

# Shape-Aware Monocular 3D Object Detection

Wei Chen, Jie Zhao, Wan-Lei Zhao\*, Song-Yuan Wu

**Abstract**—The detection of 3D objects through a single perspective camera is a challenging issue. The anchor-free and keypoint-based models receive increasing attention recently due to their effectiveness and simplicity. However, most of these methods are vulnerable to the occlusion and truncation of objects. In this paper, a single-stage monocular 3D object detection model is proposed. An instance-segmentation head is integrated into the model training, which allows the model to be aware of the visible shape of a target object. Therefore, the detection largely avoids interference from irrelevant regions surrounding the target objects. In addition, we also reveal that the popular IoU-based evaluation metrics, which were originally designed for evaluating stereo or LiDAR-based detection methods, are insensitive to the improvement achieved by the monocular 3D object detection algorithms. A novel evaluation metric, namely average depth similarity (ADS) is proposed for the monocular 3D object detection models. Our method outperforms the comparison baseline in terms of both the popular and the proposed evaluation metrics while maintaining real-time efficiency.

**Index Terms**—Monocular 3D object detection, 3D object detection, computer vision, autonomous driving.

## I. INTRODUCTION

The three-dimensional visual object detector is a core component of an autonomous driving system. The quality of the detection results has a direct impact on the subsequent tasks, such as tracking and driving planning. Recently, we have witnessed the great advances in the design of 3D object detection models. Most of them [1]–[3] rely on 3D LiDAR laser scanners. Although LiDAR point clouds allow the detectors to achieve accurate 3D scene localization, the hardware is too expensive to be equipped with normal cars. Stereo camera rigs are the alternative option in many 3D object detection models [4]–[6], where depth information is estimated by building accurate pixel correspondences between the left and the right cameras. However, binocular methods are usually computationally expensive and come with higher bandwidth costs. The stringent conditions mentioned above hinder the practical application of these methods. To circumvent these difficulties, increasing efforts have been made in developing cheaper 3D object detectors. They aim at detecting the 3D pose of objects from a single 3-channel RGB image [7]–[11] that has been captured by a monocular camera.

Despite recent progress, the performance gap between monocular 3D object detectors and stereo-based or LiDAR-based detectors remains wide. Compared with stereo-based or

LiDAR-based detectors, it is an ill-posed problem that recovers the depth information from a single image captured by a perspective camera. Due to the irreversible information loss introduced by the projection process of a perspective camera, the depth we can recover from a single image is theoretically very limited. Instead of directly regressing the object depth [10], [12], [13], most recent works [7], [9], [11] recover the depth of objects by modeling the relationship between the 3D geometric prior and the detected 2D keypoints. For instance, one can infer the depth of a target object by estimating the height of the 3D object and the height of the 2D projection, since the height of the 2D projection are inversely proportional to the depth when the 3D height is constant. These geometric constraint-based methods [7], [9], [11] usually show more strong generalization ability than the one that regresses the depth of objects directly [10], [12], [13].

The depth estimation is particularly challenging when the object is occluded or truncated in the 2D view, as the surrounding 2D keypoints of the object cannot be fully recovered. Objects are only partially visible due to either the occlusion by other objects or the truncation due to the limited FoV (Field of View) of cameras. In this case, it is difficult to provide an accurate prediction for an object since its shape is largely missing. On the one hand, information that can be extracted from the partially visible object is very limited. In addition, the available occluded training samples are over-dominated by those fully-visible ones. The training model is easily overfit on the occluded objects. On the other hand, most of the existing object detection methods are center-based, for which an object center (instead of a 3D bounding box) is assigned to a detected object. However, for occluded/truncated objects, their center points often lie on other irrelevant objects. In such cases, the features derived from the center point and its neighboring area are mixed with irrelevant contents.

In this paper, a novel auxiliary training task is designed to alleviate the over-fitting risk that is latent in the training process. Namely, an additional instance segmentation branch is added to the model. By learning to predict the visibility mask of objects, our model becomes “shape-aware”. The feature is therefore derived from a real target object region, which prevents the model from overfitting to the noisy context surrounding the occluded objects.

However, the training task of instance segmentation requires pixel-level annotation on the monocular image, which is a laborious task to undertake. To overcome this difficulty, the pixel-level object annotation is approximated by utilizing the 3D object-level annotation and the corresponding point cloud. We simply project the 3D points inside the 3D bounding box onto the image plane to form an instance mask used for training. Additionally, an uncertainty-weighted loss function is adopted to learn useful information from these noisy and

Wei Chen and Wan-Lei Zhao are with the Department of Computer Science and Technology, Xiamen University, Xiamen 361005, China (e-mail: chenwei128@stu.xmu.edu.cn; wlzhao@xmu.edu.cn).

Jie Zhao and Song-Yuan Wu are with the Ningbo Boden AI Technology Co., Ltd., Ningbo 315048, China (e-mail: jie.zhao@bodenai.com; songyuan.wu@bodenai.com).

This work has been submitted to the IEEE for possible publication. Copyright may be transferred without notice, after which this version may no longer be accessible.

sparse generated instance segmentation annotations.

In our paper, we also reveal the limitations of the current evaluation measure for 3D object detection. With the current IoU-based evaluation metrics, only the overlapping with the ground-truth objects is counted. The performance can be boosted by simply duplicating detected 3D bounding boxes and putting in different positions and layouts, which increases the chance of overlapping with the ground-truth while making no improvement in the 3D object detection. To address this issue, a new measure called average depth similarity (ADS) is proposed to evaluate the overlaps with the ground-truth in the object depth. In combination with 3D IoU-based metrics, a more objective evaluation can be made for 3D object detection methods that are based on the monocular image.

## II. RELATED WORK

In this section, we briefly review the 3D object detection methods based on 3D point clouds and binocular images first. Then the review focuses on the methods built upon a monocular image since our method also falls in this category.

