
Fully Sparse 3D Object Detection

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 As the perception range of LiDAR increases, LiDAR-based 3D object detection
2 becomes a dominant task in the long-range perception task of autonomous driving.
3 The mainstream 3D object detectors usually build dense feature maps in the
4 network backbone and prediction head. However, the computational and spatial
5 costs on the dense feature map are quadratic to the perception range, which makes
6 them hardly scale up to the long-range setting. To enable efficient long-range
7 LiDAR-based object detection, we build a fully sparse 3D object detector (FSD).
8 The computational and spatial cost of FSD is roughly linear to the number of points
9 and independent of the perception range. FSD is built upon the general sparse voxel
10 encoder and a novel sparse instance recognition (SIR) module. SIR first groups the
11 points into instances and then applies instance-wise feature extraction and predic-
12 tion. In this way, SIR resolves the issue of center feature missing, which hinders the
13 design of the fully sparse architecture for all center-based or anchor-based detectors.
14 Moreover, SIR avoids the time-consuming neighbor queries in previous point-based
15 methods by grouping points into instances. We conduct extensive experiments on
16 the large-scale Waymo Open Dataset to reveal the working mechanism of FSD, and
17 state-of-the-art performance is reported. To demonstrate the superiority of FSD in
18 long-range detection, we also conduct experiments on Argoverse 2 Dataset, which
19 has a much larger perception range (200m) than Waymo Open Dataset (75m). On
20 such a large perception range, FSD achieves state-of-the-art performance and is
21 $2.4\times$ faster than the dense counterpart. Codes will be released.

22 1 Introduction

23 Autonomous driving systems are eager for efficient
24 long-range perception for downstream tasks, espe-
25 cially in high-speed scenarios. Current 3D LiDAR-
26 based object detectors usually convert sparse features
27 into dense feature maps for further feature extrac-
28 tion and prediction. For simplicity, we name the
29 detectors utilizing dense feature maps as **dense de-**
30 **ectors**. Dense detectors perform well on current
31 popular benchmarks [28, 7, 2], where the percep-
32 tion range is relatively short (less than 75 meters).
33 However, it is impractical to scale the dense detec-
34 tors to the long-range setting (more than 200 meters,
35 Fig. 2). In such settings, the computational and spa-
36 tial complexity on dense feature maps is quadratic
37 to the perception range. Fortunately, the sparsity of
38 LiDAR point clouds also increases as the perception
39 range extends (see Fig. 2), and the calculation on the

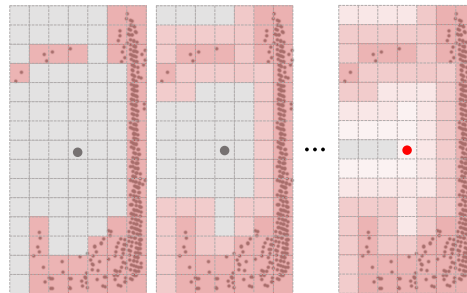


Figure 1: Illustration of **center feature missing** and **feature diffusion** on dense feature maps from Bird's Eye View. The empty instance center (red dot) is filled by the features diffused from occupied voxels (with LiDAR points), after several convolutions.

40 unoccupied area is essentially unnecessary. Given the inherent sparsity, an essential solution for
41 efficient long-range detection is to remove the dense feature maps and make the network architectures
42 *fully sparse*.

43 However, removing the dense feature map is non-trivial since it plays a critical role in current designs.
44 Commonly adopted sparse voxel encoders [35, 5, 26] only extract the features on the non-empty
45 voxels for efficiency. So without dense feature maps, the object centers are usually empty, especially
46 for large objects. We name this issue as “**Center Feature Missing (CFM)**” (Fig. 1). Almost all
47 popular voxel or pillar based detectors [25, 5, 39, 27, 35] in this field adopt center-based or anchor-
48 based assignment since the center feature is the best representation of the whole object. However,
49 CFM significantly weakens the representation power of the center voxels, even makes the center
50 feature empty in some extreme cases like super large vehicles. Given this difficulty, many previous
51 detectors [35, 39, 25, 5] have to convert sparse voxels to dense feature maps in Bird’s Eye View after
52 the sparse voxel encoder. Then they resolve the CFM issue by applying convolutions on the dense
53 feature maps to diffuse features to instance centers, which we name as **feature diffusion** (Fig. 1).

54 To properly remove the dense feature map, we then
55 investigate the purely point-based detectors because
56 they are naturally fully sparse. However, two draw-
57 backs limit the usage of point-based methods in the
58 autonomous driving scenario. (1) **Inefficiency**: The
59 time-consuming neighborhood query [22] is the long-
60 standing difficulty to apply it to large-scale point
61 cloud (more than 100K points). (2) **Coarse represen-**
62 **tation**: To reduce the computational overhead, point-
63 based methods aggressively downsample the whole
64 scene to a fixed number of points. The aggressive
65 downsampling leads to inevitable information loss
66 and insufficient recall of foreground objects [37, 40].
67 As a result, very few purely point-based detectors
68 have reached state-of-the-art performance in the re-
69 cent benchmarks with large-scale point clouds.

70 In this paper, we propose **Fully Sparse Detector**
71 **(FSD)**, which takes the advantages of both sparse
72 voxel encoder and point-based instance predictor.
73 Since the central region might be empty, the detector
74 has to predict boxes from other non-empty parts of instances. However, predicting the whole box
75 from individual parts causes a large variance on the regression targets, making the results noisy and
76 inconsistent. This motivates us to first group the points into an instance, then we further extract the
77 instance-level feature and predict a single bounding box from the instance feature. To implement this
78 principle, FSD first utilizes the sparse voxel encoder [35, 26, 5] to extract voxel features, then votes
79 object centers based on these features as in VoteNet [21]. Then the **Instance Point Grouping (IPG)**
80 module groups the voted centers into instances via Connected Components Labeling. After grouping,
81 a point-based **Sparse Instance Recognition (SIR)** module extracts instance features and predicts
82 the whole bounding boxes. As a point-based module, SIR has several desired properties. (1) Unlike
83 previous point-based modules, SIR treats instances as groups, and does not apply the time-consuming
84 neighborhood query for further grouping. (2) Similar to dynamic voxelization [43], SIR leverages
85 **dynamic broadcast/pooling** for tensor manipulation to avoid point sampling or padding. (3) Since
86 SIR covers the whole instance, it builds a sufficient receptive field regardless of the physical size of
87 the instance. We list our contributions as follows.

