

# RoIFusion: 3D Object Detection from LiDAR and Vision

CAN CHEN, LUCA ZANOTTI FRAGONARA, AND ANTONIOS TSOURDOS

Address: School of Aerospace, Transport and Manufacturing, Cranfield University, Cranfield, MK43 0AL, UK

Corresponding author: Luca Zanotti Fragonara (e-mail: l.zanottifragonara@cranfield.ac.uk).

arXiv:2009.04554v1 [cs.CV] 9 Sep 2020

**ABSTRACT** When localizing and detecting 3D objects for autonomous driving scenes, obtaining information from multiple sensor (e.g. camera, LIDAR) typically increases the robustness of 3D detectors. However, the efficient and effective fusion of different features captured from LIDAR and camera is still challenging, especially due to the sparsity and irregularity of point cloud distributions. This notwithstanding, point clouds offer useful complementary information. In this paper, we would like to leverage the advantages of LIDAR and camera sensors by proposing a deep neural network architecture for the fusion and the efficient detection of 3D objects by identifying their corresponding 3D bounding boxes with orientation. In order to achieve this task, instead of densely combining the point-wise feature of the point cloud and the related pixel features, we propose a novel fusion algorithm by projecting a set of 3D Region of Interests (RoIs) from the point clouds to the 2D RoIs of the corresponding the images. Finally, we demonstrate that our deep fusion approach achieves state-of-the-art performance on the KITTI 3D object detection challenging benchmark.

**INDEX TERMS** Sensors Fusion, 3D object detection, Region of Interests, Neural Network, Segmentation Network, Point Cloud, Image

## I. INTRODUCTION

Object detection of 3D bounding boxes is one of the fundamental challenges of situational awareness and 3D environmental perception for autonomous systems (e.g. autonomous vehicles, robots, unmanned aerial vehicles, etc.). In fact, autonomous systems need to perceive objects in their surrounding environment using different sensors (e.g. cameras, LIDAR) for navigation and obstacle avoidance. In the past few years, 2D object detection for computer vision [7], [9], [22]–[24], [31], [32] has made significant progresses, especially with the advent of Convolutional Neural Network (CNN) technology [15]. However, 3D object detection remains an open challenge, especially when multiple sensors are used to obtain a more reliable and robust information.

Recently, many researchers focused on the exploitation of LIDAR-only methods for 3D object detection due to the advantages that the point clouds provide precise depth information and dense geometric shape feature [2], [3], [25], [29], [30], [43]. *PointRCNN* [34] builds a two-stage architecture to directly process dense 3D point clouds and estimate 3D bounding boxes from all the foreground points. *VoxelNet* [53], *SECOND* [47] convert the point clouds to voxels before applying standard CNNs to achieve the same result. *Pixor* [48], *Complex-YOLO* [36], and *Birdnet* [1] operate deep CNNs on

Bird-Eye-View (BEV) maps for 3D object classification and bounding box regression.

However, a standard point cloud is incapable of offering texture information and high resolution of an object, which are actually beneficial to capture discriminative features. In contrast, the images provide rich color and texture information, but with a lack of depth and scale information without the application of complex and computationally intensive algorithms (i.e. stereography). For example, small objects (e.g. pedestrians) detected at long-distance generates only few points in the point cloud, which makes the classification or localization of these objects very difficult with only a LIDAR. Meanwhile, in the image domain, texture and color features of small objects can be still visible, due to the higher spatial resolution of images, and likely to be captured by existing mature 2D CNNs technology. As a result, the fused features, leveraging the advantages from both point clouds and images, are beneficial in exploiting more reliable representations and improving the performance of the 3D object detection architecture.

However, it is still challenging to develop an efficient and effective sensor fusion method due to the viewpoint misalignment caused by the properties of the point clouds and the images. In order to address this issue, early method *MV3D*

[5] and *AVOD* [16] perform 3D bounding box regression on fused 2D images and 2D Bird-Eye-View (BEV) feature maps, although quantization for BEV generation gives rise to a lot of geometric information losses. *Frustum Pointnets* [28] and *Pointfusion* [46] project 2D bounding boxes from the image-based 2D detector onto the point clouds to coarsely cluster potential foreground points. At this point, the PointNets are applied for 3D boxes estimation, but this procedure heavily relies on the performance of the 2D detectors. *PointPainting* [40] feeds the pixel-wise semantic features captured from image-based semantic segmentation model onto corresponding point-wise semantic features in the point cloud to boost the performance of the 3D object detection.

It can be observed that the main disadvantage of dense point-pixel fusion method [40] is that they are leading to a considerable amount of redundant computations. Meanwhile, using a BEV-image fusion method allows the deep learning-based fusion of the feature maps captured from an individual viewpoint but with geometric information losses. However, it is the authors assumption that it is not strictly necessary to densely fuse the whole point clouds with images. Conversely, it is feasible to generate a small set of potential Region of Interests (RoIs), followed by the application of a deep fusion method only on those local regions used for 3D object detection. The advantages of this fusion method are that it considerably reduces the computation cost and allows an easy alignment of the viewpoints on the local regions.

Motivated by these observations, we hereby present an efficient and lightweight deep fusion method for 3D object detection for point clouds and images. Our main contributions can be summarized as follows:

- We propose a lightweight deep fusion neural network, named *RoIFusion*, aiming at efficiently and effectively fusing the point clouds and the images for 3D object detection.
- We propose a keypoints generation layer for the estimation of the keypoints on the objects guided by the fusion of the point clouds and the images, followed by a voting layer used to generate the center points of the objects.
- We propose a *RoIFusion* layer to aggregate the 3D RoIs generated from the center points with the corresponding 2D RoIs which are obtained by projecting the 3D RoIs to the images.
- We evaluated our model on the KITTI dataset [10] and achieved state-of-the-art results compared with respect to other outstanding methods.

## II. RELATED WORK

### A. LIDAR-ONLY METHODS FOR 3D OBJECT DETECTION

Many existing methods explored the possibility of detecting the objects with 3D bounding boxes only using point clouds, as they provide accurate geometric information. It is possible to broadly classify these methods into four subcategories: projection-based methods, volumetric-based methods, pointnets-based methods, and point-voxel methods.

#### a: Projection-based methods

Several works [1], [37], [48] apply 2D CNNs directly to Bird-Eye-View (BEV) projected from the raw point clouds in order to estimate the 3D bounding box and orientation of an object. *FVNet* [52] projects the raw point clouds to the front view, which is then fed to a proposal generation network and a refinement network to estimate the parameters of the 3D bounding box (i.e. object location, size, and orientation). This method allows for building a lightweight neural network for real-time applications. However, it ignores the size and the location of the objects and suffers from lots of geometric information losses during quantization. As a result, it is unlikely to exploit sufficient discriminative features for 3D object detection.

#### b: Volumetric-based methods

Volumetric-based methods convert the raw point clouds to standard 3D grids and represents the point clouds as voxels. For instance, *VoxelNet* [53] learns discriminative voxel-wise features for 3D region proposal generation and then proceeds to solve the 3D bounding box regression problem. However, many empty voxels are generated during the voxelization process, which leads to large computational cost because of the processing of those empty cells. In order to address this problem, *SECOND* [47] improved *VoxelNet* [53] by proposing an efficient method, named *sparse convolution* [12], to ignore the empty voxels. Finally, *PointPillars* [17] converts the raw point clouds to a set of stacked pillars and then encodes the same to 2D pseudo-images, which can be used as input for 2D CNNs for 3D bounding box regression.