In the representative LiDAR-based 3D object detection method [1], 3D point clouds are encoded as 2D feature maps by projecting them into the bird-eye-view (BEV), then a 2D CNN is applied for detection. In VoxelNet [2], the neighboring cloud points are grouped into voxels. Thereafter, PointNet [14] is adopted in each voxel for feature encoding. The 3D convolution, which is used in VoxelNet for object detection, comes with an expensive computational cost. SECOND [3] improves the inference speed of VoxelNet by introducing sparse convolution based on the fact that many voxels are empty due to the point cloud sparsity. LiDAR-based detectors are able to provide accurate 3D detection results, however, the expensive sensor cost hinders the application of these methods.

Besides 3D point clouds, depth information can also be recovered from stereo cameras. Stereo R-CNN [4] extends the standard Faster R-CNN to detect and match the left and right objects, then refines the 3D detection with dense RoI (Region of Interest) alignment. In [6], pseudo-LiDAR points extracted from the disparity map are directly used as the input for 3D object detectors, achieving state-of-the-art performance. Stereo camera-based methods usually require dense disparity estimation on the full image or each RoI, which is computationally expensive. They also introduce extra bandwidth and calibration requirements compared to monocular image-based methods.

Different from above methods, monocular based 3D object detectors only take a single RGB image as input. Amongst these methods, recent models are mostly center-based [7], [8], [11], [12], where each object is represented with a representative center  $\mathbf{x}_r = (u_c, v_c)$ . Depth, dimensions, and orientation are estimated as attributes of the detected center point. Theoretically speaking, depth estimation from a monocular image is ill-posed. There are various ways to undertake depth estimation. They can be roughly divided into three categories, namely the direct depth regression, depth-map-based, and geometric constraint-based methods. Direct depth regression methods [10], [12], [13] regress the depth as the 3D attributes of the objects directly. These direct

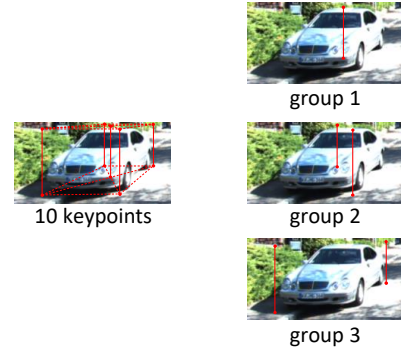


Fig. 1: Five vertical lines can be extracted from 10 keypoints. These lines are divided into three groups, each of which can produce an independent depth prediction.

depth regression methods are efficient since they do not rely on any additional depth-map estimation modules. The depth map-based methods [6], [15] are also known as the pseudo-LiDAR (PL) based method, employs an off-the-shelf dense depth map estimation Convolutional Neural Network (CNN) as a substitute for LiDAR sensors. With the resulting pseudo-LiDAR, the detection task is usually achieved by applying a state-of-the-art LiDAR-based detector. Due to biased training protocol [15], superior performance for PL-based methods is only observed on the validation set. According to [15], the training set adopted by their depth estimation modules heavily overlaps with the *val* set of the PL-based detectors.

Apart from the above two types of depth estimation, the task can also be implemented by leveraging 3D-2D geometric constraints. RTM3D [9] estimates nine predefined keypoints which introduce 18 projection constraints (9 pairs of  $x$  and  $y$ ) for each object. Then the depth of each object can be recovered by solving a nonlinear least-squares optimization problem. Due to the variety of car layouts, the estimated  $x$ -coordinate is usually less reliable than the  $y$ -coordinate for one keypoint. Therefore, only the geometric constraint on the height of keypoints is utilized in Monoflex [7]. In addition, it also divides 10 keypoints into three groups (shown in Figure 1) to produce three independent depth predictions. The final depth is a weighted average of all the depth predictions.

In the above geometric-constraint-based methods, objects are modeled as 3D bounding boxes, and keypoints are often defined as corners, centers, or top/bottom points of the bounding boxes. These keypoints mostly lie on the background, *e.g.*, the sky or the ground. Deep MANTA [16] and AutoShape [11] address this issue by annotating the keypoints of the training data with CAD templates via a semi-autonomous/autonomous approach. However, due to the mismatch between the predefined CAD model and the real object, the improvement that these methods achieve is very limited. In this paper, an auxiliary training task is introduced. Namely, instance segmentation is incorporated as another task for the detection network. As the task itself is shape-aware, the detected object regions are expected to shrink from a bounding box to the real object area. Moreover, since the annotations

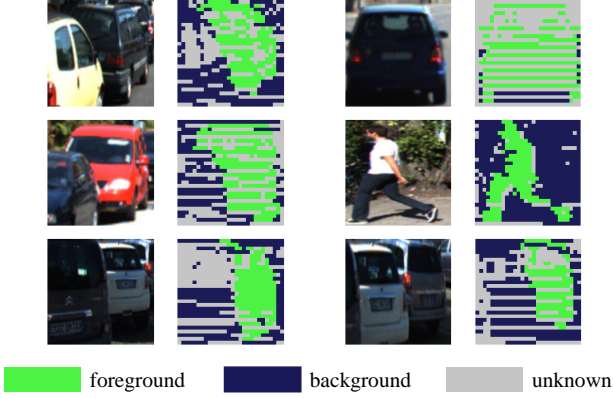


Fig. 2: Instance level masks generated by projecting LiDAR points onto the image plane and treating points inside the 3D bounding box as foreground.

on the 3D point cloud are wisely utilized to produce the segmentation mask, the introduced segmentation head boosts the overall detection performance while no laborious pixel-level annotation is required.

### III. SHAPE-AWARE 3D OBJECT DETECTION

As aforementioned, the existing monocular 3D object detection models easily get over-fit on the occluded objects. First of all, in the monocular detectors, objects are modeled as representative centers. However, for the occluded objects, their center may lie on other objects, leading to the polluted feature representation that is derived from the center. Moreover, it is challenging for the network to recover the 3D properties of an object when only a small area of the target object is visible. Furthermore, very few samples are available for the severely occluded objects in the training set.

#### A. Shape-aware Auxiliary Training Task

To address the above issue, we aim to train a model that is able to estimate the shape of a target object. Given the precise shape of a detected object is known, the detector would be able to focus on the real visible areas of an object. Therefore, more precise detection of occluded objects could be expected.