- 88 • We propose the concept of Fully Sparse Detector (FSD), which is the essential solution for
89 efficient long-range LiDAR detection. We further propose Sparse Instance Recognition (SIR)
90 to overcome the issue of Center Feature Missing in sparse feature maps. Combining SIR
91 with general sparse voxel encoders, we build an efficient and effective FSD implementation.
- 92 • FSD achieves state-of-the-art performance on the commonly used Waymo Open Dataset.
93 Besides, we further apply our method to the recently released Argoverse 2 dataset to
94 demonstrate the superiority of FSD in long-range detection. Given its challenging 200
95 meters perception range, FSD showcases state-of-the-art performance while being $2.4\times$
96 faster than state-of-the-art dense detectors.

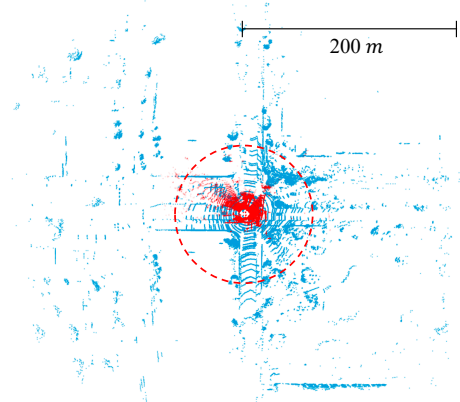


Figure 2: **Short-range point clouds (red, from KITTI [7]) v.s. long-range point clouds (blue, from Argoverse 2 [34]).** The radius of the red circle is 75 meters. The sparsity quickly increases as the range extends.

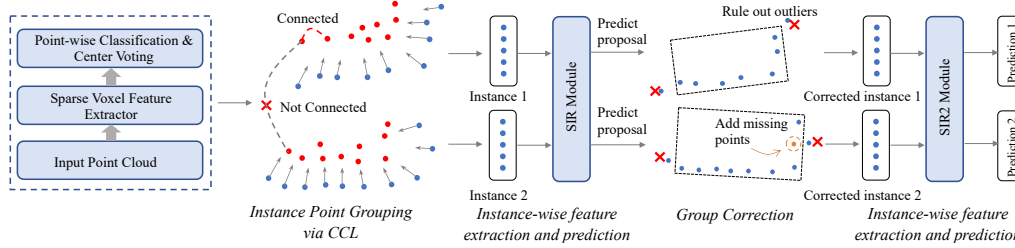


Figure 3: **Overall architecture of FSD.** For simplicity, we only use two instances to illustrate the pipeline. Red dots are the voted centers from each LiDAR point (blue dots). The SIR module and the SIR2 module all contain 3 SIR layers.

97 2 Related Work

98 **Voxel-based dense detectors** Pioneering work VoxelNet [42] uses dense convolution for voxel
 99 feature extraction. Although it achieves competitive performance, it is inefficient to apply dense
 100 convolution to 3D voxel representation. PIXOR [36] and PointPillars [12] adopt 2D dense convolution
 101 in Bird’s Eye View (BEV) feature map achieving significant efficiency improvement. We name such
 102 detectors as **dense detectors** since they convert the sparse point cloud into dense feature maps.

103 **Voxel-based semi-dense detectors** Different from the dense detectors, **semi-dense detectors**
 104 incorporate both sparse features and dense features. SECOND [35] adopts sparse convolution to
 105 extract the sparse voxel features in 3D space, which then are converted to dense feature maps in
 106 BEV to enlarge the receptive field and integrate with 2D detection head [17, 23, 41]. Based on
 107 SECOND-style semi-dense detectors, many methods attach a second stage for fine-grained feature
 108 extraction and proposal refinement [26, 25, 27, 3]. It is noteworthy that the semi-dense detector is
 109 hard to be trivially lifted to the fully sparse detector since it will face the issue of Center Feature
 110 Missing, as we discussed in Sec. 1.

111 **Point-based sparse detectors** The purely point-based detectors are born to be fully sparse. PointR-
 112 CNN [24] is the pioneering work to build the purely point-based detector. 3DSSD [37] accelerates the
 113 point-based method by removing the feature propagation layer and refinement module. VoteNet [21]
 114 first makes a center voting and then generates proposals from the voted center achieving better
 115 precision. Albeit many methods have tried to accelerate the point-based method, the time-consuming
 116 neighborhood query is still unaffordable in large-scale point clouds (more than 100k points per
 117 scene). So current benchmarks [28, 2] with large-scale point clouds are dominated by voxel-based
 118 dense/semi-dense detectors [10, 27, 13].

119 3 Methodology

120 3.1 Overall Architecture

121 Following the motivation of instances as groups, we have four steps to build the fully sparse detector
 122 (FSD): 1) We first utilize a sparse voxel encoder [5, 26, 35] to extract voxel features and vote object
 123 centers (Sec. 3.2). 2) **Instance Point Grouping** groups foreground points into instances based on the
 124 voting results (Sec. 3.2). 3) Given the grouping results, **Sparse Instance Recognition (SIR)** module
 125 extracts instance/point features and generates proposals (Sec. 3.3). 4) The proposals are utilized to
 126 correct the point grouping and refine the proposals via another SIR module (Sec. 3.4).

127 3.2 Instance Point Grouping

128 **Classification and Voting** We first extract voxel features from the point cloud with a sparse voxel
 129 encoder. Although FSD is not restricted to a certain sparse voxel encoder, we utilize sparse attention
 130 block in SST [5] due to its demonstrated effectiveness. Then we build point features by concatenating
 131 voxel features and the offsets from points to their corresponding voxel centers. These point features
 132 are passed into two heads for foreground classification and center voting. The voting is similar to
 133 VoteNet [21], where the model predicts the offsets from foreground points to corresponding object
 134 centers. L1 loss [23] and Focal Loss [16] are adopted as voting loss L_{vote} and semantic classification
 135 loss L_{sem} .

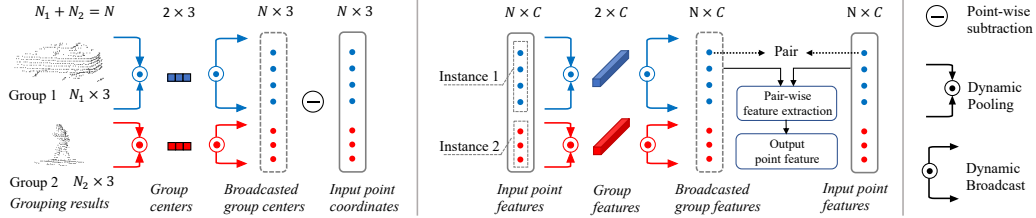


Figure 4: **Illustration of building instance-level point operators with dynamic broadcast/pooling.** Best viewed in color. Left: calculating center-to-neighbor offsets given raw point clouds. Right: updating point features. Note that the operation is parallel among all instances.