#### c: Pointnets-based methods

*PointNets* [29], [30] models are efficient in the exploitation of point cloud features. *PointRCNN* [34] sets an example in the classification and regression of 3D bounding box directly from dense point clouds. Specifically, it firstly applies *PointNet++* [30] to extract dense semantic features for all the points, and then generates 3D region proposals for all the foreground points. Successively, the second stage is applied to refine predictions. However, the dense processing leads to quite heavy computational costs.

#### d: Point-voxel methods

In order to achieve high detection performance but to also reduce the computational costs, several works [6], [26], [33], [35] introduced two-stages neural networks for 3D object detection. In the first stage, they coarsely localise the objects and estimate the parameters of the bounding box from the voxel grids generated from the raw point clouds. In the second stage, they introduce a refinement module that leverages the *PointNets* to refine the 3D bounding box. The methods leverage the sparsity of the voxel grids and the ability of *PointNets* to carry out feature extraction and to gradually detect the 3D objects starting from coarse to more refined representations.

Finally, a recent one-stage method, *3DSSD* [49] abandons the refinement stage and builds a one-stage anchor-free neural network to directly regress 3D bounding box from the estimated candidate 3D RoIs.

### B. IMAGE-ONLY METHODS FOR 3D OBJECT DETECTION

In the past few years, 2D object detection has made great progress. However, estimating 3D bounding box directly from 2D images is still quite difficult due to the lack of depth information in single camera images. *Mousavian et al.* [27] estimates the pose and 3D bounding box by learning the geometric constraints from the 2D bounding box. *Wang et al.* applies LIDAR-only 3D detectors on the Pseudo-LiDAR representations converted from the estimated image-based depth maps. *Stereo R-CNN* [20] applies *Faster R-CNN* [32] on both the left and right images and predicts 3D bounding boxes by learning the projection relations between the associated 2D left-right bounding boxes and 3D bounding box corners.

### C. MULTI-SENSOR FUSION METHODS FOR 3D OBJECT DETECTION

In order to get the best of both worlds, there are several works attempting to fuse point clouds and 2D images with various strategies. Early works such as *MV3D* [5] and *AVOD* [16] firstly used off-the-shelf 2D feature extractors to capture the feature maps from the images and the multi-view representations of the point clouds (e.g. Bird Eye View and Front View), which are then typically fused together by a sum or a concatenation operation. A Region Proposal Network (RPN) is then applied to the fused feature maps to generate 3D bounding box proposals, followed by a refinement network for final 3D bounding box prediction. The advantages of this method are that mature 2D object detector and 2D feature extractor technologies are available to be applied on the multi-view representations of the point clouds. Furthermore, the features from different sensors can interact over the stacked layers, as these features are normally obtained from similar or even the same neural networks. *Liang et al.* [21] utilizes the continuous convolution method to fuse the feature maps of the images and BEVs. Specifically, this approach proposes a continuous fusion layer that aggregates each pixel feature in the image feature maps with the features of the neighbouring points in the BEV feature maps to learn a fused local region, which allows to extract sufficient discriminative features for 3D object detection.

In order to narrow the searching space, *Frustum pointnets* [28] and *Frustum convnet* [44] introduced the concept of 3D bounding frustums. The 2D bounding boxes are obtained from mature 2D detectors, and then the 3D frustums are used to trim the point cloud data. Finally, *Pointnets* methods are applied to the trimmed point clouds for carrying out the 3D bounding box regression task. Similarly, *Pointfusion* [46] aggregates the global features of the image obtained from an off-the-shelf 2D feature extractor with the dense semantic features of the point cloud, which are captured from *Pointnet* [29].

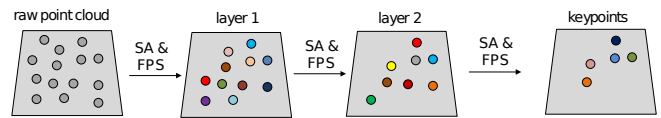


FIGURE 1. Illustration for point-guided keypoints generation

Finally, *PointPainting* [40] densely aggregates the output of the image segmentation neural network with the point clouds before applying LIDAR-only 3D detectors to boost the performance of the 3D object detection task.

### III. ROIFUSION ARCHITECTURE

In this section, we introduce and describe our RoIFusion neural network for 3D object detection as shown in Fig. 2, which uses both raw point clouds and 2D images as input. Our goal is to leverage the fusion information captured from both sensor modalities to classify and localize the objects within the oriented 3D bounding boxes. In particular, we firstly propose a fused keypoints generation layer (FKG layer) to estimate a set of 3D keypoints from the point clouds, followed by a RoIs fusion layer to fuse the 3D RoI features in the point clouds with the 2D RoI features in the images by the 3D/2D RoI pooling operation respectively. At last, a prediction layer is proposed to predict the parameters of the oriented 3D bounding box.

#### A. FUSED KEYPOINTS GENERATION (FKG) LAYER

Instead of densely generating 3D region proposals relying on all the foreground points, we only estimate a small set of 3D keypoints on the objects to generate the RoIs for deep fusion. As illustrated in part (a) of Fig. 2, our FKG layer takes the raw point clouds and the RGB images as input, and combines point-guided keypoints and pixel-guided keypoints that are generated by individually performing the LIDAR-only point cloud segmentation model and image-only segmentation network respectively. As a result, we obtain a set of keypoints on the objects leveraging both the point cloud and the image information.

