

InfoFocus: 3D Object Detection for Autonomous Driving with Dynamic Information Modeling

Jun Wang^{1*}, Shiyi Lan^{1*}, Mingfei Gao², and Larry S. Davis¹

¹University of Maryland, College Park MD 20742, USA

²Salesforce Research, Palo Alto, CA 94301, USA

{junwang,lsd}@umiacs.umd.edu, sylan@cs.umd.edu

mingfei.gao@salesforce.com

Abstract. Real-time 3D object detection is crucial for autonomous cars. Achieving promising performance with high efficiency, voxel-based approaches have received considerable attention. However, previous methods model the input space with features extracted from equally divided sub-regions without considering that point cloud is generally non-uniformly distributed over the space. To address this issue, we propose a novel 3D object detection framework with dynamic information modeling. The proposed framework is designed in a coarse-to-fine manner. Coarse predictions are generated in the first stage via a voxel-based region proposal network. We introduce InfoFocus, which improves the coarse detections by adaptively refining features guided by the information of point cloud density. Experiments are conducted on the large-scale nuScenes 3D detection benchmark. Results show that our framework achieves the state-of-the-art performance with 31 FPS and improves our baseline significantly by 9.0% mAP on the nuScenes test set.

Keywords: 3D Object Detection, Point Cloud

1 Introduction

With growing interests in autonomous vehicles, 3D object detection has received considerable attention. Due to the superior capability of modeling 3D objects, point cloud is the most popular type of data source. Most existing 3D detectors are point-based [17, 25, 11, 21, 27] and voxel-based [12, 30, 26, 28, 7]. Point-based approaches generate features from raw point cloud data directly. Although achieving promising performance, these methods suffer from high computational complexity which discourages their application in real-time scenarios. Voxel-based approaches [12, 30, 26, 28, 7] firstly convert point cloud into voxels and then employ deep convolutional neural networks (DCNN) to conduct object detection. Taking advantage of the advanced DCNN architecture, voxel-based approaches achieve the state-of-the-art performance with low computational cost.

* Equal contribution.



Fig. 1. Left: we calculate the average point density across different parts of objects in BEV of nuScenes training set. $E1-E4$ indicate four edges sorted by their normalized density scores (sum to 100%) and *others* denotes areas inside objects. We set each edge width as 10% of the length along the object size and only objects over more than 100 points are counted. Middle and right: we visualize an example of the LiDAR point cloud in 2D image and its corresponding bird’s eye view (BEV). Clearly, most of the point clouds locate on the contour of the object

Our work follows the setting of voxel-based methods for their advanced balance of efficiency and effectiveness.

Although much progress has been made in improving the performance of voxel-based detectors, an important characteristic of point cloud is not well explored: input data points are usually not uniformly distributed over the space. The density of point cloud can be affected by different factors, e.g., the distance of objects from LiDAR sensor and object self-occlusion. As illustrated in Fig. 1, the density of point cloud over objects highly depends on the relative locations of different parts. It is also intuitive that the amount of information is highly related to the point density. However, existing voxel-based detectors extract features from uniformly divided sub-regions, regardless of the actual distribution of the points. We believe that this will lead to loss of useful information and ultimately result in sub-optimal detection performance.

To fully exploit the non-uniform distribution of point cloud, we propose a novel 3D object detection framework, to adaptively model the rich feature of 3D objects according to the information density of points. Illustrated in Fig. 2, our framework contains two stages. Coarse detection results are obtained in the first stage via voxel-based region proposal network. In the second stage, we introduce InfoFocus, to model and extract the informative features from regions of interest (RoI) (formed by the coarse predictions) according to the distribution of point cloud, and the predictions are improved with the help of the refined features.

The InfoFocus is the core structure of our framework which contains three sequentially connected modules including the Point-of-interest (PoI) Pooling, the Visibility Attentive Module, and the Adaptive Point-wise Attention.

PoI Pooling. Unlike 2D objects which contain densely distributed information over the whole RoI, more of the points of 3D objects locate on the their

surfaces. Therefore, we hypothesize that most informative feature concentrates on the edge of RoI. Motivated by this intuition, we propose PoI Pooling which densely samples features on the edge and sparsely samples feature in the middle of RoI to accommodate the non-uniform information distribution of point cloud.

Visibility Attentive Module. Heavy self-occlusion is presented because of the nature of LiDAR data that no point cloud exists on the backside of object relatively to the sensor. To mitigate this issue, our proposed Visibility Attentive Module applies hard attention to emphasize the visible parts of objects and eliminate the features from invisible points.

Adaptive Point-wise Attention. PoIs may contain different amount of information, although they are all visible. We introduce Adaptive Point-wise Attention to re-weight the features to improve the modeling of 3D objects.

We conduct extensive experiments on the largest public 3D object detection benchmark, i.e, nuScenes [1]. Experimental results show that our approach significantly outperforms the baselines, achieving 39.5% mAP with 31 FPS. Results of comprehensive ablation studies demonstrate the effectiveness of our InfoFocus and that each sub-module makes considerable contributions to our framework.

2 Related Work

Point-based Detectors. Inspired by the powerful feature learning capability of PointNet [18, 19] and the advanced modeling structure of 2D object detectors [5, 4, 20], Frustum PointNets [17] extrude the 2D object proposals into frustums to generate the 3D bounding boxes from raw point cloud. Lan et al. [11] add a decomposition-aggregation module modeling local geometry to extract the global feature descriptor of point cloud. Limited by initial 2D box proposals, those methods yield low performance when objects are occluded. In contrast, PointRCNN [21] generates 3D proposals directly from point cloud instead of 2D images. The recent STD [27] attempts to refine the detection boxes in a coarse-to-fine manner. However, all those methods are computationally expensive due to the large amount of data points to be processed.

