

CenterNet3D: An Anchor free Object Detector for Autonomous Driving

Guojun Wang, Bin Tian*, Yunfeng Ai, Tong Xu, Long Chen, *Senior Member, IEEE*, Dongpu Cao, *Senior Member, IEEE*

Abstract—Accurate and fast 3D object detection from point clouds is a key task in autonomous driving. Existing one-stage 3D object detection methods can achieve real-time performance, however, they are dominated by anchor-based detectors which are inefficient and require additional post-processing. In this paper, we eliminate anchors and model an object as a single point—the center point of its bounding box. Based on the center point, we propose an anchor-free CenterNet3D Network that performs 3D object detection without anchors. Our CenterNet3D uses keypoint estimation to find center points and directly regresses 3D bounding boxes. However, because inherent sparsity of point clouds, 3D object center points are likely to be in empty space which makes it difficult to estimate accurate boundary. To solve this issue, we propose an auxiliary corner attention module to enforce the CNN backbone to pay more attention to object boundaries which is effective to obtain more accurate bounding boxes. Besides, our CenterNet3D is Non-Maximum Suppression free which makes it more efficient and simpler. On the KITTI benchmark, our proposed CenterNet3D achieves competitive performance with other one stage anchor-based methods which show the efficacy of our proposed center point representation.

I. INTRODUCTION

THE 3D object detection is the task to recognize and locate objects in 3D scene. It serves as a fundamental task for 3D scene understanding with wide applications in autonomous driving. Recent approaches utilize various types of data, including monocular images [1][2], stereo images [3][4]

This work was supported by the Key-Area Research and Development Program of Guangdong Province (2020B090921003), the National Natural Science Foundation of China (61503380), the Intel Collaborative Research Institute for Intelligent and Automated Connected Vehicles (“ICRI-IACV”). (Corresponding author: Bin Tian)

G. Wang is with the State Key Laboratory of Automotive Simulation and Control, Jilin University, Changchun 130022, China (email: 839977837wgj@gmail.com).

B. Tian is with the State Key Laboratory of Management and Control for Complex Systems, Institute of Automation, Chinese Academy of Sciences, Beijing 100190, China, with the Qingdao Academy of Intelligent Industries, Shandong, China (e-mail: bin.tian@ia.ac.cn).

Y. Ai is with the School of Artificial Intelligence, University of Chinese Academy of Sciences, Beijing 100049, China (e-mail: aiyunfeng@ucas.ac.cn).

Tong Xu is with the Waytous Inc., Beijing 100080, China. (email: tong.xu@waytous.com)

L. Chen is with School of Data and Computer Science, Sun Yat-sen University, Guangzhou, Guangdong, P.R. China (e-mail: chenl46@mail.sysu.edu.cn)

D. Cao is with the Department of Mechanical and Mechatronics Engineering, University of Waterloo, 200 University Ave West, Waterloo ON, N2L3G1 Canada (e-mail: dongpu.cao@uwaterloo.ca).

and point cloud [5][6] from LiDAR. Unlike 2D images, point cloud data has some unique properties. Each point in point clouds reflects the surface of physical objects in the real world. In bird eye view, point cloud data is scale-invariant, and objects are naturally separated, which facilitates the detection of occluded objects. Thus, LiDAR has become an indispensable sensor for autonomous driving. Due to the operation mechanism of LiDAR, point clouds are sparse, unordered which making them impossible to directly use convolution neural networks (CNNs) to parse them. Therefore, how to convert and utilize raw point cloud data has become the primary problem in 3D detection task.

In order to handle these problems and utilize the advantages of LiDAR data, many approaches have been proposed recently. Most existing methods convert point clouds from sparse points to compact representations with 2D/3D voxelization and PointNet feature extractor is applied in each voxel. We call these methods voxel-based methods, which conduct voxelization on the whole point cloud. Then the learned voxel features are converted to 2D/3D pseudo image feature map, these methods achieve both good accuracy and high inference speed. VoxelNet [5] is the pioneer of this type of method which groups point cloud data into voxels and applies PointNet to extract voxel features. Then the voxel features are processed with standard Full Convolution Network (FCN) including 2D and 3D convolutions. These voxel-based approaches such as VoxelNet [5], along with its successors, SECOND [6], PointPillar [7], SA-SSD [9] and HVNet [8] densely place pre-defined bounding boxes, called anchors, over the final feature maps and classify them directly. And the offsets from the anchors are predicted to generate bounding boxes.

But the use of anchor boxes has three drawbacks. First, a large number of anchor boxes are needed to place to overlap different aspect ratios and scales which introduces the extra burden. Thus, they require many hyper-parameters and design choices, such as anchor ranges, anchor sizes, and orientations which require prior knowledge about the statistics of objects. Second, the IOU matching threshold must be carefully determined to obtain the appropriate positive and negative samples, and the performance of the networks is very sensitive to it. Third, Non-Maximum Suppression (NMS) is necessary for anchor-based methods to suppress the overlapped high confident detection bounding boxes. So it also introduces extra computational cost which is not conducive to model

deployment in practical applications.

Recently, academic attention has been geared toward anchor-free detectors due to the emergence of FPN [10] and Focal Loss [11], such as [12][13][14][15] and [16]. These anchor-free detectors directly find objects without pre-defined anchors in two different ways. One way is to first locate several pre-defined keypoints and then regress the sizes and offset of objects. We call this type of anchor-free detector as keypoint based method [14][15][16]. Another way is to use the center region of objects to define positives and then predict the four distances from positives to the object boundary. We call this kind of anchor-free detectors as anchor point methods [12][13]. These anchor-free detectors in 2D images have eliminated those hyperparameters related to anchors and achieved better performance, making them more potential in terms of generalization ability. CenterNet is the first Non-Maximum Suppression free 2D detector based on key points, it is natural to expand its ideas for 3D point cloud, just like AFDet [17]. The differences from [17] are that we design different losses for bounding box regression. In addition, we propose an auxiliary corner classification loss to help the CNN backbone learn the structure and scale information of objects. More importantly, we have been doing related research for a long time and preparing to publish before AFDet was made public.

Inspired by the success of CenterNet [16], a keypoint based anchor-free detectors, in 2D perspective image space, we propose a CenterNet3D Network that performs 3D object detection to generate oriented 3D bounding boxes without the predefined anchors in 3D point cloud. Our method detect objects in LiDAR point clouds by representing them as center points of bounding boxes. The properties of bounding boxes such as 3D size and orientation are then regressed directly from feature maps at the center location. However, the inherent sparsity of point clouds makes this approach difficult to learn effective scene context. In images there often exists a pixel near the object center, but it is often not the case in point clouds. As LiDARs only capture surfaces of objects, 3D object centers are likely to be in empty space. Besides center point may be far from the boundary of an object so it is difficult to estimate accurate boundaries. To solve this problem, we propose an auxiliary corner attention module to predict the object corner points and category only in the training stage. The corner attention module can enforce the CNN backbone to pay more attention to object boundaries which is helpful to learn more effective corner heatmap and regress more accurate bounding boxes.

