

OBMO: One Bounding Box Multiple Objects for Monocular 3D Object Detection

Chenxi Huang, Tong He, Haidong Ren, Wenxiao Wang, Binbin Lin*, Deng Cai, *Member, IEEE*

Abstract—Compared to typical multi-sensor systems, monocular 3D object detection has attracted much attention due to its simple configuration. However, there is still a significant gap between LiDAR-based and monocular-based methods. In this paper, we find that the ill-posed nature of monocular imagery can lead to depth ambiguity. Specifically, objects with different depths can appear with the same bounding boxes and similar visual features in the 2D image. Unfortunately, the network cannot accurately distinguish different depths from such non-discriminative visual features, resulting in unstable depth training. To facilitate depth learning, we propose a simple yet effective plug-and-play module, **One Bounding Box Multiple Objects (OBMO)**. Concretely, we add a set of suitable pseudo labels by shifting the 3D bounding box along the viewing frustum. To constrain the pseudo-3D labels to be reasonable, we carefully design two label scoring strategies to represent their quality. In contrast to the original hard depth labels, such soft pseudo labels with quality scores allow the network to learn a reasonable depth range, boosting training stability and thus improving final performance. Extensive experiments on KITTI and Waymo benchmarks show that our method significantly improves state-of-the-art monocular 3D detectors by a significant margin (The improvements under the moderate setting on KITTI validation set are 1.82 ~ 10.91% mAP in BEV and 1.18 ~ 9.36% mAP in 3D). Codes have been released at <https://github.com/mrsempress/OBMO>.

Index Terms—3D object detection, Monocular images, Depth ambiguity, Camera project principles.

I. INTRODUCTION

DUE to widely deployed applications in robot navigation and autonomous driving [1]–[5], 3D object detection has become an active research area in computer vision. Although LiDAR-based 3D object detectors [6]–[8] have achieved excellent performance because of accurate depth measurements, the application of these methods is still constrained by the high cost of 3D sensors, limited working range, and sparse data representation. Monocular-based 3D detectors [9]–[14], on the other hand, have received increasing attention in autonomous driving due to their easy accessibility and rich semantic clues.

Though tremendous efforts have recently been devoted to improving the accuracy, monocular-based 3D object detection is still highly challenging, as substantiated by

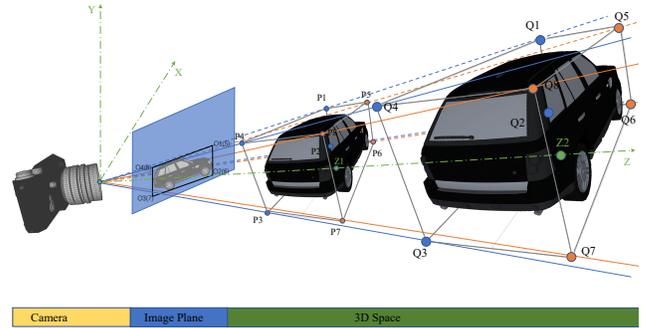


Fig. 1. Objects with different depths and dimensions in 3D space. Objects P and Q have the same bounding box and similar visual features in the 2D image, leading to depth ambiguity.

[14]–[17]. Earlier works are based on mature 2D object detection, using 2D Region of interests (ROIs) to regress 3D information. [18]–[22] follow the pipeline by either introducing geometrical priors or laying 2D-3D constraints. Promising results have been achieved; however, the gap in the accuracy between LiDAR-based and monocular-based approaches is still significant. One of the critical reasons for the less competitiveness in monocular-based methods is lacking precise knowledge of depth.

To this end, many prior works [23]–[35] focus on improving the accuracy of instance depth estimation. These methods mainly involve two strategies to encode depth prior by either building dependencies on the intermediate task of depth prediction or adding geometric constraints on the final results. For the former, the expression content of the data is enriched by initial depth prediction values from monocular depth estimation or extra-designed modules. Some methods transform the front view into other views using the predicted depth values, such as the bird’s eye view (BEV) [23], [26]. Other methods combine the predicted depth values with the corresponding RGB values into new data representations, for example, concatenating them on channels [27] or converting them into LiDAR format [28], [29]. For the latter, they add extra modules or change the objective function to assist in estimating depth. [30], [35] use the projection relationship to constrain the predictions by well-designed Volume Displacement loss and ground-aware convolution module, respectively. [32]–[34] split the depth value into a coarse value and a bias-corrected value. [32] considers the relationship between objects and regards the distance of the 3D pair as bias. [33], [34] predict the variance of depth and height, respectively, and

* Corresponding author

C. Huang and D. Cai are with the State Key Laboratory of CAD&CG, College of Computer Science, Zhejiang University, Hangzhou, Zhejiang 310058, China (emails: hcx_98@zju.edu.cn, dengcai@cad.zju.edu.cn).

T. He is with Shanghai AI Laboratory, Shanghai, China (emails: tonghe90@gmail.com).

H. Ren is with Ningbo Zhoushan Port Group Co.,Ltd., Ningbo, China. (email: rhdong@nbport.com.cn).

W. Wang and B. Lin are with College of Software, Zhejiang University, Ningbo, Zhejiang, China (emails: wenxiaowang@zju.edu.cn, binbinlin@zju.edu.cn).

use them to fine-tune the rough values.

Previous works adopt the one-to-one learning strategy, which uses a 3D object label to supervise learning from the visual features of one 2D object. However, due to the asymmetrical projection between 2D and 3D, this one-to-one learning strategy often causes depth ambiguity. For example, if moving a car along the viewing frustum from a depth of 50 m to 55 m and enlarging the size by 1.1, the visual representation remains unchanged when projected to the 2D plane. As shown in Figure 1, different objects in 3D space may have very similar bounding boxes and visual features when projected to the 2D image. Considering that the average length, width, and height of the Car are 3.88 m, 1.63 m, and 1.53 m, respectively, in KITTI, it is still maintained in a reasonable scope when the size expanded by 1.1 times. Consequently, such ambiguity often causes inferior performance the network has to distinguish different depths based on the non-discriminative visual clues.

We propose a simple yet effective plug-and-play module named **One Bounding Box Multiple Objects (OBMO)** to address the above problems. The core idea of OBMO is adding reasonable pseudo labels by shifting the depth of an object along the viewing frustum. Considering the lack of depth in 2D images, the soft pseudo labels of 3D objects play a significant role in encoding depth prior. Compared with the hard labels, such soft labels encourage the network to learn the depth distribution and stabilize the learning process, due to the less variance in the gradient between training cases [36].

Designing such soft labels is non-trivial, as the significant variation of depth often generates invalid sizes of 3D objects, making the network overwhelmed by negative samples. To this end, we design two label scoring strategies that use dimensional priors and geometric constraints to represent the quality of pseudo labels.

By introducing the OBMO module and the label scoring strategy, the one-to-many problem is addressed to some extent: the network is encouraged to learn a soft distribution of object locations rather than deterministic ones. To show the superiority of OBMO, we perform extensive experiments on KITTI and Waymo datasets. Multiple monocular 3D detectors are used, including direct regression-based detectors like RTM3D [37], Ground-aware [35], GUPNet [33] and depth-aware detectors like PatchNet [27], Pseudo-LiDAR [28]. Experimental results show that our method stabilizes the training process and improves the overall BEV and 3D detection performance, as shown in Figure 5. Concretely, on the widely used KITTI dataset, our approach significantly improves the state-of-the-art (SOTA) monocular 3D detectors by **1.82 ~ 10.91% mAP in BEV** and **1.18 ~ 9.36% mAP in 3D**. On the larger Waymo open dataset, we boost GUPNet with **3.34% mAP** gains under the LEVEL 1 (IoU = 0.5) setting.

The contributions can be summarized as follows:

- We point out that the depth ambiguity problem in monocular 3D detection has been ignored in previous methods and argue that this problem can result in unstable depth training, which undermines performance.
- To alleviate the problem of unstable depth training in monocular 3D object detection, we propose a plug-and-play module OBMO. It explicitly adds a set of suitable

pseudo labels by shifting bounding boxes along the viewing frustum for each original object.

- We design two label scoring strategies to represent the qualities of pseudo labels: IoU Label Scores and Linear Label Scores, which are inspired by the fixed dimension range of objects in the same category.
- We conduct extensive experiments on various datasets: KITTI and Waymo. The consistent improvement of the accuracy demonstrates the effectiveness of our proposed OBMO. For example, we achieved 21.41% in AP_{BEV} and 15.70% in AP_{3D} under the moderate KITTI validation set based on GUPNet, improving the state-of-the-art results substantially.

II. RELATED WORK

A. LiDAR 3D Object Detection

Due to the accurate depth measurement, most state-of-the-art 3D object detection methods are based on LiDAR [8], [38]–[41]. These methods can be roughly divided into two parts: voxel-based methods and point-based methods.

1) *voxel-based methods*: In order to tackle the irregular data format of point clouds, voxel-based methods [3], [42], [43] convert the irregular point clouds into regular voxel grids. Then, use mature convolution neural architectures to extract high-level features. However, the receptive fields are constrained by the kernel size of 2D/3D convolutions [44], [45]. Moreover, the computation and memory grow cubically with the input resolution. To this end, SECOND [46] leverages the 3D submanifold sparse convolution. In spatially sparse convolution, output points are not computed if there is no related input point, which significantly increases the speed of both training and inference. Further, PointPillars [39] is proposed to simplify the voxels to pillars. Overall, voxel-based methods can achieve good detection performance with promising efficiency. However, it is difficult to determine the optimal voxel resolution in practice since the complex geometry and various dimension objects.