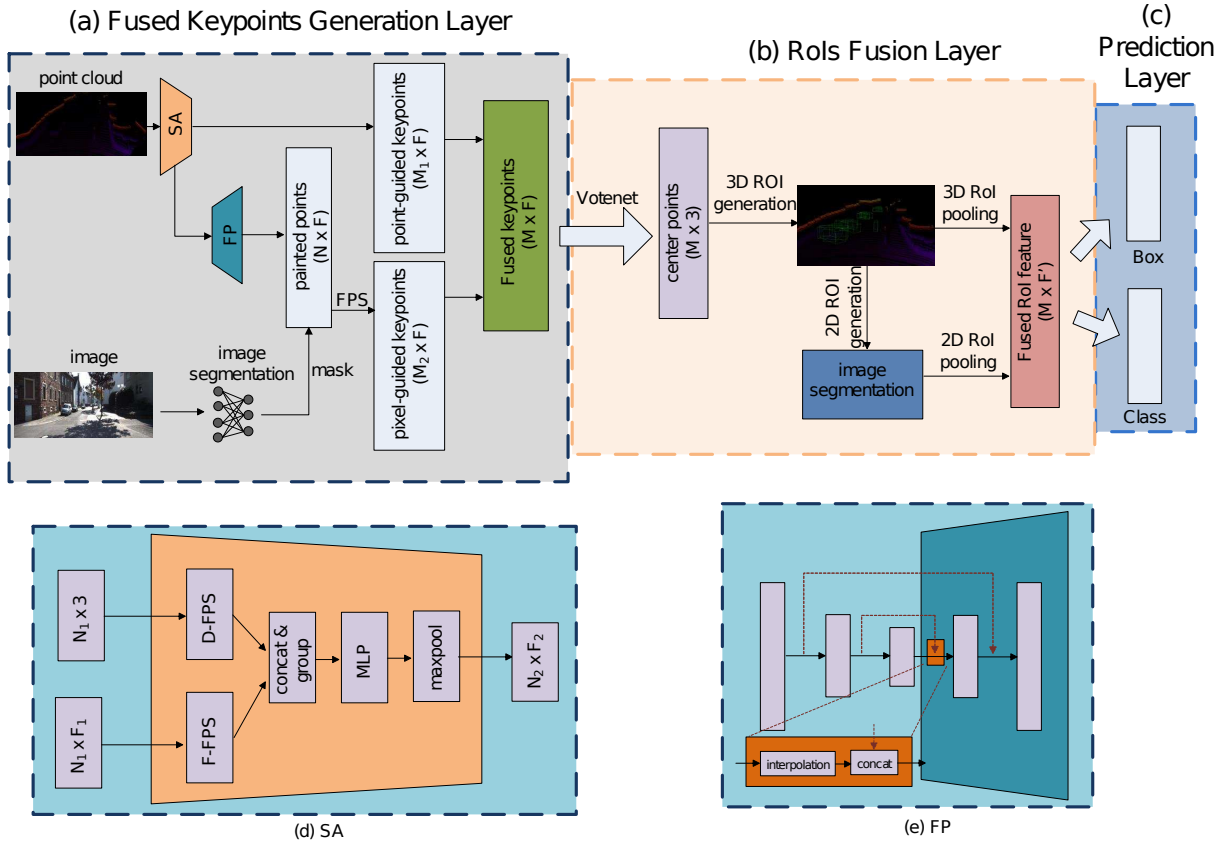
We can define the point cloud as shown in Eq. (1), where  $\mathbf{x}_i$  denotes the  $i$ -th point with the 3D space coordinates  $[x_i, y_i, z_i]$  and the measured reflectance  $r_i$ . As a result, the dimension of the point cloud data set is  $N \times 4$ .

$$\mathbf{X} = \{\mathbf{x}_i = [x_i, y_i, z_i, r_i] \in \mathbb{R}^4, i = 1, 2, \dots, N\} \quad (1)$$

a: Point-guided keypoints generation

Our point-guided keypoints  $\hat{\mathbf{X}}^{(\text{pc})} \in \mathbb{R}^{M_1, F}$ , where  $M_1, F$  are the number of keypoints and corresponding features respectively, are extracted by the set abstraction (SA) layers backbone, as illustrated in part (d) of Fig. 2. This is done as proposed by *PointNet++* [30], where a simultaneously down-sampling of the points and the extraction of the corresponding deep features are carried out.

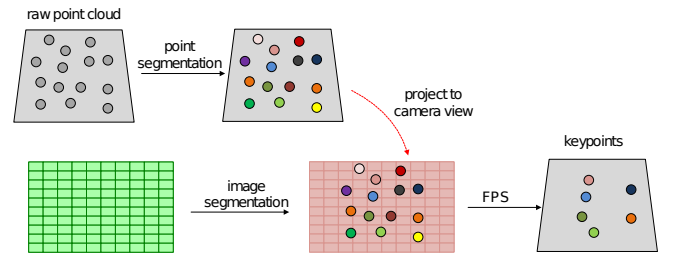
Specifically, as shown in Fig. 1, we apply a number of SA layers with downsampling operation to the raw point cloud.



**FIGURE 2. RoIFusion framework.** As a whole, the proposed architecture contains three parts: (a) Fused Keyoints Generation (FKG) layer takes the raw point clouds and the images as input and aggregates the keypoints generated from both the point cloud set abstract (SA) network and the image segmentation network respectively. (b) The keypoints are then used to estimate the center points of the potential objects and the 3D Region of Interests (Rols), which are projected to the image to obtain the 2D Rols. The 3D/2D Rol pooling layers are employed to capture the respective local features, which are finally fused together for 3D bounding box prediction as shown in part (c). As shown in (d), the input of the SA module are the 3D points  $N_1 \times 3$  and the corresponding features  $N_1 \times F_1$ . We simultaneously downsample the points and extract the corresponding deep features in the orange block, and obtain the downsampled 3D points  $N_1 \times 3$  and corresponding features  $N_2 \times F_2$ .  $N_1, F_1, N_2, F_2$  are the dimension number of the input points/features, output points/features respectively. Part (e) illustrates the Feature Propagation (FP) module: the upsampling method that takes the output of the SA module as input and generates the segmented features for all the points in the point cloud.

At each layer, a set of points are processed and a new set with higher-level but fewer points is generated. Finally, we obtain a small number of points that are treated as keypoints.

With respect to the downsampling strategy, we use an iterative farthest point sampling (FPS) method to select the points for the subset. Let us suppose an empty subset  $X_1$ , a random point is firstly picked and added to  $X_1$ , then the point having the farthest 3D geometric Euclidean distance is iteratively added to  $X_1$  until the expected  $M$  points are picked. The FPS strategy, named *D-FPS*, has a better coverage of the whole point set than random sampling. In order to preserve sufficient foreground points and filter out the background, inspired by *3DSSD* [49], we decided to also employ a specific FPS strategy, named *F-FPS*, which calculates the Euclidean distances of the semantic features for the points selection. The *F-FPS* method is beneficial to preserving foreground points (e.g. points on the objects) and removing the useless background, such as points on the ground. Finally, we follow [49] and combined both FPS strategies together to efficiently capture sufficient foreground points as the keypoints.



**FIGURE 3.** Illustration for pixel-guided keypoints generation

**b: Pixel-guided keypoints generation**

Considering the fact that the colour and texture representations are useful to localize objects within point clouds, especially for small objects that are difficult to be detected by LIDAR-only detectors, we capture the segmentation features and corresponding scores using an image segmentation network, which then is used to guide the keypoints selection as shown in Fig. 3.

The detailed procedure to extract the pixel-guided keypoints is shown in Alg. 1. Firstly, we generate the point clouds



segmentation features  $\mathbf{X}_s$  using a Feature Propagation (FP) layer as shown in part (e) of Fig. 2. In particular, we leverage the output of the SA layer as input, and upsample the points by interpolating the point features using the inverse squared Euclidean distance weighted average function as shown in Eq. (2). Furthermore, we concatenate the interpolated point features with the skip linked point features from the corresponding SA layer. As a result, our FP layer outputs the 3D geometric points and corresponding semantic features with the same number of points as the raw point cloud.

$$f(\mathbf{x}) = \frac{\sum_{i=1}^k \omega_i(\mathbf{x}) f_i}{\sum_{i=1}^k \omega_i(\mathbf{x})} \quad (2)$$

where  $\omega_i(\mathbf{x}) = \frac{1}{(\mathbf{x}-\mathbf{x}_i)^2}$  is the inverse squared Euclidean distance between a certain point  $\mathbf{x}$  and corresponding  $i$ -th neighbouring point  $\mathbf{x}_i$  of the  $k = 3$  nearest neighbours.

