

# CFDS3D: Centroid Feature Diffusion Sampling framework for 3D Object Detection

Anonymous CVPR submission

Paper ID 4265

## Abstract

3D object detection is one of the most important tasks in fields of autonomous driving and robotics. Our research focuses on the low efficiency issue of point-based methods on large-scale point clouds. Existing point-based methods adopt farthest point sampling (FPS) strategy for downsampling, which is computationally expensive in terms of inference time and memory consumption when the number of point cloud increases. In order to improve efficiency of downsampling, we propose the Centroid Instance Fusion Sampling Strategy (CIFSS) and effectively replace the first and the most complicated Set Abstraction (SA) layer. Besides, we believe that excessive background information is unnecessary for 3D object detection. Thus a local feature diffusion based background point filter (LFDBF) is constructed to exclude most invalid background points. LFDBF and CIFSS are combined to develop a highly efficient point-based 3D object detection framework, which can be embedded to any point-based models. Extensive experiments on multiple public benchmarks have demonstrated the superiority of CFDS3D. On Waymo dataset [22], the proposed framework significantly improves performance of baseline model and accelerates inference speed by 3.8 times. For the first time, **real-time detection** of point-based models in large-scale point cloud scenario is realized.

**Keywords:** 3D object detection, point-based method, farthest distance sampling.

## 1. Introduction

In recent year, as a common 3D representation, point cloud has been widely applied in 3D tasks, among which 3D object detection is crucial for autonomous driving. However, due to the orderless, sparse and irregular nature of point cloud, it is still challenging to predict 3D detection box with multiple degrees-of-free. Early works project point clouds to multi-view [1, 8, 9, 11, 20, 21, 27, 37], or voxelise point clouds [4, 10, 26, 29, 30, 34, 35, 38] and extract

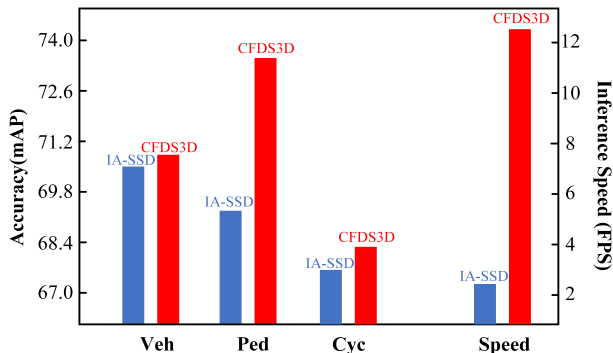


Figure 1. Comparison of performance and inference time of various models on Waymo dataset [22]. CFDS3D baseline is IA-SSD [33]. Experiment results are derived by using OpenPCDet [23] framework on a single A40 GPU. More details can be found in Table 1.

features through 3D convolution. Although these methods have achieved outstanding results, it is unavoidable to lose information when point cloud is converted to intermediate representation such as voxel. And it inevitably leads to declined model performance.

As we know, point-based methods [13–15, 19, 24, 25, 28, 33] extract point cloud features layer by layer, and rely on sophisticated downsampling strategies, such as Distance-Farthest Point Sampling (D-FPS) and Feature-Farthest Point Sampling (F-FPS) [14, 28, 33], to obtain center points. However, computational costs of these strategies are too expensive to afford when applied for large-scale 3D object detection. As shown in Table 1, the first layer of downsampling costs most of the inference time. On large datasets like Waymo [22] and Nuscenes [2], when the number of input points reaches 100k, it takes 498.1ms for D-FPS to sample 16,384 points, which severely hinders the applicability of point-based models on tackling large-scale point cloud tasks. While for 3D object detection, to the best of our knowledge, existing works have not provided an efficient and effective sampling strategy.

Moreover, as an object-oriented task, 3D detection does not need dense background context. SASA [3] employs MLP to encode point features, combining with FPS to increase the number of instance points. While IA-SSD [33] utilizes MLP to replace the last two layers of FPS for center point selection, in order to improve the recall rates of instance targets. However, current methods depend on complicated downsampling strategies or neighboring features extraction by the preceding SA layers to separate foreground and background points. And it leads to a large amount of inefficient computations in the first few SA layers.

Firstly, to reduce runtime of the downsampling module, we propose Centroid Instance Fusion Sampling Strategy (CIFSS) as an alternative to FPS. CIFSS significantly improves the inference speed of the model, thus make it possible for point-based models to achieve real-time detection in large-scale point cloud scenes. In addition, we propose a Centroid Point Offset Module to preserve the original geometry of instance targets, which helps the model to accurately regress bounding box size. Secondly, we present a Local Feature Diffusion Background Filter (LFDBF) that effectively discriminates foreground and background points and extracts voxel-based features before SA layer. It also provides multi-dimensional features of point clouds for further computations. We find that LFDBF tends to sample highly dense regions and ignores sparse foreground points at the far end. Thus we propose Density Distance Focal Loss based on density and distance of point clouds regions, to guarantee that sparse foreground points are sampled. Furthermore, we construct a 3D object detection framework CFDS3D that is compatible to almost all the point-based models. And real-time detection of point-based models in large-scale point cloud scenarios is realized.

To summarize, our contributions are listed as follows,

- We introduce a point-based 3D detection framework CFDS3D, which achieves efficient and accurate detection in large-scale point cloud scenes when inserted to existing point-based methods.
- We propose a Local Feature Diffusion Background Filter (LFDBF) to extract multi-dimensional features and exclude invalid background points in the raw point cloud input.
- We propose Centroid Instance Fusion Sampling Strategy (CIFSS) to realize real-time inference of point-based models in large-scale point cloud scenarios, as an efficient alternative of complicated downsampling strategies such as FPS.
- Extensive experiments on multiple large datasets have demonstrated effectiveness and superiority of CFDS3D.

## 2. Related Works

Due to the intricate properties of point clouds, researchers attempt to represent features by projecting into voxel grids [4, 9, 26, 29, 30, 34, 35, 38]. But these static projection methods result in information loss. Point-based methods directly use raw point cloud information as the input, and apply top-down aggregation on global features of point cloud. Moreover, combining two different forms of point cloud information can effectively improve model performance. Section 2.1 briefly introduces voxel-based methods for 3D object detection, followed by point-based methods in Section 2.2 and voxel-point based methods in Section 2.3.