Because we only have one positive “anchor” per object, and hence inference is a single network forward-pass, without Non-Maximum Suppression for post-processing. Peaks in center heatmap correspond to object centers. The main contributions of our proposed method can be summarized as follows:

1. We propose a novel 3D object detection head CenterNet3D Head for point cloud, each object is represented by a center point of its bounding box. Other properties, such as object size and orientation are then regressed directly from feature maps at each center location. Inference is a single network forward-pass, without Non-Maximum Suppression for

post-processing.

2. Based on CenterNet3D Head, we architect the CenterNet3D Network for LiDAR point clouds and this is a keypoint based anchor-free 3D detection method which achieves competitive accuracy compared with the previous one-stage detectors on the KITTI test set.

3. We propose an auxiliary corner attention module to enforce the CNN backbone to learn discriminative corner features, which makes the network aware of object size and structure information. The corner attention module is implemented with a corner classification module that enables CenterNet3D method to obtain better performance without extra cost.

II. RELATED WORK

Currently, there are four types of point cloud representations as input for 3d detectors. 1) point-based representation [18][19][20]. The raw point cloud is directly processed, and bounding boxes are predicted based on each point. 2) voxel-based representation [5][6][7][21][22]. The raw point clouds are converted to compact representations with 2D/3D voxelization. 3) Mixture of representations [23][24]. In these methods, both points and voxels are used as inputs, and their features are fused at different stages of the networks for bounding box prediction. Different methods may consume different types of point cloud representation, in this paper, we adopt voxel-based representations to obtain a tradeoff between efficiency and accuracy. Below, we will briefly review one-stage and two-stage 3D object detection, and then we emphasize possibly related anchor-free 3D object detection.

A. One-Stage 3D Detection

Similar to one-stage 2D detectors, one-stage 3D detectors are proposed to enhance the computational efficiency by processing the input data once in a fully convolutional network. First, a contiguous and regular feature representation is constructed by 2D/3D voxelization. Then a 2D/3D CNN backbone is applied to extract features to obtain a 3D bounding box. Zhou et al. proposed to extracted voxel feature by a simplified PointNet [39]. On the observation of sparsity in voxels, [6] proposes a sparse convolution algorithm and integrates it into the original framework of [5] to accelerate the calculation of convolutional layers. To improve computational efficiency further, PointPillars [7] further simplifies SECOND by implementing voxelization only in the BEV plane. To recover spatial and structure information, [9] proposes a detachable auxiliary network to learn structure information and exhibits better localization performance without extra cost. Our proposed method is also built on top of a general architecture, similar to [6][9], to reduce the complexity of computation.

B. Two-Stage 3D Detection

Unlike one-stage approaches that directly produce 3D bounding box, two-stage approaches firstly generate plausible candidate proposals, then more accurate bounding boxes are obtained by re-using the point cloud within these candidate

proposals in the second stage. Some image-driven methods [34][33] have been proposed to lift a set of 3D regions of interest (ROI) from the image and then apply a PointNet to extract ROI feature by gathering the interior points with transformed canonical coordinates. Shi et al. [18] proposed to generate 3D ROIs from raw point clouds by using a PointNet++ backbone to segment the foreground points from the scene. Its variant work [23] generates ROIs by an efficient 3D CNN. Shi et al. [35] enriched the ROI features by performing intra-object part-aware analysis and demonstrated the effectiveness of reducing the ambiguity of the bounding boxes. To reduce computational burden, [19] seeds each point with a new spherical anchor to achieve high recall with less computation. The inherent complex design of these two-stage methods makes them not applicable for autonomous driving.

C. Anchor-Free 3D Detection

Most of one-stage and two-stage approaches above have relied heavily on the design of pre-defined anchors or pre-defined object sizes which require extra computational burden and hyper-parameters. To our knowledge, there are few anchor-free 3D detectors for LiDAR-based point clouds. SGPN [36] uses a single network to predict point grouping proposals and corresponding semantic class. Then a similarity matrix was learned to group points together. This method is not scalable since the size of similarity matrix grows quadratically with the number of points. 3D-BoNet [37] directly regresses 3D bounding boxes for all instances in a point cloud, while simultaneously predicting a point-level mask for each instance. VoteNet [38] uses PointNet++ [39] to generate seed points, and independently generates votes through a shared voting module. Then the clustered points are refined to obtain box proposals. Though with some anchor-free flavor, VoteNet is not strictly anchor-free because it uses anchor boxes for the size regression, similar to Point-RCNN [18]. HotSpotNet [28] takes the voxels located in the center region of bounding boxes as positive samples and then directly regresses bounding boxes. It eliminates the anchors, but still requires additional Non-Maximum Suppression post-processing.

III. CENTERNET3D NETWORK

Our CenterNet3D represents LiDAR objects as center points of their bounding boxes. Each center point likelihood is individually predicted on each pixel of CNN feature map. For the voxelization representations, a voxel is a center sample if the center point of a bounding box resides in. Each voxel can be projected to a neuron on the feature map based on its location. During training, neurons on the feature map are assigned as centers and non-centers for each object category which are trained by a binary classifier. In inference, a neuron on a feature map is considered as a center point if it gives a local peak in a predicted heatmap.

The proposed CenterNet3D consists of a 3D feature extractor, a 2D feature extractor, and CenterNet3D head. CenterNet3D head has three modules for center classification, box regression, and corner classification. Our CenterNet3D head can be stacked

on top of any voxel-based detectors for point clouds. The whole architecture of our proposed CenterNet3D is shown in Fig. 1. The input LiDAR point clouds are voxelized into regular grids $I \in \mathbb{R}^{L \times W \times H}$ of size L, W, H along the X, Y, Z axes respectively for a frame of point cloud. The regular grids I pass through the 3D CNN to generate the 3D feature maps. Then the 3D feature maps are transformed into 2D feature maps by collapsing in Z axis and passed into 2D feature extractor. Finally, the three subnets will guide the supervision and generate the predicted 3D bounding boxes. The center points and corner points assignment happens at the last feature maps of the 2D feature extractor. The details of network architecture and the three subnets for supervision will be described below.