For what concerns the image processing, it is a common choice in literature to use mature 2D feature extractors to capture the feature maps from the RGB images. However, these feature maps are unlikely to localize the objects in the images. As a result, we use a lightweight image segmentation neural network *DeepLabv3* [4] to efficiently capture pixel-wise segmentation features  $\mathbf{I}_s$  and segmentation scores  $\mathbf{S}$ , which allows to ignore the background and conduct the keypoints selection. It is worth pointing out that our RoI-Fusion model is agnostic to the development of the image segmentation models.

After that, we project the point cloud to the image viewpoint to paint the point cloud with the segmentation features and the segmentation scores from the relevant pixels, then we mask the painted point cloud  $\mathbf{X}_{\text{img}}$  with the foreground image segmentation scores and map to the point clouds segmentation features to generate all the foreground segmentation features  $\mathbf{X}_s^{(\text{obj})}$  for the point cloud. At last, we use the F-FPS as our downsampling strategy to further select a small set of point segmentation features  $\hat{\mathbf{X}}^{(\text{img})}$  with the dimension size of  $M_2 \times F$ , where  $M_2$  and  $F$  are the number of keypoints and corresponding features respectively.

c: Keypoints fusion

We finally obtain the fused keypoints  $\hat{\mathbf{X}} \in \mathbb{R}^{M,F}$  by aggregating the point-guided keypoints with the pixel-guided keypoints along with the channel of the number of the keypoints, where  $M$  is the number of fused keypoints. We note that the points on small objects are likely to be selected due to the fact that a part of the points are captured based on the image segmentation scores.

## B. ROIS FUSION LAYER

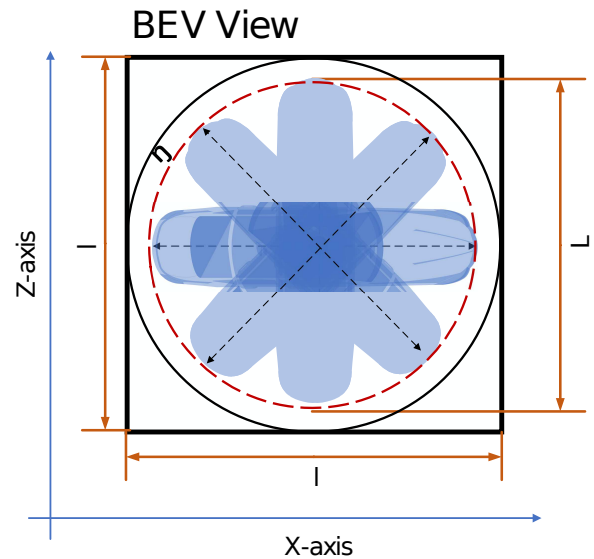
After the implementation of the FKG layer intertwined with the point cloud segmentation network and the image segmentation network, we obtain a set of keypoints scattered over the objects, which are also used to predict the center of the objects before we generate the RoIs. Considering that these keypoints are on the objects, inspired by [8], we use

### Algorithm 1. Pixel-guided keypoints generation.

**Input:** The point clouds  $\mathbf{X} \in \mathbb{R}^{N,4}$ .  
The images  $\mathbf{I} \in \mathbb{R}^{W,H,3}$ .  
Homogenous transformation matrix  $\mathbf{T} \in \mathbb{R}^{4,4}$ .  
Camera projection matrix  $\mathbf{M} \in \mathbb{R}^{3,4}$ .

**Output:** Pixel-guided keypoints features  $\hat{\mathbf{X}}^{(\text{img})} \in \mathbb{R}^{M_2,F}$ .

- 1: Apply for point cloud segmentation network to obtain segmentation features  $\mathbf{X}_s \in \mathbb{R}^{F_p}$ .
- 2: Apply for image segmentation network to obtain segmentation features  $\mathbf{I}_s \in \mathbb{R}^{W,H,F_i}$  and segmentation scores  $\mathbf{S} \in \mathbb{R}^{W,H,C}$ .
- 3:  $\mathbf{X}_{\text{img}} = \text{Projection}(\mathbf{M}, \mathbf{T}, \mathbf{X})$ .
- 4:  $\mathbf{X}_{\text{obj}}, \text{index} = \text{Mask}(\mathbf{X}_{\text{img}}, \mathbf{S})$ .
- 5:  $\mathbf{X}_s^{(\text{obj})} = \text{Mapping}(\mathbf{X}_s, \text{index})$ .
- 6:  $\hat{\mathbf{X}}^{(\text{img})} = \text{FPS}(\mathbf{X}_s^{(\text{obj})})$ .
- 7: **return**  $\hat{\mathbf{X}}^{(\text{img})}$ .

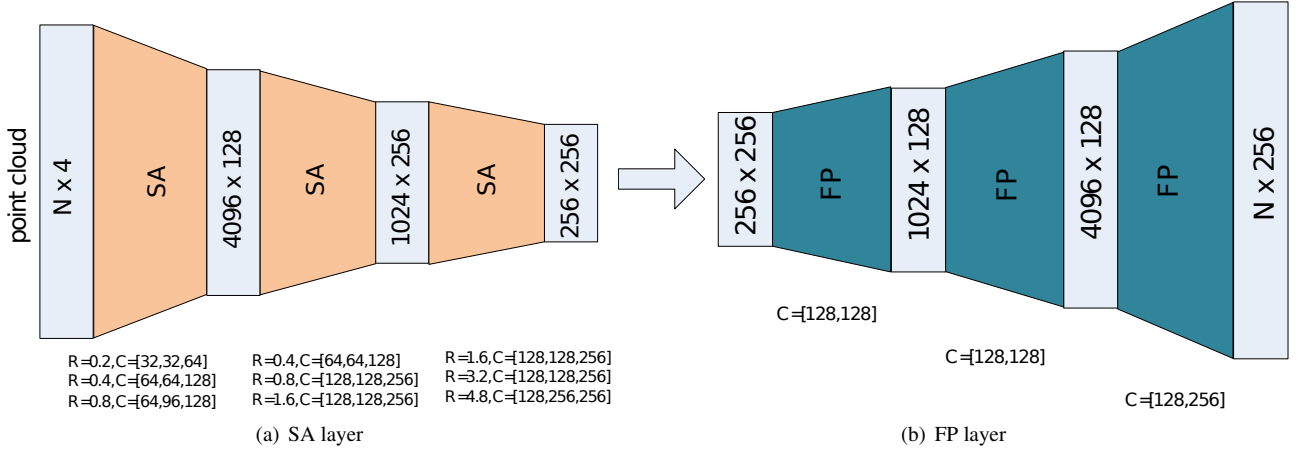


**FIGURE 4.** RoIs generation layer. The objects with all the oriented angles are illustrated in the red dot circle. The length size of the objects (e.g. cars) are defined as  $L$ .  $\eta$  is the enlarged size for the object length.  $l$  represents the size of the RoI in the Bird eye view (BEV).

