# Surround-view Fisheye BEV-Perception for Valet Parking: Dataset, Baseline and Distortion-insensitive Multi-task Framework

Zizhang  $Wu^{1\&}$  Yuanzhu  $Gan^{1\&}$  Xianzhi  $Li^{2*}$  Yunzhe  $Wu^1$  Xiaoquan  $Wang^1$  Tianhao  $Xu^3$  Fan  $Wang^1$ 

Abstract—Surround-view fisheye perception under valet parking scenes is fundamental and crucial in autonomous driving. Environmental conditions in parking lots perform differently from the common public datasets, such as imperfect light and opacity, which substantially impacts on perception performance. Most existing networks based on public datasets may generalize suboptimal results on these valet parking scenes, also affected by the fisheye distortion. In this article, we introduce a new large-scale fisheye dataset called Fisheye Parking Dataset (FPD) to promote the research in dealing with diverse realworld surround-view parking cases. Notably, our compiled FPD exhibits excellent characteristics for different surround-view perception tasks. In addition, we also propose our real-time distortion-insensitive multi-task framework Fisheye Perception Network (FPNet), which improves the surround-view fisheye BEV perception by enhancing the fisheye distortion operation and multi-task lightweight designs. Extensive experiments validate the effectiveness of our approach and the dataset's exceptional generalizability.

Index Terms—dataset, surround-view, fisheye, valet parking, multi-task

# I. INTRODUCTION

As the priority of developing an effective and safe advanced driver assistance system (ADAS) [1]–[3], valet parking attracts more attention from industry and research communities in recent years [4]–[6]. Among various driving assistance applications, valet parking exhibits an essential and challenging task. Figure 1 shows several challenging scenarios during valet parking [7]–[9]. Besides, the environmental conditions in parking scenes, such as light and opacity, significantly increase the difficulty of robust environment perception [10], [11]. Unlike the relatively clear scenarios like highway and urban areas, valet parking aims to drive the vehicle into the drop-off area such as parking slots, which faces high requirements in perception [12].

Recent advances [13], [14] demonstrate the potential to replace the LiDAR with cheap onboard cameras, which are readily available on most modern vehicles [15], [16]. Particularly, the surround-view fisheye cameras can provide a wider field-of-view (FoV) [17] than the pinhole cameras, which have grown popular in mass production. In addition, four surround-view fisheye cameras cover a 360-degree perception, which makes up for pinhole cameras' near-field perception

- 1 Zongmu Technology
- 2 Huazhong University of Science and Technology
- 3 Technical University of Braunschweig
- & These authors contributed equally to this work and should be considered co-first authors.
  - \* Corresponding author: xzli@hust.edu.cn(X. Li)

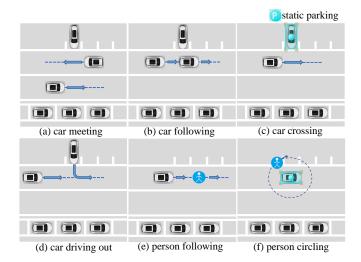


Fig. 1. Several valet parking scenes in the parking lots, including car meeting, car following, car crossing, car driving out, person following and person circling etc.

insufficiency, especially under the valet parking conditions [18]. However, the fisheye lens usually perform obvious radial distortion, which leads to substantial appearance distortion [18], [19], complicating the surrounding's recognition. In order to fully exploit the fisheye paradigm, more researchers begin to explore the fisheye surround-view perception, such as the position and pose information of the vehicles or pedestrians [12], [20].

Current datasets like KITTI [21], Cityscapes [22] etc. mostly capture images with pinhole cameras, which comfortably achieve clear and distinguishable images under the urban, rural, or highway driving scenes. There are scarce fisheye datasets for autonomous driving [23], [24], which have aided rapid progress for the fisheye surround-view perception. Woodscape [23] and KITTI360 [24] collect large-scale fisheye datasets for different perception tasks on the ground. However, these datasets place insufficient emphasis on perception with valet parking scenarios and fisheye image formation. Hence, the models trained on public pinhole datasets or fisheye datasets on the ground may reveal suboptimal performance, since lacking sufficient training samples, particularly for valet parking scenes.

To expand surround-view fisheye perception tasks' images with various occlusions and postures under valet parking scenes, we contribute the first fisheye dataset called **F**isheye **P**arking **D**ataset (**FPD**) for parking scenarios, which has the following great properties, including (i) large scale quantity

with more than 400 thousand fisheye images; (ii) high diversity with different parking lots, different periods and different parking conditions; (iii) high quality by filtering noisy and redundant images; (iv) multiple kinds of annotations for different perception tasks, like 2D object detection, 3D object detection, BEV perception, depth estimation, etc.

Different from other public autonomous driving datasets [21], [25], [26], our **FPD** dataset focuses on valet parking surround-view perception tasks and makes up the vacancy for the research in dealing with real-world parking lot scenes. Furthermore, we offer the baseline on our **FPD** and propose the real-time distortion-insensitive multi-task network **Fisheye Perception Network (FPNet)**, specifically for surround-view fisheye perception tasks, including 2D object detection, monocular 3D object detection, BEV perception, and monocular dense depth estimation. The network achieves a balance between lightweight and accuracy, additionally with a particular module for addressing fisheye distortion.