### 2.1. Voxel-based Methods

In order to process unstructured 3D point cloud, voxel-based 3D detectors convert point clouds to regular voxel grids, such that the commonly used convolutions can be applied. VoxelNet [38] voxelizes the point cloud, and employs a voxel feature encoding layer to aggregate global and local information. However, computation and storage cost of 3D convolution increase along with resolution and bring unaffordable burden. SECOND [26] alleviates this issue by introducing 3D sparse convolution to substitute traditional convolution, and effectively optimizes memory usage and computation speed. PointPillars [9] further improves detection speed. It simplifies voxel to pillar with two dimensions, projects features to bird's eye view and applies 2D convolutions to extract deep features. SE-SSD [36] utilizes teacher-student network for data augmentation. SA-SSD [5] employs segmentation and center point prediction to facilitate model for further extraction of structural information.

### 2.2. Point-based Methods

Different to voxel-based methods, point-based methods use original information as the input, and adopt top-down learning to extract unstructured features of point cloud. Existing point-based methods normally adopt architectures similar to PointNet++ [14], which aggregates features by using symmetric aggregation function. PointRCNN [17] is the first 3D object detection model based on the original point cloud. It uses foreground segmentation network to obtain valid points for detection box regression, and predicts detection box by basing on bin. 3DSSD [28] is a single-stage detection framework that combines advantages of D-FPS and F-FPS. IA-SSD [33] uses FPS and instance-aware downsampling modules to extract features point by point, and utilizes contextual clues around bounding box to predict centroids. Nevertheless, existing methods cannot fully achieve fast and accurate downsampling. Constrained by complicated strategy in SA layer, existing methods are infeasible for large-scale 3D object detection.

## 2.3. Voxel-Point Based Methods

Voxel-point based methods combines the two methods above. Voxel-based methods are able to select region proposal effectively, but their receptive fields are limited due to the size of voxel. While, point-based methods have flexible receptive fields to capture accurate contexts. PV-RCNN [16] adopts 3D convolutional networks and set abstraction of PointNet [13] to learn features. Learnt voxel features are encoded to constitute a set of keypoints that can be used to transform features to different dimensions. HVPR [12] is a single-stage detection framework that effectively integrates voxel-based and point-based features to a single 3D representation. Memory modules are used to enhance point-based features and generate voxel-based features. Although accuracy has obvious improvement when compared to voxel-based and point-based methods, voxel-point based methods require extraction of multi-modal features and extra computational resources, which makes it impracticable in real-life scenario.

The proposed CFDS3D framework not only preserves advantages of common downsampling strategies such as FPS, but also improves efficiency of downsampling modules. Besides, CFDS3D extracts both voxel-based and point-based features for subsequent models to learn multi-modal information.

## 3. The Proposed Centroid Feature Diffusion Sampling Framework

### 3.1. Overview

Compared to intensive perception tasks like 3D semantic segmentation, 3D object detection pays more attentions to target objects. Existing 3D detection models use a great amount of background points as the input, which burdens subsequent models with huge computational load. Besides, complicated downsampling strategies such as FPS are unavoidable for existing point-based methods, causing inefficient computation and long inference time. To address these problems, we propose the Centroid Feature Diffusion Sampling (CFDS3D) framework for point cloud 3D object detection. CFDS3D is comprised of Local Feature Diffusion based Background Point Filter (LFDBF), Centroid-Instance Fusion Sampling Strategy (CIFSS) and SA layer.

As shown in Figure 2, LFDBF and CIFSS are combined to classify foreground and background points, and conduct highly efficient downsampling. Then Neighborhood Feature Diffusion Module (NFDm) further enlarges diffusion range of voxle centroid features, for reducing information loss due to downsampling. CFDS3D can be inserted to any point-based 3D object detection models and perform end-to-end training.

### 3.2. Local Feature Diffusion based Background Point Filter

We propose the Local Feature Diffusion based Background Point Filter (LFDBF) to efficiently remove background points. Given a set of points  $P = \{p_i \mid i = 1, \dots, N\} \in R^{n \times c}$ ,  $n$  and  $c$  denote the number of input points and channels, respectively. In Figure 2, points are voxelized to obtain a matrix of the size  $m \times s \times c$ , where  $m$  is the number of valid voxels,  $s$  is the number of points in voxel. For extraction of local features in every voxel, we use the method proposed in PointPillars [9] to acquire relative position information in voxels. We concatenate these 6-dimensional local features to original voxel features along feature dimension to obtain relative positional features with dimension  $R^{[m, s, (c+6)]}$ . After that, we flatten these relative positional features to  $R^{[m, s \times (c+6)]}$  for further computations.

MLPs supervised by instance labels are used to encode voxel-based neighboring features of point cloud. Inspired by RandLA-Net [6], we construct a Neighborhood Feature Diffusion Module (NFDm), where multi-scale ball query replaces K nearest neighbors in RandLA-Net to accelerate processing. As shown in Figure 3, for each voxel, after the neighborhood is derived, neighboring features are diffused to other voxels in its neighborhood. Thus each voxel contains contextual information from voxels in the vicinity. Such method can effectively improve accuracy of classification network and reduce information loss of foreground points caused by downsampling.

After NFDm, each voxel is rated by voxel features. MLP is employed as the classification network to calculate classification confidence of voxels. With a classification confidence higher than  $\alpha$ , the voxel is selected and considered as the foreground voxel. And its features are expressed by  $F_v \in R^{m_s \times c_1}$ , where  $m_s$  is the number of foreground voxels and  $c_1$  is the number of channels. Points in foreground voxels are regarded as foreground points.