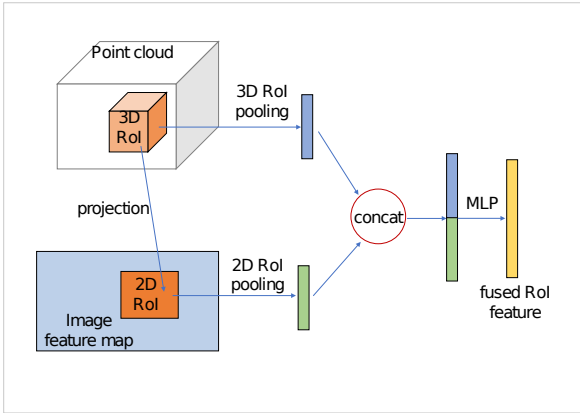
the spatial location and features of the keypoints to estimate the corresponding center of the objects. As shown in part (b) of Fig. 2, these high-level keypoints are used to generate 3D RoIs in the point cloud view and corresponding 2D RoIs in the camera view, followed by a RoI fusion operation to obtain the fused RoI features for further bounding box regression and object classification. Specifically, we build a subnetwork *VoteNet* with a single layer to learn the spatial offset between predicted center points and corresponding ground truth. We treat each center point as the centroid of the 3D bounding box of the object.

a: 3D RoIs generation and pooling

The 3D RoIs are generated for the center points we previously obtained. Successively, we apply a 3D RoI pooling layer to



**FIGURE 5.** Illustration of the SA and FP layers. The SA layer 5(a) takes  $N$  points as input, followed by stacked SA layers to downsample points and extract corresponding features. After that, the FP layer 5(b) upsamples the points to the original 16384 points using stacked FP layers.  $R$  is the radius of the ball query strategy for clustering local region points, and  $C$  is the number of the filters for the MLP layers.



**FIGURE 6.** Rols fusion layer

pool the surrounding points of each center point and learn the local features for those clustered points around each center point.

We encode our RoIs using the axis-aligned 3D bounding boxes. Specifically, the centroid of each RoI  $(x^{(c)}, y^{(c)}, z^{(c)})$  is parametrized using the obtained center point. The length  $l$  and the width  $w$  of the RoI are set to the enlarged length size of the objects to cover all the orientation scenarios as shown in Fig. 4. We finally use an enlarged height of the objects as the height  $h$  for each RoI. As a result, the dimension of the RoI is defined as  $(x_i^{(c)}, y_i^{(c)}, z_i^{(c)}, \eta + h_i, \eta + w_i, \eta + l_i)$ , where  $\eta$  is the parameter for extended size of the RoI.

After that, we shift the points inside each 3D RoI to the relative locations based on the center points for better local features learning and then apply a subnetwork equipped with stacked Multi-Layer-Perceptron (MLP) layers on the cluster the points inside the 3D RoIs to extract the local RoI pooling features.

**TABLE 1.** KITTI dataset difficulty classification levels for object detection.

Levels	Min. bounding box height	Max. occlusion level	Max. truncation
Easy	40 pixels	Fully visible	15%
Moderate	25 pixels	Partly occluded	30%
Hard	25 pixels	Difficult to see	50%

b: 2D Rols generation and pooling

Our 3D RoIs are then projected to the image to generate the corresponding 2D RoIs, followed by a 2D RoI pooling layer, inspired by [32], to learn the local texture features for the 2D RoIs.

c: RoI pooling features fusion

We finally fuse the point cloud 3D RoI and the image 2D RoI by aggregating the pooling features along with feature dimension axis as shown in Fig. 6. Specifically, we define a fusion strategy by concatenation as in Eq. (3):

$$\mathbf{F}_{\text{fuse}} = \text{MLP}(\text{concat}[\mathbf{F}_{\text{roi}}^{(\text{pc})}, \mathbf{F}_{\text{roi}}^{(\text{img})}]) \quad (3)$$

where  $\mathbf{F}_{\text{fuse}}$  is the fused feature from the 3D RoI pooling features  $\mathbf{F}_{\text{roi}}^{(\text{pc})}$  and 2D RoI pooling features  $\mathbf{F}_{\text{roi}}^{(\text{img})}$ .

### C. PREDICTION LAYER

The prediction layer, inspired by [49], use an anchor-free method to directly predict the offset between the center points and corresponding ground truth of the center of the 3D bounding box for regression. Besides, we also directly regress the 3D bounding box size from the fused RoI features. For the orientation regression, we follow the method introduced in [28] that utilizes a hybrid classification and regression

**TABLE 2.** 3D car detection results on KITTI test dataset. *Sen.* indicates involved sensors by the methods. *L* and *I* denote LIDAR and images respectively.

Method	Sen.	$AP_{Car}(\%)$		
		easy	mod.	hard
VoxelNet [53]	L	77.47	65.11	57.73
SECOND [47]		83.13	73.66	66.20
PointPillars [17]		79.05	74.99	68.30
PointRCNN [34]		85.94	75.76	68.32
Fast PointRCNN [6]		84.28	75.73	67.39
Patches [19]		87.87	77.16	68.91
Part A2 [35]		85.94	77.86	72.00
STD [53]		86.61	77.63	76.06
3DSSD [51]		88.36	79.57	74.55
MV3D [5]	L + I	71.09	62.35	55.12
AVOD [16]		81.94	71.88	66.38
F-PointNet [28]		81.20	70.39	62.19
F-ConvNet [44]		85.88	76.51	68.08
IPOD [50]		79.75	72.57	66.33
Painted PointRCNN [40]		82.11	71.70	67.08
Ours		88.32	79.54	74.47

formulations to estimate the orientation angle of the 3D bounding box. In particular, we pre-define  $H$  equally split angle bins and use the output of the RoI fusion layer to classify the angle bins, and then regress residual with respect to the classified bin.

#### IV. MODEL STRUCTURE

The model structure is presented in Fig. 2. In our experiments we randomly choose  $N = 16384$  points from the raw point cloud. We then apply the SA layer Fig. 5(a), the FP layer Fig. 5(b), and the image segmentation network *DeepLabv3* to capture  $M = 256$  keypoints. The hyper-parameters of the SA layer and the FP layer are represented in Fig. 5(a) and Fig. 5(b) respectively. Successively, we employ a single layer *Votenet* with filters (128) to estimate the center points for the 3D bounding box of the objects. The dimensions of the 3D RoIs are set to  $[h = 1.8m, w = 5.0m, l = 5.0m]$ ,  $[h = 1.8m, w = 1.0m, l = 1.0m]$ ,  $[h = 1.8m, w = 1.8m, l = 1.8m]$  for the car, pedestrians, and cyclists objects respectively. We set the constant extended value  $\eta = 1.0m$ .

#### V. EXPERIMENTS

In this section, we evaluate our deep fusion method on the widely used KITTI 3D object detection benchmark [10], [11]. We firstly introduce the KITTI dataset and explain the detailed training settings. Then, we demonstrate our results by comparison with recent state-of-the-art 3D detectors. We only test our model on the car category due to the large amount of data after preprocessing. However, we evaluate all the categories when we compare our model to the backbone model *3DSSD* [51] to present the effectiveness of our fusion method. Finally, we analyse the efficiency of our fusion

method and visualize some representative results for our 3D object detection model.