Multi-view 3D Detectors. MV3D [2] is proposed to fuse multi-view feature maps for the generation of 3D box proposals. Following [2], Ku et al. [10] explore high resolution feature maps to compensate the information loss for small objects. These methods address the feature alignment between multi-modality in a coarse level and are typically slow. Liang et al. [14] design a continuous fusion layer to deal with the continuous state of LiDAR and the discrete state of images. Later, [13, 24] leverage different strategies to jointly fuse related tasks to improve feature representation.

Voxel-based Detectors. Recently, there is a trend of using regular 3D voxel grids to represent point cloud such that the input data can be easily processed by the 3D/2D convolution networks. Among those, VoxelNet [30] is the pioneering work of performing voxelization on the raw 3D point cloud. To improve its efficiency, Second [26] adopts Sparse Convolution and speeds up detection process without compromising the detection accuracy. PointPillars [12] dynamically

converts the 3D point cloud into a 2D pseudo image, making it more suitable for the application of the existing 2D object detection techniques. In [28], Ye et al. design a new voxel generator to preserve the information loss along the vertical direction. Building upon voxel-based detectors, our model captures richer information of objects by refining their feature representations at a second stage guided by the point cloud density and ultimately improves the detection results.

There are several recent studies [3, 16] focusing on fusing the voxel-based features with PointNet-based features in order to extract more fine-grained 3D features. InfoFocus is complementary to these techniques and can be further applied on top of them. WYSIWYG [7] is the most related method to our approach since we both drive the model to encode visibility information. However, instead of using a separate branch to generate the hidden invisibility representation, our method directly aggregates the valuable point-wise features together from existing backbone network to refine the proposals in an end-to-end manner.

3 Proposed Approach

The proposed framework is illustrated in Fig. 2, which consists of a deep feature extractor followed by a two-stage architecture. The deep feature extractor containing a Pillar Feature Network and a DCNN, converts the input point cloud to representative feature maps. Specifically, the Pillar Feature Network divides the whole space into equal pillars and generates the so-called pseudo images [12]. The pseudo images are then processed by the DCNN to obtain the feature maps which are shared by the two stages, i.e., Region Proposal Network (RPN) and InfoFocus. The RPN generates the initial coarse bounding box proposals that are refined by InfoFocus, with dynamic information modeling. Note that our Deep Feature Extractor and RPN follow the setting of [12].

3.1 Deep Feature Extractor

Deep Feature Extractor is composed of two parts: 1) voxelization using Pillar Feature Network that converts the orderless point cloud into a sparse pseudo image via a simplified PointNet-like architecture and 2) feature extraction using DCNN to learn informative feature maps.

Pillar Feature Network. The Pillar Feature Network operates on the raw point cloud, and learns point-wise features for each pillar. After voxelizing raw point cloud into evenly spaced pillars, we randomly sample N points from each non-empty pillar and then obtain a dense tensor with the size of $D \times P \times N$, where D indicates the information dimension of each point, P denotes the number of non-empty pillars, and N denotes the number of points in each pillar. The Pillar Feature Network utilizes a PointNet-like block to learn a multi-dimensional feature vector for each pillar. The pillar-wise features are encoded into a 2D pseudo image with the shape of $W \times L \times C$, where W and L indicate the width and length of the pseudo image, and C is the channel of the feature map.

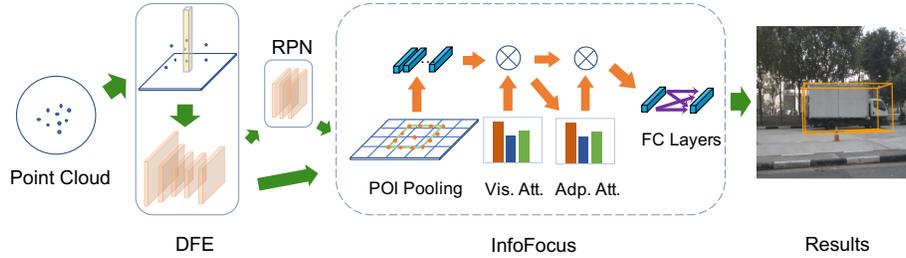


Fig. 2. The proposed 3D object detection framework. It consists of three parts: Deep Feature Extractor(DFE), Region Proposal Network, and InfoFocus. InfoFocus contains three modules: PoI Pooling, Visibility Attentive Module, and Adaptive Point-wise Attention Module

Deep Convolution Neural Network (DCNN). DCNN learns feature maps from the generated pseudo 2D image. The DCNN uses conv-deconv layers to extract features of different levels, and concatenates them to get the final features from different strides.

3.2 Region Proposal Network (RPN)

The RPN takes the feature maps provided by DCNN as inputs and produces high-quality 3D object proposals. Similar to the proposal generation in 2D object detection, anchor boxes are predefined at each position and proposals are generated by learning the offsets between anchors and the ground truths. To handle different scales of objects, a dual-head strategy is adopted. Specifically, the small-scale head takes features from the first conv-deconv phase of the DCNN, while the large-scale head takes the features from its concatenation phase.