A. CenterNet3D Head

Our CenterNet3D head network consists of three modules: 1) a center classification module that predicts the probability of center point for an object category. 2) a box regression module that generates an eight-channel feature map to regress the box properties of objects: offset regression that regresses the discretization error of center 2D locations caused by downsampling stride; z-coord regression that predicts the center location in Z axis; size regression that regresses the 3D size of objects; direction regression that regresses the rotation angle around Z axis. 3) a corner classification module that predicts the probability of corner points for an object category. For box regression, there are two alternative methods to implement. The direct way is to predict all the box properties with one common module. Alternatively, we can predict different properties with different modules to learn the specific appearance characteristics of an object instance. We term the two variants of our network according to the way to regress box properties that regressing all box properties with one module or different modules as CenterNet3D-Merge and CenterNet3D-Split.

1) Center Classification

The center classification outputs center heatmaps with several convolutional layers and each center heatmap corresponds to one category. let $\hat{Y} \in [0, 1]^{\frac{L}{R} \times \frac{W}{R} \times C}$ be the center heatmaps, where R is the downsampling stride and C is the number of center point types which equals to the number of classes. The output prediction \hat{Y} is downsampled by a factor R . A prediction $\hat{Y}_{x,y,c} = 1$ corresponds to a detected center point, while $\hat{Y}_{x,y,c} = 0$ is background. We generate ground truth center heatmaps $Y \in [0, 1]^{\frac{L}{R} \times \frac{W}{R} \times C}$ following Law and Deng [14]. For each ground truth center point $p \in \mathbb{R}^2$ of class c , we compute a low-resolution equivalent $\tilde{p} = \left\lfloor \frac{p}{R} \right\rfloor$. We then splat all ground truth center points onto the heatmaps $Y \in [0, 1]^{\frac{L}{R} \times \frac{W}{R} \times C}$ using a gaussian kernel $Y_{xyc} = \exp\left(-\frac{(x-\tilde{p}_x)^2 + (y-\tilde{p}_y)^2}{2\sigma_p^2}\right)$, where σ_p is an object size-adaptive standard deviation which is 1/6 gaussian radius in our work. We determine the gaussian radius by the size of an object by ensuring that one point within the radius would generate a bounding box with at least t IoU with the ground-truth annotation. Because those false center predictions

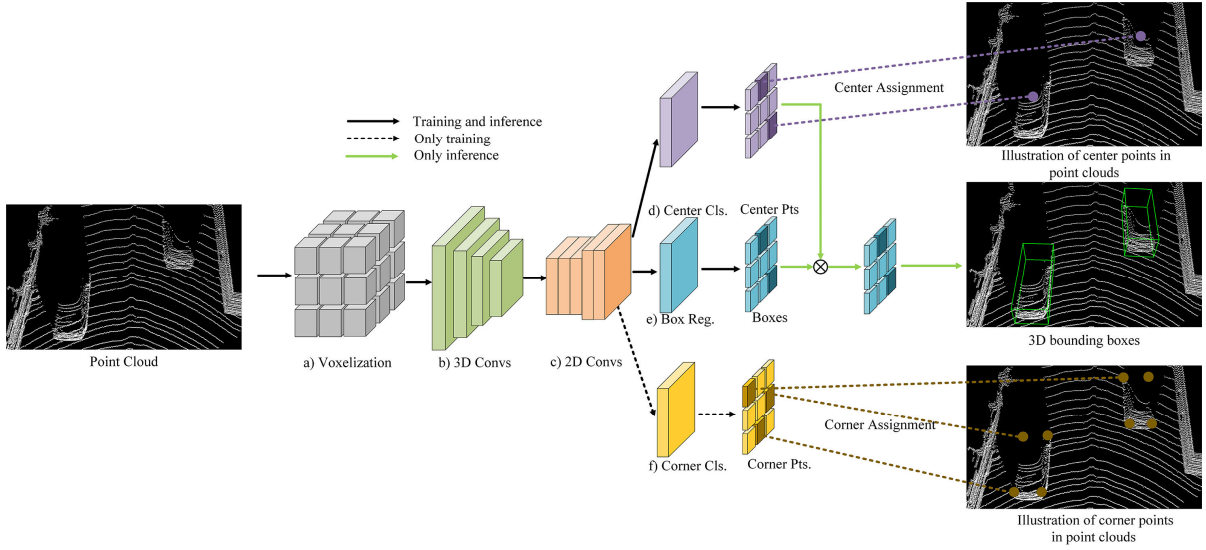


Fig. 1. Outline of CenterNet3D. The point cloud is (a) voxelization (b) 3D Convs including sparse convolution and submanifold convolution proposed in [6]. c) 2D Convs including transposed convolution and deformable convolution to produce 2D feature maps. These feature maps pass into three modules to perform d) Center Classification e) 3D Bounding Box regression (f) Corner Classification to train the network. During the inference only (d) Center Classification and (e) 3D Bounding Box Regression are performed to obtain center points and bounding boxes respectively. The center classification branch does not contribute to inference. The final bounding boxes are obtained by processing the center points prediction through max pooling, without the need of IOU-based Non-Maximum Suppression (NMS).

that are close to their respective ground truth locations can still produce a box that sufficiently overlaps the ground-truth box, a

$$L_{cls} = \frac{-1}{N} \sum_{xyc} \begin{cases} (1 - \hat{Y}_{xyc})^\alpha \log(\hat{Y}_{xyc}) & \text{if } Y_{xyc} = 1 \\ (1 - Y_{xyc})^\beta (\hat{Y}_{xyc})^\alpha \log(1 - \hat{Y}_{xyc}) & \text{otherwise} \end{cases}$$

where α and β are hyper-parameters of the focal loss, and N is the number of center points in input I . The normalization by N is chosen as to normalize all positive focal loss instances to 1. We use $\alpha = 2$ and $\beta = 4$ in all our experiments, following Law and Deng [14].

2) Box Regression

The bounding box regression module only happens on the positive center point features. For each bounding box, an eight-dimensional vector $[d_x, d_y, z, l, w, h, \cos(r), \sin(r)]$ is regressed to represent the object instance in LiDAR point clouds. d_x, d_y are the discretization deviations of center points on the last feature map, z is the absolute coordinate value in Z axis, l, w , and h are the 3D sizes: length, width and height, $\cos(r), \sin(r)$ is trigonometric function value of rotation angle r around Z axis. Thus, box regression includes offset regression, z-coord regression, size regression and direction regression.

a) Offset Regression

To recover the 2D discretization error $\frac{p}{R} - \tilde{p}$ caused by the output stride, the offset regression is used to predict an offset feature map $\hat{O} \in \mathcal{R}^{\frac{L}{R} \times \frac{W}{R} \times 2}$ for each center point. All classes c share the same offset prediction. Because the ground truth offsets fall between 0 and 1. A logistic activation is used to constrain the offset predictions to fall in this range. The offset is trained with an L1 loss:

penalty-reduced pixelwise logistic regression with focal loss is used as training objective [11].

