

MonoGround: Detecting Monocular 3D Objects from the Ground

Zequn Qin¹, Xi Li^{1,2,3*}

¹College of Computer Science, Zhejiang University

²Shanghai Institute for Advanced Study of Zhejiang University; ³Shanghai AI Lab

zequnqin@gmail.com, xilizju@zju.edu.cn

Abstract

Monocular 3D object detection has attracted great attention for its advantages in simplicity and cost. Due to the ill-posed 2D to 3D mapping essence from the monocular imaging process, monocular 3D object detection suffers from inaccurate depth estimation and thus has poor 3D detection results. To alleviate this problem, we propose to introduce the ground plane as a prior in the monocular 3D object detection. The ground plane prior serves as an additional geometric condition to the ill-posed mapping and an extra source in depth estimation. In this way, we can get a more accurate depth estimation from the ground. Meanwhile, to take full advantage of the ground plane prior, we propose a depth-align training strategy and a precise two-stage depth inference method tailored for the ground plane prior. It is worth noting that the introduced ground plane prior requires no extra data sources like LiDAR, stereo images, and depth information. Extensive experiments on the KITTI benchmark show that our method could achieve state-of-the-art results compared with other methods while maintaining a very fast speed. Our code, models, and training logs are available at <https://github.com/cfzd/MonoGround>.

1. Introduction

3D object detection is a fundamental computer vision task that aims to obtain the locations, sizes, and orientations of objects. To get the real-world 3D information, many methods adopt modalities like point clouds from LiDAR, stereo images, and depth images, which require extra sensors of data sources. Different from them, monocular 3D object detection, which only requires a single 2D image and camera calibration information, increasingly draws the community's attention for its superiorities in simplicity and cost, especially in the autonomous driving field.

Despite the superiorities of monocular 3D object detection, obtaining 3D information from a single 2D image is

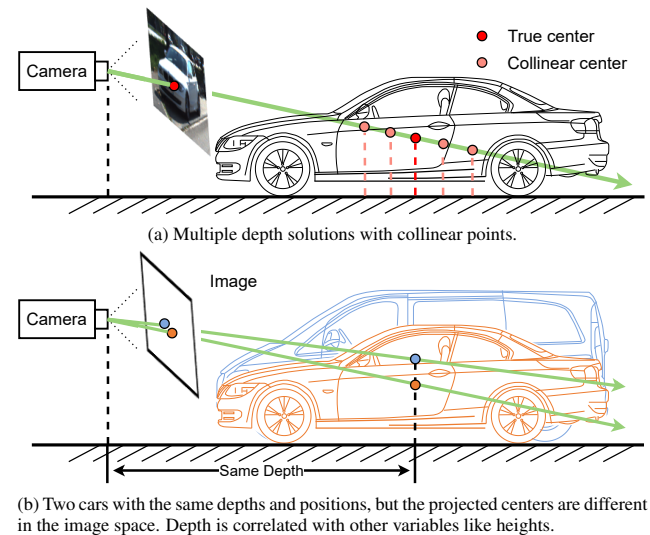


Figure 1. Difficulties in the monocular 3D object detection.

essentially hard for the following reasons: 1) the mapping from 2D to 3D is an ill-posed problem since a location in the 2D image plane corresponds to all collinear 3D positions, which are in the ray from the optical center to the 2D location. This one-to-many property makes predicting the depths of objects difficult, as shown in Fig. 1a. 2) The expression of 3D object's depth is correlated with the height. For two objects with the same depths and positions, if the heights are different, the projected center would differ, as shown in Fig. 1b. In this way, the learning of depths has to overcome the interference with irrelevant attributes of objects. 3) Mainstream monocular 3D object detection methods [19, 22, 42, 44] are commonly based on the CenterNet [43]. Under this framework, all information of each object is represented with a point, including the depth information. In this way, the supervision of depth is very sparse during training (Suppose there are two cars in an image, then only two points on the depth map would be trained, and all other points on the depth map are ignored). Such sparse supervision would lead to insufficient learning of depth estimation, which significantly differs from common monocular depth estimation tasks with dense depth supervision.

*Corresponding author.

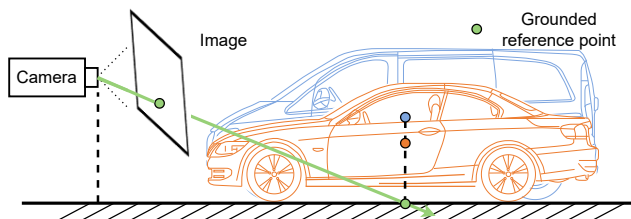


Figure 2. Illustration of the depth on the ground.

To address the above problems, we propose to introduce the ground plane prior to the monocular 3D object detection task. With the ground plane prior, the ill-posed 2D to 3D mapping becomes a well-posed problem with a unique solution. The unique solution is determined by the intersection of the camera ray and the ground plane. The expression of depth is no longer correlated with the object’s height, as shown in Fig. 2. Moreover, since the ground plane is introduced and utilized, we could expand the original sparse depth supervision to a dense depth supervision by using the dense depth from the ground plane.

With the above motivations, we propose our monocular 3D object detection model with the dense ground plane prior, termed as MonoGround. Within this formulation, we propose a depth-align method to effectively learn and predict dense grounded depth. Meanwhile, with the help of dense predicted depth, we also propose a two-stage depth inference method, which brings finer-grained depth estimation. In summary, the main contribution of this work can be summarized as follows:

- We propose to introduce the ground plane prior to the monocular 3D object detection, which could alleviate ill-posed mapping, remove irrelevant correlation of depths, and provide dense depth supervision, without any extra data like LiDAR, stereo, and depth images.
- We propose a depth-align training strategy and a fine-grained two-stage depth inference method to take advantage of the introduced ground plane prior and achieve precise depth inference.
- Experiments on the KITTI dataset show the effectiveness of introducing the ground plane prior and the proposed method, and our method achieves the SOTA performance with a very fast speed in real-time.