2) *point-based methods*: Point-based methods [47], [48] directly extract raw unstructured point cloud features via different set abstraction operations. Further, it generates specific proposals for objects of interest. These point-based methods, such as the PointNet [49] series, enable flexible receptive fields for point cloud feature learning. For example, PointRCNN [47], a two-stage 3D region proposal framework for 3D object detection, generates object proposals from segmented foreground points and exploits semantic features to regress high-quality 3D bounding boxes. PointGNN [48] generalizes graph neural networks to do 3D object detection. In conclusion, point-based methods don't need extra preprocessing steps such as voxelization. However, the main bottleneck of point-based methods is insufficient representation and inefficiency.

B. Monocular 3D Object Detection

Although LiDAR 3D object detectors present promising results, they have disadvantages of the limited working range and sparse data representation. Monocular 3D object detectors, on the other hand, enjoy the low cost and high frame

rate. Current monocular 3D object detection methods can be roughly divided into two categories: direct regression-based methods and depth-aware methods.

1) *Direct Regression-based Methods*: Direct regression-based methods [9]–[11] obtain the 3D detection results directly from RGB images without extra knowledge like depth maps, stereo images, etc.

Mono3D [20] first proposes an energy minimization approach and assumes that all vehicles are placed on the ground plane. Moreover, it scores each candidate box projected to the image plane via several intuitive potentials encoding semantic segmentation, contextual information, size and location priors, and typical object shape. Deep3DBox [19] simplifies the whole pipeline by removing extra 3D shape models and complex pre-processing operators. It is based on 2D object detection and uses geometric constraints that the 3D bounding box should fit tightly into the 2D detection bounding box. Considering geometric reasoning, MonoGRNet [22] simultaneously estimates 2D bounding boxes, instance depth, 3D location of objects, and local corners. M3D-RPN [18] proposed depth-aware convolutional layers for learning spatially-aware features to produce 3D proposals directly.

SMOKE [50] removes the 2D detection part and directly estimates 3D position by predicting projected 3D centers. RTM3D [37] adds eight corner points as keypoints so that more geometric constraints can be applied to remove false alarms. It also designs a keypoint feature pyramid, which uses soft weight by a softmax operation to denote the importance of each scale. Center3D [51] uses Linear Increasing Discretization and a combination of classification and regression branches to predict depth. MonoFlex [24] explicitly decouples the truncated objects and adaptively combines multiple approaches for object depth estimation. Specifically, it divides objects according to whether their projected centers are “inside” or “outside” the image. Furthermore, it formulates the object depth estimation as an uncertainty-guided ensemble of directly regressed object depth and solved depths from different groups of keypoints. GUPNet [33] proposes a GUP module to obtain the geometry-guided uncertainty of the inferred depth and designs a Hierarchical Task Learning strategy to reduce the instability caused by error amplification. MonoDTR [13] combines the transformer architecture and proposes a Depth-Aware Transformer module, which is used to integrate context- and depth-aware features globally.

Direct regression-based methods predict depth through a branch and employ one depth value to supervise a Region of Interest (ROI). However, we argue that this one-to-one learning strategy often suffers from depth ambiguity problems in monocular 3D object detection.

2) *Depth-aware Methods*: Depth-aware methods usually need extra depth map, which is used for 3D detection.

Pseudo-LiDAR [28] combines monocular 3D object detection task with monocular depth estimation task. It transforms RGB images to point clouds via an off-the-shelf depth estimator. Finally, effective point cloud-based 3D object detectors are employed for achieving the detection results. PatchNet [27] discovers that the data representation is not the most important one, but the coordinate transformation is. Thus it directly inte-

grates the 3D coordinates as additional channels of RGB image patches. D4LCN [52] points out that methods like Pseudo-LiDAR highly rely on the quality of depth map, and traditional 2D convolution cannot distinguish foreground pixels and background pixels. So it generates dynamic convolution kernels to extract features in different 3D locations. CaDDN [23] discretizes the range of depth and utilizes estimated categorical pixel-wise depth distribution. It changes the representation into BEV and then uses the BEV backbone to predict the 3D detection results. MonoJSG [12] reformulates the Monocular Object Depth Estimation as a progressive refinement problem and proposes a joint semantic and geometric cost volume to model the depth error.

Depth-aware methods only obtain a single depth value of the center point pixel through the Monocular Depth Estimation task. Similarly, they ignore the possibility of multiple reasonable depth values.

III. APPROACH

In this section, we first provide a detailed analysis of the widespread existence of “one bounding box with multiple objects.” Such ambiguity severely affects the training stability and accuracy of the model. Previous works ignore this problem, while we propose a simple but efficient module OBMO to lessen the impact. Since the dimension of each category has its reasonable range, we design two label scoring strategies to represent the quality of pseudo labels, making unreasonable pseudo labels ineffective.

A. Depth Ambiguity Problem

Obviously, 3D space is much larger than the projected 2D space. Using a 2D image to recover 3D space is an ill-posed task. Considering two objects with different 3D locations in 3D space, they may have similar bounding boxes and visual features in the 2D image, as shown in Figure 1. It indicates that predicting precise 3D locations from the 2D image may be impossible. We theoretically prove it in this subsection.

Without loss of generality, we assume that the camera system has been calibrated, which follows a typical pinhole imaging principle, as shown in Equation 1.

$$s \begin{bmatrix} u \\ v \\ 1 \end{bmatrix}_{3 \times 1} = \begin{bmatrix} f_x & 0 & c_x & 0 \\ 0 & f_y & c_y & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix}_{3 \times 4} \begin{bmatrix} x \\ y \\ d \\ 1 \end{bmatrix}_{4 \times 1} \quad (1)$$

In this equation, s is the scale factor, u, v represent the position of an object in image coordinates, x, y refer to its position in camera coordinates, d is the depth of the object. f_x, f_y, c_x, c_y come from intrinsic parameters of the calibrated camera. We set the scale factor $s = 1$ for notation convenience. Then, we can rewrite Equation 1 as follows:

$$\frac{u - c_x}{f_x} = \frac{x}{d}, \quad \frac{v - c_y}{f_y} = \frac{y}{d} \quad (2)$$

Regarding a point $A(u, v)$ on the image, $\frac{u - c_x}{f_x}$ and $\frac{v - c_y}{f_y}$ are fixed, as f_x, f_y, c_x, c_y are intrinsic parameters of the camera.

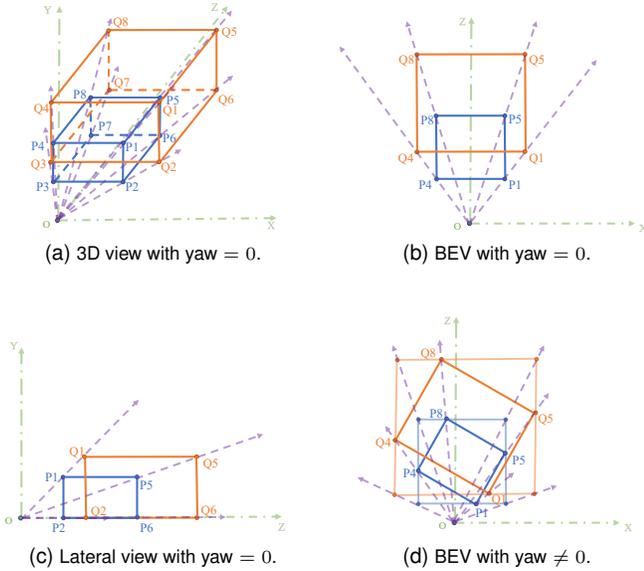


Fig. 2. Different views of object P and Q with the same 2D bounding box.

According to Equation 2, $\frac{u-c_x}{f_x}$ denotes the ratio between x and d , while $\frac{v-c_y}{f_y}$ denotes the ratio between y and d . Therefore, we call them X-Z ratio and Y-Z ratio, respectively. According to the projection relationship, we know that infinite 3D points can lead to A , as long as they have the same X-Z ratio and Y-Z ratio (along the same ray from the camera optical center to the 3D point (x, y, d)). As a result, one point on the image can correspond to multiple 3D locations, and one 2D bounding box on the image can correspond to various objects in 3D space.

Moreover, in Figure 2, we give an intuitive explanation in three different views: 3D view, bird’s eye view (BEV), and lateral view. From Figure 2b, we can see that the ratio of the widths (lengths) equals the ratio of depths. From Figure 2c, we know that the ratio of heights (lengths) equals the ratio of depths. More generally, when $\text{yaw} \neq 0$, we can get the same conclusion using bounding boxes, illustrated in Figure 2d.

In order to obtain the accurate value of depth, an intuitive solution is to estimate object dimensions both in 2D image and 3D space, then recover the depth according to geometry projection. However, the error caused by dimension estimation will amplify the depth estimation error, and it is non-trivial to predict object dimensions precisely.