136 **Connected Components Labeling (CCL)** To group points into instances, we regard all the predicted
 137 centers (red dots in Fig. 3) as vertices in a graph. Two vertices are connected if their distance is
 138 smaller than a certain threshold. Then a connected component in this graph can be viewed as an
 139 instance, and all points voted to this connected component share a group ID. Unlike the ball query in
 140 VoteNet, our CCL-based grouping greatly avoids fragmented instances. Although there are many
 141 elaborately designed instance grouping methods [11, 31, 9], we opt for the simple CCL because it is
 142 adequate in our design and can be implemented by the efficient depth-first search.

143 3.3 Sparse Instance Recognition

144 3.3.1 Preliminaries: Dynamic Broadcast/Pooling

145 Given N points belong to M groups, we define their corresponding group ID array as I in shape of $[N,]$
 146 and their feature array as F in shape of $[N, C]$, where C is the feature dimensions. $F^{(i)}$ is the feature
 147 array of points belonging to the i -th group. Dynamic pooling aggregates each $F^{(i)}$ into one group
 148 feature g_i of shape $[C,]$. Thus we have $g_i = p(F^{(i)})$, where p is a symmetrical pooling function. The
 149 dynamic pooling on all group features G of shape $[M, C]$ is formulated as $G = p(F, I)$. The dynamic
 150 broadcast can be viewed as the inverse operation to dynamic pooling, which broadcasts g_i to all
 151 the points in the i -th group. Since the broadcasting is essentially an indexing operation, we use the
 152 indexing notation $[]$ to denote it as $G[I]$, which is in shape of $[N, C]$. Dynamic broadcast/pooling is
 153 very efficient because it can be implemented with high parallelism on modern devices and well fits
 154 the sparse data with dynamic size.

155 The prerequisite of dynamic broadcast/pooling is that each point *uniquely* belongs to a group. In
 156 other words, groups should not overlap with each other. Thanks to the motivation of instances as
 157 groups, the groups in 3D space do not overlap with each other naturally.

158 3.3.2 Formulation of Sparse Instance Recognition

159 After grouping points into instances in Sec. 3.2, we can directly extract instance features by some
 160 basic point-based networks like PointNet, DGCNN, etc. There are three elements to define a basic
 161 point-based module: *group center*, *pair-wise feature* and *group feature aggregation*.

162 **Group center** The group center is the representative point of a group. For example, in the ball
 163 query, it is the local origin of the sphere. In SIR, the group center is defined as the centroid of all
 164 voted centers in a group.

165 **Pair-wise feature** defines the input for per point feature extraction. SIR adopts two kinds of features:
 166 1) relative coordinate between group center and each point, 2) feature concatenation between group
 167 and each neighbor point. Taking feature concatenation as example and using the notations in 3.3.1,
 168 the pair-wise feature can be denoted as $CAT(F, G[I])$, where CAT is channel concatenation.

169 **Group feature aggregation** In a group, a pooling function is used to aggregate neighbor features.
 170 SIR applies dynamic pooling to aggregate feature array F . Following the notations in 3.3.1, we have
 171 $G = p(F, I)$, where G is the aggregated group features.

172 **Integration** Combining the three basic elements, we could build many variants of point-based
 173 operators, such as PointNet [20], DGCNN [32], Meta-Kernel [4], etc. Fig. 4 illustrates the basic
 174 idea of how to build an instance-level point operator with dynamic broadcast/pooling. In our design,
 175 we adopt the formulation of VFE [42] as the basic structure of SIR layers, which is basically a
 176 two-layer PointNet. In the l -th layer of SIR module, given the input point-wise feature array F_l ,

177 point coordinates array X , the voted center X' and group ID array I , the output of l -th layer can be
 178 formulated as:

$$F'_l = \text{LinNormAct}(\text{CAT}(F_l, X - p_{\text{avg}}(X', I)[I])), \quad (1)$$

$$F'_{l+1} = \text{LinNormAct}(\text{CAT}(F'_l, p_{\text{max}}(F'_l, I)[I])), \quad (2)$$

179 where LinNormAct is a fully-connected layer followed by a normalization layer [30] and an activation
 180 function [8]. The p_{avg} and the p_{max} are average-pooling and max-pooling function, respectively. The
 181 output F'_{l+1} can be further used as the input of the next SIR layer, so our SIR module is a stack of a
 182 couple of basic SIR layers.

183 3.3.3 Sparse Prediction

184 With the formulation in Eqn. 1 and Eqn. 2, SIR extracts features of all instances dynamically in
 185 parallel. And then SIR makes **sparse prediction** for all groups. In contrast to two-stage sparse
 186 prediction, our proposals (i.e., groups) do not overlap with each other. Unlike one-stage dense
 187 prediction, we only generate a single prediction for a group. Sparse prediction avoids the difficulty of
 188 label assignment in dense prediction when the center feature is missing, because there is no need to
 189 attach anchors or anchor points to non-empty voxels. It is noteworthy that the fully sparse architecture
 190 may face a severe imbalance problem: short-range objects contain much more points than long-range
 191 objects. Some methods [1, 4] use hand-crafted normalization factors to mitigate the imbalance.
 192 Instead, SIR avoids the imbalance because it only generates a single prediction for a group regardless
 193 of the number of points in the group.

194 Specifically, for each SIR layer, there is a $G_l = p_{\text{max}}(F'_l, I)$ in Eqn. 2, which can be viewed as
 195 the group features. We concatenate all G_l from each SIR layer in channel dimension and use the
 196 concatenated group features to predict bounding boxes and class labels via MLPs. All the groups
 197 whose centers fall into ground-truth boxes are positive samples. For positive samples, the regression
 198 branch predicts the offsets from group centers to ground-truth centers and object sizes and orientations.
 199 L1 loss [23] and Focal Loss [16] are adopted as regression loss L_{reg} and classification loss L_{cls} ,
 200 respectively.

201 3.4 Group Correction

202 There is inevitable incorrect grouping in the Instance Point Grouping module. For example, some
 203 foreground points may be missed, or some groups may be contaminated by background clutter. So we
 204 leverage the bounding box proposals from SIR to correct the grouping. The points inside a proposal
 205 belong to a corrected group regardless of their previous group IDs. After correction, we apply an
 206 additional SIR to these new groups. To distinguish it from the first SIR module, we denote the
 207 additional SIR module as SIR2.

208 SIR2 predicts box residual from the proposal to its corresponding ground-truth box, following many
 209 two-stage detectors. The regression loss is denoted as $L_{\text{res}} = L1(\Delta_{\text{res}}, \widehat{\Delta_{\text{res}}})$, where Δ_{res} is the
 210 ground-truth residual and $\widehat{\Delta_{\text{res}}}$ is the predicted residual. Following previous methods [26, 25], the 3D
 211 Intersection over Union (IoU) between the proposal and ground-truth serves as the soft classification
 212 label in SIR2. Specifically, the soft label q is defined as $q = \min(1, \max(0, 2IoU - 0.5))$, where IoU
 213 is the IoU between proposals and corresponding ground-truth. Then cross entropy loss is adopted to
 214 train the classification branch, denoted as L_{iou} . Taking all the loss functions in grouping (Sec. 3.2)
 215 and sparse prediction into account, we have

$$L_{\text{total}} = L_{\text{sem}} + L_{\text{vote}} + L_{\text{reg}} + L_{\text{cls}} + L_{\text{res}} + L_{\text{iou}}, \quad (3)$$

216 where we omit the normalization factors. Please refer to supplementary materials for details.