2. Related Work

Monocular 3D Object Detection Monocular 3D object detection methods can be roughly divided into two groups. The first kind of method try to utilize extra available data sources to simplify the detection of 3D objects. The extra data sources can be LiDAR, depth images, and CAD models. The most commonly used extra data sources are the LiDAR [12, 32, 33] and depth information [4, 10, 23, 24, 38, 40,

45] from pre-trained depth estimation models. For example, Pseudo-LiDAR [39] generates the pseudo LiDAR information from the monocular image and depth information. CaDDN [31] projects LiDAR point clouds into the image to create depth maps, then a categorical depth distribution is learned. Besides the depth and LiDAR data, there are also methods that try to utilize the CAD models [5, 20, 26, 28] to simplify the recognition and pose estimation of objects.

The second kind of method aim to detect 3D objects without any extra data [1, 15, 16, 27, 30]. For example, M3D-RPN [1] uses depth-aware convolutional layers to detect 3D objects. MonoPair [9] proposes to use the pairwise spatial relationships to achieve better results. SMOKE [19] proposes a CenterNet-style [43] 3D detector via keypoints estimation. Then, MonoFlex [42] proposes a flexible center definition that unifies truncated objects and regular objects and an uncertainty-based depth ensemble method. To better find the bottleneck of purely monocular detectors, MonoDLE [25] examines the effects of each component in mainstream methods. GrooMeD-NMS [14] introduces differential non-maximal suppression into the monocular 3D object detection. MonoRUn [6] proposes to use the dense correspondence between 2D and 3D space to learn the reconstruction and reproject processes. To better utilize the geometry relationship in 3D and 2D space accompanying uncertainty, GUPNet [22] proposes a geometry uncertainty projection method to reduce the error in depth estimation.

Ground Plane Knowledge in the Monocular 3D Object Detection

There have been several attempts in using the ground plane knowledge in monocular 3D object detection. Mono3D [7] first tries to use the ground plane to generate 3D bounding box proposals. Besides, Ground-Aware [18] introduces the ground plane in the geometric mapping and proposes a ground-aware convolution module to enhance the detection.

In the above works, the ground plane is defined based on a fixed rule, that the camera height on KITTI is 1.65 meters. In this way, all positions at a fixed height of -1.65 meters in the 3D space would construct the ground plane. Different from the strong hypothesis of the fixed ground plane at the height of -1.65 meters, in this work, we propose a learnable ground plane prior that is based on a more reasonable hypothesis: objects should lie on the ground. As long as the objects lie on the ground, the ground plane can be substituted with the bottom surface of object’s 3D bounding box. In this way, the proposed ground plane prior with the objects’ bottom surface is object-wise adaptive and precise.

Depth Estimation in the Monocular 3D Object Detection

There are two kinds of mainstream depth estimation methods in the purely monocular 3D object detection, which are direct regression [43] and geometry depth [4] derived from

the pinhole imaging model. Further, MonoFlex [42] extends the geometry depth to the diagonal keypoint depth by averaging the geometry depth of diagonal paired keypoints. Meanwhile, many works [6,22,34,42] adopt the uncertainty along with the depth estimation to get better results.

In this work, with the help of the proposed ground plane prior, we propose a novel depth estimation method different from the above methods. It is a two-stage depth inference method that could get a more precise depth estimation compared with the geometry depth. Moreover, the proposed depth estimation can also adopt uncertainty and be extended with diagonal paired keypoints.

3. Method

In this section, we elaborate on the details of our method. First, we give a brief problem definition of monocular 3D object detection. Second, we show how to utilize the ground plane prior and generate dense grounded depths. Third, the depth-align training strategy and two-stage depth inference method based on the ground plane prior are discussed.

3.1. Problem Definition

The monocular 3D object detection task is to detect 3D objects from monocular RGB images. Besides the RGB image, calibration information (camera parameter matrix) can also be adopted. Specifically, for each object, the 3D location (x, y, z) , size (h, w, l) , and orientation θ are required. On the KITTI [11] dataset, pitch and roll angles are considered as zero, so only orientation θ is considered.

Corresponding to the above targets, mainstream methods [19, 25, 42] divide the whole task into four subtasks, which are 2D location, depth, size, and orientation estimation tasks. As discussed in Sec. 1 and pointed out in [25], depth estimation is the key bottleneck for monocular 3D object detection. In this way, we focus on precise depth estimation in this work.

3.2. Ground Plane Prior

To define the ground plane, we start from a reasonable hypothesis that objects lie on the ground. This hypothesis holds for all common objects like cars, pedestrians, and cyclists. With this hypothesis, the ground plane can be substituted with the bottom surface of objects' 3D bounding box. For each object, we first get the bottom keypoints $k_1, k_2, k_3,$ and k_4 in the 3D space. Then we conduct random sampling and interpolation on the 3D bottom surface composed of the bottom keypoints. Denote $R \in \mathbb{R}^{N \times 2}$ as a random matrix with each element $R_{ij} \in [0, 1]$. N is the number of sampling points. The sampled dense points can be written as:

$$P_{3d} = R \begin{bmatrix} k_2 - k_1 \\ k_4 - k_1 \end{bmatrix} + k_1, \quad (1)$$

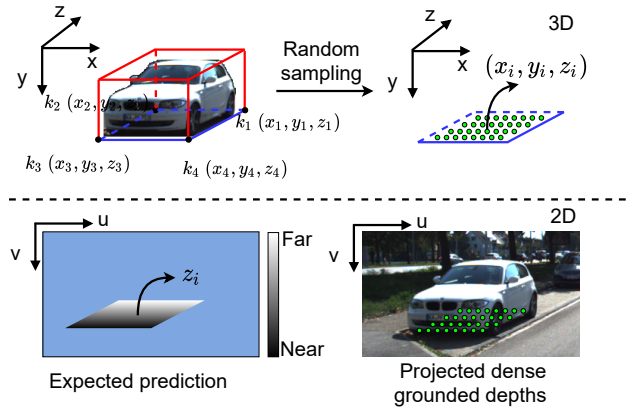


Figure 3. Illustration of the dense grounded depths. We first get the bottom of object's 3D bounding box, *i.e.*, the ground plane. A random sampling is conducted on the plane, and the sampled results are projected to the image space to get dense grounded depths.

in which $P_{3d} \in \mathbb{R}^{N \times 3}$ contains all sampled points. Suppose $K_{3 \times 3}$ is the camera parameter matrix. After the sampling, these points are projected back to image space:

$$P_{2d}^T = K_{3 \times 3} P_{3d}^T, \quad (2)$$

in which $P_{2d} \in \mathbb{R}^{N \times 3}$. The i -th row of P_{2d} is $[u_i \cdot z_i, v_i \cdot z_i, z_i]$, in which u_i and v_i are the projected coordinates in the image space, and z_i is the corresponding depth at this point. The illustration is shown in Fig. 3.