Assume the error in dimension estimation is at the centimeter level, then the depth error is $\pm 0.01 \times \text{depth}$. Taking car P and car Q in Figure 1 as an example, assume the dimension scale factor between object Q and object P is 1.02. For objects in 100-meter away, as the typical height of cars is 1.53-meter (averaged value in KITTI), 0.03-meter dimension errors ($1.53 \times (1.02 - 1) \approx 0.03$) can cause 2-meter depth errors ($100 \times (1.02 - 1) = 2$). It will significantly decrease IoU values between predictions and ground truths, which increases training difficulty and instability. It indicates that only using the dimension to resolve the depth is also infeasible.

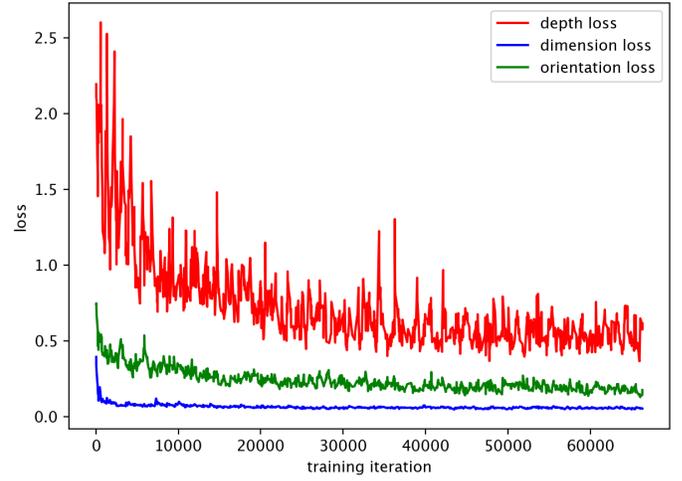


Fig. 3. The loss based on PatchNet [27]. We can see that both predictions of dimension and orientation are stable.

B. OBMO Module

The depth ambiguity causes that objects with different depths can appear very similar visual clues on the RGB image. For monocular-based methods, they have to distinguish depth from such non-discriminative features. This characteristic significantly affects the training stability. Therefore, we propose a module named OBMO to resolve the intractable depth ambiguity problem.

OBMO aims to let the network know that objects with different positions in 3D space may have similar bounding boxes and visual features in the 2D image. After looking at multiple reasonable pseudo labels, the network can give more general answers. Similar to label smoothing [53], which strengthens the network generalization ability by changing the one-hot encoding to a soft encoding that carries more information. Specifically, OBMO is a plug-and-play module capable of being applied during training to any monocular 3D detector.

To mitigate the adverse impacts of the depth ambiguity issue, we add some pseudo labels along with the viewing frustum within a reasonable range, as shown in Figure 4. This design improves the generalization ability of the network, as pseudo labels in a larger space remove the strict limitation of original hard labels. Specifically, we first calculate the X-Z ratio and Y-Z ratio for each object as defined in Equation 2. Then, we disturb the depth by a set of small offsets for each ground truth (class, $X, Y, Z, H, W, L, \text{yaw}$). The depth offsets Δ_z are determined by the dimension error tolerance and its depth Z . Taking the “Car” as an example, we consider the values of Δ_z from the set $\{-8\%, -4\%, +4\%, +8\%\}$, resulting in $Z_{\text{adjust}} = \Delta_z \cdot Z$. Then, adjust X and Y based on the X-Z ratio and Y-Z ratio, respectively. Given that the prediction of the dimensions (H, W, L) is relatively precise and consistent, as depicted in Figure 3, and considering that dimensions are inherent attributes of an object, we preserve their initial values without any alterations, directing our efforts solely towards enhancing the forecast of the 3D location. To validate our design approach, we conduct ablation studies to justify our design, as presented in Ablation Study Table X. Therefore, we get a

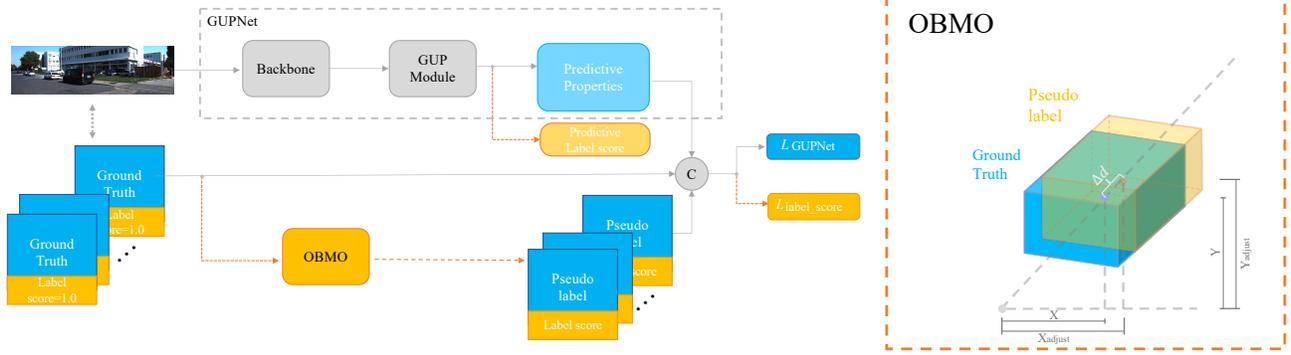


Fig. 4. The architecture of OBMO with label scoring strategies embedded on GUPNet. The differences are marked in orange. © means “compare”. The OBMO module is used to produce a set of pseudo labels and adds an extra attribute to measure their quality. The Label Score branch is inserted into GUPNet, parallel with 3D prediction branches. Moreover, the OBMO module only works in the training stage.

new pseudo label (class, X_{adjust} , Y_{adjust} , Z_{adjust} , H , W , L , yaw). To facilitate learning, we provide ground truths and pseudo labels as supervised signals to the network. It allows the model to incorporate and benefit from the knowledge encoded in the ground truths and the generated pseudo labels during training.

C. Two Label Scoring Strategies

However, the added pseudo labels are not unlimited. If the pseudo label is too far from the corresponding ground truth, then the transformation of dimensions is too heavy. Moreover, for each category of object, its dimension is limited. Therefore, to make pseudo labels reasonable, we should add constraints on depth offsets or distinguish unreasonable pseudo labels so that irrational pseudo labels will not affect training. Consequently, we design two kinds of label scores, which are used to represent the quality of pseudo labels. One is IoU Label Scores, and the other is Linear Label Scores. Both measure the similarity between pseudo labels and ground truths so they can be obtained before training.

1) *IoU Label Scores*: Since Intersection over Union (IoU) is a good measurement of how similar two bounding boxes are, we use it as the quality score. The higher the IoU value, the more significant the pseudo label. If two objects do not intersect, the 3D IoU is 0, but the 2D project IoU may not. It is common in categories with smaller lengths, such as Pedestrian. Therefore, instead of using 3D IoU, we use the IoU value of 2D project bounding boxes, defined as follows:

$$\text{IoU Label Score} = \text{IoU}(B_{\text{gt}}, B_{\text{pseudo}}), \quad (3)$$

B_{gt} refers to the original ground-truth 2D project bounding box, and B_{pseudo} is the projected project bounding 2D box of the added pseudo 3D box label.

2) *Linear Label Scores*: Furthermore, we introduce another simple yet effective scoring strategy: the Linear Label Score. It only cares about the offset of depth, and we use a simple linear function, as Equation 4 shows,

$$\text{Linear Label Score} = 1 - \frac{(|\Delta_z \cdot Z|)}{c}, \quad (4)$$

where c is a hyper-parameter, and we use it to balance the impacts of pseudo labels. The larger c is, the more enormous

impacts pseudo labels have on the training stage. Thus there is a trade-off in the choice of c . In our experiments, we choose $c = 4$ which empirically makes the score range in $[0, 1]$. This scoring strategy intuitively reflects the quality of pseudo labels. For pseudo objects too far away, Linear Label Scores are less than 0 and filter them out.

For ground truths, the quality scores under both scoring strategies are set to 1.0. In section IV-D of the ablation study, we find these two label scoring strategies have similar performance, which means that OBMO is robust to the label scoring strategies.

The quality score estimation branch is an auxiliary network that adopts the same structure as the other parallel regression heads. We use L1 loss between the ground truth Label Score and predicted Label Score, as follows:

$$\mathcal{L}_{\text{Label Score}} = |\text{Label Score}_{\text{pred}} - \text{Label Score}_{\text{gt}}|. \quad (5)$$

So, the total objective function is:

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{Baseline}} + \lambda \mathcal{L}_{\text{Label Score}}, \quad (6)$$

where λ is a trade-off between our Label Score Loss and the losses designed in the original method. If the baseline monocular detector is GUPNet, then $\mathcal{L}_{\text{Baseline}}$ is the hierarchical task loss of 2D detection (including heatmap, 2D offset and 2D size), 3D heads (containing angle, 3D offset and 3D size) and depth inference.

The whole process of the OBMO with label scoring strategies embedded on GUPNet is shown in Figure 4, which can conclude as *adding reasonable pseudo labels and adding a parallel Label Score branch*.

IV. EXPERIMENTS

A. Implementation Details

We adopt published codes from each baseline¹: PatchNet [27], Pseudo-LiDAR [28], Ground-aware [35], RTM3D [37], and GUPNet [33]. We use the same

¹The codes we referenced are: <https://github.com/xinzhuma/patchnet> (PatchNet, Pseudo-LiDAR), <https://github.com/Owen-Liuyuxuan/visualDet3D> (Ground-aware, RTM3D), and <https://github.com/SuperMHP/GUPNet> (GUPNet).