$$L_{off} = \frac{1}{N} \sum_p \sum_{i \in \{\Delta x, \Delta y\}} |\sigma(\hat{O}_{\tilde{p}, i}) - O_{p, i}|$$

b) Direction Regression

To predict the rotation angle around Z axis and solve the adversarial example problem between the cases of 0 and π radians, we encode each rotation angle r as $(\cos(r), \sin(r))$, and decode the r as $\text{atan2}(\sin(r), \cos(r))$ during inference. The rotation angle is factored into two correlated values. Thus, the direction regression predicts a direction feature map $\hat{D} \in \mathcal{R}^{\frac{L}{R} \times \frac{W}{R} \times 2}$ for each center point. The direction is also trained with an L1 loss:

$$L_{dir} = \frac{1}{N} \sum_p \sum_{i \in \{\sin(r), \cos(r)\}} |\sigma(\hat{D}_{\tilde{p}, i}) - D_{p, i}|$$

c) Z-coord Regression

The z-coord regression is used to predict the center location of the bounding box in Z axis. It outputs a z-coord feature map $\hat{Z} \in \mathcal{R}^{\frac{L}{R} \times \frac{W}{R} \times 1}$ for each center point. All classes c share the same z-coord prediction. However, owing to the unbounded regression targets z-coord, the model is sensitive to outliers. These outliers, which can be regarded as hard samples, will produce excessively large gradients that are harmful to the training process. The inliers, which can be regarded as the easy samples, contribute little gradient to the overall gradients compared with the outliers. Considering this issue, the balanced L1 loss proposed in [25] is employed for training z-coord

regression:

$$L_z = \frac{1}{N} \sum_p L_b(|\hat{Z}_{\bar{p}} - Z_p|)$$

where, L_b is balanced L1 loss:

$$L_b(x) = \begin{cases} \frac{a}{b}(b|x| + 1)\ln(b|x| + 1) - a|x| & \text{if } |x| < 1 \\ \gamma|x| + C & \text{otherwise} \end{cases}$$

in which a, b and γ are hyper-parameters of balanced L1 loss, and they are constrained by

$$a\ln(b + 1) = \gamma$$

We use $a = 0.5$ and $\gamma = 1.5$ in all our experiments, following [25].

d) Size Regression

To predict the length l , width w and height h , the size regression is employed. To limit the computational burden, we regress a single size for all classes. The size regression output a size feature map $\hat{S} \in \mathcal{R}^{\frac{L}{R} \times \frac{W}{R} \times 3}$. Like z-coord regression, the size regression also easily introduces training imbalance owing to the unbounded regression targets. So, the balanced L1 loss proposed in [25] is also employed for training size regression:

$$L_{size} = \frac{1}{N} \sum_p \sum_{i \in \{l, w, h\}} L_b|\hat{s}_{\bar{p}, i} - s_{p, i}|$$

3) Corner Classification

Unlike 2D images there often exists a pixel near the object center, as LiDARs only capture surfaces of objects, 3D object centers are likely to be in empty space. Besides center points may be far from the boundary of an object so it is difficult to estimate accurate boundaries. We want our model to learn the object shape and structure information, so we introduce another auxiliary supervision signal for corner classification. Like center classification, the corner classification module outputs corner heatmaps and each corner heatmap corresponds to one category. let $\hat{A} \in [0, 1]^{\frac{L}{R} \times \frac{W}{R} \times C}$ be the corner heatmaps, where R is the downsample stride and C is the number of classes. A prediction $\hat{A}_{x, y, c} = 1$ corresponds to a detected corner point, while $\hat{A}_{x, y, c} = 0$ is background. We generate ground truth corner heatmaps $A \in [0, 1]^{\frac{L}{R} \times \frac{W}{R} \times C}$ following the above center classification. The training corner loss L_{cor} is also a penalty-reduced focal loss similar to center classification. The corner classification branch does not contribute to inference.

B. Training and Inference

The final loss for our proposed CenterNet3D is the weighted sum of all above classification loss and regression losses:

$$L = \delta_{cls}L_{cls} + \delta_{off}L_{off} + \delta_zL_z + \delta_{size}L_{size} + \delta_{dir}L_{dir} + \delta_{cor}L_{cor}$$

where, $\delta_{cls}, \delta_{off}, \delta_z, \delta_{size}, \delta_{dir}$ and δ_{cor} are the weights to balance the center classification, offset regression, z-coord regression, size regression, direction regression and corner classification loss.

At inference time, we first extract the center predictions in center heatmaps for each class independently. We then filter all

center prediction by whether whose value is greater or equal to its 8-connected neighbors. Then we only keep those center predictions whose value is above the predefined threshold as detected center points.

Let $\hat{P}_c = \{\hat{x}_i, \hat{y}_i\}_{i=1}^n$ be the set of n detected center points of class c . (\hat{x}_i, \hat{y}_i) is the integer coordinates in heatmap \hat{Y} . $(\delta\hat{x}_i, \delta\hat{y}_i), \hat{z}_i, (\hat{l}_i, \hat{w}_i, \hat{h}_i), (\hat{s}\hat{m}_i, \hat{c}\hat{o}s_i)$ are corresponding offset, z-coord, size and direction prediction at location (\hat{x}_i, \hat{y}_i) . We use the center prediction value $\hat{Y}_{x_i, y_i, c}$ as a measure of its detection confidence, and produce a bounding box: $(\hat{x}_i + \delta\hat{x}_i, \hat{y}_i + \delta\hat{y}_i, \hat{z}_i, \hat{l}_i, \hat{w}_i, \hat{h}_i, \text{atan2}(\hat{s}\hat{m}_i, \hat{c}\hat{o}s_i))$. All outputs are produced directly from the center prediction without the need for Non-Maximum Suppression or other post-processing. The center prediction filter serves as a sufficient Non-Maximum Suppression alternative and can be implemented efficiently using a 3×3 max pooling operation.

IV. EXPERIMENTS

In this section, we summarize the dataset in Section A and introduce the implementation details of our proposed CenterNet3D in Section B. In Section C we evaluate our method on the challenging 3D detection Benchmark KITTI [41]. In Section D, we present ablation studies about our method.

A. Dataset and Evaluation

We evaluate our proposed CenterNet3D detector on the KITTI 3D/BEV object detection benchmark. The dataset contains 7,481 annotated LiDAR frames for training with 3D bounding boxes for object classes such as cars, pedestrians and cyclists. Following the common protocol, we further divide the training data into a training set with 3,712 frames and a validation set with 3,769 frames. Additionally, KITTI provides 7,518 LiDAR frames without labeling for testing. We conduct experiments on the most commonly used car category and use average precision (AP) with an (IOU) threshold 0.7 as evaluation metric. The benchmark considers three levels of difficulties: *easy*, *moderate*, and *hard* based on the object size, occlusion state, and truncation level. The average precision (AP) is calculated using 40 recall positions. In order to further compare the results with other methods on the KITTI 3D detection benchmark, we divide the KITTI training dataset into 4:1 for training and validation, and report performance on the KITTI test dataset.