Our contributions are summarized as follows:

- We build the first fisheye parking dataset FPD, concentrating on the surround-view fisheye perception including the 2D object detection, 3D object detection, BEV perception, and depth estimation. Our contributed FPD comprises over 400 thousand fisheye images and contains attractive characteristics for parking scenarios.
- We propose the baseline for our FPD: the distortioninsensitive multi-task framework FPNet for surroundview perception tasks, especially the BEV perception. FPNet utilizes the specific distortion module and lightweight designs to achieve real-time, distortioninsensitive and accurate performance.
- Comprehensive experiments validate the practicability of our collected FPD dataset and the effectiveness of FPNet.

### II. RELATED WORK

# A. Autonomous driving datasets

To meet the improving requirement for automated driving development, in the last decade, pioneer works create numerous datasets [21], [25]–[27], which cover the most autonomous driving tasks, such as object detection [21], [27], semantic segmentation [25], [26], depth estimation [21], [28], lane detection [29], motion estimation [26], [27] etc., making amazing contribution to autonomous driving.

However, these datasets are almost pinhole camera datasets, with a limited field-of-view (FoV). Actually extensive visual tasks also adopt surround-view fisheye cameras to monitor the surrounding environment, due to their large FoV [17] and sufficiently stable performance [30]. In particular, egovehicles can achieve 360-degree perception using only four fisheye cameras, which is conducive to mass production [31]. Moreover, the omnidirectional camera also captures the 360-degree field of view, which covers a full circle in the horizontal plane [32].

Valeo releases the first fisheye dataset Woodscape [23], to encourage the development of native fisheye models. But Woodscape doesn't publish the LiDAR ground truth due to data protection restrictions. KITTI360 [24] presents another

large-scale dataset that contains rich sensory information including pinhole and fisheye cameras. For omnidirectional cameras, the industry and academia provide the datasets like Stanford2D3D [33], Matterport3D [34], 360D [35], PanoSUNCG [36] etc. for the omnidirectional visual perception. However, these datasets generally focus on aboveground autonomous driving scenes, such as urban, rural, and highways. There is no publicly available benchmark dataset for the valet parking scenes. Our Fisheye Parking Dataset (FPD) could make up the vacancy, promoting research in dealing with real-world parking lot scenes.

# B. Monocular perception tasks

Autonomous driving involves various perception tasks, like object detection or depth estimation, to assist the system to cover a wider range of use cases. In this paper, we mainly deal with three tasks: 2D object detection, monocular 3D detection, and monocular depth estimation. We can derive the BEV perception from the 3D detection results.

- 1) 2D object detection: 2D object detection acts as a basic vision task. CNN-based 2D object detection frameworks consist of one-stage detection methods [38], [39] and two-stage detection methods [40]–[42]. As an end-to-end pipeline, one-stage methods achieve a significant trade-off between performance and speed, like SSD series [43]–[45], YOLO series [39], [46] and RetinaNet [38]. In addition, two-stage methods, like RCNN series [42], [47], [48], take advantage of the predefined anchors to improve the performance at the cost of speed. Furthermore, [49], [50] fuse multi-scale feature maps to improve detection with different scales.
- 2) Monocular 3D detection: Many prior works [14], [51]-[53] have tackled the inherently ill-posed problem of detecting 3D objects from monocular images. Due to the lack of depth information from images, monocular 3D detection learns harder than LiDAR-based and stereo-based counterparts. Many works [54]–[56] address this problem by utilizing 2D-3D geometric constraints to improve 3D detection performance. In addition, CenterNet [57] proposes a center-based anchorfree method but with restrained accuracy. Following this work, center-based series SMOKE [58], KM3D [59] and RTM3D [60] assist the regression of object depth by solving a Perspective-n-Point method and have achieved remarkable results. However, most existing works target the pinhole cameras instead of the fisheye ones, where fisheye cameras have a strong radial distortion and exhibit more complex projection geometry [61], which leads to appearance distortion [62].
- 3) Monocular depth estimation: CNN-based supervised methods [63]–[65] are popular in monocular depth estimation tasks due to their superior performance. As a pioneer, Eigen et al. [66] directly regress depth by employing two stacked deep networks for a coarse prediction, then refining it locally. Then Laina et al. [67] adopt the end-to-end single CNN architecture, following the residual learning. Moreover, DRO [63] introduces a deep recurrent optimizer to alternately update the depth and camera poses through iterations, to minimize the feature-metric cost. Furthermore, the recent works [68], [69] apply vision transformers to improve depth estimation.

TABLE I

THE COMPARISON OF DIFFERENT IMAGING SENSORS. 'FOV' DENOTES 'FIELD OF VIEW'.

'H' DENOTES 'HORIZONTAL FOV' AND 'V' DENOTES 'VERTICAL FOV'.