TABLE I

COMPARISONS ON KITTI TEST SET. FOR EASY TO COMPARE, WE SORT THEM ACCORDING TO THEIR 3D PERFORMANCE ON THE MODERATE LEVEL OF THE TEST SET (SAME AS THE KITTI LEADERBOARD). WE USE RED FOR THE HIGHEST ONES AND BLUE FOR THE SECOND-HIGHEST ONES.

Methods	$AP_{BEV}(\%)$			$AP_{3D}(\%)$		
	Easy	Mod.	Hard	Easy	Mod.	Hard
M3D-RPN [18]	21.02	13.67	10.23	14.76	9.71	7.42
SMOKE [50]	20.83	14.49	12.75	14.03	9.76	7.84
RTM3D [37]	19.17	14.20	11.99	14.41	10.34	8.77
PatchNet [27]	22.97	16.86	14.97	15.68	11.12	10.17
KM3D [55]	23.44	16.20	14.47	16.73	11.45	9.92
D4LCN [52]	22.51	16.02	12.55	16.65	11.72	9.51
Monodle [25]	24.79	18.89	16.00	17.23	12.26	10.29
MonoRUn [56]	27.94	17.34	15.24	19.65	12.30	10.58
GrooMeD-NMS [57]	26.19	18.27	14.05	18.10	12.32	9.65
Ground-aware [35]	29.81	17.98	13.08	21.65	13.25	9.91
CaDDN [23]	27.94	18.91	17.19	19.17	13.41	11.46
MonoEF [58]	29.03	17.26	19.70	21.29	13.87	11.71
MonoFlex [24]	28.23	19.75	16.89	19.94	13.89	12.07
GUPNet [33]	-	-	-	20.11	14.20	11.77
GUPNet (+ OBMO)	30.81	21.41	18.37	22.71	15.70	13.23
Improvements	-	-	-	+2.6	+1.50	+1.46

configuration described in their papers or projects. Take GUPNet as an example; we use DLA-34 as the backbone, train the model with the batch size of 32 for 140 epochs and adopt the initial learning rate $1.25e^{-3}$ with decay in the 90-th and the 120-th epoch. We train all models on Nvidia GTX 1080Ti GPUs with 11 GB memory.

Moreover, we set $\Delta_Z = \{-8\%, -4\%, +4\%, +8\%\}$ for all detectors and report the better one between IoU Label Score and Linear Label Score. For the monocular 3D detectors like PatchNet, which regard the scores of 2D bounding boxes as absolute confidence of objects directly, we employ the 2D-3D confidence mechanism from [54] to make the scores better describe the 3D predictions.

B. Dataset and Metrics

We conduct experiments on KITTI [59] and Waymo [60] benchmarks.

1) *KITTI*: KITTI is the widely employed dataset for monocular 3D object detection. It provides 7481 images for training and 7518 images for testing. All the scenes are pictured around Karlsruhe, Germany in clear weather and day time. To make fair comparisons, we follow previous works [33], [35], [61] to split the training images into train set (3712 images) and val set (3769 images). All experiments are performed under this dataset split. Furthermore, the detection results are evaluated under three levels of difficulty: easy, moderate and hard, which are defined according to the height of the 2D bounding box, occlusion, and truncation. We conduct experiments under two core evaluations: the average precision of 3D bounding boxes AP_{3D} and the average precision of objects in Bird’s Eye View AP_{BEV} . For the metric, we employ the recently suggested metric $AP_{BEV|R_{40}}$ and $AP_{3D|R_{40}}$ by KITTI benchmark [59]. Following common practice [27], [33], [37], we evaluate the results on the Car category under IoU threshold 0.7.

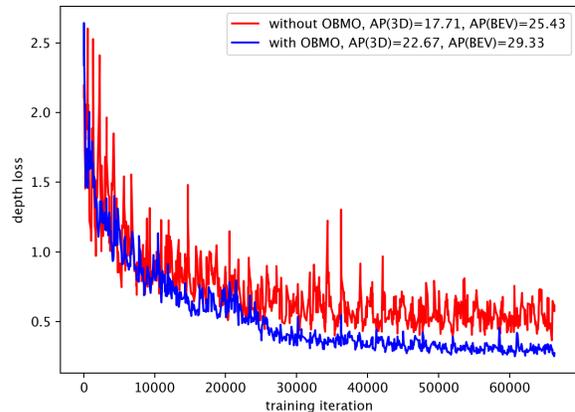


Fig. 5. The depth loss with/without the proposed module (OBMO) based on PatchNet. We can see that OBMO can stabilize depth training from the more stable loss curve.

2) *Waymo*: The Waymo dataset is a recently released large dataset for autonomous driving research. It consists of 798 training sequences and 202 validation sequences. The scenes are pictured in Phoenix, Mountain View, and San Francisco under multiple kinds of weathers and at multiple times of a day. Different from KITTI, it provides 3D box labels in the 360-degree field of view, while we only use the front view for the task of monocular 3D object detection. We use the same data processing strategy proposed in CaDDN [23]. Specifically, we sample every third frame from the training sequences to form our training set due to the large dataset size and high frame rate. We adopt the officially released evaluation to calculate the mean average precision (mAP) and the mean average precision weighted by heading (mAPH). The evaluation is separated by difficulty setting (LEVEL 1, LEVEL 2) and distance to the sensor (0 – 30 m, 30 – 50 m, and 50 m – ∞). We evaluate the Car category with IoU criteria of 0.7 and 0.5.

C. Quantitative Results

In Table I, we conduct a comprehensive comparison between our proposed method and existing state-of-the-art methods on the test sets of KITTI benchmark for Car. Without bells and whistles, our method outperforms all prior methods under $AP_{BEV|R_{40}}$ and $AP_{3D|R_{40}}$ including those with extra information. For $AP_{3D|R_{40}}$, our method is 22.71%/15.70%/13.23%, which is much higher than the baseline GUPNet on three levels of difficulty. The performance improvement is even more significant at the easy level. We suspect this is because, for foreign vehicles, the tiny depth shift needs a significant dimension transformation compared to near cars. Indeed, there is also less visual information in distant objects.

We further show the efficiency of our module OBMO embedded in other different SOTA monocular detectors in Table II. Because of the different train-val split, PatchNet* is retrained by its public code with a unified split [61]. As for RTM3D and Pseudo-LiDAR, which only report results

TABLE II

PERFORMANCE OF USING OUR OBMO METHOD, INCLUDING OBMO MODULE AND QUALITY SCORES, ON DIFFERENT SOTA MONOCULAR DETECTORS. ALL METHODS ARE EVALUATED ON KITTI VAL SET WITH METRIC $AP|_{R_{40}}$. WE REPORT BOTH THE RESULTS OF USING THE IOU LABEL SCORE AND LINEAR LABEL SCORE. AND WE ONLY COMPARE THE BEST ONE OF AP_{BEV} MOD. WITH THE BASELINE.

Methods	$AP_{2D}(\%)$			$AP_{BEV}(\%)$			$AP_{3D}(\%)$		
	Easy	Mod.	Hard	Easy	Mod.	Hard	Easy	Mod.	Hard
PatchNet [27]	97.05	94.00	86.32	43.97	25.43	20.73	32.56	17.71	13.98
(+ 3D scores [54])	97.17	94.07	87.00	41.82	28.13	24.23	32.96	21.27	17.87
(+ OBMO-Linear Label Scores)	97.20	94.00	86.96	42.91	29.33	24.53	34.16	22.67	18.22
(+OBMO-IoU Label Scores)	97.23	93.96	86.92	43.10	29.47	24.59	33.24	22.30	18.15
Improvements	+0.15	0.00	+0.64	-1.06	+3.90	+3.80	+1.60	+4.96	+4.24
PatchNet*	97.05	94.00	86.32	26.03	15.05	12.74	18.54	10.46	8.71
(+ 3D scores)	97.63	94.34	87.17	30.03	20.68	17.60	22.01	15.37	12.83
(+ OBMO-Linear Label Scores)	97.57	94.22	87.11	32.41	22.75	19.56	24.40	16.63	14.53
(+ OBMO-IoU Label Scores)	97.41	94.12	87.03	32.72	22.58	19.24	24.95	16.60	14.20
Improvements	+0.52	+0.22	+0.79	+6.38	+7.70	+6.82	+5.86	+6.17	+5.82
Pseudo-LiDAR† [28]	97.05	94.00	86.32	37.77	21.31	17.92	25.19	12.72	10.22
(+ OBMO-Linear Label Scores)	95.18	92.68	86.05	36.76	24.72	21.01	24.39	16.07	13.32
(+ OBMO-IoU Label Scores)	96.27	92.90	86.14	39.55	25.40	20.85	25.78	16.03	13.13
Improvements	-0.78	-1.10	-0.18	+1.78	+4.09	+2.93	+0.59	+3.31	+2.91
RTM3D†	97.02	91.49	83.98	21.27	15.84	13.63	15.05	11.42	9.66
(+ OBMO-Linear Label Scores)	97.03	91.44	83.93	23.25	17.60	14.93	16.50	12.74	10.59
(+ OBMO-IoU Label Scores)	97.03	91.45	83.93	22.99	17.52	14.83	16.29	12.75	10.54
Improvements	+0.01	-0.05	-0.05	+1.98	+1.76	+1.30	+1.45	+1.32	+0.93
Ground-aware [35]	-	-	-	28.95	20.11	15.51	22.80	15.41	11.43
(+ OBMO-Linear Label Scores)	96.79	84.04	64.26	31.22	21.96	16.85	23.48	16.59	12.39
(+ OBMO-IoU Label Scores)	96.83	81.72	61.90	29.58	21.75	16.73	22.64	16.40	12.22
Improvements	-	-	-	+2.27	+1.85	+1.34	+0.68	+1.18	+0.96
GUPNet [33]	-	-	-	31.07	22.94	19.75	22.76	16.46	13.72
(+ OBMO-Linear Label Scores)	96.67	88.67	78.85	33.09	23.63	20.42	24.65	17.80	15.15
(+ OBMO-IoU Label Scores)	96.58	88.64	78.92	32.20	23.88	20.67	24.48	17.94	15.26
Improvements	-	-	-	+1.13	+0.94	+0.92	+1.72	+1.48	+1.54