B. Implementation Details

1) Backbone Network

In all experiments, we employ the commonly used backbone network [5][6] as our feature extractor. We crop the point cloud based on the ground-truth distribution at $[-3, 1] \times [-40, 40] \times [0, 70]$ m along the Z, Y, X axes. We use a voxel size of $V_x = V_y = 0.05m, V_z = 0.1m$ to voxelize the cropped point cloud into a regular grid map of size $40 \times 1600 \times 1408$ along the Z, Y, X axes respectively. The car instance only occupies 30-60 voxels, which is a typical small object. In order to improve detection performance, we apply a single transposed convolution layer and deformable convolution layer to obtain

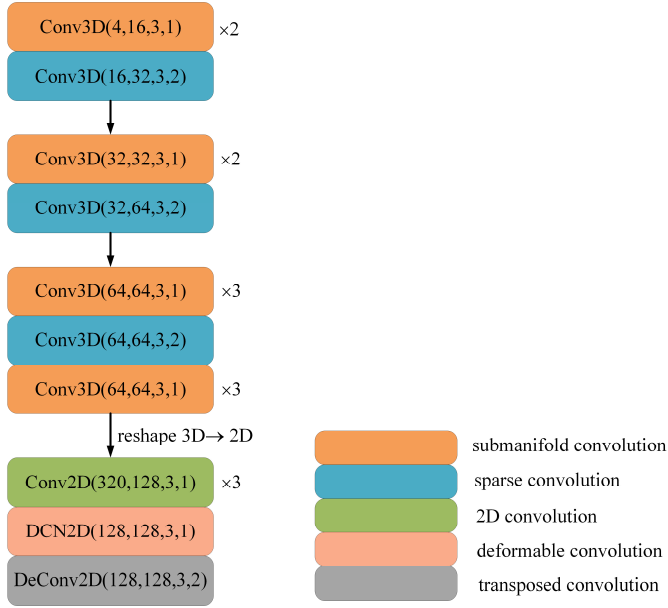


Fig. 2 The details of backbone network for CenterNet3D. Conv(cin, cout, k, s) represents a convolutional block, where cin, cout, k, s denotes input channel number, output channel number, kernel size and stride respectively. Each block consists of Convolution, BatchNorm, Leaky ReLU. The stride of each 3D sparse convolution is 2 which is used to downsample the feature maps.

high resolution and strong semantic feature maps. The details of our architecture of the backbone network are shown in Fig. 2. We use the simplified version of Voxel Feature Encoder, i.e. VFE [5], by taking the mean of a fixed amount of points sampled in a voxel. Our backbone network has 3D and 2D parts. The 3D part has 10 3D submanifold convolution blocks and 3 3D sparse convolution blocks. The 2D part has 3 2D convolution blocks, 1 deformable convolution block, and 1 transposed convolution block. We down-sample in total three times with 3D sparse convolution blocks, then the 3D feature maps are transformed to 2D by collapsing the height dimension. Finally, the 2d feature maps are processed by deformable convolution block and transposed convolution block to obtain high-resolution and strong semantic feature maps for CenterNet3D head.

2) CenterNet3D Head

Since the output feature maps of the backbone network are collapsed to bird eye view, we thus in this paper assign center points and corner points in bird eye view. For CenterNet3D-Split, the detection head has six modules (center classification, offset regression, z-coord regression, size regression, direction regression and corner classification), for CenterNet3D-Merge, the detection head has three modules (center classification, box regression and corner classification). All sub-task modules share the common backbone network. For each module, the features of the backbone are passed through a separate 3×3 convolution, Leaky ReLU, and another 1×1 convolution. For center classification, corner classification, and offset regression, another logistic activation is used to constrain the predictions in between 0 and 1. Compared with anchor-based methods, our method only has one positive sample for each ground truth which leads to an extreme imbalance of positive and negative samples. To mitigate the imbalance, we set $t = 0.01$ to obtain more non-zero pixels in ground truth center heatmap $Y \in [0, 1]^{\frac{L}{R} \times \frac{W}{R} \times C}$ and corner heatmap $A \in [0, 1]^{\frac{L}{R} \times \frac{W}{R} \times C}$. We set $\delta_{cls} = \delta_{cor} = 0.5$, $\delta_{off} = \delta_z = \delta_{size} = \delta_{dir} = 1$ in the weighted sum of losses.

3) Training and Inference

The network is trained by ADAM [26] optimizer with fixed weight decay 0.01. The learning schedule is one-cycle policy [27] with maximum learning rate $2.25e-3$, division factor 10, momentum ranges from 0.95 to 0.85. The network is trained with batch size 2 for 50 epochs. During the inference, a 3×3 max pooling operation is applied on the center heat map, then we keep 50 predictions whose center confidences are larger than 0.3.

4) Data Augmentation

We perform a common cut-and-paste strategy [6][9] for data augmentation. Specifically, we randomly adds foreground instances with their inner points from other frames to current point cloud. Each instance is followed by a collision test to

TABLE I
PERFORMANCE COMPARISON WITH PREVIOUS METHODS ON KITTI TEST SERVER. 3D OBJECT DETECTION METRICS ARE USED, REPORTED BY THE AVERAGE PRECISION (AP) WITH IOU THRESHOLD 0.7. THE BOLD VALUE INDICATES THE TOP PERFORMANCE.

Method	Modality	Stage	3D Detection (Car)			Bev Detection (Car)			FPS
			Easy	Moderate	Hard	Easy	Moderate	Hard	
MV3D [42]	L+C	Two	74.97	63.63	54.00	86.49	79.98	72.23	2.8
F-PointNet [33]	L+C	Two	82.19	69.79	60.59	91.17	84.67	74.77	5.9
AVOD-FPN [43]	L+C	Two	83.07	71.76	65.73	90.99	84.82	79.62	10
PointRCNN [18]	L	Two	86.96	75.64	70.70	92.13	87.39	82.72	10
F-ConvNet [34]	L+C	Two	87.36	76.39	66.69	91.51	85.84	76.11	2.1
Fast-PointRCNN [23]	L	Two	85.29	77.40	70.24	90.87	87.84	80.52	15.4
MMF [44]	L+C	Two	88.40	77.43	70.22	93.67	88.21	81.99	12.5
STD [19]	L	Two	87.95	79.71	75.09	94.74	89.19	86.42	10
PIXOR [21]	L	One	-	-	-	83.97	80.01	74.31	-
ComplexYOLO [22]	L	One	55.93	47.34	42.60	77.24	68.96	64.95	-
VoxelNet [5]	L	One	77.82	64.17	57.51	87.95	78.39	71.29	4.4
SECOND-V1.5 [6]	L	One	84.65	75.96	68.71	91.81	86.37	81.04	20
TANet [45]	L	One	84.39	75.94	68.82	91.58	86.54	81.19	25
PointPillars [7]	L	One	82.58	74.31	68.99	90.07	86.56	82.81	42
HotSpotNet [28]	L	One	88.12	78.34	73.49	94.06	88.09	83.24	20
CenterNet3D-Merge	L	One	86.10	77.62	72.57	90.26	88.29	83.33	26
CenterNet3D-Split	L	One	86.20	77.90	73.03	91.80	88.46	83.62	25