#### A. DATASET

The KITTI dataset [10] contains both 2D images and 3D point clouds with the corresponding annotations for the cars, pedestrians, and cyclists categories in an urban driving scenario. The sensors used for data collection are: 2 grayscale cameras, 2 color cameras, and 1 Velodyne HDL-64E LIDAR. We only used the point clouds data and images from the left color camera to train our fusion model. The dataset provides 7481 samples for training and 7518 samples for testing. As standard good practice, we further split the KITTI training dataset into 3712 samples for training and 3769 samples for validation. We evaluated our model on the validation dataset following the easy, moderate, and hard difficulty classification levels officially introduced by KITTI. Specifically, in order to align the performance of the algorithms and cover most of the traffic scene scenarios, the object detection task is divided into three levels for validation and testing with respect to the different size, occlusion, and truncation level as shown in Table 1. Besides, the average precision (AP) metric is used when we compare our results with other different models.

#### B. TRAINING SETTINGS

We used the Adam [14] algorithm as our training optimizer. The batch size was set to 4 on a NVIDIA 1080Ti GPU. The learning rate was initially set to 0.002, and then was divided by 10 at 40 epochs. Our model has been trained for a total of 50 epochs.

#### C. RESULTS

We can firstly compare our model to the backbone network *3DSSD* [51] on the validation dataset to show the effectiveness of our fusion strategy as shown in Table 3. The bottom line indicates the difference between our model and *3DSSD* for 3D car detection. It shows that our model outperforms *3DSSD* in all the categories and all the difficulty levels, which convincingly shows the efficiency of our RoI fusion method.

As shown in Table 2, our model also achieves the best performance compared to recent state-of-the-art fusion methods on the test dataset. We choose *moderate* difficulty as the main average precision (AP) metric, and compare our model to BEV-image fusion *MV3D* [5] and *AVOD* [16], our deep fusion method outperforms all others by a large margin. For the frustum method, our method outperforms *F-PointNet* [28] by 9.12%. Besides, our model significantly outperforms point-pixel-wise fusion method *PointRCNN* [34] by 7.84%. We also visualize some examples for prediction results and corresponding ground truth as shown in Fig. 9 for better representation.

#### D. ABLATION STUDY

We carried out several ablation experiments to investigate the effectiveness of extended value for RoI size and different

**TABLE 3.** 3D Car detection average precision (AP) on KITTI validation dataset compared to 3DSSD model. The *Delta* indicates the difference between our model and 3DSSD model that is our backbone for the point cloud processing. *repro* represents that the results are reproduced on our own computer.

Method	$AP_{car}\%$			$AP_{ped}\%$			$AP_{cyc}\%$		
	easy	mod.	hard	easy	mod.	hard	easy	mod.	hard
3DSSD (repro)	89.71	79.45	78.67	41.72	39.63	36.86	78.01	62.32	57.01
Ours	91.36	82.74	80.22	43.40	41.44	37.72	80.84	64.05	58.37
Delta	+1.65	+3.29	+1.55	+1.68	+1.81	+0.86	+2.83	+1.73	+1.36

**TABLE 4.** Effectiveness of the extended value  $\eta$  for RoI size on KITTI car validation dataset.

$\eta$ (m)	$AP_{easy}(\%)$	$AP_{mod.}(\%)$	$AP_{hard}(\%)$
0.0	89.33	80.61	70.47
0.5	91.12	82.53	79.85
1.0	<b>91.36</b>	<b>82.74</b>	<b>80.22</b>
1.5	90.42	82.48	79.11
2.0	89.01	81.09	78.96

**TABLE 5.** Effectiveness of the fusion strategy on KITTI car validation dataset.

fusion strategy	$AP_{easy}(\%)$	$AP_{mod.}(\%)$	$AP_{hard}(\%)$
sum	88.29	76.72	68.48
concat	<b>91.36</b>	<b>82.74</b>	<b>80.22</b>
max	87.01	74.93	67.27

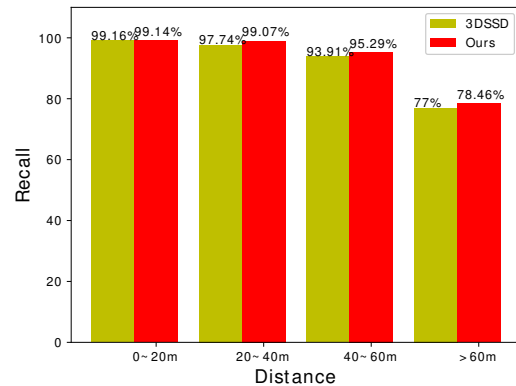
fusion methods. All the experiments are performed on the KITTI validation dataset for the 3D car detection task.

a: Effects of the extended size of RoI

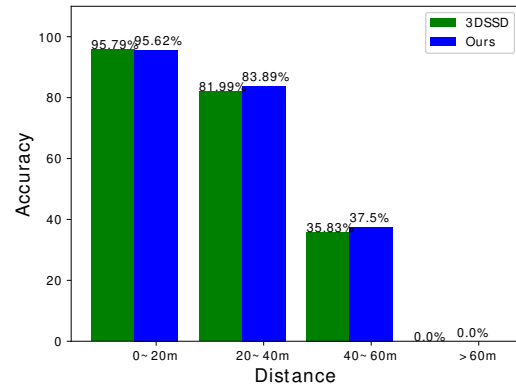
In order to cluster sufficient points around the objects, we enlarge the RoI size by an extended value *eta* for more contextual local features. Table 4 shows that the model achieves the best performance when  $\eta = 1.0$ . Besides, we notice that there is a significant drop of performance when no extended size (i.e.  $\eta = 0$ ) is used, especially for hard difficulty level of detection. It is assumed that the larger size of the box also provides sufficient information, but is likely to involve more redundant and harmful information. In contrast, smaller size only could provide part information of the cars, which is insufficient to predict the parameters of the 3D bounding box. For the hard detection level objects, they normally are occluded by other objects or far away from the sensor, which leads to very few points on the objects. As a result, involving more surrounding points is beneficial to object classification and regression.

b: Effects of the fusion strategy

We further investigate the effectiveness for the different fusion strategies. In addition to the concatenation operation as described in Section III-B0c, we also employ operations,