on $AP|_{R_{11}}$ in their paper, we evaluated them on $AP|_{R_{40}}$ by their public models (use † to represent). The improvements indicate that OBMO can be applied both in direct regression-based and depth-aware methods. The results show that the improvement in depth-aware methods is more remarkable than in direct regression-based methods. Specifically, for PatchNet, we improve the AP_{BEV}/AP_{3D} from 25.43%/17.71% to 29.33%/22.67% under the moderate setting. We think that OBMO might mitigate the influence of the worse monocular depth estimation to a certain extent. For direct regression-based methods such as RTM3D, the original detector is boosted by 1.82%/2.14% in AP_{BEV}/AP_{3D} under the moderate setting. Such significant improvements demonstrate the effectiveness and robustness of our method. We also present the 2D mAP in Table II. The 2D performances with and without OBMO are similar because the reasonable pseudo labels produced by the OBMO module are along the viewing frustum. Note the 2D detectors used in PatchNet and Pseudo-LiDAR are both Faster-RCNN, so their 2D mAPs are the same.

We further investigate the depth loss curve and the mAP curves after adding the OBMO module in training. As in Figure 5, we can easily see that the method employing OBMO has a smoother learning curve. By contrast, the curve of the original detector is unstable and contains many strong oscillations. It indicates that OBMO endows the network to steadily learn depth, stabilizing the overall learning process and thus bringing apparent improvements. As for the mAP curves shown in Figure 6, it is not difficult to find that without the OBMO module, the mAP of the training set is

TABLE III
ABLATION STUDY ON EACH COMPONENT IN OBMO. THE 3D SCORE [54] COMBINES THE 2D SCORE WITH 3D INFORMATION TO REPRESENT THE SCORE OF AN OBJECT.

3D scores	OBMO	IoU Label Scores	Linear Label Scores	$AP_{BEV}(\%)$			$AP_{3D}(\%)$		
				Easy	Mod.	Hard	Easy	Mod.	Hard
				26.03	15.05	12.74	18.54	10.46	8.71
✓				30.03	20.68	17.60	22.01	15.37	12.83
✓	✓			31.13	22.14	19.02	23.87	16.41	14.25
✓	✓	✓		32.72	22.58	19.24	24.95	16.60	14.20
✓	✓		✓	32.41	22.75	19.56	24.40	16.63	14.53

TABLE IV
ABLATION STUDY ON DIFFERENT CONSTRAINTS.

under X-Z ratio	under Y-Z ratio	$AP_{BEV}(\%)$			$AP_{3D}(\%)$		
		Easy	Mod.	Hard	Easy	Mod.	Hard
		30.03	20.68	17.60	22.01	15.37	12.83
✓		31.12	22.30	19.18	23.36	16.36	14.09
	✓	28.73	20.44	17.34	21.31	15.28	12.64
✓	✓	32.41	22.75	19.56	24.40	16.63	14.53

still rising, while there is almost no fluctuation in the mAP of the validation set or even a slight decline. It can be illustrated that OBMO overcomes overfitting to a certain extent.

D. Ablation Studies

We take PatchNet* as our baseline detector in the ablation study to save training time. By default, we set depth offset Δ_z to $\{-8\%, -4\%, +4\%, +8\%\}$ and use Linear Label Score. X-Z ratio and Y-Z ratio are both regarded as constraints.

Validity of Each Component. To study the impact brought by each component of OBMO, we investigate them through

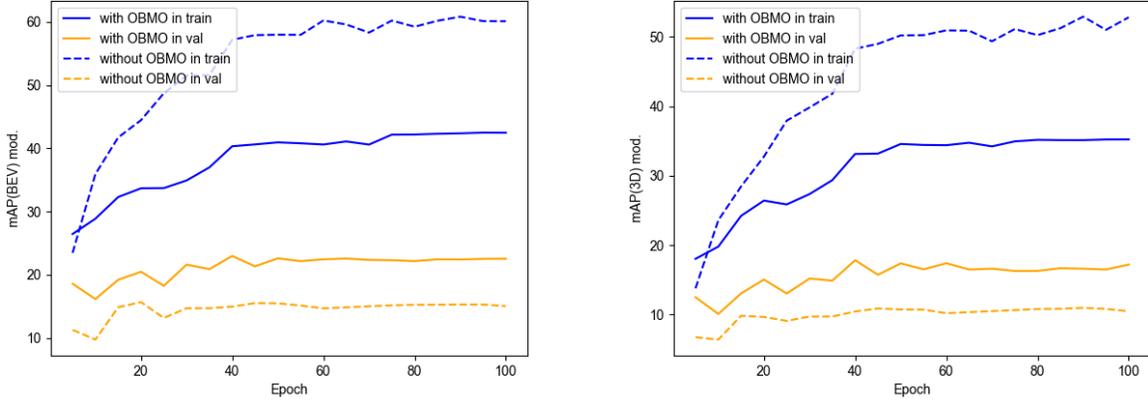


Fig. 6. The mAP (3D/BEV) with/without OBMO based on PatchNet in training and validation set. We can see that OBMO can overcome overfitting to a certain extent.

TABLE V

ABLATION STUDY ON THE DEPTH OFFSET UNDER THE SAME NUMBER OF PSEUDO LABELS.

Offset(%)	$AP_{BEV}(\%)$			$AP_{3D}(\%)$		
	Easy	Mod.	Hard	Easy	Mod.	Hard
0	30.03	20.68	17.60	22.01	15.37	12.83
2	30.78	22.09	18.84	23.78	16.45	14.04
4	32.41	22.75	19.56	24.40	16.63	14.53
6	29.35	21.85	18.72	22.11	15.98	13.86
8	29.97	21.82	18.78	22.67	15.98	13.28

TABLE VI

ABLATION STUDY ON THE NUMBER OF PSEUDO LABELS IN THE SAME DEPTH OFFSET.

the number of pseudo labels	$AP_{BEV}(\%)$			$AP_{3D}(\%)$		
	Easy	Mod.	Hard	Easy	Mod.	Hard
0	30.03	20.68	17.60	22.01	15.37	12.83
2	31.17	21.96	18.83	24.43	16.43	13.64
4	32.41	22.75	19.56	24.40	16.63	14.53
6	30.41	22.09	18.99	22.98	16.02	13.29
8	31.14	22.32	19.04	23.65	16.41	14.21

extra experiments, as shown in Table III. The results show that each component of the OBMO module is effective. As components gradually increase, the final accuracy also increases accordingly. We can see that the initial performance (AP) is boosted from 15.05%/10.46% to 22.75%/16.63% under the moderate setting, which is rather impressive. Moreover, both IoU Label Score and Linear Label Score for pseudo labels work well, suggesting that the proposed soft pseudo label strategy is robust since it is not sensitive for specifically designed manners.

Different Constraints. Also, we investigate the impact brought by different constraints, namely, the X-Z ratio and Y-Z ratio. If we do not use the X-Z ratio or Y-Z ratio as constraints, X or Y will not change in our pseudo labels. The results are reported in Table IV. Ultimately, we achieve the best performance by using them in combination.

Depth Offsets Δ_z . We further show the influences of different

TABLE VII

ABLATION STUDY ON THE NUMBER OF PSEUDO LABELS IN THE SAME DEPTH RANGE.

the number of pseudo labels	$AP_{BEV}(\%)$			$AP_{3D}(\%)$		
	Easy	Mod.	Hard	Easy	Mod.	Hard
0	30.03	20.68	17.60	22.01	15.37	12.83
2	29.17	21.37	18.49	21.01	15.35	12.78
4	32.41	22.75	19.56	24.40	16.63	14.53
6	30.69	22.16	18.98	23.64	16.40	13.52
8	31.82	22.66	19.23	23.68	16.47	14.03

depth offsets. In particular, we have to choose a suitable depth offset carefully due to discrete depth values. However, there is a dilemma in making a choice. If the depth offset is too small, we have to add multiple pseudo labels in the reasonable depth range, and the computational complexity will increase dramatically. On the contrary, if the depth offset is too large, we will lose some reasonable pseudo labels, resulting in suboptimal performance. Therefore, according to statistical dimensions information of the KITTI dataset, we try four base offset values: {2%, 4%, 6%, 8%} with 4 pseudo labels. Specially, if we choose the base of 2%, then $\Delta_z = \{-4\%, -2\%, +2\%, +4\%\}$. We report the corresponding results in Table V. It shows that the proper depth offset is indeed desired.