3.3 InfoFocus

The InfoFocus serves as the second stage of our framework, which takes the candidate proposals from RPN and extracts features of objects in a hierarchical manner from the feature maps produced by the DCNN. Specifically, given each 3D object proposal, InfoFocus dynamically focuses on the informative parts of the feature maps by gradually emphasizing the representative PoIs in the following three steps: 1) the edge points are selected out from the whole proposal region by PoI Pooling; 2) Visibility Attention module emphasizes on the informative points according to their relative visibility to the LiDAR sensor and 3) in the Adaptive Point-wise Attention module, the features of the visible points are further weighted adaptively. The re-weighted features of the visible points are then fused to form the final representation of the proposal, on top of which two fully-connected layers are utilized to predict the refined box.

PoI Pooling. When representing a 3D proposal, the most intuitive way is adopting the commonly used strategy in the two-stage 2D object detectors, i.e., RoI

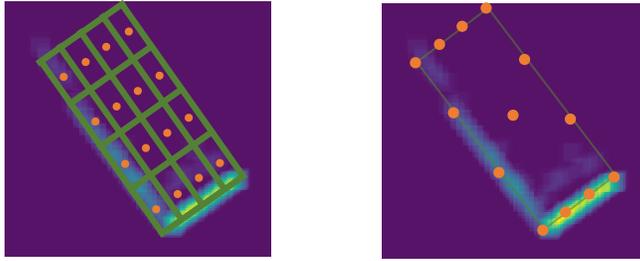


Fig. 3. RoI Pooling *vs.* PoI Pooling. The grid represents the feature map, and the dots denotes sampling points of interest. RoI Pooling samples the whole box, while PoI Pooling focuses on the key-points from edge-of-interest

Pooling (see Fig. 3 left). However, unlike the 2D images that have densely distributed information over the region proposals, the 3D point cloud mostly resides on the object surface which results in non-uniform information over the regions (most information locates on the edges of proposals).

The proposed PoI pooling is illustrated in Fig. 3 (right). Instead of equally sampling points over a region of the feature maps, we focus on sampling the points on the informative parts including four corners, the center point and key-points on the edges. Note that we consider the center position as an additional useful signal since it is likely to capture the semantic-level information.

We first project the 3D proposal to the birds’ view coordinate system. Let p_0, p_1, p_2 and p_3 represent the positions of top-left, top-right, bottom-right, and bottom-left corners of a proposal on the pseudo image, respectively and p_c denotes the center point. Along each edge, n more key-points are uniformly sampled. For example, for the top edge between p_0 and p_1 , the position of a sampled key-point $kp_j = (p_0 \frac{j}{n+1} + p_1 \frac{n+1-j}{n+1})$, where j is an integer and $1 \leq j \leq n$. To this end, $(5 + 4 * n)$ PoIs are obtained. A high-dimensional feature is extracted for each PoI according to its relative position on the feature map and then we obtain a feature set $F_{poi} = \{f_1^{poi}, f_2^{poi}, \dots, f_{N_{poi}}^{poi}\}$, where $N_{poi} = (5 + 4 * n)$ representing the number of selected PoIs within the considered region.

Visibility Attentive Module. Severe self-occlusion typically occurs in point cloud, but is ignored by most of the existing methods. The Visibility Attentive Module (see Fig. 4 left) is proposed to mitigate this issue by focusing on the information provided by the visible parts of objects. We argue that visible regions contain more useful information than the occluded ones. Formally, we propose to re-weight features of PoIs according to their corresponding visibility by exploiting the geometric relationship between the proposals and the LiDAR sensor in bird’s eye view. As shown in Eq. 1, F_{vis} denotes the updated feature set, where v_i^{poi} indicates the visibility score of the i th PoI. Different weighting strategies can be used and we use a hard attention strategy in this work for its simplicity, that is assigning $v_i^{poi} = 1$ if the i th PoI is visible and $v_i^{poi} = 0$ otherwise. In other words, we only take PoIs on the visible edges to represent the proposal.

$$F_{vis} = \{f_1^{poi} * v_1^{poi}, f_2^{poi} * v_2^{poi}, \dots, f_{N_{poi}}^{poi} * v_{N_{poi}}^{poi}\} \quad (1)$$

For the consideration of model efficiency, a simple yet effective method is used to estimate the visibility of points in the bird’s eye view. To figure out the sides of proposals facing to the sensor, we first compute the distance of each corner to the LiDAR sensor and determine the one that is closest to the sensor. Then, we consider the two edges passing this closest corner as the visible edges and the other two as the occluded ones.

Adaptive Point-wise Attention Module. PoI Pooling and Visibility Attentive Module are motivated by the nature of the non-uniform density of point cloud. However, two points may offer different amount of information even though they are all visible by the sensor. Adaptive Point-wise Attention Module provides the flexibility for the visible PoIs to contribute unequally to the prediction. Suppose $F_{vis} = \{f_1^{vis}, f_2^{vis}, \dots, f_{N_{vis}}^{vis}\}$ indicates the feature set of visible PoIs. Adaptive Point-wise Attention Module learns an attention weight, w_i , for each f_i^{vis} adaptively for the next-step feature aggregation. Specifically, a shared fully connected (FC) layer with sigmoid as the activation function is used to learn the attention weights, formally expressed as $v_i^{vis} = Sigmoid(\mathbf{W}f_i^{vis} + b)$. We use $F_{att} = \{f_1^{att}, f_2^{att}, \dots, f_{N_{vis}}^{att}\}$ to represent the re-weighted feature set of visible PoIs updated using F_{vis} and the attention weights, where $f_i^{att} = f_i^{vis} * v_i^{vis}$.