Attaching to the backbone network, a center-based segmentation branch is added, which is parallel to the detection head. Since the object detection pipeline is also center-based, these two tasks are compatible to each other. As illustrated in Figure 3, we extract a feature vector  $\mathbf{f}_i \in \mathbb{R}^{64}$  for the  $i$ th object, then it's added with a predefined position embedding map  $\mathbf{f}_{pos} \in \mathbb{R}^{s,s,64}$  with broadcasting. A feature map with size  $s \times s \times 64$  is, therefore, produced. The position embedding map  $\mathbf{f}_{pos}$  is a predefined image with two channels. Each pixel at position  $(x_p, y_p)$  of  $\mathbf{f}_{pos}$  equals to a linear transformed coordinate  $(\tilde{x}_p, \tilde{y}_p) = (\frac{x_p}{s} - 0.5, \frac{y_p}{s} - 0.5)$ .

In order to train this center-based segmentation branch, the instance masks must be ready for the training images. However, pixel-level annotation is expensive. To circumvent

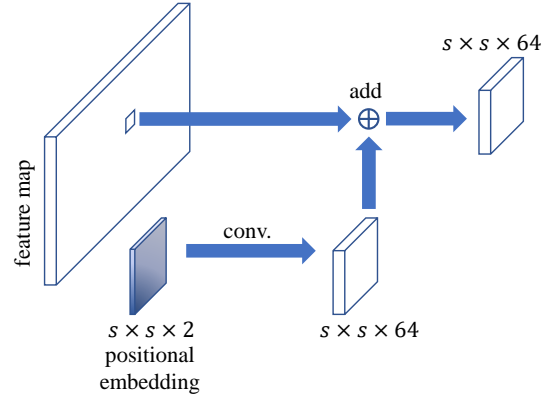


Fig. 3: The feature expanding module. In this module, the feature vector extracted from the representative center is expanded to a square feature map. This feature map is further fed into subsequent CNN modules to produce the final mask prediction.

this difficulty, we alternatively extract the instance masks from the 3D bounding boxes and the original point cloud. Inspired by the ground-truth sampling strategy used by LiDAR-based 3D detection methods [3], [17], we treat points inside the bounding boxes as instance-level foreground points. When the whole cloud of points are projected onto the image plane, the pixel-level mask label is produced for each object by checking whether a point is inside the box or not. An  $s \times s$  square image is used to represent the ground-truth mask. If there are two or more points that fall into the same pixel of the image, we keep one of them at random. As illustrated in Figure 2, the generated ground-truth masks are coarse and sparse, however, they are sufficient to make the whole detection pipeline become aware of the object shapes.

To alleviate the interference from the noisy data, an uncertainty weighted cross-entropy loss is designed for the instance segmentation head

$$L_{seg} = -\frac{1}{s^2 \sigma_{seg}^2} \sum_{(r,c) \in s \times s} (L_{seg}^+(r,c) + L_{seg}^-(r,c)) + \frac{\log \sigma_{seg}}{s^2} \sum_{(r,c) \in s \times s} \mathbf{1}(y_{rc} \neq -1), \quad (1)$$

where

$$L_{seg}^+(r,c) = \mathbf{1}(y_{rc} = 1) \log p_{rc},$$

$$L_{seg}^-(r,c) = \mathbf{1}(y_{rc} = 0) \log(1 - p_{rc}).$$

In Eqn. 1, the indicator function  $\mathbf{1}(\cdot)$  equals to 1 if the input expression is true, and 0 otherwise. The label is  $y_{rc} \in \{1, 0, -1\}$  when the pixel is foreground, background or unknown respectively. The uncertainty value  $\sigma_{seg}$  is a learnable parameter, that indicates how much noise we have in this task. In practice, we train the model to predict the  $\log \sigma_{seg}$  to avoid division by zero.

In order to increase the number of training samples as well as make the model more robust to the position error of the detected center, we assign the instance labels to each point that is inside the  $3 \times 3$  neighborhood of the object center

point. Namely, we first collect all points that belong to the  $3 \times 3$  neighborhood of representative points. These points share the same object label as the center closest to them in terms of *Euclidean* distance. The center sampling is adopted in the regression heads and the segmentation head. For regression heads, the center sampling operation is performed on regressed attributes. For the segmentation head, the center sampling is applied on original features and before the segmentation head to avoid producing abundant masks on background positions. The center sampling strategy does not affect the heatmap prediction. We still apply the penalty-reduced focal loss [12], [18], [19] to train the heatmap. At the test time, we can simply extract the object predictions by finding the maximum points of the heatmap through a max-pooling operation.

For efficiency, the segmentation head is turned on only during the training. It is already sufficient since the instance segmentation training is designed to assist the training of detection and regression tasks.

### B. Multi-Task Learning in 3D Object Detection

Similar to other center-based methods [12], [20], each target object is detected as a center point along with its attributes in our method. The detection is built upon the DLA-34 backbone network. There are several target parameters to be estimated (listed in Table I) in our model. In general, it is a multi-task deep learning framework. The classification head estimates the raw location  $(u_f, v_f)$  of the 2D center point. Other attributes are regressed by several regression heads. The segmentation head is designed to boost the performance of other heads.

TABLE I: Detection heads and the corresponding parameters to be estimated

Head	Task	Target Paras.
Classification	Heatmap	$u_f, v_f$
Segmentation	Shape-Aware Feat.	-
Regression	Depth	$z$
	Orientation	$r_y$
	2D Bounding Box	$l_b, t_b, r_b, b_b$
	3D Dimensions	$h, w, l$
	Center Offset	$\delta_u, \delta_v$
	Keypoints	$\mathbf{x}_{kpt}$
	Uncertainties	$\sigma_1, \sigma_2, \sigma_3, \sigma_4$

**2D Detection** Similar to FCOS [21], the 2D bounding box is represented as distances from the representative center to four sides of the box. GIoU loss [22] is adopted to learn 2D bounding boxes since it is robust to object scale changes.