Then we fix the base value of the depth offset to 4%, and change the number of pseudo labels. Specially, if we use four pseudo labels, then $\Delta_z = \{-8\%, -4\%, +4\%, +8\%\}$. The results are in Table VI. Adding six pseudo labels reduces performance compared to adding four pseudo labels. It means if the depth value is outside the reasonable range, the added pseudo labels do not help performance and can even be detrimental to the performance.

Intuitively, if the added pseudo labels are too dense, it will also decrease the performance. Therefore, we set the largest value of depth offset to $8\% \cdot Z$, and choose the number of pseudo labels from {2, 4, 6, 8}. Specially, if we use 4 pseudo labels, then $\Delta_z = \{-8\%, -4\%, +4\%, +8\%\}$. The results are shown in Table VII, which verifies the viewpoint.

In conclusion, the best choice of depth offset Δ_z for

TABLE VIII
ABLATION STUDY ON IMPLYING RANGE. NONE MEANS WE DON'T USE OBMO MODULE.

level \geq^*	$AP_{BEV}(\%)$			$AP_{3D}(\%)$		
	Easy	Mod.	Hard	Easy	Mod.	Hard
None	30.03	20.68	17.60	22.01	15.37	12.83
3	27.06	19.42	16.46	19.75	13.78	11.79
2	31.11	22.58	19.56	22.17	15.29	13.98
1	32.41	22.75	19.56	24.40	16.63	14.53

TABLE IX
ABLATION STUDY ON THE WEIGHT OF LABEL SCORE BRANCH.

Lambda	$AP_{BEV}(\%)$			$AP_{3D}(\%)$		
	Easy	Mod.	Hard	Easy	Mod.	Hard
0	31.13	22.14	19.02	23.87	16.41	14.25
0.5	31.98	22.46	19.31	23.05	16.21	13.47
1	32.41	22.75	19.56	24.40	16.63	14.53
2	30.07	21.15	18.72	23.19	16.20	13.46

Car category in KITTI dataset is $\{-8\%, -4\%, +4\%, +8\%\}$. We also use this setting in the Waymo dataset because the dimension of the Car category is similar.

Additionally, we evaluate the performance of applying OBMO to different difficulty levels of objects in Table VIII. The difficulty level of the object is defined according to the height of the 2D bounding box, occlusion, and truncation values. We can see that OBMO works well for all levels of objects, which means that the issue of depth ambiguity indeed exists and is widespread.

Weights of Label Score λ . We use different loss weights in the label score branch and report the results in Table IX. The model performs better when the weight of loss is set to 1.

Keep Dimension. To verify that changing dimension harms the performance, we change both the position and dimension. The results are shown in Table X. When we change the dimensions of the object, the performance drops drastically. We can't change the inherent property of objects. Otherwise, the pseudo labels are not reasonable. As for the positions that we modify, they are current state values and can be changed.

Generalization of OBMO. Furthermore, we verify the generalization of our OBMO method. On the one hand, we test on other categories: Pedestrian and Cyclist. On the other hand, we test on another larger dataset: Waymo.

For the first one, we use the same default configuration as in Car, i.e., four pseudo labels: $\Delta_z = \{-8\%, -4\%, +4\%, +8\%\}$ and IoU Label Score. The results are shown in Table XI. The IoU threshold we used is 0.5 for both of them. The improvements are apparent, proving that our method can apply to varied categories.

For the latter one, we take GUPNet [33] as our baseline and adopt the metrics with mAP and mAPH under the IoU threshold of 0.7 and 0.5, respectively. "Level 1" denotes the evaluation of the bounding boxes that contain more than 5 lidar points. "Level 2" denotes the evaluation of all bounding boxes. The results prove that our proposed OBMO method achieves consistent improvements in all settings, as shown in Table XII.

TABLE X
ABLATION STUDY OF THE INFLUENCE IN CHANGING DIMENSION.

Change	$AP_{BEV}(\%)$			$AP_{3D}(\%)$		
	Easy	Mod.	Hard	Easy	Mod.	Hard
None	30.03	20.68	17.60	22.01	15.37	12.83
Dimension	29.06	16.90	14.15	21.95	11.85	9.58
Position	32.41	22.75	19.56	24.40	16.63	14.53
Both	7.07	6.86	5.92	3.28	3.40	3.00

TABLE XI
 $AP_{3D|40}$ ON PEDESTRIAN AND CYCLIST ON KITTI VALIDATION SET.

Categories	Methods	$AP_{BEV}(\%)$			$AP_{3D}(\%)$		
		Easy	Mod.	Hard	Easy	Mod.	Hard
Pedestrian	PatchNet	10.55	8.23	6.48	8.82	6.82	5.14
	(+ OBMO)	16.17	11.93	9.38	12.80	9.55	7.40
	Improvements	+5.62	+3.70	+2.89	+3.98	+2.73	+2.25
Cyclist	PatchNet	7.47	3.64	3.03	5.83	2.87	2.60
	(+ OBMO)	9.30	4.72	4.45	7.81	4.06	3.60
	Improvements	+1.83	+1.08	+1.42	+1.98	+1.19	+1.00

E. Qualitative Results

To visually evaluate the performance of our method based on GUPNet, we illustrate some examples in Figure 7. To clearly show the position of objects in the 3D world space, we also visualize the LiDAR signals and the ground-truth 3D boxes. We can observe that our outputs are remarkably accurate for the cases at a reasonable distance. Unfortunately, it remains a challenge for occluded and truncated objects, a common dilemma for most monocular 3D detectors.

V. LIMITATIONS AND FUTURE WORK

Although our work tries to alleviate the effect of the depth ambiguity problem, the prediction of depth in monocular images is still an ill-posed problem. The occluded and truncated objects which drop some pixel information are even more challenging to be detected. Our OBMO module only allows the network to learn a reasonable depth range, making depth prediction more flexible. And it can not improve the confidence of an object. If an object with a low 3D object score, it is still difficult to know its depth range. We will consider the above situations in future work.

VI. CONCLUSION

In this paper, we point out that it is hard to predict depth accurately due to the enormous 3D space. According to this discovery, we design a simple but elegant plug-and-play module OBMO. We add pseudo labels under the X-Z ratio and Y-Z ratio, and design two kinds of label scores: IoU Label Score and Linear Label Score. Compared with existing monocular 3D object detection methods, OBMO achieves better performance on challenging KITTI and Waymo benchmarks.

ACKNOWLEDGMENTS

This work was supported in part by The National Nature Science Foundation of China (Grant Nos: 62036009, 62273302, 62273303, 62303406, 61936006), in part by Ningbo Key R&D Program (No.2023Z231, 2023Z229), in

TABLE XII
EXPERIMENTAL RESULTS OF THE CAR CATEGORY ON THE WAYMO OPEN DATASET VALIDATION SET.

Difficulty	Threshold	Method	Overall	3D mAP / 3D mAPH		
				0 – 30 m	30 – 50 m	50 – ∞
LEVEL 1	IoU=0.7	PatchNet	0.39/0.37	1.67/1.63	0.13/0.12	0.03/0.03
		CADDN [23]	5.03/4.99	14.54/14.43	1.47/1.45	0.10/0.10
		PCT [62]	0.89/0.88	3.18/3.15	0.27/0.27	0.07/0.07
		MonoJSG [12]	0.97/0.95	4.65/4.59	0.55/0.53	0.10/0.09
		GUPNet [33] (+ OBMO)	8.00/7.94	22.71/22.54	3.17/3.15	0.39/0.38
	IoU=0.5	PatchNet	2.92/2.74	10.03/9.75	1.09/0.96	0.23/0.18
		CADDN [23]	17.54/17.31	45.00/44.46	9.24/9.11	0.64/0.62
		PCT [62]	4.20/4.15	14.70/14.54	1.78/1.75	0.39/0.39
		MonoJSG [12]	5.65/5.47	20.86/20.26	3.91/3.79	0.97/0.92
		GUPNet [33] (+ OBMO)	17.52/17.37	43.95/43.59	11.33/11.24	1.04/1.03
LEVEL 2	IoU=0.7	PatchNet	0.38/0.36	1.67/1.63	0.13/0.11	0.03/0.03
		CADDN [23]	4.49/4.45	14.50/14.38	1.42/1.41	0.09/0.09
		PCT [62]	0.66/0.66	3.18/3.15	0.27/0.26	0.07/0.07
		MonoJSG [12]	0.91/0.89	4.64/4.65	0.55/0.53	0.09/0.09
		GUPNet [33] (+ OBMO)	7.57/7.51	22.64/22.47	3.10/3.08	0.36/0.36
	IoU=0.5	PatchNet	2.42/2.28	10.01/9.73	1.07/0.94	0.22/0.16
		CADDN [23]	16.51/16.28	44.87/44.33	8.99/8.86	0.58/0.55
		PCT [62]	4.03/3.99	14.67/14.51	1.74/1.71	0.36/0.35
		MonoJSG [12]	5.34/5.17	20.79/20.19	3.79/3.67	0.85/0.82
		GUPNet [33] (+ OBMO)	16.41/16.28	43.80/43.44	10.99/10.91	0.90/0.90
			19.41/19.23	47.91/47.47	16.46/16.35	0.59/0.59