Since dense regions in the vicinity are more likely to be chosen, small and sparse regions in the distant foreground are normally neglected. To solve this problem, we propose the Density Distance Focal Loss  $L_{DDFL}$  based on normal distribution. It prevents object instances at a distance from exclusion. Density constraint  $M_{Den}$  assigns various weights according to point density in voxels, and is written as follows,

$$M_{Den} = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(\frac{N_v}{N_{max}} - \mu)^2}{2\sigma^2}\right) \times \sqrt{2\pi}\sigma \quad (1)$$

where  $\mu$  and  $\sigma$  are position and scale parameters of normal distribution.  $N_v$  and  $N_{max}$  represent the number of valid points in voxel and the maximum value, respectively. Distance constraint  $M_{Dis}$  assigns more weights on objects at a

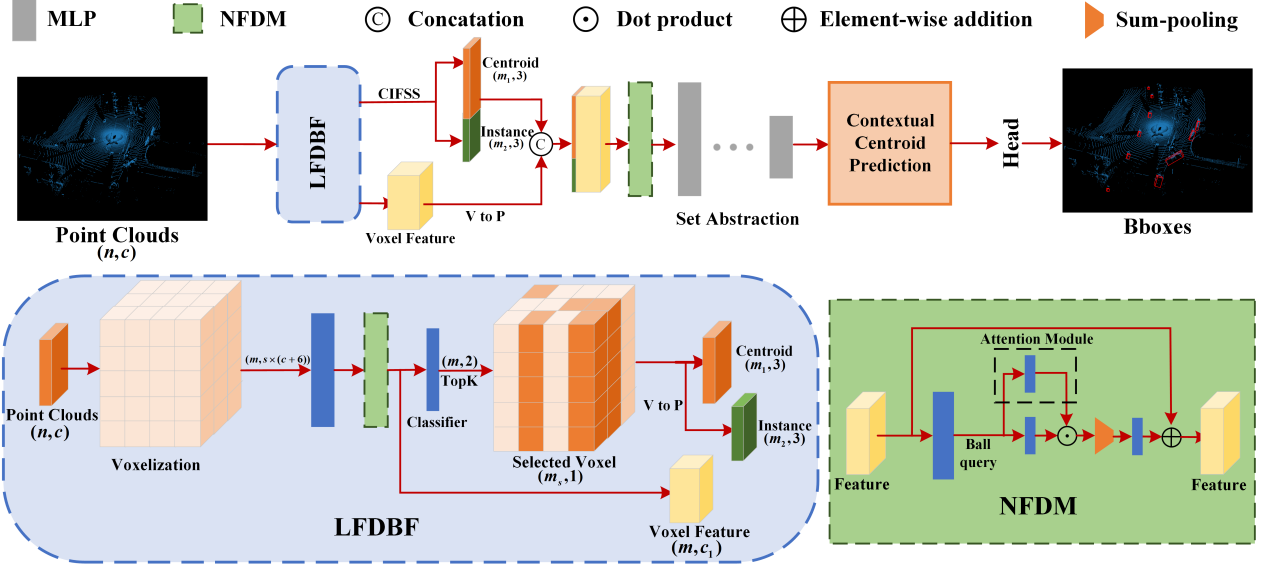


Figure 2. Diagram of CFDS3D framework. V to P represents the conversion from voxel to point.  $n$ ,  $m$ ,  $ms$ ,  $m1$ , and  $m2$  represent the input number of point cloud, the number of valid voxels, the number of selected voxels, the number of selected centroid points and the number of selected Instance points, respectively. Voxels colored in orange are selected as the foreground voxels.  $c$  and  $c1$  represent the number of input channel and feature channel respectively.

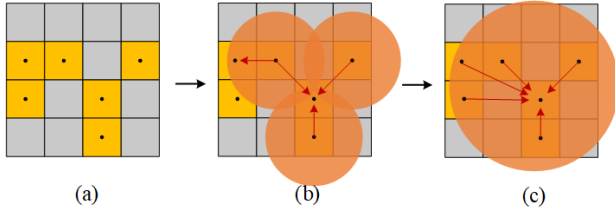


Figure 3. Diagram of Neighborhood Feature Diffusion Module. (a) represents neighbourhood features in voxel. (b) represents the diffusion process. (c) demonstrates the diffusion process after second NFD module. Yellow regions are valid voxels, orange regions are the diffusion coverage of voxels. Each voxel uses ball query to obtain its neighbourhood, and diffuses its neighbourhood features to other valid voxels. After the first NFD module, a voxel can obtain information from neighboring voxels. After the second NFD module, the voxel can perceive information outside the neighborhood.

distance, and is expressed as follows,

$$M_{Dis} = \frac{\exp\left(\frac{D}{M_D}\right)}{e} \quad (2)$$

where  $D$  is the distance between the point and the origin in the coordinate system,  $M_D$  is the distance from the farthest point to the origin. And density distance focal loss  $L_{DDFL}$  is written as follows,

$$L_{DDFL} = (1 - p_t)^\gamma \log(p_t) \times M_{Den} \times M_{Dis} \quad (3)$$

where  $p_t$  represents the difference between predicted value and ground truth value,  $\gamma$  is an adjustable factor.

### 3.3. Centroid-Instance Fusion Sampling Strategy

CIFSS aims to sample center points with high efficiency, and replace FPS in the first SA layer. Reasons for the success of FPS are twofold. One is that FPS adaptively selects center points according to point density, i.e. more center points are selected in regions of high density. The other is that FPS samples from the original point cloud and preserves geometric structural information, which is beneficial to improve accuracy of the following detection box regression. Therefore we add part of foreground points to center points, in order to increase sample density of instance objects. A centroid point offset module is constructed to restore the original geometric structure of point cloud.

**Raw Instance Points Sampling.** After foreground voxels are derived, centroid positions of voxels are calculated and noted as  $D_i \in R^{ms \times 3}$ . We sort centroid points according to classification confidence from the highest to the lowest, and select the highest  $m_1$  points. In addition, we select  $m_2$  points that have shortest distances to the origin. Both instance points and centroid points are regarded as the center points for the first SA layer in subsequent models. Algorithm 1 is the detailed procedure of CIFSS.

**Voxel Centroid Points Offset.** Adding positional information of the original point cloud to the following network layer by layer helps adapting to the original point cloud structure. However there exist deviations between centroid



and actual point cloud position. Direct use of centroids might cause losing the original positional information, and the model to fail in predicting accurate size of bounding box during regression. Consequently, we propose a centroid point offset module that moves centroid to the nearest instance point, for effective restoration on original size of targets. And the centroid point offset loss function  $L_{CB}$  is written as follows,

$$L_{CB} = \text{SmoothL1} \left( b \times \text{Mask}, \tilde{b} \right) \quad (4)$$

where  $b$  is the predicted offset between the centroid and its nearest instance point,  $\tilde{b}$  represents the actual offset of the centroid point to the nearest instance point.  $\text{Mask}$  is a matrix to denote whether centroid is present in the instance box or not. If true, the corresponding value in the matrix is set to 1, otherwise set to 0.