Name	FOV	Distortion	Dataset	Characteristics
pinhole camera	H:<180° V:<180°	small	KITTI [21] nuScenes [26] Waymo [27] Cityscapes [25]	universal in most public perception methods; caring little about the distortion
fisheye camera	H:>180° V:<180°	large	Woodscape [23] KITTI 360 [24]	common in surround-view system; designing specific module to operate distortion
omnidirectional camera	H:360° V:180°	large	Stanford2D3D [33] Matterport3D [34] 360D [35] PanoSUNCG [36]	higher complexities and distortions implicated in 360° full-view panoramas [32], [37]



Fig. 2. Illustration of our collected real valet parking scenes, with different parking conditions, different lighting conditions, and different occlusions.

#### C. Visual perception for different imaging sensors

Most public perception methods design modules for the pinhole images (subsection B), which conduct the perspective projection and care little about the camera distortion. However, when utilizing wide-angle cameras (the fisheye camera or omnidirectional camera), researchers must concern about the large distortion and other issues. The surround-view autonomous driving system usually adopts the fisheye cameras to achieve the surrounding perception [70]–[75], which builds specific components to operate the distortion [76]-[79]. The omnidirectional camera reveals the higher complexities and distortions implicated in 360° full-view panoramas [32], [37], where the omnidirectional images' perception also deals with the distortion and the problems caused from the large resolution, image preprocessing and knowledge transformation, etc. [35], [80]–[83]. Table I shows more comparison of the three imaging sensors.

# D. Multi-task visual perception

Numerous autonomous driving investigations address complex real-world scenarios by joint learning of different subtasks [84]–[89]. MultiNet [86] accomplishes road segmentation, detection, and classification tasks via a universal efficient architecture. NeurAll [88] brings unified CNN architecture including object recognition, motion, depth estimation, and facilitating visual SLAM. Li et al. [89], in like manner, propose a unified end-to-end framework (MJPNet) that shares predictions among multiple subtasks. These studies validate

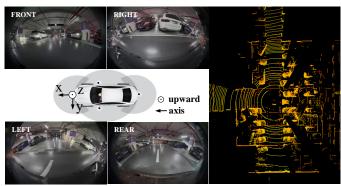


Fig. 3. Sensors' installation positions and produced images or point cloud's visualization, where the center of front bumper acts as the origin of the egovehicle coordinate system.

the benefit of joint multi-task learning in different autonomous driving scenes.

However, most studies focus on undistorted pinhole visual conditions and infrequent exploration in wide-angle fisheye camera-based perception. OmniDet [12] forms an encoder-shared framework with six primary tasks on the fisheye dataset WoodScape [23]. On the whole, fisheye models reveal a significant growth potential with rich scenario distributions.

#### III. FISHEYE PARKING DATASET

In this section, we introduce our Fisheye Parking Dataset (**FPD**) in detail, including the process of data collection and annotation, dataset description and outstanding characteristics.

#### A. Data Collection

To ensure the diversity of autonomous driving scenarios, we collect a total of three cities, over one hundred parking lots, two periods (daytime and nighttime), and capture more than four hundred videos together with the point cloud sequences from LiDAR. Figure 2 displays several our collected real fisheye images of valet parking scenes.

Specifically, our master LiDAR adopts the RoboSense RS-Ruby, which has 128 beams, 10Hz capture frequency,  $360^{\circ}$  horizontal FOV and  $-25^{\circ}$  to  $+15^{\circ}$  vertical FOV. Besides, we choose four ZongMu fisheye RGB cameras, with  $1920\times1280$  resolution and 20Hz capture frequency. Figure 3 demonstrates these sensors' installation positions and their produced images or point cloud's performance. During the data recording process, the system aligns the timestamp between the videos from cameras and point cloud sequences from LiDAR, so qualified for the following annotation. Furthermore, we conduct the

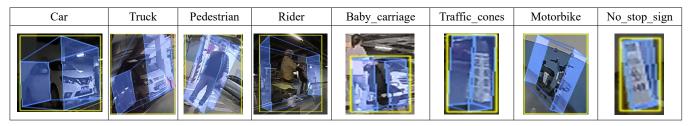


Fig. 4. Our data annotation contains eight categories, including the car, truck, pedestrian, rider, baby carriage, traffic cones, motorbike and no-stop sign.

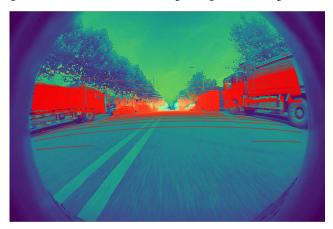


Fig. 5. We project the LiDAR points to the camera image plane to rectify the camera's extrinsics where the color denotes the points' different distances.