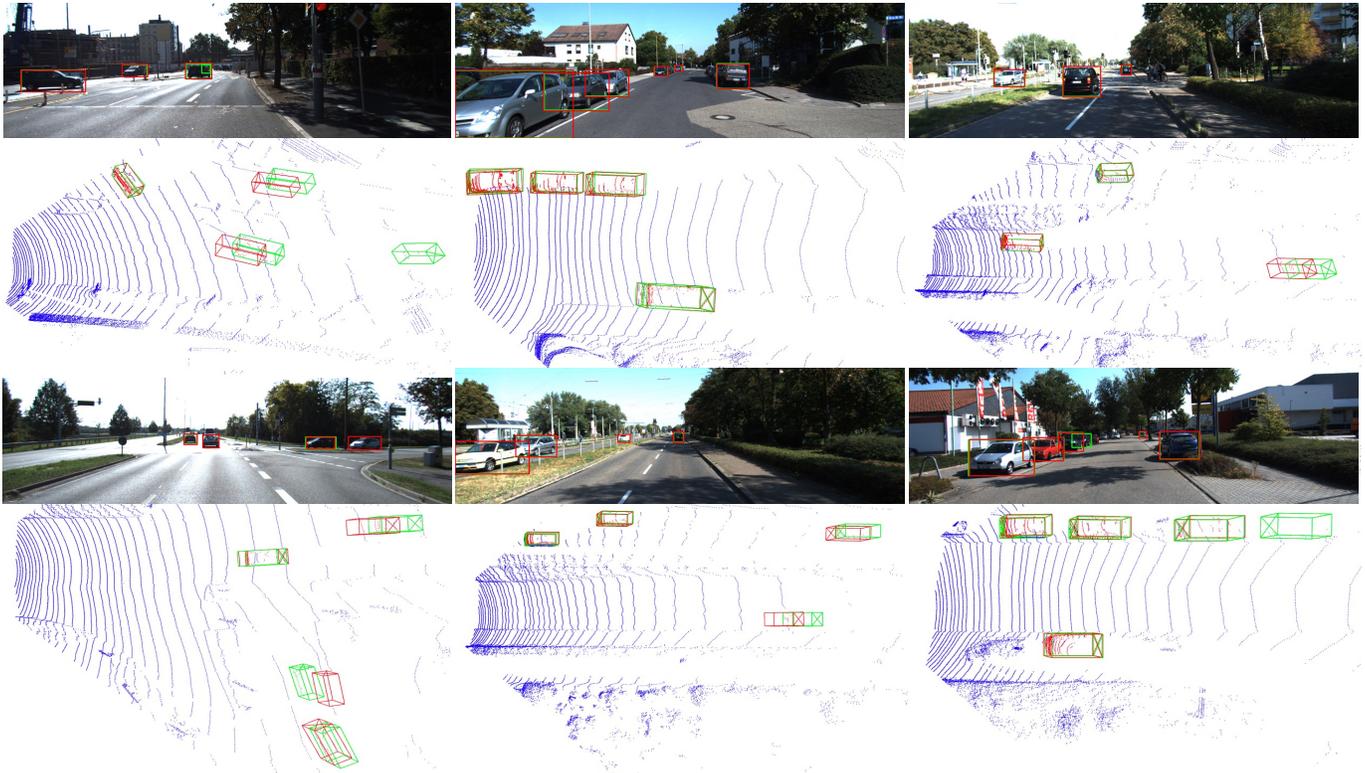


Fig. 7. Qualitative results on the KITTI val set. LiDAR point clouds are plotted for reference but not used in our method. We use green and red to denote predictions and ground truths, respectively.

part by the Key R&D Program of Zhejiang Province, China (2023C01135), in part by Yongjiang Talent Introduction Programme (Grant No: 2022A-240-G, 2023A-194-G). We sincerely thank Liang Peng, Minghao Chen, Menghao Guo for their suggestion in writing and Zhou Yang for the color scheme on the figures.

REFERENCES

- [1] M. Bertozzi, A. Broggi, and A. Fascioli, "Vision-based intelligent vehicles: State of the art and perspectives," *Robotics and Autonomous systems*, vol. 32, no. 1, pp. 1–16, 2000.
- [2] J. Janai, F. Güney, A. Behl, A. Geiger *et al.*, "Computer vision for autonomous vehicles: Problems, datasets and state of the art," *Foundations and Trends® in Computer Graphics and Vision*, vol. 12, no. 1–3, pp. 1–308, 2020.
- [3] X. Chen, H. Ma, J. Wan, B. Li, and T. Xia, "Multi-view 3d object detection network for autonomous driving," in *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, 2017, pp. 1907–1915.
- [4] D. Rey, G. Subsol, H. Delingette, and N. Ayache, "Automatic detection and segmentation of evolving processes in 3d medical images: Application to multiple sclerosis," *Medical image analysis*, vol. 6, no. 2, pp. 163–179, 2002.
- [5] C. Chen, A. Seff, A. Kornhauser, and J. Xiao, "Deepdriving: Learning affordance for direct perception in autonomous driving," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 2722–2730.
- [6] Q. Xu, Y. Zhou, W. Wang, C. R. Qi, and D. Anguelov, "Spg: Unsupervised domain adaptation for 3d object detection via semantic point generation," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 15 446–15 456.
- [7] Q. Xu, Y. Zhong, and U. Neumann, "Behind the curtain: Learning occluded shapes for 3d object detection," *arXiv preprint arXiv:2112.02205*, 2021.
- [8] W. Zheng, W. Tang, L. Jiang, and C.-W. Fu, "Se-ssd: Self-ensembling single-stage object detector from point cloud," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 14 494–14 503.
- [9] W. Bao, B. Xu, and Z. Chen, "Monofenet: Monocular 3d object detection with feature enhancement networks," *IEEE Transactions on Image Processing*, vol. 29, pp. 2753–2765, 2019.
- [10] 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*, 2021, pp. 3142–3152.
- [11] X. Liu, N. Xue, and T. Wu, "Learning auxiliary monocular contexts helps monocular 3d object detection," *arXiv preprint arXiv:2112.04628*, 2021.
- [12] Q. Lian, P. Li, and X. Chen, "Monojsjg: Joint semantic and geometric cost volume for monocular 3d object detection," *CoRR*, vol. abs/2203.08563, 2022. [Online]. Available: <https://doi.org/10.48550/arXiv.2203.08563>
- [13] K. Huang, T. Wu, H. Su, and W. H. Hsu, "Monodtr: Monocular 3d object detection with depth-aware transformer," *CoRR*, vol. abs/2203.10981, 2022. [Online]. Available: <https://doi.org/10.48550/arXiv.2203.10981>
- [14] M.-Q. V. Bui, D. T. Ngo, H.-A. Pham, and D. D. Nguyen, "Gac3d: improving monocular 3d object detection with ground-guide model and adaptive convolution," *PeerJ Computer Science*, vol. 7, p. e686, 2021.
- [15] D. Beker, H. Kato, M. A. Morariu, T. Ando, T. Matsuoka, W. Kehl, and A. Gaidon, "Monocular differentiable rendering for self-supervised 3d object detection," in *European Conference on Computer Vision*. Springer, 2020, pp. 514–529.
- [16] L. Wang, L. Zhang, Y. Zhu, Z. Zhang, T. He, M. Li, and X. Xue, "Progressive coordinate transforms for monocular 3d object detection," *Advances in Neural Information Processing Systems*, vol. 34, 2021.
- [17] F. Nobis, F. Brunhuber, S. Janssen, J. Betz, and M. Lienkamp, "Exploring the capabilities and limits of 3d monocular object detection—a study on simulation and real world data," in *2020 IEEE 23rd International Conference on Intelligent Transportation Systems (ITSC)*. IEEE, 2020, pp. 1–8.
- [18] G. Brazil and X. Liu, "M3d-rpn: Monocular 3d region proposal network for object detection," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 9287–9296.
- [19] 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*, 2017, pp. 7074–7082.
- [20] X. Chen, K. Kundu, Z. Zhang, H. Ma, S. Fidler, and R. Urtasun, "Monocular 3d object detection for autonomous driving," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 2147–2156.
- [21] T. He and S. Soatto, "Mono3d++: Monocular 3d vehicle detection with two-scale 3d hypotheses and task priors," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, 2019, pp. 8409–8416.
- [22] Z. Qin, J. Wang, and Y. Lu, "Monogrnet: A geometric reasoning network for monocular 3d object localization," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, 2019, pp. 8851–8858.
- [23] C. Reading, A. Harakeh, J. Chae, and S. L. Waslander, "Categorical depth distribution network for monocular 3d object detection," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 8555–8564.
- [24] 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*, 2021, pp. 3289–3298.
- [25] 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*, 2021, pp. 4721–4730.
- [26] J. Gu, B. Wu, L. Fan, J. Huang, S. Cao, Z. Xiang, and X. Hua, "Homography loss for monocular 3d object detection," *CoRR*, vol. abs/2204.00754, 2022. [Online]. Available: <https://doi.org/10.48550/arXiv.2204.00754>
- [27] X. Ma, S. Liu, Z. Xia, H. Zhang, X. Zeng, and W. Ouyang, "Rethinking pseudo-lidar representation," in *European Conference on Computer Vision*. Springer, 2020, pp. 311–327.
- [28] 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*, 2019, pp. 8445–8453.
- [29] Y. You, Y. Wang, W.-L. Chao, D. Garg, G. Pleiss, B. Hariharan, M. Campbell, and K. Q. Weinberger, "Pseudo-lidar++: Accurate depth for 3d object detection in autonomous driving," *arXiv preprint arXiv:1906.06310*, 2019.
- [30] A. Naiden, V. Paunescu, G. Kim, B. Jeon, and M. Leordeanu, "Shift r-cnn: Deep monocular 3d object detection with closed-form geometric constraints," in *2019 IEEE International Conference on Image Processing (ICIP)*, 2019, pp. 61–65.
- [31] 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 *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020*. AAAI Press, 2020, pp. 10 478–10 485. [Online]. Available: <https://ojs.aaai.org/index.php/AAAI/article/view/6618>
- [32] Y. Chen, L. Tai, K. Sun, and M. Li, "Monopair: Monocular 3d object detection using pairwise spatial relationships," in *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA, June 13-19, 2020*. Computer Vision Foundation / IEEE, 2020, pp. 12 090–12 099. [Online]. Available: https://openaccess.thecvf.com/content_CVPR_2020/html/Chen_MonoPair_Monocular_3D_Object_Detection_Using_Pairwise_Spatial_Relationships_CVPR_2020_paper.html
- [33] 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*, 2021, pp. 3111–3121.
- [34] X. Shi, Q. Ye, X. Chen, C. Chen, Z. Chen, and T. Kim, "Geometry-based distance decomposition for monocular 3d object detection," in *2021 IEEE/CVF International Conference on Computer Vision, ICCV 2021, Montreal, QC, Canada, October 10-17, 2021*. IEEE, 2021, pp. 15 152–15 161. [Online]. Available: <https://doi.org/10.1109/ICCV48922.2021.01489>
- [35] 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.