---

#### Algorithm 1 CIFSS

---

**Input:** foreground voxel index  $i$ , foreground voxel feature  $F_v \in R^{m_s \times c_1}$ , classification confidence  $\sigma$ , raw point cloud position  $P \in R^{n \times 3}$ , hashtable that maps voxel to point  $\text{HashMap}()$ .

**Output:** center point feature  $F_C \in R^{(m_1+m_2) \times (3+c_1)}$ .

- 1:  $m_1$  = the number of selected centroid points
  - 2:  $m_2$  = the number of selected instance points
  - 3: Calculate centroid point position of each foreground voxel  $P_{ctr} \in R^{m_s \times 3}$  and centroid point feature  $F_{ctr} \in R^{m_s \times (3+c_1)} \leftarrow \text{Concat}(P_{ctr}, F_v)$ .
  - 4: Calculate foreground point position  $P_{ins} \in R^{n_s \times 3} \leftarrow \text{HashMap}(P, i)$  and foreground point feature  $F_{ins} \in R^{n_s \times (3+c_1)} \leftarrow \text{Concat}(P_{ins}, \text{HashMap}(F_v, i))$ .
  - 5: Sort centroid points according to classification confidence  $\sigma$  in descending order, and select the centroid feature  $F_c \in R^{m_1 \times (3+c_1)}$  from  $F_{ctr}$ .
  - 6: Select the instance feature  $F_i \in R^{m_2 \times (3+c_1)}$  from  $F_{ins}$  by distance.
  - 7: Derive center point feature by merging centroid feature and instance feature :  $F_C \in R^{(m_1+m_2) \times (3+c_1)} \leftarrow \text{Concat}(F_c, F_i)$ .
- 

### 3.4. Centroid Feature Diffusion Sampling framework

As a plug-and-play framework, CFDS3D is compatible to any point-based 3D object detection models. We apply another Neighborhood Feature Diffusion Module to enlarge diffusion range of each centroid feature. Through stacking two NFDMS, CFDS3D can effectively reduce information loss due to downsampling and replace the sophisticated and first SA layer.

**Total Loss.** The proposed CFDS3D framework can be integrated with other models for end-to-end training. Multiple

loss functions are combined for optimization. Total loss includes density distance focal loss  $L_{DDFL}$ , centroid point offset loss  $L_{CB}$ , classification loss  $L_{cls}$  and bounding box generation loss  $L_{box}$ , as given in Equation 5.

$$L_{total} = L_{DDFL} + L_{CB} + L_{cls} + L_{box} \quad (5)$$

## 4. Experiments and Discussions

### 4.1. Implementation Details

We choose IA-SSD [33] as the baseline and construct our model. Extensive experiments are conducted to demonstrate the effectiveness of our method. First of all, we voxelize the input point cloud, and set voxel size to  $[0.075, 0.075, 1]$ . In LFDBF module, after dimensional expansion, point clouds are fed into three MLP layers of size  $(16, 16, 32)$ . Diffusion radius of the first and the second NFDMS are set to 0.4 and  $[0.2, 0.8]$ . The maximum number of diffusion points of both NFDMS are set to 16. Confidence threshold is set to 0.45 in voxel classification module. We set  $\mu = 0.5$  and  $\sigma = 0.5$  in the density distance focal loss function (Equation 3). Centroid offset module in CIFSS contains two MLP layers with size of  $(16, 3)$ .

Following the configuration of IA-SSD [33], multi-scale grouping is applied in the last three SA layers. Multi-scale receptive fields with radius  $[0.8, 1.6]$ ,  $[1.6, 4.8]$  and  $[4.8, 6.4]$  are provided. In contextual centroid prediction module, three MLP layers of size  $(256, 256, 3)$  are employed to predict offsets of instance centroids. All experiments are conducted on NVIDIA A40 GPU and AMD EPYC 7402 CPU with a batch size of 8. We apply Adam [7] optimizer with learning rate 0.01 and weight decay 0.01.

### 4.2. State-of-the-art Comparison

**3D Detection on Waymo.** To validate performance of CFDS3D on large-scale point cloud scenes, we conduct quantitative experiments on Waymo [22] dataset. Waymo dataset [22] contains 160K samples in training set and 40K samples in validation set with two difficulty levels of challenge, *Level1* and *Level2*. We construct two CFDS3D frameworks in two different settings. The maximum number of sampled centroid points and instance points are set to 16384/2048 and 30720/8197, respectively. The number of sampling points of the last three SA layers are set to 4096, 2048 and 1024. IoU threshold is set to 0.25.

Table 1 demonstrates that the proposed CFDS3D significantly improves the inference speed and accuracy of baseline model. Compared to the baseline model IA-SSD [33], on *Level1* CFDS3D(16384/2048) improves 0.35%/0.24%, 4.01%/5.59%, 0.53%/1.17% in terms of mAP/mAPH. Inference speed is improved by 3.8 times. For CFDS3D(30720+8192) with more input points, on *Level2* mAP/mAPH are improved by (1.14%/1.45%, 4.42%/5.36%, 1.04%/1.49%). Inference

Table 1. Comparison of different works on the Waymo *val* set. 20% of training set (32K) is used for training. Evaluation metrics are mean average precision (mAP) and mAP weighted by heading accuracy (mAPH). L1 and L2 denote *Level1* and *Level2*. FPS represents the inference speed of the model when batch size is set to 1. Bold texts are our results and the best results are underlined. For fair comparison, inference time and performance evaluation metrics are derived by reproducing methods under OpenPCDet framework [23].

Method	Type	Veh.(L1) mAP/mAPH	Veh.(L2) mAP/mAPH	Ped.(L1) mAP/mAPH	Ped.(L2) mAP/mAPH	Cyc.(L1) mAP/mAPH	Cyc.(L2) mAP/mAPH	FPS
SECOND [26]	1-stage	70.96/70.34	62.58/62.02	65.23/54.24	57.22/47.49	57.13/55.62	54.97/53.53	33.4
Pointpillars [9]	1-stage	70.43/69.83	62.18/61.64	66.21/46.32	58.18/40.64	55.26/51.75	53.18/49.80	26.1
CenterPoint [31]	1-stage	71.33/70.76	63.16/62.65	72.09/65.49	64.27/58.23	68.68/67.39	66.11/64.87	22.9
PV-RCNN [16]	2-stage	75.41/74.74	67.44/66.80	71.98/61.24	63.70/63.95	65.88/64.25	63.39/61.82	2.8
Part - A <sup>2</sup> [18]	2-stage	74.66/74.74	65.82/65.32	71.71/62.24	62.46/54.06	66.53/65.18	64.05/62.75	8.1
IA-SSD [33]	1-stage	70.53/69.67	61.55/60.80	69.38/58.47	60.30/50.73	67.67/65.30	64.98/62.71	2.7
CFDS3D(16384/2048)	1-stage	<b>70.88/69.91</b>	<b>67.75/60.96</b>	<b>73.39/64.06</b>	<b>64.27/56.01</b>	<b>68.20/66.47</b>	<b>65.79/64.02</b>	<b>12.9</b>
CFDS3D(30720/8192)	1-stage	<b>71.47/70.81</b>	<b>62.84/62.25</b>	<b>73.80/64.18</b>	<b>64.80/56.09</b>	<b>68.71/66.65</b>	<b>66.17/64.20</b>	<b>8.5</b>