217 3.5 Discussion

218 The center voting in FSD is inspired by VoteNet [21], while FSD has two essential differences from
 219 VoteNet.

- 220 • After voting, VoteNet simply aggregates features around the voted centers without further
 221 feature extraction. Instead, FSD builds a highly efficient SIR module taking advantage of

222 dynamic broadcast/pooling for further instance-level feature extraction. Thus, FSD extracts
223 more powerful instance features than VoteNet, which is experimentally demonstrated in
224 Sec. 4.6.

- 225 • VoteNet is a typical point-based method. As we discussed in Sec. 1, it aggressively down-
226 samples the whole scene to a fixed number of points for efficiency, causing inevitable
227 information loss. Instead, the dynamic characteristic and efficiency of SIR enable fine-
228 grained point feature extraction from any number of input points without any downsampling.
229 In Sec. 4.6, we showcase the efficiency of our design in processing large-scale point clouds
230 and the benefits from fine-grained point representation.

231 4 Experiments

232 4.1 Setup

233 **Dataset: Waymo Open Dataset (WOD)** We conduct our main experiments on WOD [28]. WOD
234 is currently the largest and most trustworthy benchmark for LiDAR-based 3D object detection. WOD
235 contains 1150 sequences (more than 200K frames), 798 for training, 202 for validation, and 150 for
236 test. The detection range in WOD is 75 meters (cover area of $150m \times 150m$).

237 **Dataset: Argoverse 2 (AV2)** We further conduct long-range experiments on the recently released
238 Argoverse 2 dataset [34] to demonstrate the superiority of FSD in long-range detection. AV2 has a
239 similar scale to WOD, and it contains 1000 sequences in total, 700 for training, 150 for validation,
240 and 150 for test. In addition to *average precision* (AP), AV2 adopts a *composite score* as evaluation
241 metric, which takes both AP and localization errors into account. The perception range in AV2 is
242 200 meters (cover area of $400m \times 400m$), which is much larger than WOD. Such a large perception
243 range leads to a huge memory footprint for dense detectors.

244 **Implementation Details** We use 4 sparse regional attention blocks [5] in SST as our voxel feature
245 extractor. To demonstrate the generality of SIR, we also tried to adopt sparse convolution based
246 U-Net in PartA2 [26] as the voxel feature extractor, and the details are presented in supplementary
247 materials. Unless otherwise specified, we use SST-based sparse voxel encoder. Both the SIR module
248 and SIR2 module consist of 3 SIR layers. A SIR layer is defined by Eqn. 1 and Eqn. 2. Our model
249 converges much faster than SST, so we train our models for 6 epochs instead of the $2\times$ schedule
250 (24 epochs) in SST. All other hyper-parameters and training schemes are the same with SST. All the
251 models in our experiments use single-frame point clouds. We present more details in supplementary
252 materials.

253 4.2 Comparison to State-of-the-art Methods

254 We first compare FSD with state-of-the-art detectors and our baseline in Table 1. Since we adopt
255 the training scheme of SST and its SRA blocks as our voxel feature extractor, SST and SST_TS
256 are our baseline methods. FSD achieves the state-of-the-art performance among all the detectors in
257 Table 1 and surpasses the SST baseline by a large margin. It is noteworthy that FSD achieves exciting
258 performance on pedestrian class with single-frame point clouds.

259 4.3 Study of Treatments to Center Feature Missing

260 In what follows, we conduct experiments on WOD to elaborate the issue of **Center Feature Missing**
261 **(CFM)**. We first build several models with different characteristics. Note that all the following models
262 adopt the same voxelization resolution, so they face the same degree of CFM at the beginning.

- 263 • **FSD_{plain}**: After the sparse voxel encoder, FSD_{plain} directly predicts the box from each voxel.
264 The voxels inside ground-truth boxes are assigned *positive*. Although FSD_{plain} uses the most
265 straightforward solution for CFM, it suffers from the large variance of regression targets and
266 low-quality predictions from hard voxels.
- 267 • **SST_{center}**: It replaces the anchor-based head in SST with CenterHead [41, 39]. Based on sparse
268 voxel encoder, SST_{center} converts sparse voxels into dense feature maps and applies several

Table 1: **Performances of FSD on the Waymo Open Dataset validation split.** We mark the best result in bold. *: Two-stage SST. ‡: from [27]. †: Using sparse convolution based U-Net as sparse voxel encoder, referring to supplementary material for details.

Methods	mAP/mAPH	Vehicle 3D AP/APH		Pedestrian 3D AP/APH		Cyclist 3D AP/APH	
	L2	L1	L2	L1	L2	L1	L2
SECOND ‡ [35]	61.0/57.2	72.3/71.7	63.9/63.3	68.7/58.2	60.7/51.3	60.6/59.3	58.3/57.0
MVF [43]	-/-	62.9/-	-/-	65.3/-	-/-	-/-	-/-
AFDet [6]	-/-	63.7/-	-/-	-/-	-/-	-/-	-/-
Pillar-OD [33]	-/-	69.8/-	-/-	72.5/-	-/-	-/-	-/-
RangeDet [4]	65.0/63.2	72.9/72.3	64.0/63.6	75.9/71.9	67.6/63.9	65.7/64.4	63.3/62.1
PointPillars [12]	62.8/57.8	72.1/71.5	63.6/63.1	70.6/56.7	62.8/50.3	64.4/62.3	61.9/59.9
Voxel RCNN [3]	-/-	75.6/-	66.6/-	-/-	-/-	-/-	-/-
RCD [1]	-/-	69.0/68.5	-/-	-/-	-/-	-/-	-/-
VoTr-TSD [19]	-/-	74.9/74.3	65.9/65.3	-/-	-/-	-/-	-/-
LiDAR-RCNN [14]	65.8/61.3	76.0/75.5	68.3/67.9	71.2/58.7	63.1/51.7	68.6/66.9	66.1/64.4
Pyramid RCNN [18]	-/-	76.3/75.7	67.2/66.7	-/-	-/-	-/-	-/-
3D-MAN [38]	-/-	74.5/74.0	67.6/67.1	71.7/67.7	62.6/59.0	-/-	-/-
Part-A2-Net ‡ [26]	66.9/63.8	77.1/76.5	68.5/68.0	75.2/66.9	66.2/58.6	68.6/67.4	66.1/64.9
CenterNet-Pillar [39]	-/-	76.1/75.5	68.0/67.5	76.1/65.1	68.1/57.9	-/-	-/-
CenterPoint-Voxel [39]	69.8/67.6	76.6/76.0	68.9/68.4	79.0/73.4	71.0/65.8	72.1/71.0	69.5/68.5
IA-SSD [40]	62.3/58.1	70.5/69.7	61.6/61.0	69.4/58.5	60.3/50.7	67.7/65.3	65.0/62.7
PV-RCNN [25]	66.8/63.3	77.5/76.9	69.0/68.4	75.0/65.6	66.0/57.6	67.8/66.4	65.4/64.0
RSN [29]	-/-	75.1/74.6	66.0/65.5	77.8/72.7	68.3/63.7	-/-	-/-
PV-RCNN++ [27]	68.4/64.9	78.8/78.2	70.3/69.7	76.7/67.2	68.5/59.7	69.0/67.6	66.5/65.2
SST_TS* [5]	-/-	76.2/75.8	68.0/67.6	81.4/74.0	72.8/65.9	-/-	-/-
SST [5]	67.8/64.6	74.2/73.8	65.5/65.1	78.7/69.6	70.0/61.7	70.7/69.6	68.0/66.9
FSD _{sparse} (ours) †	70.7/68.4	76.7/76.3	67.5/67.1	82.8/77.2	73.7/68.6	73.1/71.9	70.8/69.5
FSD (ours)	71.7/69.4	77.3/76.9	69.8/69.3	83.3/77.7	74.4/69.3	73.2/71.9	70.8/69.6