**3D Detection** As illustrated in Figure 1, the keypoint regression head predicts 10 keypoints for each object, from which we can extract five vertical lines. The depth value of vertical line  $l$  is estimated by utilizing the relative proportion between pixel height and estimated object height [7], [23]

$$z_l = \frac{f \times h}{h_l}, \quad (2)$$

where  $f$  is the focal length.  $h_l$  and  $h$  are the height of line  $l$  in pixel and the 3D height of the object respectively. Following with [7], five vertical lines are divided into three groups, and the depth value of each vertical line is averaged within each

group, resulting in three independent center depth values. In addition, another depth value is also estimated directly by regression. The resulting 4 depth values are weighted by their uncertainties (given in Eqn. 3).

$$z = \left( \sum_{i=1}^4 \frac{z_i}{\sigma_i} \right) / \left( \sum_{i=1}^4 \frac{1}{\sigma_i} \right), \quad (3)$$

where the uncertainty estimation branch is trained under the laplacian uncertainty loss [24]–[27].

When an object is detected on the heatmap at position  $(u_f, v_f)$ , its center  $\mathbf{x}_c$  is given by  $(s_0 \cdot u_f + \delta_u, s_0 \cdot v_f + \delta_v)$ , where  $s_0$  is the downsampling factor. With the predicted projected 3D center  $\mathbf{x}_c = (u_c, v_c)$ , The object location can be decoded as

$$(x, y, z) = \left( \frac{(u_c - c_u) \cdot z}{f}, \frac{(v_c - c_v) \cdot z}{f}, z \right), \quad (4)$$

where  $(c_u, c_v)$  is the principle point.

For dimension estimation, the relative changes with respect to the average dimension is predicted. For each class  $c$ , the average dimension is denoted as  $(\bar{h}_c, \bar{w}_c, \bar{l}_c)$ , the L1 loss for dimension regress is defined as

$$L_{dim} = \sum_{k \in \{h, w, l\}} |\bar{k}_c \cdot e^{\delta_k} - k^*|, \quad (5)$$

where  $k^*$  is the ground-truth dimension,  $\delta_k$  is the relative offset to be regressed. The dimensions are estimated by scaling the average dimensions, namely  $(h, w, l) = (\bar{h}_c \cdot e^{\delta_h}, \bar{w}_c \cdot e^{\delta_w}, \bar{l}_c \cdot e^{\delta_l})$ .

MultiBin loss [28] is used to estimate the local orientation  $\alpha$ . The global orientation  $r_y$  can be obtained by calculating

$$r_y = \alpha + \arctan(x/z). \quad (6)$$

The final 3D bounding box is then encoded as  $(x, y, z, w, h, l, r_y)$ . We refer readers to [7] for more details of the 3D detection head.

### C. Implementation Details

Following [7], [10], [29], a modified DLA-34 [12] is adopted as our backbone. For each input image<sup>1</sup>, the backbone produces a feature map with the down-sampling ratio 4. Different headers are responsible to estimated different parameters. The parameters to be estimated and the corresponding header are shown in Table I. Every regression head consists of one  $3 \times 3 \times 256$  convolution layer, BatchNorm [30], ReLU, and another  $1 \times 1 \times c_o$  convolution layer, where  $c_o$  is the output channels. For heatmap prediction, we use the same structure except that there is a sigmoid function padded after the final Conv layer. For center offset and heatmap prediction, the edge fusion [7] module is applied after the  $3 \times 3$  Conv layer to decouple the feature learning of truncated objects. For the instance segmentation head, a two-layer CNN block with a ReLU activation function is used to process the predefined position embedding map. The processed position embedding map is then added with extracted features vectors with broadcasting. The resulted tensor is then fed into another two-layer

<sup>1</sup>All the input images are padded into the same size of  $384 \times 1280$ .

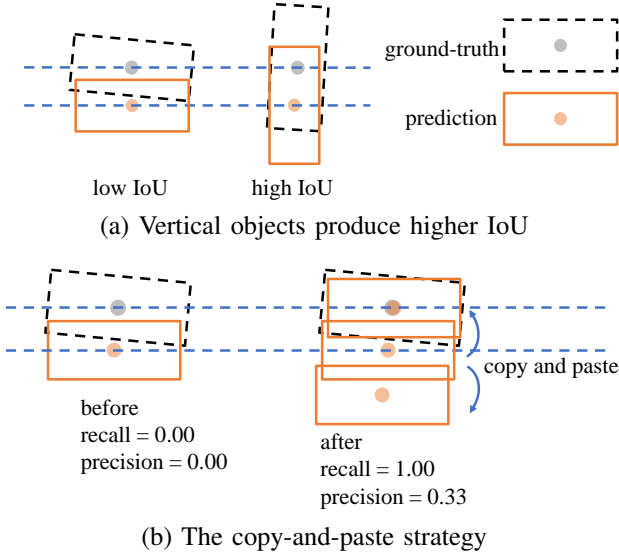


Fig. 4: Two scenarios where  $AP_{3D}$  cannot reflect the real performance of a monocular 3D detector.

CNN, which has a GroupNorm [31] layer embedded before each ReLU. The final mask is obtained by further applying one convolution layer with a sigmoid activation function.

The model is trained using AdamW optimizer with an initial learning rate as  $3 \times 10^{-4}$  and decay rate as  $1 \times 10^{-5}$ . We train the model for 200 epochs with a batch size of 18 on three 1080-Ti GPUs. The learning rate is divided by 10 at the 190th epoch and the 195th epoch.

#### IV. AVERAGE DEPTH SIMILARITY

In the literature, KITTI [32], [33] is one of the most popularly used evaluation benchmarks. The performance of monocular 3D object detection is measured by the 3D IoU-based Average Precision (AP) on the 3D space ( $AP_{3D}$ ) and the BEV ( $AP_{BEV}$ ). A detection is considered positive if it overlaps a ground-truth with an IoU larger than a threshold. Such metrics are reasonable to evaluate LiDAR-based or stereo-based methods. However, we find they could not reflect the real performance of a monocular object detector.