- [36] G. E. Hinton, O. Vinyals, and J. Dean, “Distilling the knowledge in a neural network,” *CoRR*, vol. abs/1503.02531, 2015. [Online]. Available: <https://arxiv.org/abs/1503.02531>
- [37] 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*. Springer, 2020, pp. 644–660.
- [38] Z. Yang, Y. Sun, S. Liu, and J. Jia, “3dssd: Point-based 3d single stage object detector,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 11 040–11 048.
- [39] A. H. Lang, S. Vora, H. Caesar, L. Zhou, J. Yang, and O. Beijbom, “Pointpillars: Fast encoders for object detection from point clouds,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 12 697–12 705.
- [40] S. Shi, C. Guo, L. Jiang, Z. Wang, J. Shi, X. Wang, and H. Li, “Pv-rnn: Point-voxel feature set abstraction for 3d object detection,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 10 529–10 538.
- [41] M.-H. Guo, J.-X. Cai, Z.-N. Liu, T.-J. Mu, R. R. Martin, and S.-M. Hu, “Pct: Point cloud transformer,” *Computational Visual Media*, vol. 7, no. 2, pp. 187–199, 2021.
- [42] Y. Zhou and O. Tuzel, “Voxelnet: End-to-end learning for point cloud based 3d object detection,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 4490–4499.
- [43] J. Deng, S. Shi, P. Li, W. Zhou, Y. Zhang, and H. Li, “Voxel rnn: Towards high performance voxel-based 3d object detection,” *arXiv preprint arXiv:2012.15712*, 2020.
- [44] C. Peng, X. Zhang, G. Yu, G. Luo, and J. Sun, “Large kernel matters—improve semantic segmentation by global convolutional network,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 4353–4361.
- [45] M.-H. Guo, C.-Z. Lu, Z.-N. Liu, M.-M. Cheng, and S.-M. Hu, “Visual attention network,” *arXiv preprint arXiv:2202.09741*, 2022.
- [46] Y. Yan, Y. Mao, and B. Li, “SECOND: sparsely embedded convolutional detection,” *Sensors*, vol. 18, no. 10, p. 3337, 2018. [Online]. Available: <https://doi.org/10.3390/s18103337>
- [47] S. Shi, X. Wang, and H. Li, “Pointnncnn: 3d object proposal generation and detection from point cloud,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 770–779.
- [48] W. Shi and R. Rajkumar, “Point-gnn: Graph neural network for 3d object detection in a point cloud,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 1711–1719.
- [49] 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*, 2017, pp. 652–660.
- [50] Z. Liu, Z. Wu, and R. Tóth, “Smoke: Single-stage monocular 3d object detection via keypoint estimation,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, 2020, pp. 996–997.
- [51] Y. Tang, S. Dorn, and C. Savani, “Center3d: Center-based monocular 3d object detection with joint depth understanding,” in *DAGM German Conference on Pattern Recognition*. Springer, 2020, pp. 289–302.
- [52] M. Ding, Y. Huo, H. Yi, Z. Wang, J. Shi, Z. Lu, and P. Luo, “Learning depth-guided convolutions for monocular 3d object detection,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, 2020, pp. 1000–1001.
- [53] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna, “Rethinking the inception architecture for computer vision,” in *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27–30, 2016*. IEEE Computer Society, 2016, pp. 2818–2826. [Online]. Available: <https://doi.org/10.1109/CVPR.2016.308>
- [54] L. Peng, F. Liu, S. Yan, X. He, and D. Cai, “Ocm3d: Object-centric monocular 3d object detection,” *arXiv preprint arXiv:2104.06041*, 2021.
- [55] P. Li and H. Zhao, “Monocular 3d detection with geometric constraint embedding and semi-supervised training,” *IEEE Robotics and Automation Letters*, vol. 6, no. 3, pp. 5565–5572, 2021.
- [56] 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*, 2021, pp. 10 379–10 388.
- [57] A. Kumar, G. Brazil, and X. Liu, “Groomed-nms: Grouped mathematically differentiable nms for monocular 3d object detection,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 8973–8983.
- [58] Y. Zhou, Y. He, H. Zhu, C. Wang, H. Li, and Q. Jiang, “Monocular 3d object detection: An extrinsic parameter free approach,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 7556–7566.
- [59] A. Geiger, P. Lenz, and R. Urtasun, “Are we ready for autonomous driving? the kitti vision benchmark suite,” in *2012 IEEE conference on computer vision and pattern recognition*. IEEE, 2012, pp. 3354–3361.
- [60] P. Sun, H. Kretzschmar, X. Dotiwalla, A. Chouard, V. Patnaik, P. Tsui, J. Guo, Y. Zhou, Y. Chai, B. Caine, V. Vasudevan, W. Han, J. Ngiam, H. Zhao, A. Timofeev, S. Ettinger, M. Krivokon, A. Gao, A. Joshi, Y. Zhang, J. Shlens, Z. Chen, and D. Anguelov, “Scalability in perception for autonomous driving: Waymo open dataset,” in *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA, June 13–19, 2020*. Computer Vision Foundation / IEEE, 2020, pp. 2443–2451. [Online]. Available: https://openaccess.thecvf.com/content_CVPR_2020/html/Sun_Scalability_in_Perception_for_Autonomous_Driving_Waymo_Open_Dataset_CVPR_2020_paper.html
- [61] 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*. Citeseer, 2015, pp. 424–432.
- [62] L. Wang, L. Zhang, Y. Zhu, Z. Zhang, T. He, M. Li, and X. Xue, “Progressive coordinate transforms for monocular 3d object detection,” *CoRR*, vol. abs/2108.05793, 2021. [Online]. Available: <https://arxiv.org/abs/2108.05793>



Chenxi Huang is a Ph.D. candidate supervised by Prof. Deng Cai in the State Key Lab of CAD&CG, College of Computer Science at Zhejiang University, China. Her research interests include computer vision and 3D scene understanding.



Deng Cai is a Professor in the State Key Lab of CAD&CG, College of Computer Science at Zhejiang University, China. He received a Ph.D. degree in computer science from the University of Illinois at Urbana Champaign in 2009. His research interests include machine learning, data mining and information retrieval.



Tong He received the Ph.D. degree in computer science from the University of Adelaide, Australia, in 2020. He is currently a researcher at Shanghai AI Laboratory. His research interests include computer vision and machine learning.



Haidong Ren is the deputy chief engineer of Ningbo Zhoushan Port Group Co.,Ltd, Ningbo, China. He is the director of Technology and Information Management Department, the director of the Science and Technology Center, a senior engineer, a post-doctoral supervisor, and a member of the Digital Reform Expert Group of the Zhejiang SASAC. He has long been engaged in the research and practice of port equipment and port informatization.



Wenxiao Wang is an assistant professor, School of Software Technology at Zhejiang University, China. He received the Ph.D. degree in computer science and technology from Zhejiang University in 2022. His research interests include deep learning and computer vision.



BinBin Lin is an assistant professor in the School of Software Technology at Zhejiang University, China. He received a Ph.D degree in computer science from Zhejiang University in 2012. His research interests include machine learning and decision making.