Table 2. Quantitative results on Waymo *val* set for 3D object detection. † indicates that the whole scene is divided into four parts in the first SA layer to accelerate FPS.

Method	Veh.(L1) mAP/mAPH	Veh.(L2) mAP/mAPH	Ped.(L1) mAP/mAPH	Ped.(L2) mAP/mAPH	Cyc.(L1) mAP/mAPH	Cyc.(L2) mAP/mAPH	FPS
IA-SSD† [26]	70.38/69.82	61.33/60.19	68.23/58.44	60.11/50.33	66.84/65.06	64.32/62.61	9.5
CFDS3D† (16384/2048)	<b>69.29/69.37</b>	<b>60.88/60.06</b>	<b>72.50/63.08</b>	<b>63.35/55.05</b>	<b>66.91/65.39</b>	<b>64.66/62.89</b>	<b>17.9</b>
CFDS3D† (30720/8192)	<b>71.19/70.48</b>	<b>62.62/61.99</b>	<b>72.93/63.11</b>	<b>63.72/55.21</b>	<b>67.22/65.49</b>	<b>64.75/63.09</b>	<b>12.6</b>

speed is improved by 2.1 times. Since CIFSS circumvents heavy computational burden in the first SA layer, the proposed CFDS3D is enabled to become the first point-based framework that realizes real-time detection on large-scale point cloud dataset. The visualization results are shown in Figure 4. For more visualization results, please refer to Figure 5, 6 and 7. In Table 2, we divide the whole scene into four parts to accelerate FPS at the first SA layer, and compare accuracy and inference speed before and after inserting CFDS3D to IA-SSD [33]. Note that this acceleration strategy improves the inference speed but degrades the model performance to some extent. Nevertheless, after inserting CFDS3D models approximate and even surpass baseline performances. The two CFDS3D frameworks accelerate inference by 88.4% and 32.6%.

**3D Detection on Nuscenes.** We further conduct comparison experiments on Nuscenes dataset [2] to verify the robustness of CFDS3D. Nuscenes dataset [2] contains 40K annotated keyframes with 23 object categories. mAP and NDS denote mean Average Precision and Nuscenes detection score. The maximum number of sampled centroid points and instance points are set to 30720 and 8197. The number of sampling points for the last three SA layers are set to 4096, 2048 and 1024. IoU threshold is set to 0.25. Visualization results are shown in Figure 8.

We report our results on Nuscenes dataset [2] in Table 3. Our CFDS3D outperforms baseline model by 0.3% and 0.2% in NDS and mAP while accelerating inference by 1.8 times. Results on Nuscenes [2] and Waymo [22] datasets

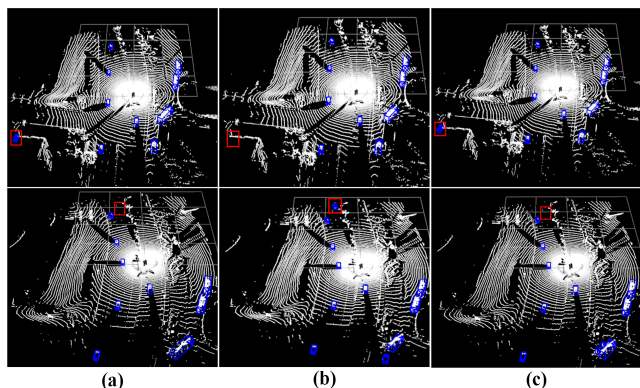


Figure 4. Visualization results of Waymo dataset *Vehicle*. (a) ground-truth, (b) IA-SSD, (c) CFDS3D+IA-SSD. Red boxes represent differences of various results.

demonstrate robustness and good performance of CFDS3D in a variety of large-scale point cloud scenarios.

### 4.3. Ablation Study

We conduct various ablation experiments on Waymo dataset [22]. As shown in Table 4, we construct three different CFDS3D variants by ablating each proposed module. Experiment results indicate that: (1) NFDM alleviates the information loss caused by downsampling while enhancing model performance. (2) LFDBF can separate foreground and background points better by adding distance and density constraints, especially for small objects in the far dis-

Table 3. Comparison results of various methods on Nuscenes *val* set. Bold texts are our results and best results are underlined. For fair comparison, inference time and performance evaluation metrics are derived by reproducing methods under OpenPCDet framework [23].

Method	NDS	mAP	Car	Truck	Bus	Tra.	C.V.	Ped.	Motor	Bicy.	T.C.	Barrier	FPS
Pointpillars [9]	46.8	28.2	75.5	31.6	44.9	23.7	4.0	49.6	14.6	0.4	8.0	30.0	<u>37.4</u>
3D-CVF [32]	49.8	42.2	79.7	37.9	55.0	36.3	-	<u>71.3</u>	37.2	-	<u>40.8</u>	47.1	-
SASA [3]	<u>61.0</u>	<u>45.0</u>	76.8	45.0	66.2	<u>36.5</u>	16.1	69.1	39.6	<u>16.9</u>	29.9	<u>53.6</u>	1.9
3DSSD [28]	56.4	42.6	<u>81.2</u>	<u>47.2</u>	61.4	30.5	12.6	70.2	36.0	8.6	31.1	47.9	2.0
IA-SSD [33]	48.8	44.0	73.8	45.1	67.0	29.7	<u>17.0</u>	66.9	40.6	14.6	32.0	53.2	2.4
<b>CFDS3D + IA-SSD</b>	<b>49.1</b>	<b>44.2</b>	<b>74.8</b>	<b>45.7</b>	<b>66.2</b>	<b>29.8</b>	<b>16.8</b>	<b>70.5</b>	<b>40.7</b>	<b>14.8</b>	<b>29.4</b>	<b>53.6</b>	<b>6.8</b>

Table 4. Ablation study on CFDS3D. We report mAP value of Waymo [22] dataset on *LEVEL1*, where NFDM stands for Neighbourhood Feature Diffusion Module, DDFL denotes Density Distance Focal Loss, and Ctr-bias represents Centroid Point Offset Module.