avoid the violation of the physical rule. All ground-truth boxes are individually augmented. Each box is randomly rotated and translated. The noise for the rotation is uniformly drawn from $[-\frac{\pi}{4}, \frac{\pi}{4}]$ and the noise for the translation is drawn from $\mathcal{N}(0,1)$, $\mathcal{N}(0,1)$ and $\mathcal{N}(0,1)$ for X, Y, Z respectively. In addition, we apply random flipping, global rotation, and global scaling to the whole point cloud. The noise for global rotation is uniformly drawn from $[-\frac{\pi}{4}, \frac{\pi}{4}]$ and the scaling factor is uniformly drawn from $[0.95, 1.05]$.

C. Experimental results on KITTI benchmark

We compare our CenterNet3d 3D point cloud detector with other state-of-the-art approaches by submitting the detection results to the KITTI server for evaluation. As shown in Table I, we evaluate our method on the 3D detection benchmark and the bird's eye view detection benchmark of the KITTI test dataset. As we can see, our proposed CenterNet3D-Split and CenterNet3D-Merge can both achieve competitive performance with anchor-based methods. But our method does not need predefined anchors and not have a complex post-processing process such as IOU-based NMS to filter out results. In particular, CenterNet3D-Split shows better performance on hard level where objects are usually far away, occluded and truncated. The proposed CenterNet3D-Merge performs slightly

worse than CenterNet3D-Split. It is demonstrated that different modules for different sub-tasks can be more sperate and more specialized, they can learn the specific appearance characteristics of an object instance. In the rest of the paper, without further emphasizing, we will adopt CenterNet3D-Split as our method for quantitative evaluations. Note that all other methods listed in Table I are anchor-based except PIXOR [21] and HotSpotNet [28]. PIXOR is a 2D detector which cannot output height predictions. HotSpotNet is an anchor-point based 3D detectors. Although it does not require predefined anchors, it assigns multiple anchor points to one object. So it also requires heavy NMS post-processing which is hard to differentiate and train. The proposed CenterNet3D assigns the “anchor” based solely on location, So our method does not need Non-Maximum Suppression. The inspiring results show the success of representing 3D objects as center points as well as potentials of keypoint based anchor-free detectors in 3D detection. Our approach also beats some classic 3D two-stage detectors, even when they fuse LiDAR and RGB images information. We can conclude that our proposed one-stage anchor free and NMS free CenterNet3D achieve competitive or even better results and the potentials of keypoint based anchor-free 3D detectors are demonstrated.

We also visualize some prediction results on the validation

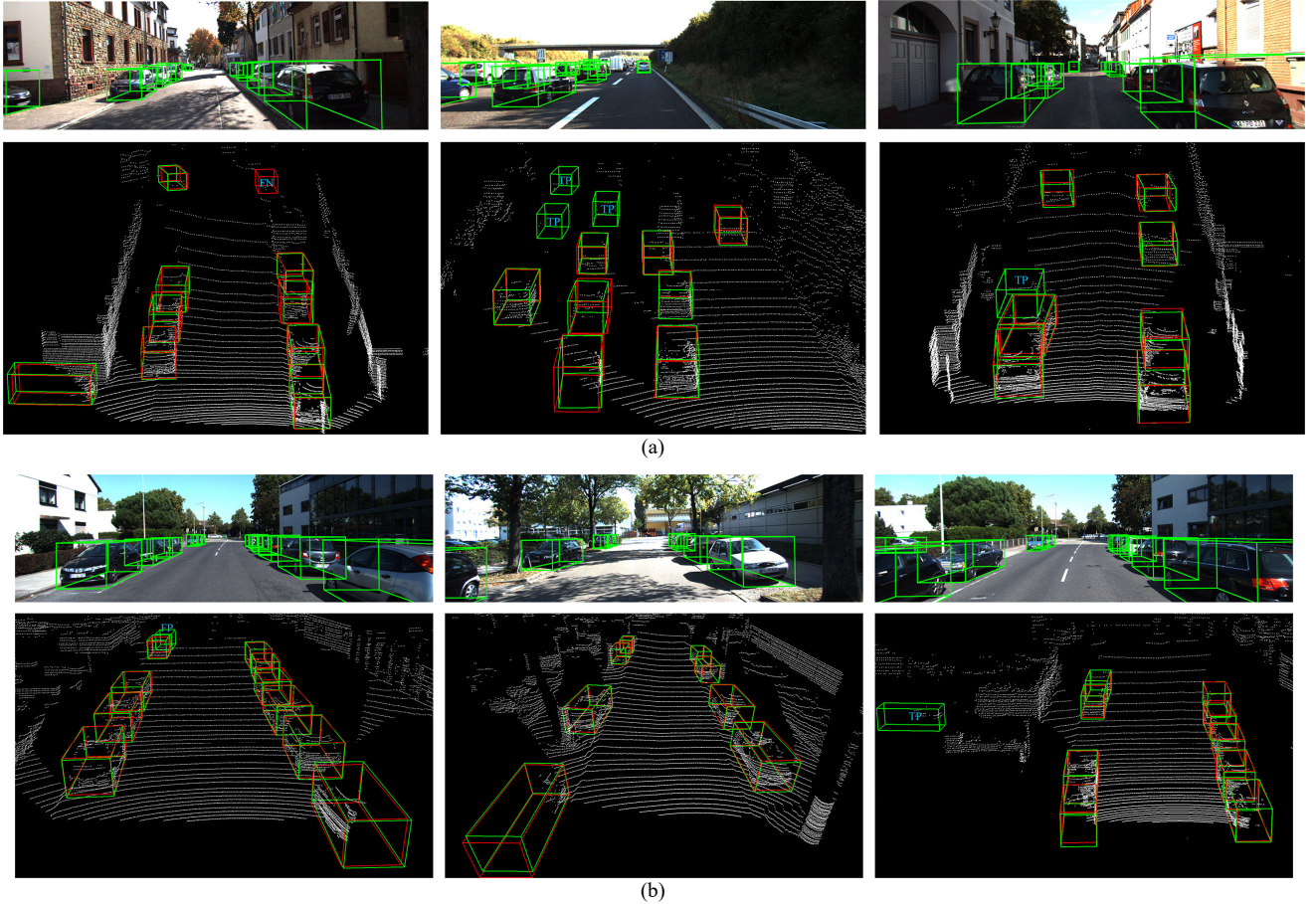


Fig. 3 Qualitative results on KITTI validation set. The predicted bounding boxes are shown in green. The ground truth bounding boxes are shown in red. The predictions are projected onto the RGB images (upper row) for better visualization. The word “FN” inside each box represents a false negative, the word “TP” represents a true positive, but the corresponding ground truth annotation is not provided because of heavy occlusion or invisibility in the camera view. The word “FP” represents a false positive.