With the above formulation, we can get a dense and grounded depth representation, which addresses the problem of sparse depth supervision in mainstream methods. For the ill-posed 2D to 3D mapping problem, the dense grounded depth provides the ground plane constraint and makes it a well-posed problem, and the unique solution to the mapping is the intersection of camera ray and ground plane. Moreover, the irrelevant correlation problem is naturally avoided since our formulation has no relation to the object's height. More importantly, the introduced ground plane prior and dense grounded depth are derived from the 3D bounding box annotation, which can be easily accessed and requires no other data sources like LiDAR and depth.

3.3. Detecting Objects with Dense Grounded Depth

With the dense grounded depth, a natural question is how to use the ground plane prior to help the detecting of 3D objects. To answer this question, we start from three aspects in this section, which are network design, depth-align training, and two-stage depth inference.

Network Design The framework of our method is shown in Fig. 4. We use DLA-34 [41] as our backbone network. Then two branches are adopted, which are the ground branch and the 3D branch. The ground branch is to predict

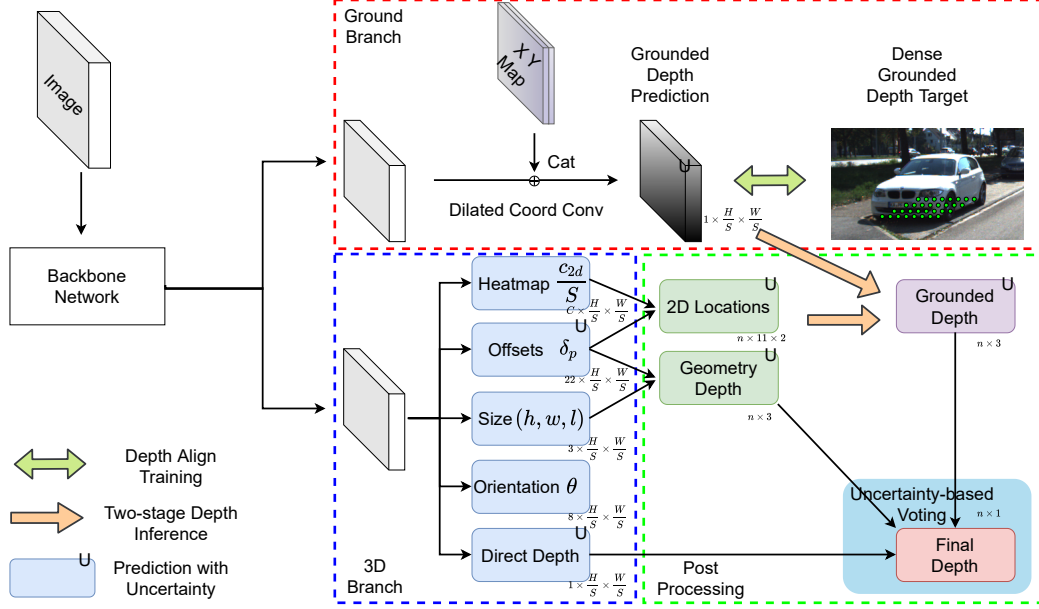


Figure 4. Framework design of MonoGround. There are three main parts, which are ground branch, 3D branch and post processing. In the post processing, n is the number of objects.

the grounded dense depth. The 3D branch is based on CenterNet [43] and is to predict attributes like location, size, orientation, *etc.*

In the 3D branch, we use MonoFlex [42] as our baseline and simultaneously predict five kinds of outputs. The first one is the heatmap, which is used to represent and distinguish each object using its center. We follow the same definition of object center c^{2d} as MonoFlex [42]. Heatmap can only predict rough locations due to quantization error. To get precise locations, we use offset maps to regress the fine-grained locations. Besides regressing the object center, we also use the offset map to regress other ten points, which are eight projected 3D bounding box vertices $\{k_1^{2d}, k_2^{2d}, \dots, k_8^{2d}\}$ and two projected bottom b^{2d} and top t^{2d} centers. The regression target of the offset map is shown in Fig. 5, and we have:

$$\delta_p = \frac{p}{S} - \lfloor \frac{c^{2d}}{S} \rfloor, \quad (3)$$

$$s.t. \quad p \in \{c^{2d}, k_1^{2d}, \dots, k_8^{2d}, b^{2d}, t^{2d}\},$$

in which S is the output stride of the heatmap, and $\lfloor \cdot \rfloor$ is the floor function. Suppose δ_p^* is the ground truth, the loss function of the offset map can be written as:

$$L_{offset} = \sum_{p \in \{c^{2d}, k_1^{2d}, \dots, k_8^{2d}, b^{2d}, t^{2d}\}} |\delta_p - \delta_p^*|. \quad (4)$$

Moreover, we also predict the size, orientation, and direct depth maps in the 3D branch as the baseline MonoFlex.