NFDM	DDFL	Ctr-bias	Veh.(L1) mAP	ped.(L1) mAP	Cyc.(L1) mAP
×	×	×	64.89	59.71	61.49
✓	×	×	67.02	60.75	63.66
✓	✓	×	68.15	68.33	67.18
✓	✓	✓	<b>71.47</b>	<b>73.80</b>	<b>68.71</b>
<b>Improvement</b>			<b>+6.58</b>	<b>+14.09</b>	<b>+7.22</b>

Table 5. Ablation study of CFDS3D variants using different downsampling strategies in the first layer. mAP values of Waymo [22] dataset in L1 difficulty are reported. Ins. Points represent using instance points as center points. Ctr. Points represent using centroid points as center points. FPS represents using the farthest distance sampling strategy to calculate center points. *FPS* represents the model inference speed when batch size is set to 1.

Ins. Points	Ctr. Points	FPS	Veh. (L1) mAP	ped. (L1) mAP	Cyc. (L1) mAP	<i>FPS</i>
×	×	✓	70.53	69.38	67.67	2.7
×	✓	×	70.17	70.51	67.08	13.3
✓	✓	×	<b>70.88</b>	<b>73.39</b>	<b>68.20</b>	<b>12.9</b>
<b>Improvement</b>			<b>+0.35</b>	<b>+4.01</b>	<b>+0.53</b>	<b>+10.2</b>

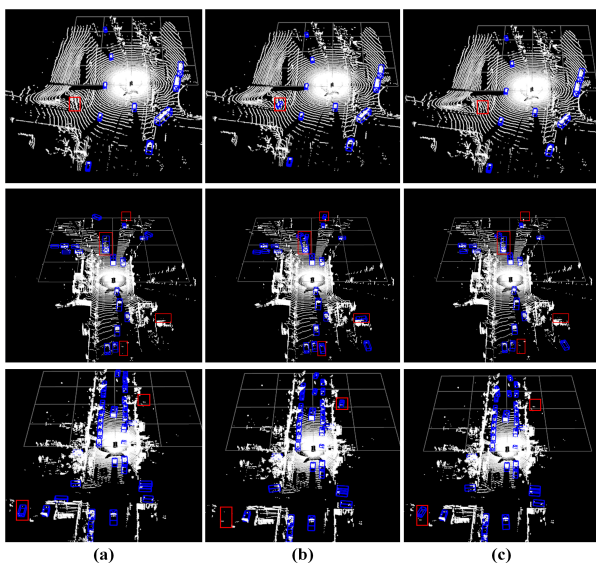


Figure 5. Complementary visualization of Waymo Dataset Vehicle, (a) ground truth, (b) IA-SSD, (c) CFDS3D + IA-SSD. Bounding boxes in red are the difference, blue are the Vehicle.

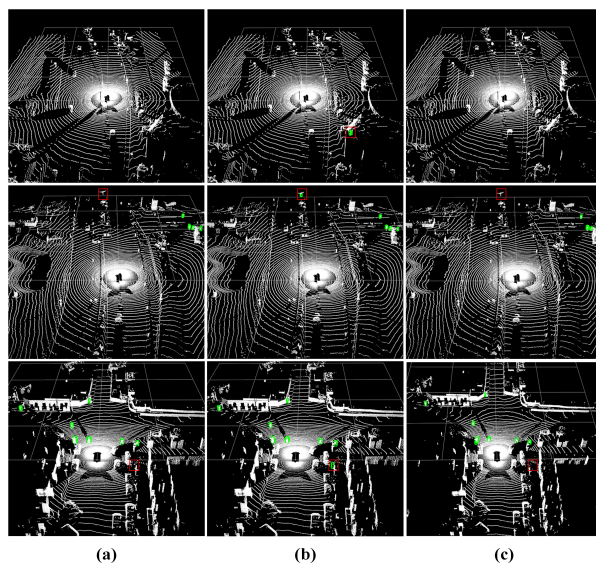


Figure 6. Complementary visualization of Waymo Dataset Pedestrian, (a) ground truth, (b) IA-SSD, (c) CFDS3D + IA-SSD. Bounding boxes in red are the difference, green are the Pedestrian.



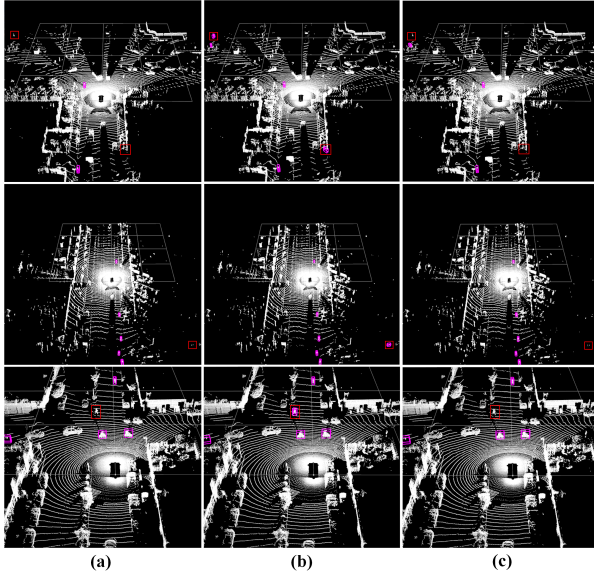


Figure 7. Complementary visualization of Waymo Dataset Cyclist, (a) ground truth, (b) IA-SSD, (c) CFDS3D + IA-SSD. Bounding boxes in red are the difference, purple are the cyclist.