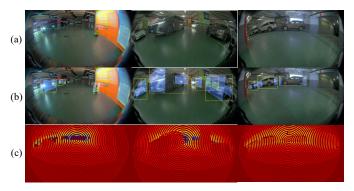


Fig. 6. Visualization of our annotated labels. (a) our captured monocular fisheye images; (b) associated 2D object labels (yellow bbox) and projected 3D object labels (blue bbox); (c) associated depth ground truth after depth completion.

calibration process of the sensors in the following three steps. Firstly, we can directly calculate the camera intrinsics from the lens' distortion lookup table provided by the initial factory settings. Secondly, we achieve the extrinsics (x, y, z, pitch, yaw, roll) of the LiDAR and cameras based on the ego-vehicle coordinate system via the measurement equipment. Thirdly, we further rectify the camera's extrinsics by the alignment of projected LiDAR points and images' semantics. As shown in Figure 5, we project the LiDAR points to the image plane via the above calibrations, then we manually adjust the camera's extrinsics to match the projected points and the semantic contents.

To cover various realistic scenes of the parking lot, we artificially arrange kinds of driving scenes to collect data, like car meeting, car crossing and person circling, etc., as shown in Figure 1, which are common but critical cases for autonomous parking task.

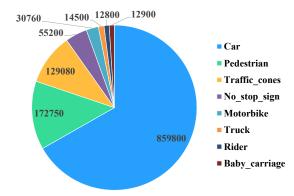


Fig. 7. The detailed data distribution of FPD in different categories.

TABLE II STATISTICS OF DATA ANNOTATIONS WITH OUR  $\mbox{\bf FPD}.$  We obtain more than  $400,\!000$  annotations totally.

	Training	Validation	Testing	Sum
# images	210,000	126,000	84,000	420,000

#### B. Data Annotation

We annotate the dataset in the same way as KITTI dataset [21], by drawing a tight bounding box around each object's complete point cloud body. We don't cover all objects' continuous moving process for redundant annotations. Instead, we remove similar segments and annotate data with intervals of three to five frames. In addition, we restrict the visible range (within 15 meters), so we abandon too far-away objects. For occluded objects, we reserve the 3D bounding box if the occlusion rate performs less than 80%, by way of imagining the full 3D bounding box with annotators' experience. Figure 6 shows several examples of our annotated labels.

Our annotation contains eight categories, including the car, truck, pedestrian, rider, baby carriage, traffic cones, motorbike and no-stop sign. Figure 4 demonstrates annotation demos of the eight categories, where the blue bounding boxes mean the 3D annotation's 2D visualization, projected from point cloud plane to image plane. Yellow bounding boxes indicate the outer bounding rectangle of blue projected points, as our 2D object detection ground truth, also shown in Figure 6 (b).

Furthermore, we project the point cloud to the monocular images with the assistance of calibration and distortion parameters for the sparse depth map. Then we adopt the depth completion method IP-Basic [90] to create more robust depth ground truth, as shown in Figure 6 (c).

# C. Dataset Description

Table II and Figure 7 illustrate the statistics of our **FPD**. Totally, we obtain more than 400,000 data, where one data con-

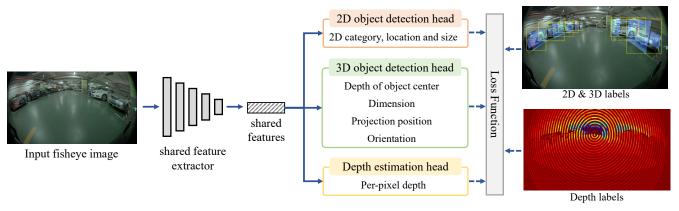


Fig. 8. The framework of our FPNet, which mainly consists of a shared feature extractor and a multi-task perception module, including a 2D object detection head, a 3D object detection head and a depth estimation head.

tains four fisheye images and one point cloud with annotation. In addition, one data accompanies one intrinsic parameter, one extrinsic parameter and one fisheye distortion parameter. We can project the point cloud annotation to the image via the intrinsic, extrinsic and fisheye distortion parameters to obtain the 2D object bounding box and depth ground truth.

Moreover, we split the **FPD** into training, validation and testing sets with the ratio of 5:3:2 and the amount of 210,000, 126,000, and 84,000. The ratio of daytime and nighttime scenes reveals 2:1. In addition, an average of 4 thousand data per parking lot composes more than 400,000 annotations, where the most frequent categories indicate the cars, the pedestrians, and the traffic cones, as shown in Figure 7.

# D. Dataset Characteristics

As the first large-scale real-world fisheye dataset, our compiled **FPD** exhibits the following nice properties:

- 1) First fisheye dataset for parking scenarios: We provide the first fisheye dataset **FPD**, which focuses on multiple autonomous driving tasks in parking scenarios, also distinct from natural scenes of public datasets. Environmental conditions in parking scenarios, such as light and opacity, significantly increase the detection difficulty. Concerning a variety of tough parking scenarios, **FPD** can promote research in dealing with real-world parking problems.
- 2) Great quantity: So far, our **FPD** contains more than 400 thousand data from over 200 hours of parking scene videos and point cloud sequences. In the future, we will continue to collect diverse parking to enrich the existing dataset.
- 3) High quality and diversity: Our **FPD** covers three cities, over one hundred parking lots from different periods, and different parking cases. Besides, we carefully pick out high-quality images and point clouds with high resolution, ensuring our dataset's superiority.
- 4) Multi-purpose: As a point-cloud-based dataset, our FPD's potential not only lies in the three tasks (i.e., 2D object detection, monocular 3D object detection and depth estimation) but also in other vision tasks, such as point cloud 3D object detection, 2D or 3D semantic segmentation, video object detection. Therefore, the FPD is multi-purpose for diverse tasks.