First of all, as pointed out in [34], the depth estimation error increases with respect to the object’s distance to the camera. The detected objects that are far away from the camera are largely ignored due to low 3D IoU. Figure 4 further shows another two scenarios in which we cannot have an honest observation about the detector with 3D IoU-based evaluation. In Figure 4(a), we can see the evaluation favors object in vertical layout even both of them are estimated with the same depth accuracy. In Figure 4(b), we show a trick that  $AP_{3D}$  can be cheated.  $AP_{3D}$  score is boosted by simply duplicating and shifting several distant predictions. As we will see later in the experiment,  $AP_{3D}$  can be boosted by 25% by such simple “result sampling”. Although this trick is effective in producing better  $AP_{3D}$ , it is essentially a trick that makes a better trade-off between precision and recall. It is non-pragmatic since it

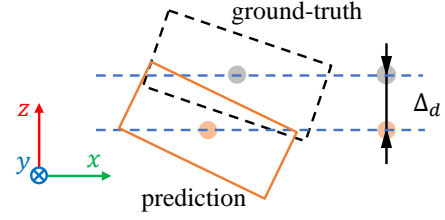


Fig. 5: The depth similarity is measured by calculating the absolute distance error along the z-axis.

neither improves the quality of the results nor enhances the ability of the model.

To address this issue, average depth similarity (ADS) is proposed as a complementary to  $AP_{3D}$ . Compared to  $AP_{3D}$  which is based on 3D-IoU, ADS matches the 2D predictions with the ground-truth. The depth similarity is defined as

$$s_d(r) = \frac{1}{|\mathcal{D}(r)|} \sum_{i \in \mathcal{D}(r)} \exp(-|\Delta_d^{(i)}|) \delta_i, \quad (7)$$

where  $\mathcal{D}(r)$  are the set of all detected objects at recall rate  $r$  and  $\Delta_d^i$  is the difference of the depth of detection  $i$ .  $\delta_i$  is 1 if detection  $i$  has been associated with a ground-truth object and 0 otherwise. The definition of ADS follows with the original design of AOS [32], [33], while the difference lies in the similarity definition. The error in depth  $\Delta_d^{(i)} \in [0, +\infty)$  is normalized into the range of  $[0, 1]$ .

As illustrated in Figure 5, the calculation of depth similarity solely depends on the absolute depth error along the z-axis. Unlike the 3D IoU, the calculation of depth similarity is independent of object attributes such as size, orientation, and position. In addition, ADS avoids penalizing predictions with small IoUs. Therefore, it enables a comprehensive observation of the behavior of a monocular detector when combined with traditional evaluation metrics  $AP_{3D}$  and  $AP_{BEV}$ .

#### V. EXPERIMENTS

Our method is evaluated on the popular KITTI dataset [32], [33] in comparison to state-of-the-art methods. The dataset consists of 7,481 training images and 7,518 testing images. Due to the restrictions on the *test* set submission to the KITTI official site, the training images are split into *train* (3,712) and *val* (3,769) sets following [40]. All our models jointly predict three categories, including Car, Pedestrian, and Cyclist.

The methods we consider in the experiment include one-stage center-based methods such as AutoShape [11], LPCG-MonoFlex [35], MonoCon [8], MonoFlex [7] and GUP-Net [24]. The comparison also covers one-stage anchor-based methods such as Ground-Aware [37], DLE [38], and LPCG-M3D-RPN [35] and one-stage FCOS-like method DD3D [36]. A two-stage method MonoRun [25] is incorporated in our comparative study as well. In the evaluation, MonoFlex [7] that is retrained with the same settings as our method is treated as the comparison baseline. For clarity, this run is given as “MonoFlex\*”, while the run that is loyal to the original paper is given as “MonoFlex”.

TABLE II: Ablation study on the KITTI *val* set for three categories, *i.e.*, the Car, Pedestrian and Cyclist

Methods	Car						Pedestrian						Cyclist					
	Easy		Moderate		Hard		Easy		Moderate		Hard		Easy		Moderate		Hard	
	AP <sub>3D</sub>	ADS	AP <sub>3D</sub>	ADS	AP <sub>3D</sub>	ADS	AP <sub>3D</sub>	ADS	AP <sub>3D</sub>	ADS	AP <sub>3D</sub>	ADS	AP <sub>3D</sub>	ADS	AP <sub>3D</sub>	ADS	AP <sub>3D</sub>	ADS
MonoFlex*	21.74	69.59	15.79	60.16	13.32	53.02	5.47	46.23	4.40	36.84	3.47	31.12	5.24	36.41	2.69	21.28	2.59	19.94
+ CS	23.25	70.00	16.78	60.79	14.81	54.96	7.54	52.37	5.54	43.77	4.52	36.45	5.62	37.82	3.17	23.57	2.72	22.15
+ CS + SA	<b>24.92</b>	<b>73.53</b>	<b>18.39</b>	<b>62.68</b>	<b>15.56</b>	<b>56.61</b>	<b>8.80</b>	<b>53.34</b>	<b>6.88</b>	<b>44.56</b>	<b>5.78</b>	<b>37.13</b>	<b>6.77</b>	<b>44.18</b>	<b>3.24</b>	<b>26.86</b>	<b>3.12</b>	<b>24.93</b>

TABLE III: Comparisons for Car category on the KITTI benchmark. For models marked with \*, we evaluate their performance on the *val* set by training with official codes. For other models, their results on the *val* set are either cited from their papers or by evaluating the released pre-trained model from the authors. Methods are ordered by their AP<sub>3D</sub> values under the moderate setting. For clarity, the best results are in red, while the second is highlighted in blue

Name	Publication	Extra Dataset	AP <sub>3D</sub> (test/val)			ADS (val)		
			Easy	Moderate	Hard	Easy	Moderate	Hard
AutoShape [11]	ICCV 2021	KINS	22.47/20.09	14.17/14.65	11.36/12.07	-	-	-
LPCG-M3D-RPN [35]	arxiv preprint	KITTI raw	22.73/26.17	14.82/19.61	12.88/16.80	-	-	-
DD3D [36]	ICCV 2021	DDAD15M	23.22/ -	16.34/ -	14.20/ -	-	-	-
LPCG-MonoFlex [35]	arxiv preprint	KITTI raw	25.56/31.15	17.80/23.42	15.38/20.60	-	-	-
MonoRUn [25]	CVPR 2021		19.65/21.63	12.30/15.22	10.58/12.83	67.95	58.91	52.94
Ground-Aware* [37]	IEEE RA-L 2021		21.65/23.14	13.25/15.71	9.91/11.76	69.42	54.47	42.51
MonoFlex [7]	CVPR 2021		19.94/24.22	13.89/17.34	12.07/15.14	71.59	<b>62.19</b>	<b>56.10</b>
DLE* [38]	BMVC 2021	None	<b>24.23/25.58</b>	14.33/16.50	10.30/12.27	70.26	56.44	43.01
GUPNet [39]	ICCV 2021		22.26/23.19	15.02/16.23	13.12/13.57	71.17	59.32	52.14
MonoCon [8]	AAAI 2022		22.50/ <b>26.33</b>	<b>16.46/19.03</b>	<b>13.95/16.00</b>	<b>72.15</b>	60.29	54.40
Ours	-	None	<b>23.84/24.92</b>	<b>16.52/18.39</b>	<b>13.88/15.56</b>	<b>73.53</b>	<b>62.68</b>	<b>56.61</b>

