate that the fused confidences outperform the cases which use the estimated centers or foreground points only. It is reasonable to set the fusion ratio to 0.5 for a balanced fusion between the foreground points and the centers. Although a large fusion ratio setting like 1.0 can further improve the recall rate, the accuracy will be affected since too many candidate points are introduced.

4.4.3. Effect of Foreground Point Compensation

By comparing the first row and the third row of Table 4, we can observe that our FPC brings a remarkable improvement by 2.49% on 3D mAP for small objects. It is noted that the FPC significantly boosts the 3D AP for “pedestrian” by 2.97%, 3.72%, 4.81% on three difficulties, where the improvement on “moderate” and “hard” pedestrians is greater than the “easy” pedestrians. As the difficulty reflects the occlusion level, this result suggests that the proposed FPC can recognize the occluded objects and thus significantly improves the 3D AP on “moderate” and “hard” difficulties.

5. Conclusion

In this work, we propose a point-based network PSA-Det3D to alleviate the limitations of small object detection in 3D point cloud. The pillar set abstraction (PSA) and the foreground point compensation (FPC) are the core parts of PSA-Det3D. The PSA improves the query approach of SA with pillar query to learn informative features for small objects. The FPC module fuses the foreground points as well as the centroid information to recognize more occluded small objects during the proposal generation.

The experiments on the KITTI 3D detection benchmark show that the proposed method outperforms existing point-based algorithms for small object detection.

References

- [1] Arnold, E., Al-Jarrah, O.Y., Dianati, M., Fallah, S., Oxtoby, D., Mouzakitis, A., 2019. A survey on 3d object detection methods for autonomous driving applications. *IEEE Transactions on Intelligent Transportation Systems* 20, 3782–3795.
- [2] Charles, R.Q., Su, H., Kaichun, M., Guibas, L.J., 2017. Pointnet: Deep learning on point sets for 3d classification and segmentation, in: 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 77–85.
- [3] Chen, C., Chen, Z., Zhang, J., Tao, D., 2022. Sasa: Semantics-augmented set abstraction for point-based 3d object detection. *Proceedings of the AAAI Conference on Artificial Intelligence* 36, 221–229. URL: <https://ojs.aaai.org/index.php/AAAI/article/view/19897>, doi:10.1609/aaai.v36i1.19897.
- [4] Chen, Q., Sun, L., Wang, Z., Jia, K., Yuille, A., 2020. Object as hotspots: An anchor-free 3d object detection approach via firing of hotspots, in: *Computer Vision – ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXI*, Springer-Verlag. p. 68–84.
- [5] Chen, Y., Liu, S., Shen, X., Jia, J., 2019. Fast point r-cnn, in: 2019 IEEE/CVF International Conference on Computer Vision (ICCV), pp. 9774–9783.
- [6] Cui, Y., Chen, R., Chu, W., Chen, L., Tian, D., Li, Y., Cao, D., 2022. Deep learning for image and point cloud fusion in autonomous driving: A review. *IEEE Transactions on Intelligent Transportation Systems* 23, 722–739.
- [7] Fan, L., Pang, Z., Zhang, T., Wang, Y.X., Zhao, H., Wang, F., Wang, N., Zhang, Z., 2022. Embracing single stride 3d object detector with sparse transformer, in: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 8458–8468.
- [8] Geiger, A., Lenz, P., Stiller, C., Urtasun, R., 2013. Vision meets robotics: The kitti dataset. *International Journal of Robotics Research* 32, 1231–1237.
- [9] Guan, T., Wang, J., Lan, S., Chandra, R., Wu, Z., Davis, L., Manocha, D., 2022. M3detr: Multi-representation, multi-scale, mutual-relation 3d object detection with transformers, in: 2022 IEEE/CVF Winter Conference on Applications of Computer Vision (WACV), pp. 2293–2303. doi:10.1109/WACV51458.2022.00235.
- [10] Guo, Y., Wang, H., Hu, Q., Liu, H., Liu, L., Bennamoun, M., 2021. Deep learning for 3d point clouds: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 43, 4338–4364.