In the ground branch, we use dilated coord convolutions to get the final prediction of grounded depth. The rea-

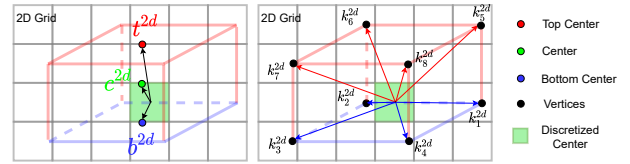


Figure 5. Illustration of the regression targets of the offset map.

son why we use dilated convolution is that the predicted grounded depth is mostly occluded by the object, as shown in Fig. 3. So we use dilated convolution to expand the receptive field of the network, and thus the surrounding ground plane can be seen by the network. In this way, the grounded depth can be better estimated. Moreover, coord convolutions [17] are also adopted for the following two reasons. First, the estimation of grounded depth is related to both the visual and positional features. Plain convolution can only provide visual features, while coord convolution can integrate both visual and positional features. Second, from Eq. (2) we can see that the 3D coordinates (x, y, z) (z is the depth) is a linear transformation of 2D coordinates. The coord convolutions can explicitly concatenate the 2D coordinates in the convolution and make it an external hint to realize a better estimation of depth.

In the post processing, we first get the 2D locations of the 11 keypoints as Eq. (3), then we get the geometry depth [4]:

$$z = \frac{f \cdot h}{h_{2d}}, \quad (5)$$

where f is the focal length, h is the predicted object height, and h_{2d} is the pixel height obtained from 2D locations b^{2d}

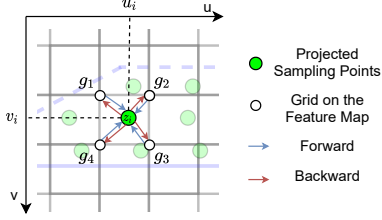


Figure 6. Demonstration of the depth-align training.

and t^{2d} . We also extend the geometry depth by averaging the geometry depth of diagonal keypoints as [42]. At last the final depth is obtained by an uncertainty-based voting from seven depths z_{pred} (1 direct depth, 3 geometry depths, 3 grounded depths that will be discussed in Fig. 7 and the two-stage depth inference paragraph):

$$z_{final} = \left(\sum_{i=1}^7 \frac{z_{pred}^i}{\sigma_i} \right) / \left(\sum_{i=1}^7 \frac{1}{\sigma_i} \right), \quad (6)$$

where σ_i is the predicted uncertainty along with the depth.

Depth-Align Training After obtaining the prediction map of grounded depth, how to train this map is a major problem. There exists an issue of misalignment between the prediction map and projected dense grounded depths. The projected dense grounded points are scattered across the prediction map with arbitrary locations, while the prediction map is composed of uniform grids. In other words, the prediction map can only work with integer locations, while the projected dense grounded points are distributed in non-integer locations. To solve this problem, we propose a depth-align training method, as shown in Fig. 6.

The forward calculation of depth-align training is the same as the bilinear interpolation. Suppose the projected grounded point is (u_i, v_i) with a depth of z_i , and $\{g_1, g_2, g_3, g_4\}$ are the four clockwise surrounding grid points from the left upper corner. Then the forward calculation can be written as:

$$\begin{aligned} \lambda_1 &= u_i - \lfloor u_i \rfloor, \quad \lambda_2 = \lceil u_i \rceil - u_i, \\ \lambda_3 &= v_i - \lfloor v_i \rfloor, \quad \lambda_4 = \lceil v_i \rceil - v_i, \\ pred_i &= \lambda_2 \lambda_4 g_1 + \lambda_1 \lambda_4 g_2 + \lambda_1 \lambda_3 g_3 + \lambda_2 \lambda_3 g_4, \end{aligned} \quad (7)$$

in which $\lfloor \cdot \rfloor$ and $\lceil \cdot \rceil$ are the floor and ceiling functions, and $pred_i$ is the prediction of depth. Then the loss can be:

$$L_{da} = \sum_{i=1}^N |z_i - pred_i|. \quad (8)$$

The variables z_i , u_i , and v_i are obtained from the P_{2d} in Eq. (2). In this way, we can directly optimize the dense depths with non-integer locations, and the backward calculation is to pass the gradient to the surrounding grids according to the weights of the bilinear interpolation.

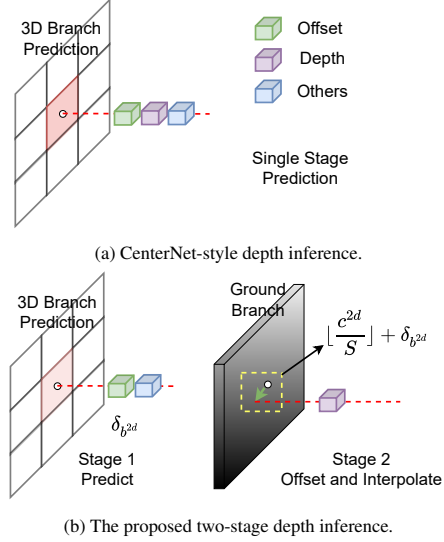


Figure 7. Illustration and comparison of different depth inference methods.

Two-stage Depth Inference In CenterNet-style [43] inference, objects are represented by coarse grid points on the feature map, and all other attributes are predicted from the features at the points at once. In this work, we call this one-stage inference. However, the coarse grid points can only reflect the rough location of objects, and this is why we need to further regress the difference between ground truth locations and rough grid locations.

With the help of the densely trained grounded depth map, we can achieve more precise depth estimation compared with the original CenterNet-style depth inference by using the regressed fine-grained locations instead of rough grid locations. Suppose $\lfloor c^{2d}/S \rfloor$ is the center predicted from the heatmap. We first apply the regressed offset to the center to obtain the fine location $b^{2d}/S = \lfloor c^{2d}/S \rfloor + \delta_{b^{2d}}$, then we use the interpolation as Eq. (7) on the prediction of ground branch to get the precise depth. Since b^{2d} is the bottom center of the object, the obtained grounded depth of b^{2d} is exactly the same as the depth of c^{2d} , i.e., the depth of object center. The illustration and comparison of different depth inference methods are shown in Fig. 7.

Similarly, we can obtain precise grounded depth for the keypoints $\{k_1^{2d}, k_2^{2d}, k_3^{2d}, k_4^{2d}\}$. According to the 3D geometry, the center depth equals to the average depth of the diagonal keypoints. In this way, we can get another two depth estimations from keypoints k_1^{2d}, k_3^{2d} and k_2^{2d}, k_4^{2d} by averaging the obtained diagonal results. Since we utilize the regression results before conducting depth estimation, we call this two-stage depth inference.