# IV. DISTORTION-INSENSITIVE MULTI-TASK FRAMEWORK

In this section, we introduce our surround-view fisheye monocular distortion-insensitive multi-task framework FPNet.

#### A. Overview

Figure 8 demonstrates the framework of our FPNet, which mainly consists of a shared feature extractor and a multi-task perception module, including a 2D object detection head, a 3D object detection head and a depth estimation head.

Given a fisheye monocular image  $I \in R^{3 \times H_I \times W_I}$ , we directly adopt our shared feature extractor to obtain features  $F^{2d} \in R^{C \times H_F \times W_F}$ , prepared for the following 2D object detection, 3D monocular object detection (BEV perception) and dense depth estimation. Through the multi-task perception head, we complete the predictions of three tasks (Section IV-B). During training, we project 3D ground truth to the monocular image plane to create the predictions' supervision, where the fisheye distortion module operates the distortion of fisheye camera projection. During testing, we utilize the post-processing module to decode the network's prediction, together with the fisheye distortion module's operation (Section IV-C and Section IV-D). Furthermore, we deploy our model to the embedded system with Qualcomm 8155 chip, to achieve real-time and excellent perception performance (Section IV-E).

#### B. Network

- 1) Shared feature extractor: To balance the trade-off between performance and speed, we select DLA34 [91] as our shared feature extractor. Additionally, we apply some improvements to achieve the lightweight requirement. Firstly, we downsample the input image with radio 8 instead of the usual 4, which saves lots of time but maintains the accuracy (see Section V-C3). Secondly, we get rid of some redundant layers according to their inference time on the embedded device, still with great performance (see Section V-C3).
- 2) Multi-task perception head: After extracting features from fisheye monocular images, we make predictions for our three perception tasks via multi-task perception head.
- **2D** object detection and monocular **3D** object detection. 2D object detection mainly focuses on searching for objects' bounding box in the image, while monocular **3D** object detection tries to locate objects' **3D** position, and regress objects' dimension and orientation, which obviously performs

more difficult compared to 2D object detection. For shared feature extractor design and real-time requirements, the centerbased framework exhibits more comfortable for our 2D and 3D object detection. Specifically, inspired by MonoCon [92], we build our 3D object detection head, whose input remains the shared features. Different from MonoCon [92], we directly predict the projected 3D bbox center heatmap, instead of the 2D bbox center heatmap and offset from projected 3D center to 2D center, which enhances the detection accuracy (see Section V-C4). Then we reserve the 3D-related predictions of object depth and uncertainty, shape dimensions and heading angle (MultiBin [93] regression). For 2D object detection, to reduce computation, we predict the offset from 2D center to the projected 3D center, then we can obtain the 2D center prediction by adding the above projected 3D center prediction and this offset. Besides, we also predict 2D bbox's height and width. At last, we abandon other auxiliary monocular contexts [92], also remaining satisfying performance (see Section V-C3).

Monocular depth estimation. Different from the above object centers' depth estimation, we also estimate the dense depth in this task. We refer DRO [63] to our depth estimation model, which indicates a gated recurrent network and iteratively updates the depth map between two images by minimizing a feature-metric cost [63]. So we are essential to cache the front frame image, to meet the tuple input of the depth estimation model. In addition, we replace the original feature backbone with our shared feature extractor, to fulfill multitask architecture. Furthermore, the depth estimation model also suits the self-supervised task. But in this work, we focus on supervised depth estimation.

#### C. Training and testing

- 1) Training and loss function: For 2D object detection and 3D object detection, we settle on different loss function for different subtasks as follows:
  - (i) The focal loss for projected 3D bbox center heatmap:

$$\mathcal{L}(\mathcal{H}, \mathcal{H}^*) = \frac{-1}{N} \sum_{(x,y)} \begin{cases} (1 - \mathcal{H}_{xy})^{\gamma} \log(\mathcal{H}_{xy}), & if \ \mathcal{H}_{xy}^* = 1, \\ (1 - \mathcal{H}_{xy}^*)^{\beta} (\mathcal{H}_{xy})^{\gamma} \log(1 - \mathcal{H}_{xy}), else \end{cases}$$

$$\tag{1}$$

To acquire the ground truth of projected 3D bbox center, we project the 3D bbox labels to the image plane via the fisheye distortion module. Then we follow CenterNet [94] to generate the ground-truth heatmap  $\mathcal{H}^* \in R^{1 \times H_F \times W_F}$ .  $\mathcal{H} \in R^{1 \times H_F \times W_F}$  means the predicted heatmap. We adopt the focal loss (Equ. 1) where the  $\alpha$  and  $\beta$  are hyper-parameters ( $\alpha$  = 4.0 and  $\beta$  = 2.0). We will detail the fisheye distortion module in section IV-D.