The final representation of each proposal aggregates the features of its visible PoIs. Let e_0, e_1, e_2 and e_3 denote the top, right, down, left edges of a proposal, respectively. We first compute f_i^e by applying max pooling to all the visible points on e_i . Then, the final representation is obtained by $f_0^e || f_1^e || f_2^e || f_3^e || f_c^p$, where f_c^p indicates the feature of the center point and $||$ indicates concatenation.

3.4 Loss Function

Given the output PoI feature representation from the aforementioned three modules topped by fully-connected layers, the head network consists of three branches predicting the box class, localization and direction. The ground truth and anchor boxes are parameterized as $(x, y, z, w, l, h, \theta)$, where (x, y, z) is the center of box, (w, l, h) is the dimension of box, and θ is the heading along the z-axis in the LiDAR coordinate system. The box regression target is computed as the residuals between the ground truth and the anchors as:

$$\begin{aligned} \Delta x &= \frac{x^{gt} - x^a}{d^a}, \Delta y = \frac{y^{gt} - y^a}{d^a}, \Delta z = \frac{z^{gt} - z^a}{h^a}, \\ \Delta w &= \log\left(\frac{w^{gt}}{w^a}\right), \Delta l = \log\left(\frac{l^{gt}}{l^a}\right), \Delta h = \log\left(\frac{h^{gt}}{h^a}\right), \\ \Delta \theta &= \theta_{gt} - \theta_a \end{aligned} \quad (2)$$

where x^{gt} and x^a refer to ground truth and anchor box respectively, and $d^a = \sqrt{(w^a)^2 + (l^a)^2}$. To deal with severe class imbalance problem in the dataset,

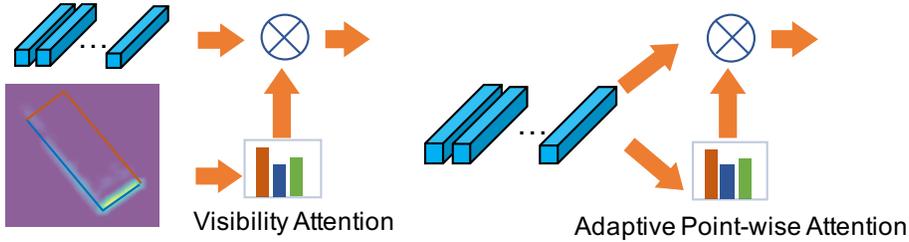


Fig. 4. Left: illustration of the Visibility Attentive Module. We compute hard attention for each sampled point depending on whether it is visible to the sensor. We also show the visibility map on the bottom left. Points on the blue line are visible while points on the orange line are invisible. Right: the architecture of the Adaptive Point-wise Attention Module. The point-wise attention is generated using a fully connected (FC) layer followed by a Sigmoid function. The input of FC layer is the feature of each point

we adopt the focal loss [15] for the classification loss. Smooth L1 loss [5] is used for the regression loss. In addition, to compensate for direction prediction missing in the regression, we adopt a softmax classification loss on orientation prediction.

Similar with that of the vanilla PointPillars network [12], we formally define a multi-task loss for both stages as threefold,

$$L_{stage.i} = \frac{1}{N_{pos}} (\beta_{cls} L_{cls.i} + \beta_{reg} L_{reg.i} + \beta_{dir} L_{dir.i}), \quad (3)$$

where i could be either RPN or InfoFocus stage, N_{pos} refers to the number of positive anchors and $\beta_{cls}, \beta_{reg}, \beta_{dir}$ are chosen to balance the weights among classification loss, regression loss and direction loss.

3.5 Comparisons with Existing Approaches

Point-based Approaches. Our framework uses PointNet to extract features from equally divided sub-grids and employs a DCNN to generate 2D feature maps while point-based techniques [17, 25, 11, 21] only use PointNet as its backbone. Both our approach and point-based approaches apply two-stage architecture to infer objects. Meanwhile, we both sample features considering the distribution of point cloud. However, compared to PointNet, InfoFocus is more computationally efficient without performance degradation.

Fusion-based Approaches. Fusion-based detectors [3, 16] make use of both RGB images and point cloud data for 3D object detection. InfoFocus is much faster than fusion-based approaches, since they contain two backbones to process multi-view sources and are heavily engineered. On the other hand, InfoFocus also achieves competitive results compared to fusion-based approaches.

Traditional Voxel-based Approaches. Our method shares the similar backbone as the existing voxel-based architectures [12, 30, 26, 28]. However, previous voxel-based detectors pay less attention to the distribution of LiDAR data that

most 3D point cloud locates on the surface of the objects. Our proposed PoI Pooling, Visibility Attentive Module, and Adaptive Point-wise Attention model the non-uniform point cloud using dynamic information focus. First, the PoI Pooling decreases the sampling from the inside of objects where few points locate. Next, the Visibility Attentive Module eliminates the noise from the back of objects where points are occluded. Last, we apply the Adaptive Point-wise Attention to learn the focus on each sampled points. Jointly, these modules contribute significantly to the superior performance of InfoFocus.