TABLE IV: Validation Test on ADS Measure

Methods	AP <sub>3D</sub>			ADS		
	Easy	Mod.	Hard	Easy	Mod.	Hard
Ours	24.92	18.39	15.56	<b>73.53</b>	<b>62.68</b>	<b>56.61</b>
+ sampling	<b>27.86</b>	<b>22.31</b>	<b>19.42</b>	69.22	52.29	46.45

Before we proceed with the comprehensive empirical study, the proposed evaluation measure, namely average depth similarity (ADS) is validated. A simple sampling [41] is conducted on the results produced by our method. As shown in Table IV, the AP<sub>3D</sub> score of our method becomes considerably higher with simple sampling. This is where ADS comes to complement. The performance of “+ sampling” run is much poorer when it is measured by ADS. This is because the AP<sub>3D</sub> of the “+ sampling” increases at the cost of producing more blind predictions. It, therefore, makes no contribution to the quality of 2D bounding boxes nor depth predictions. In the following experiments, both AP<sub>3D</sub> and ADS scores are reported on the *val* for all the methods.

#### A. Ablation Study

In our first experiment, an ablation study is conducted to show the contributions of center sampling (CS) and shape-aware training (SA), both of which are introduced by us. We observed that, when applying center sampling alone, 6.27% and 1.05% improvement in terms of AP<sub>3D</sub> and ADS respectively are observed on the moderate objects. When the shape-aware auxiliary training head is further integrated, extra 9.59% and 3.11% respectively are observed. As shown in Table II, this improvement is consistent on the Pedestrian and Cyclist categories.

In our second experiment, we further confirm the choice of feature expanding (illustrated in Figure 3) as the feature

TABLE V: Comparison between different RoI feature extractors. Experiments are conducted on the KITTI *val* set

Methods	AP <sub>3D</sub>			ADS		
	Easy	Mod.	Hard	Easy	Mod.	Hard
MonoFlex*	21.74	15.79	13.32	69.59	60.16	53.02
+ RoIAlign	19.97	14.65	12.91	69.19	59.79	52.86
+ Ours	<b>24.92</b>	<b>18.39</b>	<b>15.56</b>	<b>73.53</b>	<b>62.68</b>	<b>56.61</b>

extractor in the segmentation head. In this ablation study, the performance of using feature expanding is compared to the configuration that uses the RoIAlign [42] in the segmentation head. As shown in Table V, the performance of the detector drops when RoIAlign is used as the feature extractor in the segmentation head. In contrast, considerable improvement (in comparison to MonoFlex\*) is observed when the RoIAlign is replaced with feature expanding in the feature extractor. As discussed earlier, the misalignment between the RoI-based and center-based representation leads to potential conflicts between two training heads. The segmentation becomes harmful to the detection task, which leads to even poorer performance than the case it is, otherwise, not integrated.

#### B. Performance Analysis

In this section, the performance of our shape-aware detector is studied on both the *test* and *val* set of KITTI in comparison to state-of-the-art methods. The 3D average precision AP<sub>3D</sub> and the average depth similarity ADS are reported for all the methods. Note that we only evaluate ADS on the *val* set since the ground-truth of the *test* set is not available. Evaluation metrics are divided into easy, moderate, and hard settings according to the height, occlusion, and truncation level of objects. All evaluation metrics use 40 recall points instead of 11 recall points as recommended in [43]. As shown in