It is worth mentioning that this kind of precise two-stage depth inference can only be carried out with the help of the proposed ground plane prior because this process needs the

ability to obtain depths from any given non-integer location, *i.e.*, interpolating across densely trained **grounded depth** maps. Previous works are mainly designed to cope with the sparse and non-grounded depth supervision, which is hard to adopt a similar two-stage depth inference.

4. Experiments

4.1. Implementation Details

Dataset¹ In this work, we use KITTI [11] vision benchmark to evaluate and test our method. KITTI vision benchmark is a multimodal and multi-task dataset. We use the data from the 3D object detection track, which contains 7481 training images and 7518 testing images. Besides the RGB images and camera calibration matrices, no other data is utilized in our work. We also follow the split from [8] which divide the whole 7481 training images into a *train* set (3712) and a *val* set (3769).

Metrics We use the 3D average precision with 40 recall positions [35] $AP_{3D|R40}$ under three different difficulties (easy, moderate, hard) as our main evaluation metrics. Bird-eye view (BEV) detection is also adopted as our evaluation metrics. To better compare the performance of depth estimation, we also use the mean percentage error (MPE) as our metric:

$$MPE = \sum_i \left| \frac{pred_i - gt_i}{gt_i} \right|. \quad (9)$$

Training Details The size of the input image is padded to 384×1280 , and random horizontal flip augmentation is utilized. In the ground branch, we use two successive 3×3 Conv-BN-ReLU layers and one output convolution layer. All convolutions in the ground branch are coord convolutions [17], and their dilations are set to 2. For each object, the number of sampling points of dense grounded depth is the same as the polygon area of the projected bottom quadrilateral. The biggest number of sampling points does not exceed 5500 on KITTI.

We use PyTorch [29] to implement our method and train it with RTX 2080Ti GPU. AdamW [21] optimizer is adopted with a learning rate of $3e-4$ and a weight decay of $1e-5$. The batch size is set to 8, and the number of training epochs is set to 100. The learning rate will be decayed by a factor of 0.1 at the 80th and the 90th epochs.

4.2. Ablation Study

In this section, we discuss our method from two aspects, which are the effectiveness of the grounded depth, the ablation of each component in the ground branch.

¹We also add the experiments on NuScenes [3] in our open-source code repository.

Effectiveness of the Grounded Depth To examine the effectiveness of the grounded depth, we compare the MPE of our method and the widely used geometry depth estimation methods. The results are shown in Fig. 8. “Geo” means using the geometry depth estimation, and “Gnd” means using the grounded depth. The postfixes “13” and “24” means using the average depths from diagonal keypoints k_1, k_3 and k_2, k_4 , respectively.

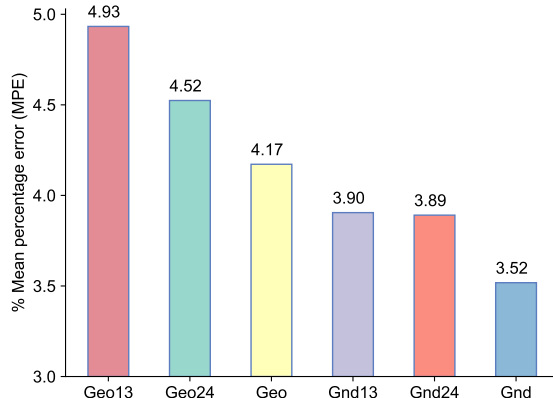


Figure 8. The mean percentage errors (MPEs) of the geometry depths and the grounded depths (lower is better).

From Fig. 8, we can see that the proposed grounded depth outperforms the widely used geometry depth in all settings. This shows the effectiveness of the grounded depth and the idea of introducing the ground plane prior.

Ablation Study on the Ground Branch To verify the effectiveness of the design of the ground branch, we show the ablation study in this section, as shown in Tab. 1.

Table 1. Ablation study on the KITTI *val* set. We use the average $AP_{3D|R40}$ with a IOU threshold of 0.7 is used as the evaluation metric. “sp.” means adding the grounded depth branch, but with the sparse supervision as the conventional depth training. “de.” means adding the dense grounded depth branch. “co.” means the coord convolution is adopted. “di.” means using the dilated convolution.

Setting	Easy	Moderate	Hard
Baseline	22.18	16.37	13.98
Baseline+sp.	19.23	14.14	11.84
Baseline+de.	20.54	15.54	13.34
Baseline+de.+co.	21.54	16.23	13.84
Baseline+de.+di.	22.65	16.67	14.37
Baseline+de.+co.+di.	22.98	17.37	14.78

We can see that the introduced grounded depth with dense training outperforms that with the sparse supervision. Moreover, the coord and dilated convolution also gains better performance as discussed in Sec. 3.3.

Table 2. Comparison on the KITTI *test* set for the car class. The results are tested on the KITTI testing server.