4 Experiments

Our method is mainly evaluated on the nuScenes dataset [1] which is considered as the most challenging 3D object detection benchmark. We first present our implementation details. We compare with the existing approaches both quantitatively and qualitatively. Then, extensive ablation studies are conducted to demonstrate the effectiveness of each designed module. Last, we analyze the inference time and the desired speed accuracy trade-off provided by our method.

4.1 Dataset and Evaluation Metric

NuScenes [1] is one of the largest datasets for autonomous driving. There are 1000 scenes of 20s duration each, including 23 object classes annotated with 28,130 training, and 6,019 validation samples. We use the LiDAR point cloud as the only input to our method and all the experiments follow the standard protocol on the training and validation sets. Officially, nuScenes evaluates the detection accuracy across different classes, based on the average precision metric (AP) which is computed based on 2D center distance between ground truth and the detection box on the ground plane. In details, the AP score is determined as the normalized area under the precision recall curve above 10%. The final mean AP (mAP) is the average among the set of ten classes over matching thresholds of $\mathbb{D} = \{0.5, 1, 2, 4\}$ meters.

4.2 Implementation Details

We integrate InfoFocus into a state-of-the-art real-time 3D object detector [12] to improve the detection performance without largely compromising speed. Closely following the codebase¹ recommended by the authors of PointPillars [12], we use PyTorch to implement our InfoFocus modules and integrate it into vanilla PointPillars network. More details will be introduced in the supplementary materials. **RPN.** For each class of objects, the RPN anchor size is set by calculating the average of all objects from the corresponding class in training set of nuScenes. In addition, the matching thresholds are based on the custom configuration following the suggested codebase. 1,000 proposals are obtained from RPN, on

¹ <https://github.com/traveller59/second.pytorch>.

which NMS with a threshold of 0.5 is applied to remove the overlapping proposals for both training and inference. The final top-ranked 300 proposals are kept for the InfoFocus stage to simultaneously predict the category, location and direction of objects during both the training and inference.

InfoFocus. The second stage is our proposed InfoFocus. The three novel modules process object-centric feature sequentially based on the initial bounding box proposals from RPN. The number of sampled key-points for each edge, n , is set to be 2. Thus, the total number of PoIs, N_{poi} , is 13, including a center, 4 corners and 2 key-points on each edge. Similar to RoIAlign [6], bi-linear interpolation is used to compute the deep feature from four neighboring regular locations of each point.

As mentioned before, we apply a max-pool layer to summarize the features of points along each edge, resulting in 5 features for each proposal, including features from top, right, down, left edges and the center. When concatenating these features, we always treat the edge that is closest to the sensor as the top edge. A fully connected layer with a single node is used to generate point-wise attention weight for each point.

The feature of each proposal is transformed by two consecutive FC layers with 512 nodes each and passed to three sibling linear layers, a box-regression branch, a box-classification branch and a box-direction branch. For the regression target assignment, anchors having Intersection over Union (IoU) bigger than 0.6 with the ground truth are considered positive, and smaller than 0.55 are assigned negative labels.

Training Parameters. Experiments are conducted on a single NVIDIA 1080Ti GPU. The weight decay is set to be 0.01. We adopt the Adam optimizer [9], and use a one-cycle scheduler proposed in [23]. We train our model with a total of 20 epochs as a default choice, taking about 40 hours from scratch. For the first 8 epochs, the learning rate progressively increases from 3×10^{-4} to 3×10^{-3} with decreasing momentum from 0.95 to 0.85, while in the remaining 12 epochs learning rate decreases from 3×10^{-3} to 3×10^{-6} with increasing momentum from 0.85 to 0.95. The focal loss [15] with $\alpha = 0.25$ and $\gamma = 2.0$ is adopted for the classification loss. The balancing weights for the classification, box regression, and direction loss $\beta_{cls}, \beta_{reg}, \beta_{dir}$ of both stages are 1, 2 and 0.2, respectively.

4.3 Main Results

First, we compare our framework with the state-of-the-art methods on the nuScenes validation set, including the vanilla PointPillars [12] as our baseline, and recently published WYSIWYG [7]. As can be seen from Table. 1, the baseline has an mAP of 29.5% with a single stage, while InfoFocus improves it by a massive 6.9%. This demonstrates the effectiveness of InfoFocus. We also visualize the detection results of our framework on 2D and 3D BEV images in Fig. 5. As shown in Fig. 6, compared to the vanilla PointPillars qualitatively, InfoFocus helps remove the false positives significantly and obtains better results.

In addition, we submit the detection results of test set on the nuScenes test server. The results show that our method achieves the state-of-the-art perfor-

Table 1. Object detection results (%) on nuScenes validation set

Method	car	ped.	barri.	traff.	truck	bus	trail.	const.	motor.	bicyc.	mAP
PointPillars [12]	70.5	59.9	33.2	29.6	25.0	34.3	16.7	4.5	20.0	1.6	29.5
WYSIWYG [7]	80.0	66.9	34.5	27.9	35.8	54.1	28.5	7.5	18.5	0	35.4
Ours	77.6	61.7	43.4	33.4	35.4	50.5	25.6	8.3	25.2	2.5	36.4

mance with inference speed of 31 FPS, improving the baseline performance by 7% mAP. Note that all methods listed in Table. 2 are LiDAR-based except that MonoDIS [22] and CenterNet[29] are camera-based methods. Without bells and whistles, our approach works better than WYSIWYG [7]. Considering that our model contains more parameters than the vanilla PointPillars, we empirically increase the number of the training epoch by 2 times. With all the others settings the same, our method is improved by 2% mAP on the nuScenes test set as shown in Table. 2 (Ours 2×). In total, our method outperforms WYSIWYG[7] by 4.5% mAP on the nuScenes test set. In the rest of paper, the default setting of training epochs is adopted. To the best of our knowledge, our framework is superior than all the published real-time methods with respect to mAP.