269 convolutions to diffuse features to the empty object centers as in Fig. 1. Then it makes predictions
 270 from the diffused center feature.

- 271 • **FSD_{nogc}**: It removes the group correction and SIR2 module in FSD.
- 272 • **CenterPoint-PP**: It does not resort to any sparse voxel encoders. Instead, it applies multiple
 273 dense convolutions soon after voxelization for feature diffusion, greatly eliminating CFM. It also is
 274 equipped with CenterHead avoiding large variance of regression targets.

275 Experiments and analyses

276 There is usually a quite large
 277 unoccupied area around the centers
 278 of large vehicles. Thus the
 279 performance of large vehicles
 280 is an appropriate indicator that
 281 reveals the effect of CFM. So we
 282 build a customized evaluation
 283 tool, which breaks down the ob-
 284 ject length following the COCO
 285 evaluation [15]. Then we use it
 286 to evaluate the performance of
 287 vehicles with different lengths. Table 2 shows the results, and we list our findings as follows.

Table 2: **Vehicle detection with vehicle length breakdown.** †: re-implemented by ourselves. *: official Waymo L2 overall metric. ‡: use CenterHead [39] instead of anchor-based head. Arrows indicate the performance changes from SST_{center}.

Methods	Vehicle length (m)				Official*
	[0, 4)	[4, 8)	[8, 12)	[12, +∞)	
CenterPoint-PP†	34.3	69.3	42.0	43.6	66.2
FSD _{plain}	32.2	64.6	41.3	42.2	62.3
SST _{center} ‡ [5]	36.0	69.4	33.7	30.5	66.3
FSD _{nogc}	33.5 ↓ 2.5	68.2 ↓ 1.2	47.7 ↑ 14.0	47.9 ↑ 17.4	65.2 ↓ 1.1
FSD	36.7 ↑ 0.7	71.0 ↑ 1.6	51.3 ↑ 17.6	53.7 ↑ 23.2	69.3 ↑ 3.0

- 288 • Comparing FSD_{plain} with SST_{center}, they share the same attention-based sparse voxel encoder.
 289 However, the trend is totally opposite w.r.t vehicle size. With feature diffusion, SST_{center} attains
 290 much worse performance than FSD_{plain} on large vehicles. It suggests feature diffusion is a sub-
 291 optimal solution for CFM in the case of large objects. For those large objects, the features may not
 292 be diffused to the centers or the diffused features are too weak to make accurate predictions.
- 293 • However, FSD_{plain} obtains the worst performance among all detectors on vehicles with normal
 294 sizes. Note that the CFM issue is minor for the normal size vehicles. So, in this case, the center-
 295 based assignment in SST_{center} shows its superiority to the assignment in FSD_{plain}. It suggests the
 296 solution for CFM in FSD_{plain} is also sub-optimal, even if it achieves better performance in large
 297 objects.

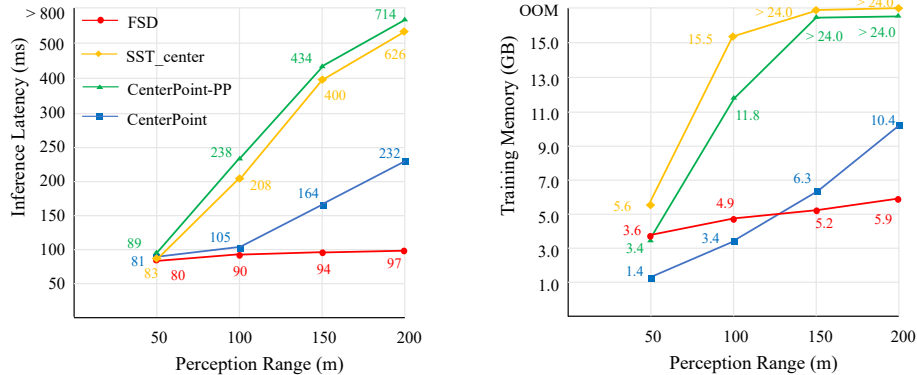


Figure 5: **Memory footprints and inference latency in different perception ranges.** All the statistics are obtained on a single 3090 GPU with batch size 1. Inference latency is evaluated by the standard benchmark script in MMDetection3D without any test-time optimization. CenterPoint-PP and SST_{center} are defined in Sec. 4.3. Best viewed in color.

- 298 • Comparing FSD_{nogc} with SST_{center} , they share the same sparse voxel encoder while FSD_{nogc}
 299 replaces the dense part in SST_{center} with SIR. The huge improvements of FSD_{nogc} on large
 300 vehicles fairly reveal that SIR effectively resolves CFM and is better than feature diffusion.
- 301 • CenterPoint-PP suffers much less from CFM because it leverages dense feature maps from very
 302 beginning of the network. It is also equipped with the advanced center-based assignment. Even so,
 303 FSD_{nogc} and FSD still outperform CenterPoint-PP, especially on large vehicles.

Table 3: **Performance in Argoverse 2 validation split.** †: provided by authors of AV2 dataset. *: re-implemented by ourselves. C-Barrel: construction barrel. MPC-Sign: mobile pedestrian crossing sign. A-Bus: articulated bus. C-Cone: construction cone. V-Trailer: vehicular trailer. We omit the results of dog, wheelchair and message board trailer because these categories contain very few instances. The average results take all categories into account, including the omitted categories. We mark the categories attaining notable improvements in **bold**.