Method	Year	Extra Data	Runtime (ms)	AP _{3D R40}			BEV		
				easy	moderate	hard	easy	moderate	hard
Decoupled-3D [4]	AAAI20	Depth	-	11.08	7.02	5.63	23.16	14.82	11.25
MonoPSR [13]	CVPR19	LiDAR	200	10.76	7.25	5.85	18.33	12.58	9.91
AM3D [24]	ICCV19	Depth	400	16.50	10.74	9.52	25.03	17.32	14.91
PatchNet [23]	ECCV20	Depth	400	15.68	11.12	10.17	22.97	16.86	14.97
DA-3Ddet [40]	ECCV20	Depth	-	16.80	11.50	8.90	-	-	-
D4LCN [10]	CVPR20	Depth	-	16.65	11.72	9.51	22.51	16.02	12.55
Kinem3D [2]	ECCV20	Multi-frames	120	19.07	12.72	9.17	26.69	17.52	13.10
PCT [38]	NeurIPS21	Depth	45	21.00	13.37	11.31	29.65	19.03	15.92
CaDDN [31]	CVPR21	LiDAR	63	19.17	13.41	11.46	27.94	18.91	17.19
DFR-Net [45]	ICCV21	Depth	180	19.40	13.63	10.35	28.17	19.17	14.84
AutoShape [20]	ICCV21	CAD Models	50	22.47	14.17	11.36	30.66	20.08	15.59
M3D-RPN [1]	ICCV19	No	160	14.76	9.71	7.42	21.02	13.67	10.23
SMOKE [19]	CVPRW20	No	30	14.03	9.76	7.84	20.83	14.49	12.75
MonoPair [9]	CVPR20	No	57	13.04	9.99	8.65	19.28	14.83	12.89
MonoDLE [25]	CVPR21	No	40	17.23	12.26	10.29	24.79	18.89	16.00
MonoRUn [6]	CVPR21	No	70	19.65	12.30	10.58	27.94	17.34	15.24
GrooMeD [14]	CVPR21	No	120	18.10	12.32	9.65	26.19	18.27	14.05
MonoRCNN [34]	ICCV21	No	70	18.36	12.65	10.03	25.48	18.11	14.10
DDMP-3D [37]	CVPR21	No	180	19.71	12.78	9.80	28.08	17.89	13.44
MonoEF [44]	CVPR21	No	30	21.29	13.87	11.71	29.03	19.70	17.26
MonoFlex [42]	CVPR21	No	30	19.94	13.89	12.07	28.23	19.75	16.89
GUPNet [22]	ICCV21	No	30	20.11	14.20	11.77	-	-	-
MonoGround		No	30	21.37	14.36	12.62	30.07	20.47	17.74

Table 3. Comparison on the KITTI *val* set for the car class. * means we use the results from the official open-source code, which are slightly different from the reported results.

Method	AP _{3D R40} @IOU=0.7			BEV@IOU=0.7			AP _{3D R40} @IOU=0.5			BEV@IOU=0.5		
	easy	moderate	hard	easy	moderate	hard	easy	moderate	hard	easy	moderate	hard
CenterNet [43]	0.60	0.66	0.77	3.46	3.31	3.21	20.00	17.50	15.57	34.36	27.91	24.65
MonoGRNet [30]	11.90	7.46	5.76	19.72	12.81	10.15	47.59	32.28	25.50	48.53	35.94	28.59
MonoDIS [35]	11.06	7.60	6.37	18.45	12.58	10.66	-	-	-	-	-	-
M3D-RPN [1]	14.53	11.07	8.65	20.85	15.62	11.88	48.53	35.94	28.59	53.35	39.60	31.76
MoVi-3D [36]	14.28	11.13	9.68	22.36	17.87	15.73	-	-	-	-	-	-
MonoPari [9]	16.28	12.30	10.42	24.12	18.17	15.76	55.38	42.39	37.99	61.06	47.63	41.92
MonoDLE [25]	17.45	13.66	11.68	24.97	19.33	17.01	55.41	43.42	37.81	60.73	46.87	41.89
GrooMeD [14]	19.67	14.32	11.27	27.38	19.75	15.92	55.62	41.07	32.89	61.83	44.98	36.29
GUPNet [22]	22.76	16.46	13.72	31.07	22.94	19.75	57.62	42.33	37.59	61.78	47.06	40.88
MonoFlex* [42]	24.22	17.34	15.13	31.65	23.29	20.02	60.70	45.65	39.91	66.26	49.30	44.42
MonoGround	25.24	18.69	15.58	32.68	24.79	20.56	62.60	47.85	41.97	67.36	51.83	45.65

4.3. Performance Evaluation

We report the evaluation results on the both KITTI *test* and *val* sets. The results for the car class are shown in Tabs. 2 and 3. From Tab. 2 we can see that our method outperforms all other purely monocular methods while maintaining a fast speed in real-time. When compared with the methods using extra data sources, our method still gets the highest performance in the “moderate” and “hard” settings. Similar results can also be seen in Tab. 3. Moreover, we can see that our method has a relatively large performance

gap compared with other methods. When the intersection over union (IOU) is set to 0.5, the performance gap becomes larger (~2%).

We also show the results for the pedestrian and cyclist classes on the KITTI *test* set. As shown in Tab. 4, our method still performs well on these two classes. Our method obtains the best results for the cyclist class and the second-best results for the pedestrian class. A potential reason for the relatively low performance for the pedestrian class is that the pedestrians’ sizes (width and length) are rel-



Figure 9. Visualization of the detection results on the KITTI *test* set.

Table 4. The performance of the $AP_{3D|R40}$ on the KITTI *test* set for the pedestrian and cyclist classes. “mod.” means moderate.

Method	Pedestrian			Cyclist		
	easy	mod.	hard	easy	mod.	hard
M3D-RPN [1]	4.92	3.48	2.94	0.94	0.65	0.47
DDMP-3D [37]	4.93	3.55	3.01	4.18	2.50	2.32
Movi3D [36]	8.99	5.44	4.57	1.08	0.63	0.70
MonoPair [9]	10.02	6.68	5.53	3.79	2.12	1.83
MonoFlex [42]	9.43	6.31	5.26	4.17	2.35	2.04
MonoDLE [25]	9.64	6.55	5.44	4.59	2.66	2.45
MonoRUn [6]	10.88	6.78	5.83	1.01	0.61	0.48
GUPNet [22]	14.72	9.53	7.87	4.18	2.65	2.09
MonoGround	12.37	7.89	7.13	4.62	2.68	2.53

atively small compared with cars and cyclists. In this way, the bottom rectangle of the 3D bounding box is small and the projected area on the ground is small. It would be hard to conduct random sampling on this small area as shown in Fig. 3. So the performance increases for the pedestrian class is not as large as other classes. However, our method still outperforms other methods (except GUPNet) by a relatively big margin for the pedestrian class ($\sim 2\%$, $\sim 1\%$, and $\sim 1.5\%$ for the “easy”, “moderate”, and “hard” settings).