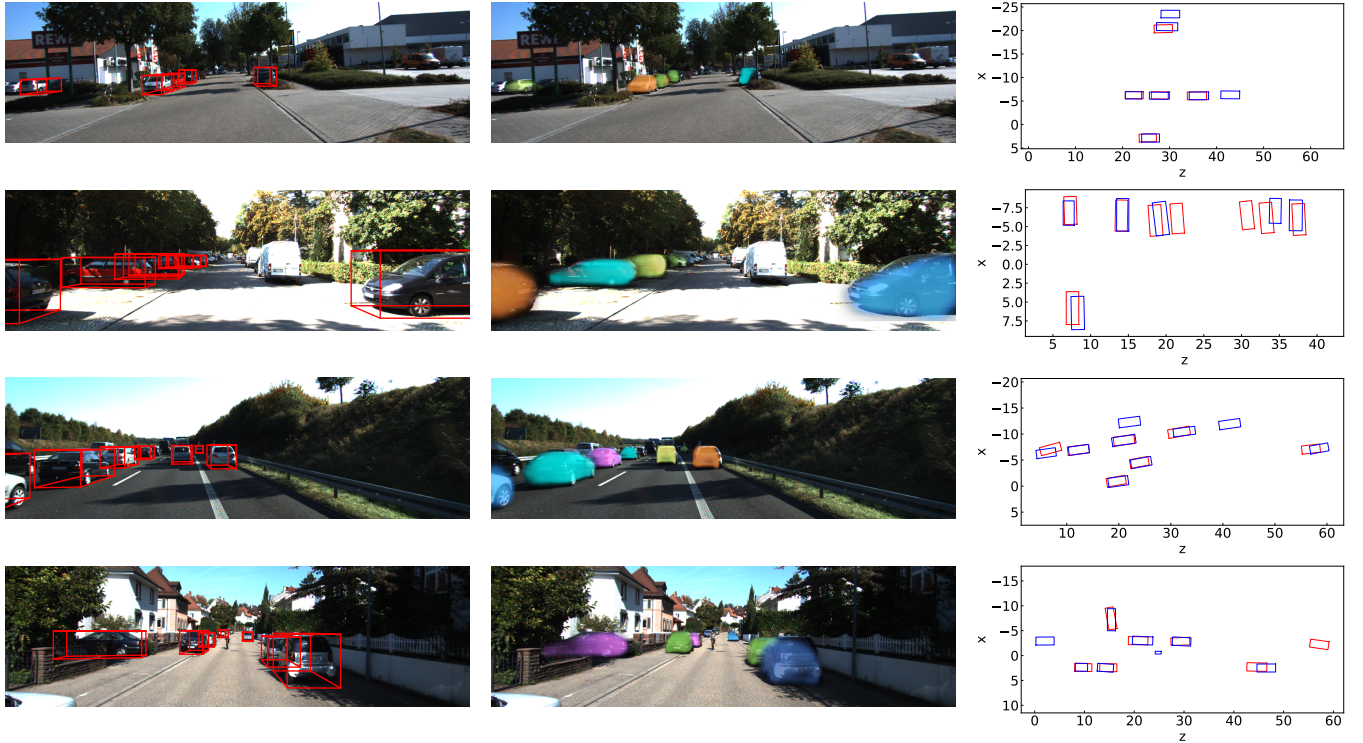


Fig. 6: Qualitative results on KITTI *val* set. The *1st* column shows the detected 3D boxes along with the objects. The *2nd* column shows the visible area of objects predicted by the segmentation head. On the *3rd* column, the detection results are shown along with the ground-truth in the BEV. The detected boxes and the ground-truth are in red and blue respectively.

Table III, our method in general outperforms the most recent method in the literature. In particular it outperforms MonoFlex on the *test* set considerably. Moreover, our method shows consistently better accuracy in depth estimation than the rest of methods. Figure 6 shows the qualitative results on four monocular street views. The second column shows the results output from our segmentation head. The produced mask for a target object indicates the possibility of a pixel being on the object and the visibility of the object. For occluded objects, the masks highlight their visible areas, which enables the model to be free from the interference of the context noises.

### C. Performance Analysis on Occluded Objects

TABLE VI: Performance on the occluded and fully visible objects. The evaluation is made on Cars of the KITTI *val* set

Approaches	AP <sub>3D</sub>		ADS	
	Fully Visible	Occluded	Fully Visible	Occluded
MonoFlex*	13.34	6.29	62.35	43.59
Ours	<b>15.21</b>	<b>8.92</b>	<b>64.68</b>	<b>50.55</b>

To further investigate the improvement that our method brings to the detection of occluded objects, we divide objects into two categories, the fully visible and the occluded, according to their occlusion state and truncation score. In KITTI, objects with occlusion state and truncation score equal to 0 are marked as fully visible. The rest that their occlusion states are 1 or 2, or truncation scores are larger than 0 are marked as

occluded. The objects whose occlusion state is 3 (unknown) are ignored in the evaluation.

The evaluation results are reported in Table VI. As seen from the table, the occluded objects are harder to be precisely located than the fully visible objects in the 3D space. Compared to MonoFlex, our method boosts the performance on the occluded objects by 41.81% and 15.97% in terms of the AP<sub>3D</sub> and the ADS respectively. Moreover, the shape-aware detector is also beneficial to the fully visible objects, for which 14.02% and 3.74% improvements are observed for the AP<sub>3D</sub> and the ADS respectively.

## VI. CONCLUSION

We have presented a simple but effective auxiliary training task named shape-aware feature learning, which aims to improve the model performance on the occluded objects. With uncertainty weighted loss function, our method is able to learn from sparse and noisy segmentation labels, relieving of the laborious manual mask annotation. Considerable improvements are observed on the full-view objects as well as the occluded objects in particular. In addition, a metric called average depth similarity that is complementary to the current popular evaluation protocol is proposed to measure the performance of a monocular 3D detector. It allows a more comprehensive understanding about the monocular 3D object detection models.

## ACKNOWLEDGMENT

This work is fully supported by Ningbo Boden AI Technology Co., Ltd. from Ningbo, China. It is also supported by National Natural Science Foundation of China under grants 61572408 and 61972326.