set in Fig. 3, and we project the 3D bounding boxes detected from LiDAR to the RGB images for better visualization. The predicted bounding boxes are shown in green. The ground truth bounding boxes are shown in red. The word “FN” inside each box represents a false negative, the word “TP” represents a true positive, but the corresponding ground truth annotation is not provided because of heavy occlusion or invisibility in camera view. However, these occluded or invisible objects are easily detected with 3D LiDAR. The word “FP” represents a false positive which is far away from the sensors and appears with sparse point clouds. In future work, we will focus on investigating ways to incorporate appearance cues from RGB images to prevent false positives.

D. Ablation Studies

1) Corner Attention module

To prove the effectiveness of the corner attention module, we show the results of our CenterNet3D with and without the corner attention module on KITTI validation split for cars in Table II. We can see that when our algorithm trained with the proposed corner attention module, the overall performance is boosted. Especially, great improvement can be observed in hard level. It is demonstrated that the proposed corner attention module is effective, and the corner features are beneficial for hard samples.

TABLE II
EFFECT OF CORNER CLASSIFICATION MODULE ON KITTI VALIDATION SET FOR “CAR” DETECTION.

Method	3D Detection			Bev Detection		
	Easy	Moderate	Hard	Easy	Moderate	Hard
w/o corner classification	88.34	78.30	77.20	89.97	87.81	85.85
w corner classification	89.22/+0.88	79.12/+0.82	78.63/+1.43	90.93/+0.96	88.54/+0.73	87.12/+1.27

2) Box Size Loss

In this part, we evaluate the performance of our approach when using different types of box size loss. which include L1

loss, smooth L1 [40] and balanced L1 loss [25]. Our experiments in Table III show that balanced L1 is considerably better than Smooth L1 and L1 loss.

TABLE III
A COMPARISON OF THE PERFORMANCES OF DIFFERENT BOX SIZE REGRESSION LOSS ON THE KITTI VALIDATION SET FOR “CAR” DETECTION.

Method	3D Detection			Bev Detection		
	Easy	Moderate	Hard	Easy	Moderate	Hard
Smooth L1 Loss	87.92	76.84	75.74	89.97	86.81	85.85
L1 Loss	88.56/+0.64	78.34/+1.50	77.32/+1.58	89.92/-0.05	88.23/+1.42	87.23/+1.38
Balanced L1 Loss	89.22/+1.30	79.12/+2.28	78.63/+2.89	90.93/+0.96	88.54/+1.73	87.12/+1.27

3) Direction Loss

We also compare the performances when using different types of direction loss, which include sine-error Loss [6], bin-based loss [18], residual-based loss [5] and sin-cos loss [21]

used in this paper. Our experiments in Table IV show that sin-cos loss and bin-based loss achieve approximatively excellent performance. However, bin-based loss contains complex encoding and decoding for bin classification and residual regression which is computationally inefficient.

TABLE IV
A COMPARISON OF THE PERFORMANCES OF ANGLE REGRESSION LOSS ON THE KITTI VALIDATION SET FOR “CAR” DETECTION.

Method	3D Detection			Bev Detection		
	Easy	Moderate	Hard	Easy	Moderate	Hard
Residual-based Loss	86.53	77.12	75.23	88.96	86.74	85.28
Sine-error Loss	88.56/+2.03	78.34/+1.22	77.27/+2.04	89.97/+1.01	87.76/+1.02	86.52/+1.22
Bin-based Loss	89.25/+2.72	79.18/+2.06	77.52/+2.29	90.98/+2.02	88.23/+1.49	87.01/+1.73
Sin-cos loss	89.22/+2.69	79.12/+2.00	78.63/+3.40	90.93/+1.97	88.54/+1.80	87.12/+1.84

V. CONCLUSION

We propose a novel representation, object as center points, and a one-stage anchor-free and NMS-free 3D object detector, CenterNet3D, for autonomous driving. Our CenterNet3D is a keypoint based 3D object detector which finds the center points of objects, and directly regresses bounding boxes. To solve the sparsity of point clouds and the hollowness near the object centers, we propose an auxiliary corner classification module to enforce the CNN backbone to learn discriminative corner features, which makes the network aware of object size and structure information. The proposed CenterNet3D is simple, fast, accurate, and end-to-end differentiable without any NMS

postprocessing. On the KITTI benchmark, our proposed CenterNet3D achieves competitive performance with other one stage anchor-based methods which show the potentials of keypoint based anchor-free detectors in 3D detection.

VI. REFERENCES

- [1] A. Simonelli, S. R. Bulo, L. Porzi, M. López-Antequera, and P. Kotschieder, “Disentangling Monocular 3D Object Detection,” in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 1991-1999.
- [2] B. Xu, and Z. Chen, “Multi-level fusion based 3d object detection from monocular images,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 2345-2353.