Table 2. Object detection results (%) on nuScenes test set. Note that MonoDIS and CenterNet are camera based methods, and the rest are LiDAR based. Ours 2× indicates 2× training time with other settings being the same with Ours

Method	car	ped.	barri.	traff.	truck	bus	trail.	const.	motor.	bicyc.	mAP
MonoDIS [22]	47.8	37.0	51.1	48.7	22.0	18.8	17.6	7.4	29.0	24.5	30.4
PointPillars [12]	68.4	59.7	38.9	30.8	23.0	28.2	23.4	4.1	27.4	1.1	30.5
SARPNET [28]	59.9	69.4	38.3	44.6	18.7	19.4	18.0	11.6	29.8	14.2	32.4
CenterNet [29]	53.6	37.5	53.3	58.3	27.0	24.8	25.1	8.6	29.1	20.7	33.8
WYSIWYG [7]	79.1	65.0	34.7	28.8	30.4	46.6	40.1	7.1	18.2	0.1	35.0
Ours	77.2	61.5	45.3	40.4	31.5	44.1	35.9	9.8	25.1	4.0	37.5
Ours 2×	77.9	63.4	47.8	46.5	31.4	44.8	37.3	10.7	29.0	6.1	39.5

4.4 Ablation Studies

To understand the contribution of our major component to the success of InfoFocus, Table. 3 summarizes the performance of our framework when a certain module is disabled, including PoI Pooling, Visibility Attention Module and Adaptive Attention Module.

PoI Pooling. To investigate the effect of PoI Pooling, we simply add the PoI Pooling on top of the vanilla PointPillar. This baseline introduces 3.0% mAP improvement. However, when we vary the number of pooling key-points on each edge, we see that our framework with four key-points ($n = 4$) on each edge degrades slightly by 0.8% mAP than that of two key-points ($n = 2$). A possible

Table 3. Ablation studies on nuScenes validation set. "Vis. Att" and "Adp. Att." refer to Visibility Attention Module and Adaptive Attention Module, respectively

PoIPool	Vis. Att.	Adp. Att.	mAP
			29.5
✓			32.5
✓	✓		34.8
✓		✓	34.8
✓	✓	✓	36.4

reason is that the higher number of samples along each edge might bring more noise which harms the detection performance.

Visibility Attention. We further add the Visibility Attention module to filter out invisible edges before PoI pooling. Table. 3 shows that when using the features from two visible edges, the mAP result is improved by 2.3% mAP compared to *baseline+PoIPool*. Generally, the visible parts of objects correspond to their sides closer to the LiDAR sensor, thus they may capture richer information. By applying visibility attention, our method focuses more on the representative information which results in better performance.

Adaptive Point-wise Attention. Without the Adaptive Point-wise Attention module, the framework naturally allows the same weight for each PoI feature. As we can see in Table. 3, when adding this module, the result of *baseline+PoIPool* improves by 2.3% mAP and that of *baseline+PoIPool+Vis.Att.* improves by 1.6%. These results suggest that the Adaptive Point-wise Attention module helps emphasize on useful points which leads to a better performance.

Table 4. Inference time of 3D object detectors. Note that inference time for the baseline here is the network reproduced by ourselves

Method	Input Format	mAP	Inference Time (ms)
Baseline [12]	LiDAR	30.5	26.9
MonoDIS [22]	RGB	30.4	29.0
SARPNET [28]	LiDAR	32.4	70.0
Ours	LiDAR	37.5	32.9

Rotated RoIAlign Comparison. One widely considered way to extract the region-wise features in the two-stage architecture is RoIAlign [6]. So, it is intuitive to compare with this strategy under the setting of 3D object detection. We implement rotated RoIAlign (RRoI) operation [8] to compensate for the rotated bounding box, since in our case they are often not axis-aligned. We conduct experiments exploring two different pooling sizes, 4×4 (pooled length and pooled width), and 8×4 with 4 sampled points in each bin. One of the reasons that we use 8×4 is that most of the objects like car and bus’s length is larger than their width. With all other implementations the same as InfoFocus, Table. 5 presents

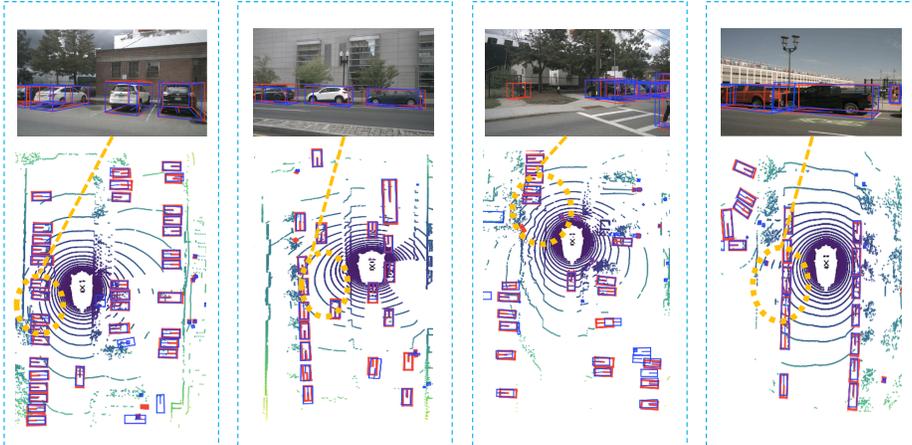


Fig. 5. We visualize the detection results on nuScenes with 2D and 3D BEV images. On the top, we demonstrate the 2D images with the 3D bounding box annotated, while the BEV of LiDAR with ground truth (red) and detection (blue) box are shown on the bottom. Note that the line in the frame denotes the direction of the object