(ii) The L1 loss for the center offset, 2D bounding box's width and height, 3D bounding box's width, height and length, the intra-bin angle residual in heading angles, the dense depth estimation:

$$\mathcal{L}(\mathcal{S}, \mathcal{S}^*) = \lambda \cdot ||S - S^*||_1, \tag{2}$$

where  $\mathcal{S}^*$  indicates subtasks' ground truth and  $\mathcal{S}$  means the predicted values.

(iii) The Laplacian aleatoric uncertainty loss function for the object depth estimation. Following [92], we use the Laplace distribution to model the uncertainty and optimize the depth and uncertainty at the same time.

- (iv) The cross-entropy loss function for the bin index in heading angles. We assign the bin index task as a classification task, so we adopt the cross-entropy loss function.
- 2) Testing and post-processing: During testing, except for the dense depth estimation directly predicted from the network, we ought to utilize the post-processing to produce final 2D and 3D detection results. Specifically, we first infer the network to output the projected 3D center's heatmap. Then we calculate the local maximum of the heatmap and select the top-k's positions as the predicted 3D centers. According to 3D centers' positions, we can obtain from the predictions: the offset from the 2D center to the projected 3D center; 2D bbox's width and height; 3D bbox's width, height, and length; 3D object's depth, uncertainty, heading angles' bin index and intra-bin angle residual. Through the fisheye distortion module, the projected 3D center, together with the predicted object depth, can resume the 3D position. In addition, the object confidence depends on both the projected 3D center's score in the heatmap and the uncertainty, as follows:

$$C_{obj} = C_{proj\_center} * e^{(-\sigma)}$$
 (3)

where  $\sigma$  denotes the uncertainty. For heading angles, we recover the rotation with the predicted bin index and intra-bin residual, referring to MultiBin [93]. Furthermore, we project the 3D results to the Bird's-Eye-View (BEV) plane to obtain the BEV perception. Actually, when given the front, left, right and rear fisheye images at the same time, we can generate the 360-degree BEV perception after some fusion procedures like ReID or object filtering.

# D. Fisheye distortion module

We achieve the distortion-insensitive function by our Fisheye Distortion Module (FDM), which correctly builds fisheye projection between 3D space and 2D image plane to exclude the distortion's disturbing. This module has two main functions: (i) producing 2D labels from projecting 3D ground truth; (ii) restoring 3D positions from 2D image points. Compared with pinhole model projection, fisheye model projection needs to consider the influence of fisheye distortion. To achieve these purposes, we summarize our procedures of fisheye projection.

**Task (i).** Given the 3D point (x, y, z) in the camera coordinate system and camera parameters, we want to solve out the 2D point (u, v) in the image plane. According to the [18], for pinhole projection, we demand to calculate the field-angle  $\theta$  of the imaged point (the field angle of the projected ray [18]):

$$\theta = a \tan(r) \tag{4}$$

The pinhole projection coordinates (a,b) denote a=x/z, b=y/z and  $r=\sqrt{a^2+b^2}$ . Then for fisheye distortion projection, we utilize the fisheye distortion parameters to compute the rectified angle  $\theta_d$ :

$$\theta_d = \theta(1 + k_1\theta^2 + k_2\theta^4 + k_3\theta^6 + k_4\theta^8) \tag{5}$$

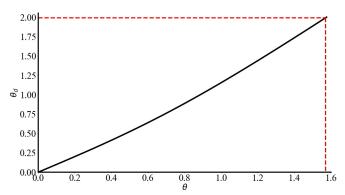


Fig. 9. The transformation between  $\theta$  and  $\theta_d$ :  $\theta_d$  is monotonically increasing with  $\theta$ , where  $\theta$  ranges from 0 to  $\pi/2$ .

Afterwards, we obtain the distorted point coordinates (x', y'):

$$x' = (\theta_d/r)a \tag{6}$$

$$y' = (\theta_d/r)b \tag{7}$$

Finally, we achieve the final pixel coordinates vector (u, v):

$$\begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = \begin{bmatrix} f_u & 0 & c_u \\ 0 & f_v & c_v \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x' \\ y' \\ 1 \end{bmatrix}$$
 (8)

where  $f_u, f_v, c_u, c_v$  are camera's intrinsic parameters. So according to Equ.(4)-(8), we successfully construct the projected 2D ground-truth coordinates (u, v) based on distorted ground-truth 3D coordinates (x, y, z).

**Task (ii).** Given the 2D image point (u, v), point's depth z and camera parameters, we need to recovery the 3D position (x, y, z). Firstly, we compute the distorted point coordinates (x', y'), according to the inverse operation:

$$\begin{bmatrix} x' \\ y' \\ 1 \end{bmatrix} = \begin{bmatrix} f_u & 0 & c_u \\ 0 & f_v & c_v \\ 0 & 0 & 1 \end{bmatrix}^{-1} \begin{bmatrix} u \\ v \\ 1 \end{bmatrix}$$
 (9)

Then combining the Equ.(6) and (7), we can acquire the  $\theta_d$ , as follows:

$$(x')^2 + (y')^2 = (\theta_d)^2 (a^2 + b^2)/r^2 = (\theta_d)^2$$
 (10)

$$\theta_d = \sqrt{(x')^2 + (y')^2}$$
 (11)

So in Equ. (5), we have to solve out the  $\theta$  with the given  $\theta_d$ , with a mathematical tool's help. Additionally, we restrict the solution  $\theta$  under the filter of the unique value ranging from 0 to  $\pi/2$  in reality, and finally we achieve the unique value  $\theta$ . Afterwards, we can calculate a and b from two equations:

$$a = (x' \cdot r)/\theta_d \tag{12}$$

$$b = (y' \cdot r)/\theta_d \tag{13}$$

where  $r = \sqrt{a^2 + b^2}$ .