Methods	Average	Vehicle	Bus	Pedestrian	Stop Sign	Box Truck	Bollard	C-Barrel	Motorcyclist	MPC-Sign	Motorcycle	Bicycle	A-Bus	School Bus	Truck Cab	C-Cone	V-Trailer	Sign	Large Vehicle	Stroller	Bicyclist
<i>Precision</i>																					
CenterPoint† [39]	13.5	61.0	36.0	33.0	28.0	26.0	25.0	22.5	16.0	16.0	12.5	9.5	8.5	7.5	8.0	8.0	7.0	6.5	3.0	2.0	14
CenterPoint*	22.0	67.6	38.9	46.5	16.9	37.4	40.1	32.2	28.6	27.4	33.4	24.5	8.7	25.8	22.6	29.5	22.4	6.3	3.9	0.5	20.1
FSD (Ours)	24.0	67.1	39.8	57.4	21.3	38.3	38.3	38.1	30.0	23.6	38.1	25.5	15.6	30.0	20.1	38.9	23.9	7.9	5.1	5.7	27.0
<i>Composite Score</i>																					
CenterPoint*	17.6	57.2	32.0	35.7	13.2	31.0	28.9	25.6	22.2	19.1	28.2	19.6	6.8	22.5	17.4	22.4	17.2	4.8	3.0	0.4	16.7
FSD (Ours)	19.1	56.0	33.0	45.7	16.7	31.6	27.7	30.4	23.8	16.4	31.9	20.5	12.0	25.6	15.9	29.2	18.1	6.4	3.8	4.5	22.1

304 4.4 Long-range Detection

305 Several widely adopted 3D detection benchmarks [28, 7, 2] have relatively short perception range. To
 306 unleash the potential of FSD, we conduct long-range detection experiments on the recently released
 307 Argoverse 2 dataset (AV2), with a perception range of 200 meters. In addition, AV2 contains objects
 308 in 30 classes, facing the challenging long-tail issue.

309 **Main results** We first list the main results of FSD on AV2 in Table 3. The authors of AV2 provide
 310 a baseline CenterPoint model, but the results are mediocre. To make a fair comparison, we re-
 311 implement a stronger CenterPoint model on the AV2 dataset. The re-implemented CenterPoint adopts
 312 the same training scheme with FSD, including ground-truth sampling to alleviate the long-tail issue.
 313 Please refer to supplementary materials for more details. FSD outperforms CenterPoint in the average
 314 metric. It is noteworthy that FSD significantly outperforms CenterPoint in some tiny objects (e.g.,
 315 Pedestrian, Construction Cone) as well as some objects with extremely large sizes (e.g., Articulated
 316 Bus, School Bus). We owe this to the virtue of instance-level fine-grained feature extraction in SIR.

317 **Range Scaling** To demonstrate the efficiency of FSD in long-range detection, we depict the trend
 318 of training memory and inference latency of three detectors when the perception range increases

319 in Fig. 5. Fig. 5 shows dramatic latency/memory increase when applying dense detectors to larger
 320 perception ranges. Designed to be fully sparse, the resource needed for FSD is roughly linear to the
 321 number of input points, so its memory and latency only slightly increase as the perception range
 322 extends.

323 4.5 More Sparse Scenes

324 Argoverse 2 dataset provides a highly reliable HD map, which could be utilized as a prior to
 325 remove uninterested regions making the scene more sparse. Thus we proceed with experiments
 326 removing some uninterested regions to show the advantages of FSD in more sparse scenarios.
 327 The results are summarized in Table 4. FSD has a significantly lower memory footprint and latency
 328 with an acceptable precision loss after removing the uninterested regions. On the contrary, the
 329 efficiency improvement of CenterPoint is minor. It reveals that FSD benefits more from the increase
 330 of data sparsity, which is another advantage of the fully sparse architecture.
 331
 332
 333
 334

Table 4: **Performance with different detection areas.** †: Region of Interest is defined by the HD map in AV2 dataset. *: mean AP of all classes.

	FSD			CenterPoint		
	Mem.	Latency(ms)	mAP*	Mem.	Latency(ms)	mAP*
all	5.9	97	24.0	10.4	232	22.0
only RoI†	3.2↓ 45.8%	81↓ 16.5%	23.2	9.9↓ 4.8%	227↓ 2.2%	21.5
w/o ground	2.3↓ 61.0%	74↓ 25.8%	21.0	9.7↓ 6.7%	217↓ 6.4%	19.8

335 4.6 Ablation Study

336 **Effectiveness of Components** In addition to FSD_{plain} and FSD_{nogc} (Sec. 4.3), we also de-
 337 grade FSD to FSD_{agg} to understand the mechanism of FSD. In FSD_{agg}, we aggregate grouped
 338 point features by dynamic pooling and then directly make predictions from the pooled features,
 339 after Instance Point Grouping. FSD_{agg} is similar to the way in VoteNet [21] as we discussed
 340 in Sec. 3.5. Thus, FSD_{agg} can explicitly leverage instance-level features other than the point-level
 341 features in FSD_{plain}. However, FSD_{agg} can not take advantage of further point feature extraction in
 342 SIR. As can be seen in Table 5, the improvement is limited if we only apply grouping without SIR.
 343 The combination of grouping and SIR attain notable improvements.
 344
 345
 346
 347

Table 5: **Ablation of design factors in SIR.** Performances are evaluated on Waymo validation split.

	Grouping	SIR	Group Correction	L2 3D APH		
				Vehicle	Pedestrian	Cyclist
FSD _{plain}				62.29	64.31	64.49
FSD _{agg}	✓			63.13	65.13	64.52
FSD _{nogc}		✓		65.20	67.39	67.78
FSD	✓	✓	✓	69.30	69.30	69.60

348 **Downsampling in SIR** The efficiency of SIR makes it feasible to extract fine-grained point
 349 features without any point downsampling. This is another notable difference between FSD and
 350 VoteNet. To demonstrate the superiority, we apply voxelization on the raw points before
 351 SIR module and treat the centroids of voxels as downsampled points. We conduct experiments
 352 on AV2 dataset because it contains a couple of categories in a tiny size, which may be sensitive to downsampling. As expected, small objects have
 353 notable performance loss when adopting downsampling, and we list some of them in Table 6. We
 354 also evaluate the inference latency of the SIR module on 3090Ti GPU. As can be seen, compared
 355 with the overall latency (97ms, Fig. 5), the SIR module is highly efficient.
 356
 357
 358
 359
 360

Table 6: **Performances with different representation granularity.** †: Latency of SIR module.