detection results utilizing the rotated-ROI with different pooling sizes. Compared with the vanilla PointPillars [12], adding the RoIAlign layer with size of 4×4 increases the mAP performance by 4.4%. However, InfoFocus still outperforms RoIAlign by 2.5% with the better information modeling scheme.

Table 5. Comparison with rotated RoIAlign feature extraction results (%) on the nuScenes validation set

Method	car	peds.	barri.	traff.	truck	bus	trail.	const.	motor.	bicyc.	mAP
RRoI 4x4	76.9	60.1	37.6	29.5	32.4	50.6	22.4	5.0	20.8	3.8	33.9
RRoI 8x4	77.0	59.5	36.7	29.2	33.2	51.5	25.4	4.5	24.0	1.8	34.3
Ours	77.6	61.7	43.4	33.4	35.4	50.5	25.6	8.3	25.2	2.5	36.4

4.5 Real-time Inference Analysis

As indicated in Table. 4, our framework takes about 32.9 ms to perform detection on an example of point cloud in the nuScenes, compared with 26.9 ms of the vanilla PointPillars when both are evaluated on a single NVIDIA 1080Ti GPU. In details, the pillar feature extraction time is 12.6 ms, the DCNN costs 1.1 ms, RPN takes 7.3 ms to generate proposals, and the InfoFocus stage takes 11.9 ms. Specifically, the proposal generation for the InfoFocus stage including NMS is 5.1 ms, the PoI feature extraction time is 3.1 ms, and the second stage including three branches takes 0.7 ms. We also note that WYSIWYG [7] provides the

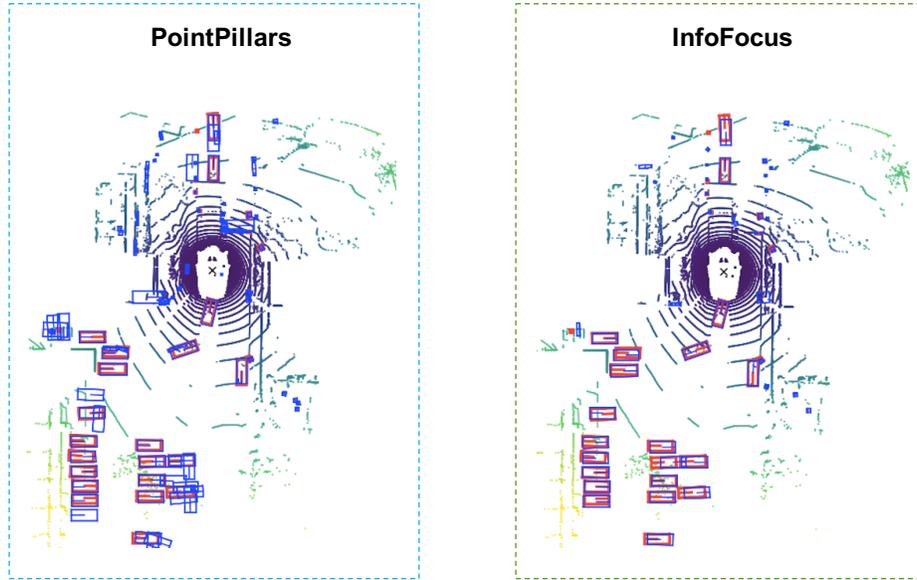


Fig. 6. We visualize the BEV detection results for the same point cloud sample on nuScenes with the vanillar PointPillars (left) and InfoFocus (right)

overhead of computing visibility over a 32-beam LiDAR point to be 24.4 ± 3.5 ms on average and InfoFocus is faster than WYSIWYG [7] since we share the similar backbone network. The framework with RROIAlign has an inference time of 32.2 ms. Further, compared with other point-based methods [21, 27], InfoFocus is considerably faster and conceptually simpler.

5 Conclusions

Non-uniform distribution of point cloud causes varying amount of information at different locations. Therefore, we argue that this imbalance distribution of information may result in degradation on previous 3D voxel-based detectors when modeling 3D objects. To address this issue, we propose a 3D object detection framework with InfoFocus to dynamically conduct information modeling. InfoFocus contain three effective modules including PoI Pooling, the Visibility Attentive Module, and the Adaptive Point-wise Attention. Demonstrated by the comprehensive experiments, our framework achieves the state-of-art performance among all the real-time detectors on the challenging nuScenes dataset.

Acknowledgement. This work was supported by the Intelligence Advanced Research Projects Activity (IARPA) via DOI/IBC contract numbers D17PC00345 and D17PC00287. The U.S. Government is authorized to reproduce and distribute reprints for Governmental purposes not withstanding any copyright annotation thereon. The authors would like to thank Zuxuan Wu and Xingyi Zhou for proofreading the manuscript.