## REFERENCES

- [1] B. Yang, W. Luo, and R. Urtasun, "Pixor: Real-time 3d object detection from point clouds," in *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pp. 7652–7660, 2018.
- [2] Y. Zhou and O. Tuzel, "Voxelnet: End-to-end learning for point cloud based 3d object detection," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 4490–4499, 2018.
- [3] Y. Yan, Y. Mao, and B. Li, "Second: Sparsely embedded convolutional detection," *Sensors*, vol. 18, no. 10, 2018.
- [4] P. Li, X. Chen, and S. Shen, "Stereo r-cnn based 3d object detection for autonomous driving," in *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 7636–7644, 2019.
- [5] J. Sun, L. Chen, Y. Xie, S. Zhang, Q. Jiang, X. Zhou, and H. Bao, "Disp r-cnn: Stereo 3d object detection via shape prior guided instance disparity estimation," in *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 10545–10554, 2020.
- [6] Y. Wang, W.-L. Chao, D. Garg, B. Hariharan, M. Campbell, and K. Q. Weinberger, "Pseudo-lidar from visual depth estimation: Bridging the gap in 3d object detection for autonomous driving," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 8445–8453, 2019.
- [7] Y. Zhang, J. Lu, and J. Zhou, "Objects are different: Flexible monocular 3d object detection," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 3289–3298, 2021.
- [8] T. W. Xianpeng Liu, Nan Xue, "Learning auxiliary monocular contexts helps monocular 3d object detection," in *36th AAAI Conference on Artificial Intelligence (AAAI)*, February 2022.
- [9] P. Li, H. Zhao, P. Liu, and F. Cao, "Rtm3d: Real-time monocular 3d detection from object keypoints for autonomous driving," in *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part III 16*, pp. 644–660, Springer, 2020.
- [10] Z. Liu, Z. Wu, and R. Toth, "Smoke: Single-stage monocular 3d object detection via keypoint estimation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, June 2020.
- [11] Z. Liu, D. Zhou, F. Lu, J. Fang, and L. Zhang, "Autoshape: Real-time shape-aware monocular 3d object detection," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 15641–15650, 2021.
- [12] X. Zhou, D. Wang, and P. Krähenbühl, "Objects as points," *arXiv preprint arXiv:1904.07850*, 2019.
- [13] S. Luo, H. Dai, L. Shao, and Y. Ding, "M3dssd: Monocular 3d single stage object detector," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 6145–6154, 2021.
- [14] C. R. Qi, H. Su, K. Mo, and L. 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*, pp. 652–660, 2017.
- [15] A. Simonelli, S. R. Bulo, L. Porzi, P. Kotschieder, and E. Ricci, "Are we missing confidence in pseudo-lidar methods for monocular 3d object detection?," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 3225–3233, 2021.
- [16] F. Chabot, M. Chaouch, J. Rabarisoa, C. Teulière, and T. Chateau, "Deep manta: A coarse-to-fine many-task network for joint 2d and 3d vehicle analysis from monocular image," in *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 1827–1836, 2017.
- [17] B. Zhu, Z. Jiang, X. Zhou, Z. Li, and G. Yu, "Class-balanced grouping and sampling for point cloud 3d object detection," *arXiv preprint arXiv:1908.09492*, 2019.
- [18] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, "Focal loss for dense object detection," in *Proceedings of the IEEE international conference on computer vision*, pp. 2980–2988, 2017.
- [19] H. Law and J. Deng, "Cornernet: Detecting objects as paired keypoints," in *Proceedings of the European conference on computer vision (ECCV)*, pp. 734–750, 2018.
- [20] X. Zhou, V. Koltun, and P. Krähenbühl, "Probabilistic two-stage detection," *arXiv preprint arXiv:2103.07461*, 2021.
- [21] Z. Tian, C. Shen, H. Chen, and T. He, "Fcos: Fully convolutional one-stage object detection," in *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 9627–9636, 2019.
- [22] H. Rezaatoughi, N. Tsoi, J. Gwak, A. Sadeghian, I. Reid, and S. Savarese, "Generalized intersection over union: A metric and a loss for bounding box regression," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
- [23] Y. Cai, B. Li, Z. Jiao, H. Li, X. Zeng, and X. Wang, "Monocular 3d object detection with decoupled structured polygon estimation and height-guided depth estimation," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, pp. 10478–10485, 2020.
- [24] A. Kendall, *Geometry and uncertainty in deep learning for computer vision*. PhD thesis, University of Cambridge, UK, 2019.
- [25] H. Chen, Y. Huang, W. Tian, Z. Gao, and L. Xiong, "Monorun: Monocular 3d object detection by reconstruction and uncertainty propagation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10379–10388, 2021.
- [26] A. Kendall and Y. Gal, "What uncertainties do we need in bayesian deep learning for computer vision?," *arXiv preprint arXiv:1703.04977*, 2017.
- [27] A. Kendall, Y. Gal, and R. Cipolla, "Multi-task learning using uncertainty to weigh losses for scene geometry and semantics," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 7482–7491, 2018.
- [28] A. Mousavian, D. Anguelov, J. Flynn, and J. Kosecka, "3d bounding box estimation using deep learning and geometry," in *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pp. 7074–7082, 2017.
- [29] Y. Chen, L. Tai, K. Sun, and M. Li, "Monopair: Monocular 3d object detection using pairwise spatial relationships," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 12093–12102, 2020.
- [30] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," in *International conference on machine learning*, pp. 448–456, PMLR, 2015.
- [31] Y. Wu and K. He, "Group normalization," in *Proceedings of the European conference on computer vision (ECCV)*, pp. 3–19, 2018.
- [32] A. Geiger, P. Lenz, and R. Urtasun, "Are we ready for autonomous driving? the kitti vision benchmark suite," in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2012.
- [33] A. Geiger, P. Lenz, C. Stiller, and R. Urtasun, "Vision meets robotics: The kitti dataset," *International Journal of Robotics Research (IJRR)*, 2013.
- [34] X. Ma, Y. Zhang, D. Xu, D. Zhou, S. Yi, H. Li, and W. Ouyang, "Delving into localization errors for monocular 3d object detection," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 4721–4730, 2021.
- [35] L. Peng, F. Liu, Z. Yu, S. Yan, D. Deng, Z. Yang, H. Liu, and D. Cai, "Lidar point cloud guided monocular 3d object detection," *arXiv preprint arXiv:2104.09035*, 2021.
- [36] D. Park, R. Ambrus, V. Guizilini, J. Li, and A. Gaidon, "Is pseudo-lidar needed for monocular 3d object detection?," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 3142–3152, 2021.
- [37] Y. Liu, Y. Yixuan, and M. Liu, "Ground-aware monocular 3d object detection for autonomous driving," *IEEE Robotics and Automation Letters*, vol. 6, no. 2, pp. 919–926, 2021.
- [38] C. Liu, S. Gu, L. Van Gool, and R. Timofte, "Deep line encoding for monocular 3d object detection and depth prediction," in *32nd British Machine Vision Conference (BMVC 2021)*, p. 354, 2021.
- [39] Y. Lu, X. Ma, L. Yang, T. Zhang, Y. Liu, Q. Chu, J. Yan, and W. Ouyang, "Geometry uncertainty projection network for monocular 3d object detection," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 3111–3121, 2021.
- [40] X. Chen, K. Kundu, Y. Zhu, A. G. Berneshawi, H. Ma, S. Fidler, and R. Urtasun, "3d object proposals for accurate object class detection," in *Advances in Neural Information Processing Systems*, pp. 424–432, Citeseer, 2015.
- [41] L. Peng, S. Yan, C. Huang, X. He, and D. Cai, "Digging into output representation for monocular 3d object detection," 2022.
- [42] K. He, G. Gkioxari, P. Dollár, and R. Girshick, "Mask r-cnn. facebook ai research (fair)," *arXiv preprint arXiv:1703.06870*, 2018.
- [43] A. Simonelli, S. R. Bulo, L. Porzi, M. López-Antequera, and P. Kotschieder, "Disentangling monocular 3d object detection," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 1991–1999, 2019.