Voxel size	AP				Latency (ms)†
	CC	Bollard	Bicyclist	Stop Sign	
30cm	35.4	36.5	24.6	18.3	3.5
20cm	37.3	37.3	26.4	20.0	4.1
10cm	38.9	38.3	27.0	21.3	4.5
Point	39.3	38.6	27.1	21.5	6.3

361 5 Conclusion

362 This paper proposes FSD, a fully sparse 3D object detector, aiming for efficient long-range object
 363 detection. FSD utilizes a highly efficient point-based Sparse Instance Recognition module to solve the
 364 center feature missing in fully sparse architecture. FSD achieves not only competitive performance
 365 on the widely-used Waymo Open Dataset, but also state-of-the-art performance in the long-range
 366 Argoverse 2 dataset with a much faster inference speed than previous detectors.

367 **Limitation** A more elaborately designed grouping strategy may help with performance improve-
 368 ments. However, it is beyond our design goal in this paper, and we will pursue it in future work.

369 **Societal Impacts** Our method can be deployed in the autonomous driving system. Performance
 370 loss caused by improper usage may increase security risks.

371 References

- 372 [1] Alex Bewley, Pei Sun, Thomas Mensink, Dragomir Anguelov, and Cristian Sminchisescu. Range Conditioned Dilated Convolutions for Scale Invariant 3D Object Detection. *arXiv preprint arXiv:2005.09927*,
373 2020.
374
- 375 [2] Holger Caesar, Varun Bankiti, Alex H Lang, Sourabh Vora, Venice Erin Liong, Qiang Xu, Anush Krishnan,
376 Yu Pan, Giancarlo Baldan, and Oscar Beijbom. nuscenes: A multimodal dataset for autonomous driving.
377 In *CVPR*, 2020.
- 378 [3] Jiajun Deng, Shaoshuai Shi, Peiwei Li, Wengang Zhou, Yanyong Zhang, and Houqiang Li. Voxel R-CNN:
379 Towards High Performance Voxel-based 3D Object Detection. In *AAAI*, 2021.
- 380 [4] Lue Fan, Xuan Xiong, Feng Wang, Naiyan Wang, and ZhaoXiang Zhang. RangeDet: In Defense of Range
381 View for LiDAR-Based 3D Object Detection. In *ICCV*, 2021.
- 382 [5] Lue Fan, Ziqi Pang, Tianyuan Zhang, Yu-Xiong Wang, Hang Zhao, Feng Wang, Naiyan Wang, and
383 Zhaoxiang Zhang. Embracing Single Stride 3D Object Detector with Sparse Transformer. In *CVPR*, 2022.
- 384 [6] Runzhou Ge, Zhuangzhuang Ding, Yihan Hu, Yu Wang, Sijia Chen, Li Huang, and Yuan Li. AFDet:
385 Anchor Free One Stage 3D Object Detection. *arXiv preprint arXiv:2006.12671*, 2020.
- 386 [7] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? the kitti vision
387 benchmark suite. In *2012 IEEE conference on computer vision and pattern recognition*. IEEE, 2012.
- 388 [8] Dan Hendrycks and Kevin Gimpel. Gaussian Error Linear Units (GELUs). *arXiv preprint*
389 *arXiv:1606.08415*, 2016.
- 390 [9] Fangzhou Hong, Hui Zhou, Xinge Zhu, Hongsheng Li, and Ziwei Liu. Lidar-based panoptic segmentation
391 via dynamic shifting network. In *CVPR*, 2021.
- 392 [10] Yihan Hu, Zhuangzhuang Ding, Runzhou Ge, Wenxin Shao, Li Huang, Kun Li, and Qiang Liu. AFDetV2:
393 Rethinking the Necessity of the Second Stage for Object Detection from Point Clouds. *arXiv preprint*
394 *arXiv:2112.09205*, 2021.
- 395 [11] Li Jiang, Hengshuang Zhao, Shaoshuai Shi, Shu Liu, Chi-Wing Fu, and Jiaya Jia. Pointgroup: Dual-set
396 point grouping for 3d instance segmentation. In *CVPR*, 2020.
- 397 [12] Alex H Lang, Sourabh Vora, Holger Caesar, Lubing Zhou, Jiong Yang, and Oscar Beijbom. PointPillars:
398 Fast Encoders for Object Detection from Point Clouds. In *CVPR*, 2019.
- 399 [13] Yingwei Li, Adams Wei Yu, Tianjian Meng, Ben Caine, Jiquan Ngiam, Daiyi Peng, Junyang Shen, Bo Wu,
400 Yifeng Lu, Denny Zhou, et al. DeepFusion: Lidar-Camera Deep Fusion for Multi-Modal 3D Object
401 Detection. In *CVPR*, 2022.
- 402 [14] Zhichao Li, Feng Wang, and Naiyan Wang. LiDAR R-CNN: An Efficient and Universal 3D Object
403 Detector. In *CVPR*, 2021.
- 404 [15] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár,
405 and C Lawrence Zitnick. Microsoft COCO: Common Objects in Context. In *European conference on*
406 *computer vision*. Springer, 2014.
- 407 [16] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal Loss for Dense Object
408 Detection. In *ICCV*, 2017.
- 409 [17] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and
410 Alexander C Berg. SSD: Single Shot Multibox Detector. In *ECCV*, 2016.
- 411 [18] Jiageng Mao, Minzhe Niu, Haoyue Bai, Xiaodan Liang, Hang Xu, and Chunjing Xu. Pyramid R-CNN:
412 Towards Better Performance and Adaptability for 3D Object Detection. In *ICCV*, 2021.
- 413 [19] Jiageng Mao, Yujing Xue, Minzhe Niu, Haoyue Bai, Jiashi Feng, Xiaodan Liang, Hang Xu, and Chunjing
414 Xu. Voxel Transformer for 3D Object Detection. In *ICCV*, 2021.
- 415 [20] Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. PointNet: Deep Learning on Point Sets for
416 3D Classification and Segmentation. In *CVPR*, 2017.
- 417 [21] Charles R Qi, Or Litany, Kaiming He, and Leonidas J Guibas. Deep Hough Voting for 3D Object Detection
418 in Point Clouds. In *ICCV*, 2019.