At last, we acquire the 3D position (x, y, z) by:

$$x = a \cdot z \tag{14}$$

$$y = b \cdot z \tag{15}$$

So according to Equ. (9)-(15), we complete the precess of our task (ii) for the fisheye distortion module.

**Speeding up.** Reviewing the above procedures, we detect that it is time-consuming for Equ. (5), especially when solving

out the  $\theta$  given  $\theta_d$ . However, on closer observation (see Figure 9),  $\theta_d$  is monotonically increasing with  $\theta$ , where  $\theta$  ranges from 0 to  $\pi/2$ .

So for speeding up our framework, we further introduce a look-up table indexing approach. Specifically, we build a look-up table consisting of 900 grids in which  $\theta_d$  can be found with corresponded  $\theta$  or  $\theta$  with corresponded  $\theta_d$ . The procedure of looking up is economical of time, enhancing the real-time performance (see Section V-C3).

# E. Deployment

We deploy our model to the embedded system with Qualcomm 8155 chip, which only accepts the Qualcomm Deep Learning Container (DLC) file. So we transform our PyTorch model to the ONNX model, then convert the ONNX model to the DLC file with the help of Qualcomm conversion tools. In addition, we run our post-processing on the chip instead of the embedded system's CPU, to save CPU load (from 30ms to 3ms). Furthermore, we artificially quantize the DLC file to satisfy 8155 chip's AIP mode and reduce inference errors compared with the automatically quantized DSP mode.

#### V. EXPERIMENTS

This section exhibits our fisheye multi-task network FPNet on our **FPD**. Then we discuss **FPD**'s generalization across datasets. Furthermore, we analyze the effects of lightweight design and fisheye distortion module FDM by ablation study and comparison. Finally, we reveal several qualitative perception results by our approach.

# A. Evaluation Metrics

For 2D object detection, we select the IoU criterion of 0.7 for object detection metrics: Average Precision (AP) and Average Recall (AR), indicating  $AP_{2D}$  and  $AR_{2D}$ . For 3D object detection, the 3D Average Precision ( $AP_{3D}$ ) and BEV Average Precision ( $AP_{BEV}$ ) are two vital evaluation metrics. We utilize the IoU threshold of 0.5 for AP<sub>40</sub> with all categories. We adopt the depth metric for dense depth estimation, absolute relative error (abs rel.) Specially, less absolute relative error (abs rel.) means better performance, which acts different from the  $AP_{2D}$ ,  $AR_{2D}$ ,  $AP_{3D}$  and  $AP_{BEV}$ . We evaluate the hardware platform's inference speed with time consumption (ms).

#### B. Implementation Details

We conduct experiments with framework Pytorch on the Ubuntu system and employ eight NVIDIA RTX A6000s. The initial learning rate is set to be 0.1, and the momentum and learning decay rates are 0.9 and 0.01. We adopt stochastic gradient descent (SGD) solver and 16 batch size settings. For preprocessing of initial images (width: 1920 pixels and height: 1280 pixels), we firstly crop the height (upper 200 pixels and lower 210 pixels) to remove the blank object space, then resize them to the final size: 640 pixels width and 480 pixels height.

For 2D object detection, baseline detectors contain CenterNet [94] and RetinaNet [38]. For 3D object detection, we choose KM3D [59] and MonoCon [92] as our baseline detectors. For dense depth estimation, we pick DRO [63] as the baseline. All the baselines have occupied the related field

TABLE III

The performance of baselines on our FPD dataset with three visual tasks: 2D object detection, monocular 3D object detection, and monocular dense depth estimation. Highest result is marked with Red and the second highest is marked with blue. Different from the  $AP_{2D}$ ,  $AR_{2D}$ ,  $AP_{3D}$  and  $AP_{BEV}$ , less abs rel. means better performance.

Method	$AP_{2D}$	$AR_{2D}$	$AP_{3D}$	$AP_{BEV}$	abs rel.
RetinaNet [38]	49.3	44.6	-	-	-
CenterNet [94]	45.4	53.7	-	-	-
KM3D [59]	46.0	55.1	58.25	73.41	-
MonoCon [92]	48.8	54.3	63.71	76.49	-
DRO [63]	-	-	-	-	0.095
FPNet	53.4	57.2	66.38	80.74	0.088
Improvement	+4.1	+2.1	+2.67	+4.25	+0.007

TABLE IV

Evaluation of the generalizability of the FPD dataset to public datasets based on CenterNet [94] and DRO [63].  $A \Rightarrow B$  means adopting the model pre-trained on dataset A to finetune on dataset B.