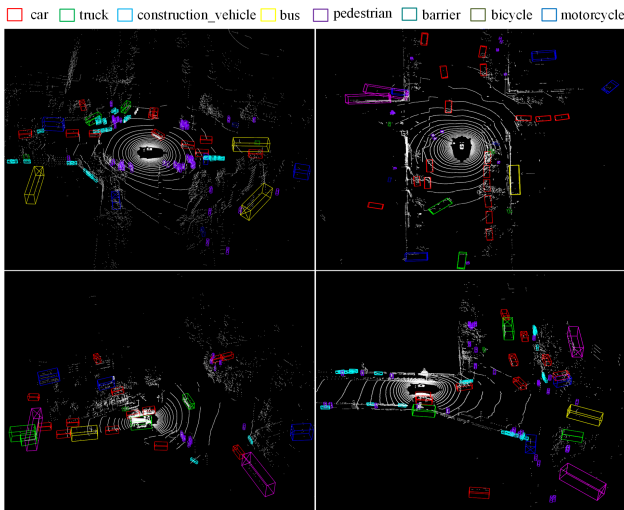


Figure 8. Complementary visualization of our methods on Nuscenes Dataset.

tance. (3) Retrieving raw point cloud geometry from voxels helps subsequent network to better capture scale information of instance object. Figure 9 compares results before and after adding DDFL. We can find that distant targets are sampled, which effectively improves model recall rate.

Table 5 reports ablation study of CFDS3D on various downsampling strategies. It is shown that centroid-point-based downsampling significantly improves the inference speed of CFDS3D. But it does not refer to the density distribution of instance target, resulting in decreasing accuracy.

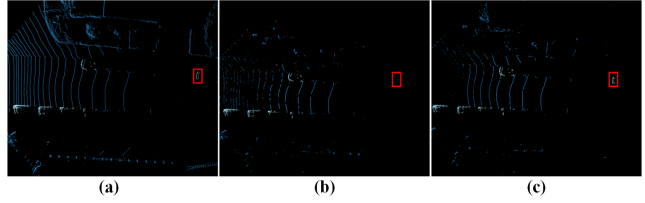


Figure 9. Visualisation results before and after DDFL. (a) original point cloud input, (b) Without DDFL, (c) With DDFL. White dots are foreground points. Bounding boxes in red are the difference.

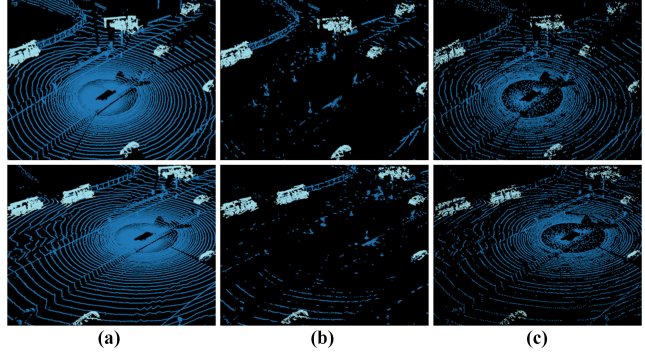


Figure 10. Distribution of sample points by various sampling strategies, (a) original point cloud input, (b) CIFSS, (c) FPS. White dots are real instance points, blue dots are background points.

Combining instance points with centroid points as sampling center not only improves the performance of baseline by 0.35%, 4.01% and 0.53%, but also still realizes real-time detection (12.9 FPS). As given in 10, compared to FPS, CIFSS samples more foreground points while excluding most of the background points, which allows the model to better focus on the instance target.

## 5. Conclusions

In this paper, an efficient framework is proposed for accelerating point-based 3D object detection model, termed as CFDS3D. It contains a fast center sampling strategy CIFSS, which effectively avoids huge computational burdens brought by FPS strategy in the first SA layer. To extract sufficient information from foreground points and exclude invalid background points, LFDBF is built to reduce information loss during downsampling. Experiment results on Waymo and Nuscenes datasets have demonstrated that our proposed CFDS3D framework is capable of surpassing baseline model performance and increasing inference speed significantly. Moreover, it is the first time that real-time detection of point-based model is realized in large-scale point cloud scenario. And it indicates a promising future for point-based 3D object detection methods to be applied in autonomous driving and other related areas.