- 419 [22] Charles Ruizhongtai Qi, Li Yi, Hao Su, and Leonidas J Guibas. PointNet++: Deep Hierarchical Feature
420 Learning on Point Sets in a Metric Space. In *NeurIPS*, 2017.
- 421 [23] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster R-CNN: Towards Real-time Object
422 Detection with Region Proposal Networks. *NeurIPS*, 28, 2015.
- 423 [24] Shaoshuai Shi, Xiaogang Wang, and Hongsheng Li. PointRCNN: 3D Object Proposal Generation and
424 Detection from Point Cloud. In *CVPR*, 2019.
- 425 [25] Shaoshuai Shi, Chaoxu Guo, Li Jiang, Zhe Wang, Jianping Shi, Xiaogang Wang, and Hongsheng Li.
426 PV-RCNN: Point-Voxel Feature Set Abstraction for 3D Object Detection. In *CVPR*, 2020.
- 427 [26] Shaoshuai Shi, Zhe Wang, Jianping Shi, Xiaogang Wang, and Hongsheng Li. From Points to Parts: 3D
428 Object Detection from Point Cloud with Part-aware and Part-aggregation Network. *IEEE Transactions on*
429 *Pattern Analysis and Machine Intelligence*, 2020.
- 430 [27] Shaoshuai Shi, Li Jiang, Jiajun Deng, Zhe Wang, Chaoxu Guo, Jianping Shi, Xiaogang Wang, and
431 Hongsheng Li. PV-RCNN++: Point-Voxel Feature Set Abstraction With Local Vector Representation for
432 3D Object Detection. *arXiv preprint arXiv:2102.00463*, 2021.
- 433 [28] Pei Sun, Henrik Kretschmar, Xerxes Dotiwalla, Aurelien Chouard, Vijaysai Patnaik, Paul Tsui, James
434 Guo, Yin Zhou, Yuning Chai, Benjamin Caine, et al. Scalability in Perception for Autonomous Driving:
435 Waymo Open Dataset. In *CVPR*, 2020.
- 436 [29] Pei Sun, Weiyue Wang, Yuning Chai, Gamaleldin Elsayed, Alex Bewley, Xiao Zhang, Cristian Smin-
437 chisescu, and Dragomir Anguelov. RSN: Range Sparse Net for Efficient, Accurate LiDAR 3D Object
438 Detection. In *CVPR*, 2021.
- 439 [30] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz
440 Kaiser, and Illia Polosukhin. Attention Is All You Need. In *NeurIPS*, 2017.
- 441 [31] Weiyue Wang, Ronald Yu, Qiangui Huang, and Ulrich Neumann. SGP: Similarity Group Proposal
442 Network for 3d Point Cloud Instance Segmentation. In *CVPR*, 2018.
- 443 [32] Yue Wang, Yongbin Sun, Ziwei Liu, Sanjay E Sarma, Michael M Bronstein, and Justin M Solomon.
444 Dynamic graph cnn for learning on point clouds. *Acm Transactions On Graphics (tog)*, 38(5):1–12, 2019.
- 445 [33] Yue Wang, Alireza Fathi, Abhijit Kundu, David Ross, Caroline Pantofaru, Tom Funkhouser, and Justin
446 Solomon. Pillar-based Object Detection for Autonomous Driving. In *ECCV*, 2020.
- 447 [34] Benjamin Wilson, William Qi, Tanmay Agarwal, John Lambert, Jagjeet Singh, Siddhesh Khandelwal,
448 Bowen Pan, Ratnesh Kumar, Andrew Hartnett, Jhony Kaesemodel Pontes, Deva Ramanan, Peter Carr,
449 and James Hays. Argoverse 2: Next Generation Datasets for Self-Driving Perception and Forecasting. In
450 *NeurIPS Datasets and Benchmarks 2021*, 2021.
- 451 [35] Yan Yan, Yuxing Mao, and Bo Li. SECOND: Sparsely Embedded Convolutional Detection. *Sensors*, 18
452 (10), 2018.
- 453 [36] Bin Yang, Wenjie Luo, and Raquel Urtasun. PIXOR: Real-time 3D Object Detection from Point Clouds.
454 In *CVPR*, 2018.
- 455 [37] Zetong Yang, Yanan Sun, Shu Liu, and Jiaya Jia. 3DSSD: Point-based 3D Single Stage Object Detector.
456 In *CVPR*, 2020.
- 457 [38] Zetong Yang, Yin Zhou, Zhifeng Chen, and Jiquan Ngiam. 3D-MAN: 3D Multi-Frame Attention Network
458 for Object Detection. In *CVPR*, 2021.
- 459 [39] Tianwei Yin, Xingyi Zhou, and Philipp Krähenbühl. Center-based 3D Object Detection and Tracking.
460 *arXiv preprint arXiv:2006.11275*, 2020.
- 461 [40] Yifan Zhang, Qingyong Hu, Guoquan Xu, Yanxin Ma, Jianwei Wan, and Yulan Guo. Not All Points Are
462 Equal: Learning Highly Efficient Point-based Detectors for 3D LiDAR Point Clouds. In *CVPR*, 2022.
- 463 [41] Xingyi Zhou, Dequan Wang, and Philipp Krähenbühl. Objects as Points. *arXiv preprint arXiv:1904.07850*,
464 2019.
- 465 [42] Yin Zhou and Oncel Tuzel. VoxelNet: End-to-End Learning for Point Cloud Based 3D Object Detection.
466 In *CVPR*, 2018.
- 467 [43] Yin Zhou, Pei Sun, Yu Zhang, Dragomir Anguelov, Jiyang Gao, Tom Ouyang, James Guo, Jiquan Ngiam,
468 and Vijay Vasudevan. End-to-End Multi-View Fusion for 3D Object Detection in LiDAR Point Clouds. In
469 *CoRL*, 2020.

470 **Checklist**

- 471 1. For all authors...
- 472 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
473 contributions and scope? [Yes]
- 474 (b) Did you describe the limitations of your work? [Yes]
- 475 (c) Did you discuss any potential negative societal impacts of your work? [Yes]
- 476 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
477 them? [Yes]
- 478 2. If you are including theoretical results...
- 479 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 480 (b) Did you include complete proofs of all theoretical results? [N/A]
- 481 3. If you ran experiments...
- 482 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
483 mental results (either in the supplemental material or as a URL)? [No] The code will
484 be anonymously public soon, after refactoring.
- 485 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
486 were chosen)? [Yes]
- 487 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
488 ments multiple times)? [Yes]
- 489 (d) Did you include the total amount of compute and the type of resources used (e.g., type
490 of GPUs, internal cluster, or cloud provider)? [Yes]
- 491 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 492 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 493 (b) Did you mention the license of the assets? [Yes]
- 494 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
495 We are not releasing new assets.
- 496 (d) Did you discuss whether and how consent was obtained from people whose data you're
497 using/curating? [Yes]
- 498 (e) Did you discuss whether the data you are using/curating contains personally identifiable
499 information or offensive content? [Yes]
- 500 5. If you used crowdsourcing or conducted research with human subjects...
- 501 (a) Did you include the full text of instructions given to participants and screenshots, if
502 applicable? [N/A]
- 503 (b) Did you describe any potential participant risks, with links to Institutional Review
504 Board (IRB) approvals, if applicable? [N/A]
- 505 (c) Did you include the estimated hourly wage paid to participants and the total amount
506 spent on participant compensation? [N/A]