- [3] P. Li, X. Chen, and S. Shen, "Stereo r-cnn based 3d object detection for autonomous driving," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 7644-7652.
- [4] Z. Qin, J. Wang, and Y. Lu, "Triangulation learning network: from monocular to stereo 3d object detection," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 7607-7615.
- [5] Y. Zhou and O. 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*, 2018, pp. 4490-4499.
- [6] Y. Yan, Y. Mao and B. Li, "SECOND Sparsely Embedded Convolutional Detection," *Sensors*, vol. 18, no. 10, pp. 3337, Oct. 2018.
- [7] A. H. Lang, S. Vora, H. Caesar, L. Zhou, J. Yang and O. Beijbom, "PointPillars: Fast Encoders for Object Detection From Point Clouds," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 12697-12705.
- [8] M. Ye, S. Xu, and T. Cao, "HVNet: Hybrid Voxel Network for LiDAR Based 3D Object Detection," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 1631-1640.
- [9] C. He, H. Zeng, J. Huang, X. S. Hua, and L. Zhang, "Structure Aware Single-stage 3D Object Detection from Point Cloud," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 11873-11882.
- [10] T. Y. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, and S. Belongie, "Feature pyramid networks for object detection," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 2117-2125.
- [11] T. Y. Lin, P. Goyal, R. Girshick, K. He and P. Dollar, "Focal loss for dense object detection," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2980-2988.
- [12] Z. Tian, C. Shen, H. Chen, and T. He, "Fcos: Fully convolutional one-stage object detection," in *Proceedings of the IEEE international conference on computer vision*, 2019, pp. 9627-9636.
- [13] T. Kong, F. Sun, H. Liu, Y. Jiang, and J. Shi, "Foveabox: Beyond anchor-based object detector," *arXiv preprint arXiv:1904.03797*.
- [14] H. Law and J. Deng, "Cornernet: Detecting objects as paired keypoints," in *Proceedings of the European Conference on Computer Vision*, 2018, pp. 734-750.
- [15] X. Zhou, J. Zhuo, and P. Krahenbuhl, "Bottom-up object detection by grouping extreme and center points," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 850-859.
- [16] X. Zhou, D. Wang, and P. Krähenbühl, "Objects as points," *arXiv preprint arXiv:1904.07850*.
- [17] R. Ge, Z. Ding, Y. Hu, Y. Wang, S. Chen, L. Huang, and Y. Li, "AFDet: Anchor free one stage 3d object detection," *arXiv preprint arXiv:2006.12671*.
- [18] S. Shi, X. Wang, and H. Li, "Pointtrcn: 3d object proposal generation and detection from point cloud," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 770-779.
- [19] Z. Yang, Y. Sun, S. Liu, X. Shen, and J. Jia, "Std: Sparse-to-dense 3d object detector for point cloud," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 1951-1960.
- [20] Z. Yang, Y. Sun, Y. S. Liu, and J. Jia, "3dssd: Point-based 3d single stage object detector," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 11040-11048.
- [21] B. Yang, W. J. Luo and R. Urtasun, "Pixor: Real-time 3d object detection from point clouds," in *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, 2018, pp. 7652-7660.
- [22] M. Simon, S. Milz, K. Amende and H. M. Gross, "Complex-yolo: An euler-regionproposal for real-time 3d object detection on point clouds," in *Proceedings of the European Conference on Computer Vision*, 2018, pp. 0-0.
- [23] Y. Chen, S. Liu, X. Shen, and J. Jia, "Fast point r-cnn," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 9775-9784.
- [24] S. Shi, C. Guo, L. Jiang, Z. Wang, J. Shi, X. Wang, and H. 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*, 2020, pp. 10529-10538.
- [25] J. M. Pang, K. Chen, J. P. Shi, H. J. Feng, W. L. Ouyang and D. Lin, "Libra R-CNN: Towards Balanced Learning for Object Detection," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2019, pp. 821-830.
- [26] I. Loshchilov and F. Hutter, "Fixing weight decay regularization in adam," *arXiv preprint arXiv:1711.05101*, 2017.
- [27] L. N. Smith and N. Topin, "Superconvergence: Very fast training of neural networks using large learning rates," *Artificial Intelligence and Machine Learning for Multi-Domain Operations Applications. International Society for Optics and Photonics*, vol. 11006, no. 1100612, 2019.
- [28] Q. Chen, L. Sun, Z. Wang, K. Jia and A. Yuille, "Object as Hotspots: An Anchor-Free 3D Object Detection Approach via Firing of Hotspots," *arXiv preprint arXiv:1912.12791*.
- [29] H. Kuang, B. Wang, J. An, M. Zhang, and Z. Zhang, "Voxel-FPN: Multi-Scale Voxel Feature Aggregation for 3D Object Detection from LIDAR Point Clouds," *Sensors*, vol. 20, no. 3, pp. 704, 2020.
- [30] A. Newell, Z. Huang, and J. Deng, "Associative embedding: End-to-end learning for joint detection and grouping," in *Advances in neural information processing systems*, 2017, pp. 2277-2287.
- [31] X. Zhou, A. Karpur, L. Luo, and Q. Huang, "Starmap for category-agnostic keypoint and viewpoint estimation," in *Proceedings of the European Conference on Computer Vision*, 2018 pp. 318-334.
- [32] B. Wu, A. Wan, X. Yue, and K. Keutzer, "SqueezeSeg: Convolutional neural nets with recurrent crf for real-time road-object segmentation from 3d lidar point cloud," in *IEEE International Conference on Robotics and Automation*, 2018, pp. 1887-1893.
- [33] C. R. Qi, W. Liu, C. Wu, H. Su, and L. J. Guibas, "Frustum pointnets for 3d object detection from rgb-d data," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 918-927.
- [34] Z. Wang, and K. Jia, "Frustum convnet: Sliding frustums to aggregate local point-wise features for amodal 3d object detection," *arXiv preprint arXiv:1903.01864*.
- [35] S. Shi, Z. Wang, J. Shi, X. Wang, and H. Li, "From points to parts: 3d object detection from point cloud with part-aware and part-aggregation network," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020.
- [36] W. Wang, R. Yu, Q. Huang, and U. Neumann, "Sgpn: Similarity group proposal network for 3d point cloud instance segmentation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 2569-2578.
- [37] B. Yang, J. Wang, R. Clark, Q. Hu, S. Wang, A. Markham, and N. Trigoni, "Learning object bounding boxes for 3d instance segmentation on point clouds," in *Advances in Neural Information Processing Systems*, 2019, pp. 6740-6749.
- [38] C. R. Qi, O. Litany, K. He, and L. J. Guibas, "Deep hough voting for 3d object detection in point clouds," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 9277-9286.
- [39] C. R. Qi, L. Yi, H. Su, and L. J. Guibas, "Pointnet++: Deep hierarchical feature learning on point sets in a metric space," in *Advances in neural information processing systems*, 2017, pp. 5099-5108.
- [40] R. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," in *Advances in neural information processing systems*, 2015, pp. 91-99.
- [41] A. Geiger, P. Lenz, C. Stiller, and R. Urtasun, "Vision meets robotics: The kitti dataset," *The International Journal of Robotics Research*, vol. 32, no. 11, pp. 1231-1237.
- [42] X. Chen, H. Ma, J. Wan, B. Li, and T. Xia, "Multi-view 3d object detection network for autonomous driving," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 1907-1915.
- [43] J. Ku, M. Mozifian, J. Lee, A. Harakeh, and S. L. Waslander, "Joint 3d proposal generation and object detection from view aggregation," in *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2018, pp. 1-8.
- [44] M. Liang, B. Yang, Y. Chen, R. Hu, and R. Urtasun, "Multi-task multi-sensor fusion for 3d object detection," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 7345-7353.
- [45] Z. Liu, X. Zhao, T. Huang, R. Hu, Y. Zhou, and X. Bai, "TANet: Robust 3D Object Detection from Point Clouds with Triple Attention," in *AAAI*, 2020, pp. 11677-11684.