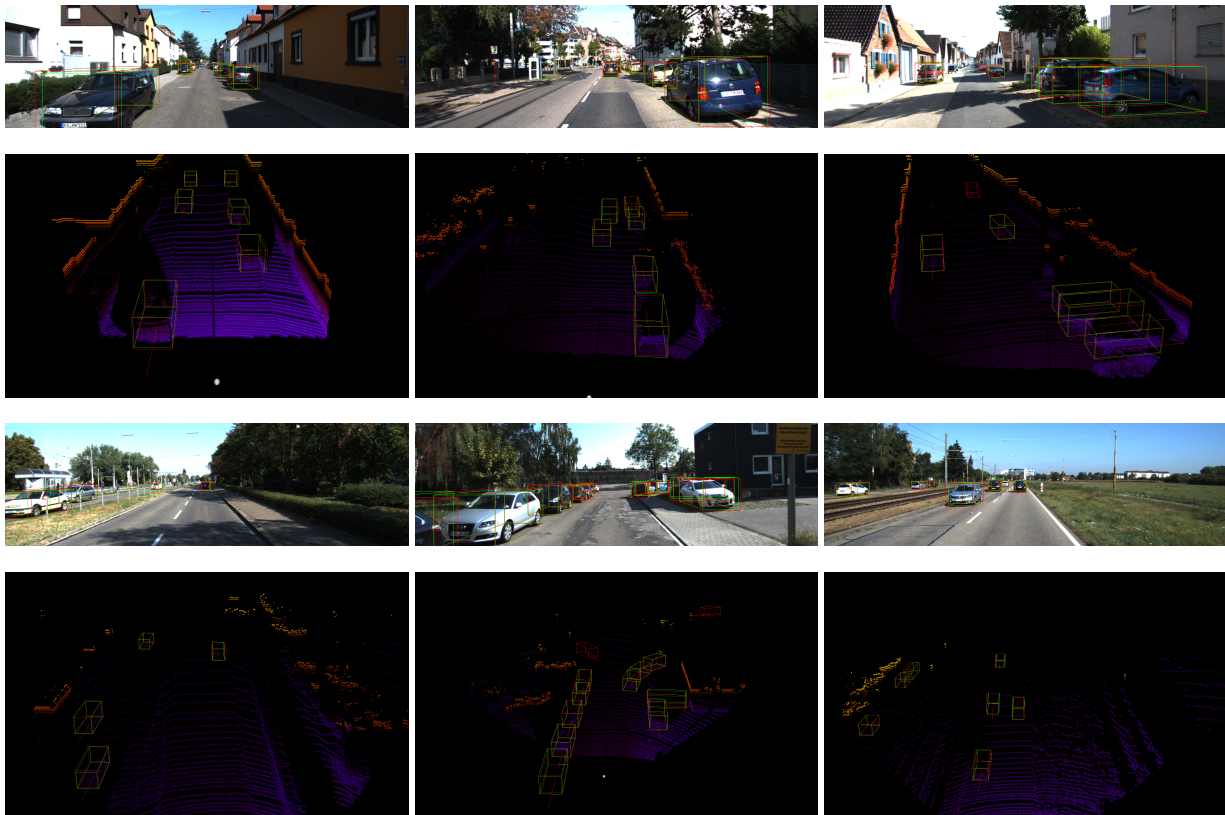
**FIGURE 7.** The comparison of the recall for the detected objects in the various distance range.



**FIGURE 8.** The comparison of the accuracy for the objects in the various distance range.

such as sum, max operation, and compare the results for different choices. As shown in Table 5, the concatenation operation for RoI features fusion achieves 91.36% 82.74% 80.22% performance for easy, moderate, and hard difficulties respectively. The results show that the concatenation operation could fuse more discriminative features from the 3D RoIs and corresponding 2D RoIs. This can be linked to the fact that both the sum operation and the max operation could obtain signature features, but the concatenation operation allows to keep all the features from different sensors, which then is likely to allow capturing more useful features for classification and regression.





**FIGURE 9.** Qualitative results on the KITTI validation dataset. The predicted objects and the ground truth objects are shown in red and green bounding boxes respectively. We also project the bounding boxes to the RGB images for better visualization.

**TABLE 6.** Inference time and accuracy on the *moderate* level for different fusion methods.

	AVOD [16]	F-PointNet [28]	PointPainting [40]	Ours
time	100 ms	170 ms	400 ms	220 ms
AP	71.48%	70.39%	71.70%	79.54%

#### c: Effects of the distance of the objects

As we know that LIDAR-only methods are not efficient to detect the objects that include few points, such as small objects and objects in the long distance. In order to better compare our fusion method to the LIDAR-only model *3DSSD* [51] which only provide car detection model, we detect the cars in the various distance range for the different size of object, and then compare the recall and the accuracy to the *3DSSD*. As illustrated in the Fig. 7 and the Fig. 8, when the objects are within 20 meters, our results are nearly no different to the *3DSSD*. However, the performance becomes worse when the objects are beyond 20 meters for both models, but our fusion model performs much better than *3DSSD* when detect the objects located in the long distance. The results convincingly prove that our fusion method successfully predicts more true positive objects in the long distance range by learning sufficient colour and texture information for the point clouds.

#### d: Inference time

We tested the inference time on KITTI validation dataset with a NVIDIA 1080Ti GPU, and then compare to existing fusion methods in Table 6. Our model achieves the best trade-off compared to BEV-image fusion method *AVOD* [16] and frustum method *F-PointNet* [28]. We also note that our RoI fusion method is much better than point-pixel fusion method *PointPainting* [40] in terms of both accuracy and inference time.

### E. QUANTITATIVE ANALYSIS

The comparison of our model with other multi-sensor methods as shown in Table 2, shows that our model is more efficient and effective. RoI fusion enables to globally fuse the local area from the point clouds and the images, which make it easier to align the viewpoint when we concatenate the 3D/2D RoI pooling features. In contrast, the point-pixel fusion is unlikely to obtain discriminative features due to the fact that the point clouds are sparse and irregular, but the images distribute as the standard grid. Compared to BEV-image fusion, our model outperforms others by a large margin due to the information losses during BEV generation.

### VI. CONCLUSION

In this paper, we propose a novel deep fusion method, named RoIFusion, to efficiently fuse the point clouds and the images

for 3D object detection. We build a lightweight neural network to generate 3D RoIs from the point clouds and 2D RoIs from the images, and then employ a 3D RoI pooling layer and a 2D RoI pooling layer to obtain the geometric features and the texture features respectively for a potential local area in both point clouds and the images. Finally, we fuse them together to predict the oriented 3D bounding box for the detected object. Our fusion method is flexible and could combine any other LIDAR-only segmentation networks and image segmentation networks. The state-of-the-art performance of our model convincingly show that the fusion method proposed can successfully boost the performance of 3D object detection.