4.4. Visualization

In this section, we show the visualization results of the proposed method on the KITTI *test* set, as shown in Fig. 9. We can see that our method gives good visualization results in detecting various 3D objects.

5. Conclusion and Limitation

In this work, we have proposed a monocular 3D object detection method termed as MonoGround, which introduces the ground plane prior in monocular 3D object detection. The introduced dense grounded depth with the ground plane prior requires no extra data sources like LiDAR, stereo images, and depth images. Moreover, we also have proposed a depth-align training strategy and a two-stage precise depth inference method to cope with the introduced dense grounded depth in the network. The proposed depth inference method can be easily extended with uncertainty and has outperformed the conventional geometry depth inference method. With these components, MonoGround has achieved state-of-the-art performance on KITTI while maintaining a very fast real-time speed.

However, there is a limitation in this work. As discussed in Sec. 4.3, the introduced ground plane prior needs a random dense sampling process, which is less friendly for objects with a small bottom surface like pedestrians, since the sampled points and depths could be rare. Although we still get the second-best result for the pedestrian class, this could be a direction for the future work.

Acknowledgements

This work is supported in part by National Key Research and Development Program of China under Grant 2020AAA0107400, Zhejiang Provincial Natural Science Foundation of China under Grant LR19F020004, National Natural Science Foundation of China under Grant U20A20222, and Zhejiang University K.P.Chao’s High Technology Development Foundation.

References

- [1] Garrick Brazil and Xiaoming Liu. M3d-rpn: Monocular 3d region proposal network for object detection. In *ICCV*, pages 9287–9296, 2019.
- [2] Garrick Brazil, Gerard Pons-Moll, Xiaoming Liu, and Bernt Schiele. Kinematic 3d object detection in monocular video. In *ECCV*, pages 135–152. Springer, 2020.
- [3] Holger Caesar, Varun Bankiti, Alex H. Lang, Sourabh Vora, Venice Erin Liong, Qiang Xu, Anush Krishnan, Yu Pan, Giancarlo Baldan, and Oscar Beijbom. nuscenes: A multimodal dataset for autonomous driving. In *CVPR*, pages 11621–11631, 2020.
- [4] Yingjie Cai, Buyu Li, Zeyu Jiao, Hongsheng Li, Xingyu Zeng, and Xiaogang Wang. Monocular 3d object detection with decoupled structured polygon estimation and height-guided depth estimation. In *AAAI*, volume 34, pages 10478–10485, 2020.
- [5] Florian Chabot, Mohamed Chaouch, Jaonary Rabarisoa, Céline Teulière, and Thierry Chateau. Deep manta: A coarse-to-fine many-task network for joint 2d and 3d vehicle analysis from monocular image. In *CVPR*, pages 2040–2049, 2017.
- [6] Hansheng Chen, Yuyao Huang, Wei Tian, Zhong Gao, and Lu Xiong. Monorun: Monocular 3d object detection by reconstruction and uncertainty propagation. In *CVPR*, pages 10379–10388, 2021.
- [7] Xiaozhi Chen, Kaustav Kundu, Ziyu Zhang, Huimin Ma, Sanja Fidler, and Raquel Urtasun. Monocular 3d object detection for autonomous driving. In *CVPR*, pages 2147–2156, 2016.
- [8] Xiaozhi Chen, Kaustav Kundu, Yukun Zhu, Andrew G Berneshawi, Huimin Ma, Sanja Fidler, and Raquel Urtasun. 3d object proposals for accurate object class detection. In *NeurIPS*, pages 424–432. Citeseer, 2015.
- [9] Yongjian Chen, Lei Tai, Kai Sun, and Mingyang Li. Monopair: Monocular 3d object detection using pairwise spatial relationships. In *CVPR*, pages 12093–12102, 2020.
- [10] Mingyu Ding, Yuqi Huo, Hongwei Yi, Zhe Wang, Jianping Shi, Zhiwu Lu, and Ping Luo. Learning depth-guided convolutions for monocular 3d object detection. In *CVPR*, pages 11672–11681, 2020.
- [11] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? the kitti vision benchmark suite. In *CVPR*, 2012.
- [12] Jason Ku, Melissa Mozifian, Jungwook Lee, Ali Harakeh, and Steven L Waslander. Joint 3d proposal generation and object detection from view aggregation. In *IEEE Int. Conf. Intell. Robots Syst.*, pages 1–8. IEEE, 2018.
- [13] Jason Ku, Alex D Pon, and Steven L Waslander. Monocular 3d object detection leveraging accurate proposals and shape reconstruction. In *CVPR*, pages 11867–11876, 2019.
- [14] Abhinav Kumar, Garrick Brazil, and Xiaoming Liu. Groomed-nms: Grouped mathematically differentiable nms for monocular 3d object detection. In *CVPR*, pages 8973–8983, 2021.
- [15] Buyu Li, Wanli Ouyang, Lu Sheng, Xingyu Zeng, and Xiaogang Wang. Gs3d: An efficient 3d object detection framework for autonomous driving. In *CVPR*, pages 1019–1028, 2019.
- [16] Lijie Liu, Jiwen Lu, Chunjing Xu, Qi Tian, and Jie Zhou. Deep fitting degree scoring network for monocular 3d object detection. In *CVPR*, pages 1057–1066, 2019.
- [17] Rosanne Liu, Joel Lehman, Piero Molino, Felipe Petroski Such, Eric Frank, Alex Sergeev, and Jason Yosinski. An intriguing failing of convolutional neural networks and the coordconv solution. In *NeurIPS*, 2018.
- [18] Yuxuan Liu, Yuan Yixuan, and Ming Liu. Ground-aware monocular 3d object detection for autonomous driving. *IEEE Robot. Autom. Lett.*, 6(2):919–926, 2021.
- [19] Zechen Liu, Zizhang Wu, and Roland Tóth. Smoke: Single-stage monocular 3d object detection via keypoint estimation. In *CVPRW*, pages 996–997, 2020.
- [20] Zongdai Liu, Dingfu Zhou, Feixiang Lu, Jin Fang, and Liangjun Zhang. Autoshape: Real-time shape-aware monocular 3d object detection. In *ICCV*, pages 15641–15650, 2021.
- [21] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *ICLR*, 2019.
- [22] Yan Lu, Xinzhu Ma, Lei Yang, Tianzhu Zhang, Yating Liu, Qi Chu, Junjie Yan, and Wanli Ouyang. Geometry uncertainty projection network for monocular 3d object detection. In *ICCV*, pages 3111–3121, 2021.
- [23] Xinzhu Ma, Shinan Liu, Zhiyi Xia, Hongwen Zhang, Xingyu Zeng, and Wanli Ouyang. Rethinking pseudo-lidar representation. In *ECCV*, pages 311–327. Springer, 2020.
- [24] Xinzhu Ma, Zhihui Wang, Haojie Li, Pengbo Zhang, Wanli Ouyang, and Xin Fan. Accurate monocular 3d object detection via color-embedded 3d reconstruction for autonomous driving. In *ICCV*, pages 6851–6860, 2019.
- [25] Xinzhu Ma, Yinmin Zhang, Dan Xu, Dongzhan Zhou, Shuai Yi, Haojie Li, and Wanli Ouyang. Delving into localization errors for monocular 3d object detection. In *CVPR*, pages 4721–4730, 2021.
- [26] Fabian Manhardt, Wadim Kehl, and Adrien Gaidon. Roi-10d: Monocular lifting of 2d detection to 6d pose and metric shape. In *CVPR*, pages 2069–2078, 2019.
- [27] Arsalan Mousavian, Dragomir Anguelov, John Flynn, and Jana Kosecka. 3d bounding box estimation using deep learning and geometry. In *CVPR*, pages 7074–7082, 2017.
- [28] J Krishna Murthy, GV Sai Krishna, Falak Chhaya, and K Madhava Krishna. Reconstructing vehicles from a single image: Shape priors for road scene understanding. In *ICRA*, pages 724–731. IEEE, 2017.
- [29] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In *NeurIPS*, pages 8024–8035. Curran Associates, Inc., 2019.
- [30] Zengyi Qin, Jinglu Wang, and Yan Lu. Monogmet: A geometric reasoning network for monocular 3d object localization. In *AAAI*, volume 33, pages 8851–8858, 2019.