## References

- [1] Waleed Ali, Sherif Abdelkarim, Mahmoud Zidan, Mohamed Zahran, and Ahmad El Sallab. Yolo3d: End-to-end real-time 3d oriented object bounding box detection from lidar point cloud. In *Proceedings of the European conference on computer vision (ECCV) workshops*, pages 0–0, 2018. 1
- [2] Holger Caesar, Varun Bankiti, Alex H Lang, Sourabh Vora, Venice Erin Liong, Qiang Xu, Anush Krishnan, Yu Pan, Giancarlo Baldan, and Oscar Beijbom. nuscenes: A multi-modal dataset for autonomous driving. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11621–11631, 2020. 1, 6
- [3] Chen Chen, Zhe Chen, Jing Zhang, and Dacheng Tao. Sasa: Semantics-augmented set abstraction for point-based 3d object detection. In *AAAI Conference on Artificial Intelligence*, volume 1, 2022. 2, 7
- [4] Jiajun Deng, Shaoshuai Shi, Peiwei Li, Wengang Zhou, Yanyong Zhang, and Houqiang Li. Voxel r-cnn: Towards high performance voxel-based 3d object detection. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 1201–1209, 2021. 1, 2
- [5] Chenhong He, Hui Zeng, Jianqiang Huang, Xian-Sheng Hua, and Lei Zhang. Structure aware single-stage 3d object detection from point cloud. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11873–11882, 2020. 2
- [6] Qingyong Hu, Bo Yang, Linhai Xie, Stefano Rosa, Yulan Guo, Zhihua Wang, Niki Trigoni, and Andrew Markham. Randla-net: Efficient semantic segmentation of large-scale point clouds. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11108–11117, 2020. 3
- [7] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 5
- [8] Jason Ku, Melissa Mozifian, Jungwook Lee, Ali Harakeh, and Steven L Waslander. Joint 3d proposal generation and object detection from view aggregation. In *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 1–8. IEEE, 2018. 1
- [9] Alex H Lang, Sourabh Vora, Holger Caesar, Lubing Zhou, Jiong Yang, and Oscar Beijbom. Pointpillars: Fast encoders for object detection from point clouds. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 12697–12705, 2019. 1, 2, 3, 6, 7
- [10] Ming Liang, Bin Yang, Yun Chen, Rui Hu, and Raquel Urtasun. Multi-task multi-sensor fusion for 3d object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7345–7353, 2019. 1
- [11] Haihua Lu, Xuesong Chen, Guiying Zhang, Qiuha Zhou, Yanbo Ma, and Yong Zhao. Scanet: Spatial-channel attention network for 3d object detection. In *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1992–1996. IEEE, 2019. 1
- [12] Jongyoun Noh, Sanghoon Lee, and Bumsub Ham. Hvpr: Hybrid voxel-point representation for single-stage 3d object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14605–14614, 2021. 3
- [13] Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 652–660, 2017. 1, 3
- [14] Charles Ruizhongtai Qi, Li Yi, Hao Su, and Leonidas J Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. *Advances in neural information processing systems*, 30, 2017. 1, 2
- [15] Rui Qian, Divyansh Garg, Yan Wang, Yurong You, Serge Belongie, Bharath Hariharan, Mark Campbell, Kilian Q Weinberger, and Wei-Lun Chao. End-to-end pseudo-lidar for image-based 3d object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5881–5890, 2020. 1
- [16] Shaoshuai Shi, Chaoxu Guo, Li Jiang, Zhe Wang, Jianping Shi, Xiaogang Wang, and Hongsheng Li. Pv-rcnn: Point-voxel feature set abstraction for 3d object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10529–10538, 2020. 3, 6
- [17] Shaoshuai Shi, Xiaogang Wang, and Hongsheng Li. Point-rcnn: 3d object proposal generation and detection from point cloud. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 770–779, 2019. 2
- [18] Shaoshuai Shi, Zhe Wang, Jianping Shi, Xiaogang Wang, and Hongsheng Li. From points to parts: 3d object detection from point cloud with part-aware and part-aggregation network. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020. 6
- [19] Weijing Shi and Raj Rajkumar. Point-gnn: Graph neural network for 3d object detection in a point cloud. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 1711–1719, 2020. 1
- [20] Martin Simon, Karl Amende, Andrea Kraus, Jens Honer, Timo Samann, Hauke Kaulbersch, Stefan Milz, and Horst Michael Gross. Complexer-yolo: Real-time 3d object detection and tracking on semantic point clouds. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 0–0, 2019. 1
- [21] Martin Simony, Stefan Milz, Karl Amende, and Horst-Michael Gross. Complex-yolo: An euler-region-proposal for real-time 3d object detection on point clouds. In *Proceedings of the European Conference on Computer Vision (ECCV) Workshops*, pages 0–0, 2018. 1
- [22] Pei Sun, Henrik Kretzschmar, Xerxes Dotiwalla, Aurelien Chouard, Vijaysai Patnaik, Paul Tsui, James Guo, Yin Zhou, Yuning Chai, Benjamin Caine, et al. Scalability in perception for autonomous driving: Waymo open dataset. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2446–2454, 2020. 1, 5, 6, 7
- [23] OpenPCDet Development Team. Openpcdet: An open-source toolbox for 3d object detection from point clouds. <https://github.com/open-mmlab/OpenPCDet>, 2020. 1, 6, 7

[24] Sumesh Thakur and Jiju Peethambaran. Dynamic edge weights in graph neural networks for 3d object detection. *arXiv preprint arXiv:2009.08253*, 2020. 1

[25] Yue Wang, Yongbin Sun, Ziwei Liu, Sanjay E Sarma, Michael M Bronstein, and Justin M Solomon. Dynamic graph cnn for learning on point clouds. *Acm Transactions On Graphics (tog)*, 38(5):1–12, 2019. 1

[26] Yan Yan, Yuxing Mao, and Bo Li. Second: Sparsely embedded convolutional detection. *Sensors*, 18(10):3337, 2018. 1, 2, 6

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

[28] Zetong Yang, Yanan Sun, Shu Liu, and Jiaya Jia. 3dssd: Point-based 3d single stage object detector. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11040–11048, 2020. 1, 2, 7

[29] Maosheng Ye, Shuangjie Xu, and Tongyi Cao. Hvnet: Hybrid voxel network for lidar based 3d object detection. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 1631–1640, 2020. 1, 2

[30] Hongwei Yi, Shaoshuai Shi, Mingyu Ding, Jiankai Sun, Kui Xu, Hui Zhou, Zhe Wang, Sheng Li, and Guoping Wang. Segvoxelnet: Exploring semantic context and depth-aware features for 3d vehicle detection from point cloud. In *2020 IEEE International Conference on Robotics and Automation (ICRA)*, pages 2274–2280. IEEE, 2020. 1, 2

[31] Tianwei Yin, Zhou Xingyi, and Philipp Krahenbuhl. Centerbased 3d object detection and tracking. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11784–11793, 2021. 6

[32] J. H. Yoo, Y. Kim, J. S. Kim, and J. W. Choi. Generating joint camera and lidar features using cross-view spatial feature fusion for 3d object detection. *arXiv preprint arXiv:2004.12636*, 3, 2020. 7

[33] Yifan Zhang, Qingyong Hu, Guoquan Xu, Yanxin Ma, Jianwei Wan, and Yulan Guo. Not all points are equal: Learning highly efficient point-based detectors for 3d lidar point clouds. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18953–18962, 2022. 1, 2, 5, 6, 7

[34] Na Zhao, Tat-Seng Chua, and Gim Hee Lee. Sess: Self-ensembling semi-supervised 3d object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11079–11087, 2020. 1, 2

[35] Wu Zheng, Weiliang Tang, Sijin Chen, Li Jiang, and Chi-Wing Fu. Cia-ssd: Confident iou-aware single-stage object detector from point cloud. In *Proceedings of the AAAI conference on artificial intelligence*, volume 35, pages 3555–3562, 2021. 1, 2

[36] Wu Zheng, Weiliang Tang, Li Jiang, and Chi-Wing Fu. Sessd: Self-ensembling single-stage object detector from point cloud. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14494–14503, 2021. 2

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

[38] Yin Zhou and Oncel Tuzel. Voxelnet: End-to-end learning for point cloud based 3d object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4490–4499, 2018. 1, 2

1026  
1027  
1028  
1029  
1030  
1031  
1032  
1033  
1034  
1035  
1036  
1037  
1038  
1039  
1040  
1041  
1042  
1043  
1044  
1045  
1046  
1047  
1048  
1049  
1050  
1051  
1052  
1053  
1054  
1055  
1056  
1057  
1058  
1059  
1060  
1061  
1062  
1063  
1064  
1065  
1066  
1067  
1068  
1069  
1070  
1071  
1072  
1073  
1074  
1075  
1076  
1077  
1078  
1079