- [11] He, C., Zeng, H., Huang, J., Hua, X.S., Zhang, L., 2020. Structure aware single-stage 3d object detection from point cloud, in: 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 11870–11879.
- [12] Lang, A.H., Vora, S., Caesar, H., Zhou, L., Yang, J., Beijbom, O., 2019. Pointpillars: Fast encoders for object detection from point clouds, in: 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 12689–12697.
- [13] Le, D.T., Shi, H., Rezatofghi, H., Cai, J., 2021. PiFeNet: Pillar-Feature Network for Real-Time 3D Pedestrian Detection from Point Cloud. arXiv e-prints , arXiv:2112.15458arXiv:2112.15458.
- [14] Li, Z., Yao, Y., Quan, Z., Xie, J., Yang, W., 2022. Spatial information enhancement network for 3d object detection from point cloud. Pattern Recognition 128, 108684.
- [15] Liang, Z., Zhang, Z., Zhang, M., Zhao, X., Pu, S., 2021. RangeiouDET: Range image based real-time 3d object detector optimized by intersection over union, in: 2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 7136–7145.
- [16] Liu, Z., Zhao, X., Huang, T., Hu, R., Zhou, Y., Bai, X., 2020. Tanet: Robust 3d object detection from point clouds with triple attention. Proceedings of the AAAI Conference on Artificial Intelligence 34, 11677–11684.
- [17] Mao, J., Shi, S., Wang, X., Li, H., 2022. 3d object detection for autonomous driving: A review and new outlooks doi:10.48550/arXiv.2206.09474.
- [18] Pan, X., Xia, Z., Song, S., Li, L.E., Huang, G., 2021. 3d object detection with pointformer, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 7463–7472.
- [19] Qi, C.R., Litany, O., He, K., Guibas, L., 2019. Deep hough voting for 3d object detection in point clouds, in: 2019 IEEE/CVF International Conference on Computer Vision (ICCV), pp. 9276–9285.
- [20] Qi, C.R., Yi, L., Su, H., Guibas, L.J., 2017. Pointnet++: Deep hierarchical feature learning on point sets in a metric space, in: Proceedings of the 31st International Conference on Neural Information Processing Systems, Curran Associates Inc., Red Hook, NY, USA. p. 5105–5114.
- [21] Qian, R., Lai, X., Li, X., 2022. Badet: Boundary-aware 3d object detection from point clouds. Pattern Recognition 125, 108524.
- [22] Shi, S., Guo, C., Jiang, L., Wang, Z., Shi, J., Wang, X., Li, H., 2020. Pv-rcnn: Point-voxel feature set abstraction for 3d object detection, in: 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 10526–10535.

- [23] Shi, S., Wang, X., Li, H., 2019. Pointcnn: 3d object proposal generation and detection from point cloud, in: 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 770–779.
- [24] Shi, W., Rajkumar, R., 2020. Point-gnn: Graph neural network for 3d object detection in a point cloud, in: 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 1708–1716.
- [25] Smith, L.N., Topin, N., 2019. Super-convergence: very fast training of neural networks using large learning rates, in: Pham, T. (Ed.), *Artificial Intelligence and Machine Learning for Multi-Domain Operations Applications*, International Society for Optics and Photonics. SPIE. pp. 369 – 386.
- [26] Team, O.D., 2020. Openpcdet: An open-source toolbox for 3d object detection from point clouds. <https://github.com/open-mmlab/OpenPCDet>.
- [27] Vora, S., Lang, A.H., Helou, B., Beijbom, O., 2020. Pointpainting: Sequential fusion for 3d object detection, in: 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 4603–4611.
- [28] Wang, Q., Chen, J., Deng, J., Zhang, X., 2021. 3d-centernet: 3d object detection network for point clouds with center estimation priority. *Pattern Recognition* 115, 107884.
- [29] Wang, Y., Fathi, A., Kundu, A., Ross, D.A., Pantofaru, C., Funkhouser, T., Solomon, J., 2020. Pillar-based object detection for autonomous driving, in: Vedaldi, A., Bischof, H., Brox, T., Frahm, J.M. (Eds.), *Computer Vision – ECCV 2020*, Springer International Publishing, Cham. pp. 18–34.
- [30] Wang, Z., Jia, K., 2019. Frustum convnet: Sliding frustums to aggregate local point-wise features for amodal 3d object detection, in: 2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), pp. 1742–1749.
- [31] Yan, Y., Mao, Y., Li, B., 2018. Second: Sparsely embedded convolutional detection. *Sensors* 18, 3337.
- [32] Yang, Z., Sun, Y., Liu, S., Jia, J., 2020. 3dssd: Point-based 3d single stage object detector, in: 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 11037–11045.
- [33] Zhang, Y., Hu, Q., Xu, G., Ma, Y., Wan, J., Guo, Y., 2022. Not all points are equal: Learning highly efficient point-based detectors for 3d lidar point clouds, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.
- [34] Zhou, Y., Tuzel, O., 2018. Voxelnet: End-to-end learning for point cloud based 3d object detection, in: 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 4490–4499.