References

1. Caesar, H., Bankiti, V., Lang, A.H., Vora, S., Liong, V.E., Xu, Q., Krishnan, A., Pan, Y., Baldan, G., Beijbom, O.: nuscenes: A multimodal dataset for autonomous driving. arXiv preprint arXiv:1903.11027 (2019)
2. Chen, X., Ma, H., Wan, J., Li, B., Xia, T.: Multi-view 3d object detection network for autonomous driving. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 1907–1915 (2017)
3. Chen, Y., Liu, S., Shen, X., Jia, J.: Fast point r-cnn. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 9775–9784 (2019)
4. Girshick, R.: Fast r-cnn. In: The IEEE International Conference on Computer Vision (ICCV) (December 2015)
5. Girshick, R., Donahue, J., Darrell, T., Malik, J.: Rich feature hierarchies for accurate object detection and semantic segmentation. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2014)
6. He, K., Gkioxari, G., Dollár, P., Girshick, R.: Mask r-cnn. In: Proceedings of the IEEE international conference on computer vision. pp. 2961–2969 (2017)
7. Hu, P., Ziglar, J., Held, D., Ramanan, D.: What you see is what you get: Exploiting visibility for 3d object detection. arXiv preprint arXiv:1912.04986 (2019)
8. Huang, J., Sivakumar, V., Mnatsakanyan, M., Pang, G.: Improving rotated text detection with rotation region proposal networks. arXiv preprint arXiv:1811.07031 (2018)
9. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980 (2014)
10. Ku, J., Mozifian, M., Lee, J., Harakeh, A., Waslander, S.L.: Joint 3d proposal generation and object detection from view aggregation. In: 2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS). pp. 1–8. IEEE (2018)
11. Lan, S., Yu, R., Yu, G., Davis, L.S.: Modeling local geometric structure of 3d point clouds using geo-cnn. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2019)
12. Lang, A.H., Vora, S., Caesar, H., Zhou, L., Yang, J., Beijbom, O.: Pointpillars: Fast encoders for object detection from point clouds. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 12697–12705 (2019)
13. Liang, M., Yang, B., Chen, Y., Hu, R., Urtasun, R.: Multi-task multi-sensor fusion for 3d object detection. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 7345–7353 (2019)
14. Liang, M., Yang, B., Wang, S., Urtasun, R.: Deep continuous fusion for multi-sensor 3d object detection. In: The European Conference on Computer Vision (ECCV) (September 2018)
15. Lin, T.Y., Goyal, P., Girshick, R., He, K., Dollár, P.: Focal loss for dense object detection. In: Proceedings of the IEEE international conference on computer vision. pp. 2980–2988 (2017)
16. Liu, Z., Tang, H., Lin, Y., Han, S.: Point-voxel cnn for efficient 3d deep learning. In: Advances in Neural Information Processing Systems. pp. 963–973 (2019)
17. Qi, C.R., Liu, W., Wu, C., Su, H., Guibas, L.J.: Frustum pointnets for 3d object detection from rgb-d data. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 918–927 (2018)
18. Qi, C.R., Su, H., Mo, K., Guibas, L.J.: Pointnet: Deep learning on point sets for 3d classification and segmentation. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 652–660 (2017)

19. Qi, C.R., Yi, L., Su, H., Guibas, L.J.: Pointnet++: Deep hierarchical feature learning on point sets in a metric space. In: *Advances in neural information processing systems*. pp. 5099–5108 (2017)
20. Ren, S., He, K., Girshick, R., Sun, J.: Faster r-cnn: Towards real-time object detection with region proposal networks. In: *Advances in neural information processing systems*. pp. 91–99 (2015)
21. Shi, S., Wang, X., Li, H.: Pointcnn: 3d object proposal generation and detection from point cloud. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. pp. 770–779 (2019)
22. Simonelli, A., Bulo, S.R., Porzi, L., López-Antequera, M., Kotschieder, P.: Disentangling monocular 3d object detection. In: *Proceedings of the IEEE International Conference on Computer Vision*. pp. 1991–1999 (2019)
23. Smith, L.N.: A disciplined approach to neural network hyper-parameters: Part 1—learning rate, batch size, momentum, and weight decay. *arXiv preprint arXiv:1803.09820* (2018)
24. Vora, S., Lang, A.H., Helou, B., Beijbom, O.: Pointpainting: Sequential fusion for 3d object detection. *arXiv preprint arXiv:1911.10150* (2019)
25. Wang, Z., Jia, K.: Frustum convnet: Sliding frustums to aggregate local point-wise features for amodal 3d object detection. *arXiv preprint arXiv:1903.01864* (2019)
26. Yan, Y., Mao, Y., Li, B.: Second: Sparsely embedded convolutional detection. *Sensors* **18**(10), 3337 (2018)
27. Yang, Z., Sun, Y., Liu, S., Shen, X., Jia, J.: Std: Sparse-to-dense 3d object detector for point cloud. In: *Proceedings of the IEEE International Conference on Computer Vision*. pp. 1951–1960 (2019)
28. Ye, Y., Chen, H., Zhang, C., Hao, X., Zhang, Z.: Sarpnet: Shape attention regional proposal network for lidar-based 3d object detection. *Neurocomputing* **379**, 53–63 (2020)
29. Zhou, X., Wang, D., Krähenbühl, P.: Objects as points. *arXiv preprint arXiv:1904.07850* (2019)
30. Zhou, Y., Tuzel, O.: Voxelnets: End-to-end learning for point cloud based 3d object detection. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. pp. 4490–4499 (2018)