- [31] Cody Reading, Ali Harakeh, Julia Chae, and Steven L Waslander. Categorical depth distribution network for monocular 3d object detection. In *CVPR*, pages 8555–8564, 2021.
- [32] Shaoshuai Shi, Chaoxu Guo, Li Jiang, Zhe Wang, Jianping Shi, Xiaogang Wang, and Hongsheng Li. Pv-rcnn: Point-voxel feature set abstraction for 3d object detection. In *CVPR*, pages 10529–10538, 2020.
- [33] Shaoshuai Shi, Xiaogang Wang, and Hongsheng Li. Point-rcnn: 3d object proposal generation and detection from point cloud. In *CVPR*, pages 770–779, 2019.
- [34] Xuepeng Shi, Qi Ye, Xiaozhi Chen, Chuangrong Chen, Zhixiang Chen, and Tae-Kyun Kim. Geometry-based distance decomposition for monocular 3d object detection. In *ICCV*, pages 15172–15181, October 2021.
- [35] Andrea Simonelli, Samuel Rota Buló, Lorenzo Porzi, Manuel López-Antequera, and Peter Kotschieder. Disentangling monocular 3d object detection. In *ICCV*, pages 1991–1999, 2019.
- [36] Andrea Simonelli, Samuel Rota Buló, Lorenzo Porzi, Elisa Ricci, and Peter Kotschieder. Towards generalization across depth for monocular 3d object detection. In *ECCV*, pages 767–782. Springer, 2020.
- [37] Li Wang, Liang Du, Xiaoqing Ye, Yanwei Fu, Guodong Guo, Xiangyang Xue, Jianfeng Feng, and Li Zhang. Depth-conditioned dynamic message propagation for monocular 3d object detection. In *CVPR*, pages 454–463, 2021.
- [38] Li Wang, Li Zhang, Yi Zhu, Zhi Zhang, Tong He, Mu Li, and Xiangyang Xue. Progressive coordinate transforms for monocular 3d object detection. *arXiv preprint arXiv:2108.05793*, 2021.
- [39] Yan Wang, Wei-Lun Chao, Divyansh Garg, Bharath Hariharan, Mark Campbell, and Kilian Q Weinberger. Pseudo-lidar from visual depth estimation: Bridging the gap in 3d object detection for autonomous driving. In *CVPR*, pages 8445–8453, 2019.
- [40] Xiaoqing Ye, Liang Du, Yifeng Shi, Yingying Li, Xiao Tan, Jianfeng Feng, Errui Ding, and Shilei Wen. Monocular 3d object detection via feature domain adaptation. In *ECCV*, pages 17–34. Springer, 2020.
- [41] Fisher Yu, Dequan Wang, Evan Shelhamer, and Trevor Darrell. Deep layer aggregation. In *CVPR*, 2018.
- [42] Yunpeng Zhang, Jiwen Lu, and Jie Zhou. Objects are different: Flexible monocular 3d object detection. In *CVPR*, pages 3289–3298, June 2021.
- [43] Xingyi Zhou, Dequan Wang, and Philipp Krähenbühl. Objects as points. *arXiv preprint arXiv:1904.07850*, 2019.
- [44] Yunsong Zhou, Yuan He, Hongzi Zhu, Cheng Wang, Hongyang Li, and Qinhong Jiang. Monocular 3d object detection: An extrinsic parameter free approach. In *CVPR*, pages 7556–7566, 2021.
- [45] Zhikang Zou, Xiaoqing Ye, Liang Du, Xianhui Cheng, Xiao Tan, Li Zhang, Jianfeng Feng, Xiangyang Xue, and Errui Ding. The devil is in the task: Exploiting reciprocal appearance-localization features for monocular 3d object detection. In *ICCV*, pages 2713–2722, 2021.