## REFERENCES

- [1] Jorge Beltrán, Carlos Guindel, Francisco Miguel Moreno, Daniel Cruzado, Fernando Garcia, and Arturo De La Escalera. Birdnet: a 3d object detection framework from lidar information. In 2018 21st International Conference on Intelligent Transportation Systems (ITSC), pages 3517–3523. IEEE, 2018.
- [2] Can Chen, Luca Zanotti Fragonara, and Antonios Tsourdos. Gapnet: Graph attention based point neural network for exploiting local feature of point cloud. arXiv preprint arXiv:1905.08705, 2019.
- [3] Can Chen, Luca Zanotti Fragonara, and Antonios Tsourdos. Go wider: An efficient neural network for point cloud analysis via group convolutions. Applied Sciences, 10(7):2391, 2020.
- [4] Liang-Chieh Chen, George Papandreou, Florian Schroff, and Hartwig Adam. Rethinking atrous convolution for semantic image segmentation. arXiv preprint arXiv:1706.05587, 2017.
- [5] Xiaozhi Chen, Huimin Ma, Ji Wan, Bo Li, and Tian Xia. Multi-view 3d object detection network for autonomous driving. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pages 1907–1915, 2017.
- [6] Yilun Chen, Shu Liu, Xiaoyong Shen, and Jiaya Jia. Fast point r-cnn. In Proceedings of the IEEE International Conference on Computer Vision, pages 9775–9784, 2019.
- [7] Jifeng Dai, Yi Li, Kaiming He, and Jian Sun. R-fcn: Object detection via region-based fully convolutional networks. In Advances in neural information processing systems, pages 379–387, 2016.
- [8] Zhipeng Ding, Xu Han, and Marc Niethammer. Votenet: A deep learning label fusion method for multi-atlas segmentation. In International Conference on Medical Image Computing and Computer-Assisted Intervention, pages 202–210. Springer, 2019.
- [9] Kaiwen Duan, Song Bai, Lingxi Xie, Honggang Qi, Qingming Huang, and Qi Tian. Centernet: Keypoint triplets for object detection. In Proceedings of the IEEE International Conference on Computer Vision, pages 6569–6578, 2019.
- [10] Andreas Geiger, Philip Lenz, Christoph Stiller, and Raquel Urtasun. Vision meets robotics: The kitti dataset. The International Journal of Robotics Research, 32(11):1231–1237, 2013.
- [11] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? the kitti vision benchmark suite. In 2012 IEEE Conference on Computer Vision and Pattern Recognition, pages 3354–3361. IEEE, 2012.
- [12] Benjamin Graham and Laurens van der Maaten. Submanifold sparse convolutional networks. arXiv preprint arXiv:1706.01307, 2017.
- [13] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In Proceedings of the IEEE conference on computer vision and pattern recognition, pages 770–778, 2016.
- [14] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980, 2014.
- [15] Alex Krizhevsky, I Sutskever, and G Hinton. Imagenet classification with deep convolutional neural. In Neural Information Processing Systems, pages 1–9, 2014.
- [16] 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.
- [17] Alex H Lang, Sourabh Vora, Holger Caesar, Luning Zhou, Jiong Yang, and Oscar Beijbom. Pointpillars: Fast encoders for object detection from point clouds. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pages 12697–12705, 2019.
- [18] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. nature, 521(7553):436–444, 2015.
- [19] Johannes Lehner, Andreas Mitterecker, Thomas Adler, Markus Hofmarcher, Bernhard Nessler, and Sepp Hochreiter. Patch refinement-localized 3d object detection. arXiv preprint arXiv:1910.04093, 2019.
- [20] Peiliang Li, Xiaozhi Chen, and Shaojie Shen. Stereo r-cnn based 3d object detection for autonomous driving. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pages 7644–7652, 2019.
- [21] Ming Liang, Bin Yang, Shenlong Wang, and Raquel Urtasun. Deep continuous fusion for multi-sensor 3d object detection. In Proceedings of the European Conference on Computer Vision (ECCV), pages 641–656, 2018.
- [22] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature pyramid networks for object detection. In Proceedings of the IEEE conference on computer vision and pattern recognition, pages 2117–2125, 2017.
- [23] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In Proceedings of the IEEE international conference on computer vision, pages 2980–2988, 2017.
- [24] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and Alexander C Berg. Ssd: Single shot multibox detector. In European conference on computer vision, pages 21–37. Springer, 2016.
- [25] Yongcheng Liu, Bin Fan, Shiming Xiang, and Chunhong Pan. Relation-shape convolutional neural network for point cloud analysis. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pages 8895–8904, 2019.
- [26] Zhe Liu, Xin Zhao, Tengting Huang, Ruolan Hu, Yu Zhou, and Xiang Bai. Tanet: Robust 3d object detection from point clouds with triple attention. arXiv preprint arXiv:1912.05163, 2019.
- [27] Arsalan Mousavian, Dragomir Anguelov, John Flynn, and Jana Kosecka. 3d bounding box estimation using deep learning and geometry. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pages 7074–7082, 2017.
- [28] Charles R Qi, Wei Liu, Chenxia Wu, Hao Su, and Leonidas J Guibas. Frustum pointnets for 3d object detection from rgb-d data. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pages 918–927, 2018.
- [29] 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.
- [30] Charles Ruizhongtai Qi, Li Yi, Hao Su, and Leonidas J Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. In Advances in neural information processing systems, pages 5099–5108, 2017.
- [31] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You only look once: Unified, real-time object detection. In Proceedings of the IEEE conference on computer vision and pattern recognition, pages 779–788, 2016.
- [32] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In Advances in neural information processing systems, pages 91–99, 2015.
- [33] 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. arXiv preprint arXiv:1912.13192, 2019.
- [34] Shaoshuai Shi, Xiaogang Wang, and Hongsheng Li. Pointtrcn: 3d object proposal generation and detection from point cloud. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pages 770–779, 2019.
- [35] 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.
- [36] Martin Simon, Stefan Milz, Karl Amende, and Horst-Michael Gross. Complex-yolo: An euler-region-proposal for real-time 3d object detection on point clouds. In European Conference on Computer Vision, pages 197–209. Springer, 2018.
- [37] 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), pages 0–0, 2018.

- [38] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. arXiv preprint arXiv:1409.1556, 2014.
- [39] Zhi Tian, Chunhua Shen, Hao Chen, and Tong He. Fcos: Fully convolutional one-stage object detection. In Proceedings of the IEEE International Conference on Computer Vision, pages 9627–9636, 2019.
- [40] Sourabh Vora, Alex H Lang, Bassam Helou, and Oscar Beijbom. Point-painting: Sequential fusion for 3d object detection. arXiv preprint arXiv:1911.10150, 2019.
- [41] Dominic Zeng Wang and Ingmar Posner. Voting for voting in online point cloud object detection. In Robotics: Science and Systems, volume 1, pages 10–15607, 2015.
- [42] Yan Wang, Wei-Lun Chao, Divyansh Garg, Bharath Hariharan, Mark Campbell, and Kilian Q Weinberger. Pseudo-lidar from visual depth estimation: Bridging the gap in 3d object detection for autonomous driving. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pages 8445–8453, 2019.
- [43] 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.
- [44] Zhixin Wang and Kui Jia. Frustum convnet: Sliding frustums to aggregate local point-wise features for amodal 3d object detection. arXiv preprint arXiv:1903.01864, 2019.
- [45] Xinshuo Weng and Kris Kitani. Monocular 3d object detection with pseudo-lidar point cloud. arXiv preprint arXiv:1903.09847, 2019.
- [46] Danfei Xu, Dragomir Anguelov, and Ashesh Jain. Pointfusion: Deep sensor fusion for 3d bounding box estimation. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pages 244–253, 2018.
- [47] Yan Yan, Yuxing Mao, and Bo Li. Second: Sparsely embedded convolutional detection. *Sensors*, 18(10):3337, 2018.
- [48] 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.
- [49] Zetong Yang, Yanan Sun, Shu Liu, and Jiaya Jia. 3dssd: Point-based 3d single stage object detector. arXiv preprint arXiv:2002.10187, 2020.
- [50] Zetong Yang, Yanan Sun, Shu Liu, Xiaoyong Shen, and Jiaya Jia. Ipod: Intensive point-based object detector for point cloud. arXiv preprint arXiv:1812.05276, 2018.
- [51] Zetong Yang, Yanan Sun, Shu Liu, Xiaoyong Shen, and Jiaya Jia. Std: Sparse-to-dense 3d object detector for point cloud. In Proceedings of the IEEE International Conference on Computer Vision, pages 1951–1960, 2019.
- [52] Jie Zhou, Xuequan Lu, Xin Tan, Zhiwei Shao, Shouhong Ding, and Lizhuang Ma. Fvnet: 3d front-view proposal generation for real-time object detection from point clouds. arXiv preprint arXiv:1903.10750, 2019.
- [53] 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.

...