Dataset	$AP_{2D}$	$AR_{2D}$	$AP_{3D}$	$AP_{BEV}$	abs rel.
COCO [95] FPD⇒COCO	45.2 <b>47.8</b>	39.5 <b>43.2</b>	-	-	-
KITTI [21] FPD⇒KITTI	80.6 <b>84.1</b>	82.2 <b>84.7</b>	48.01 <b>51.25</b>	53.39 <b>56.13</b>	0.056 <b>0.053</b>

in recent years. To ensure comparability, all baselines utilize the same experimental settings as their release. Additionally, we insert our fisheye distortion module to make them suitable for fisheye images.

Moreover, we choose two public datasets for our cross-dataset evaluation: COCO [95] and KITTI [21]. COCO [95] dataset only contains 2D object detection task, while KITTI [21] covers the 3D object detection, 2D object detection and depth estimation.

# C. Results and Analysis

1) **Results on FPD**: To demonstrate the effectiveness of our FPNet method, we conduct three tasks' baseline evaluation based on **FPD**, as shown in Table III.

For 2D object detection on FPD, our FPNet outperforms both anchor-based pipeline RetinaNet [38] and the anchorfree pipeline CenterNet [94], KM3D [59] and MonoCon [92]. For anchor-based RetinaNet [38], our FPD contains abundant objects with various and large-angle poses, so presetting anchors maybe not fit with our dataset. Our FPNet also follows the center-based structure [59], [92], [94], but these baselines [59], [92], [94] directly predict 2D bbox center, which reveals greater challenge, because our 2D labels come from projected 3D labels instead of directly labeling on 2D images, so our 2D bbox centers may not indicate the semantics of 2D object centers. Projected 3D center prediction with offset estimation from 3D center to 2D center assists FPNet to create a more correct 2D bbox center, resulting in better  $AP_{2D}$  and  $AR_{2D}$ , also greater  $AP_{3D}$  and  $AP_{BEV}$ , for precise 2D centers strengthen the other subtasks, like object center's depth and the heading angle. Consequently, our FPNet also makes considerable 3D metrics improvement, compared with KM3D [59]

TABLE V
ABLATION EXPERIMENTS WITH TWO SETTINGS FOR THE SHARED FEATURE EXTRACTOR: DOWN-SAMPLING RADIO (DR) AND LAYER REMOVAL (LR).

	DR	LR	$AP_{2D}$	$AP_{3D}$	abs rel.	time
(a)	4	-	53.7	66.47	0.076	72
(b)	4	$\checkmark$	53.7	66.46	0.076	64
(c)	8	-	53.5	66.40	0.078	40
(d)	8	$\checkmark$	53.4	66.38	0.078	23

#### TABLE VI

Ablation experiments with two settings for auxiliary monocular contexts: eight projected keypoints heatmap and offsets vectors estimation (8 Key.); and quantization residual estimation (Quan.).

	8 Key.	Quan.	$AP_{2D}$	$AP_{3D}$	abs rel.	time
(a)	✓	✓	53.6	66.50	0.079	40
(b)	-	$\checkmark$	53.5	66.41	0.078	25
(c)	✓	-	53.6	66.43	0.079	34
(d)	-	-	53.4	66.38	0.078	23

and MonoCon [92]. For dense depth estimation, the multi-task structure shares the perception knowledge, boosting the single depth estimation model, so our FPNet gains advancement.

- 2) Cross-dataset Evaluation: We validate our FPD's generalization to other public datasets. We train CenterNet [94] object detection models and DRO [63] depth estimation model on two general datasets, COCO [95] and KITTI [21], as shown in Table IV. Firstly we train the models on our FPD (first row), then adopt the pre-trained models to finetune on another dataset (second row). After finetuning, our FPD gains about 2% to 4% enhancement to the baselines directly trained on public datasets. FPD's various and diverse cases cover abundant autonomous driving scenes, which considerably lifts generalization ability.
- 3) Comparative Study for Lightweight Design: We conduct several comparative experiments for our targeted lightweight designs, where we manage to reduce the model's parameters, while still preserving the performance.

**Shared feature extractor.** Firstly, we explore the two changes with the shared feature extractor, as shown in Table V. From (a), the initial setting of normal down-sampling radio 4 and complete layers performs best. But from experiments (b) and (c), we find that only a tolerable and small performance drops when the input images are smaller (c) or remove some redundant layers (b). So we finally adopt the settings (d) to achieve the most lightweight shared feature extractor.

Other auxiliary monocular contexts. Secondly, we investigate the effects of leaving out other auxiliary monocular contexts. The other auxiliary monocular predictions consist of eight projected keypoints heatmap and offset vectors estimation, the quantization residual estimation with bbox center. The experiment is exhibited in Table VI. From settings (a) and (b), we see the additional keypoints tasks assist the 2D and 3D object detection, but consume more time. From settings (a) and (c), quantization residual estimation seems not much essential with little performance improvement. We discard both settings in (d), and the operation makes the progress in real-time aspect, still with great results.

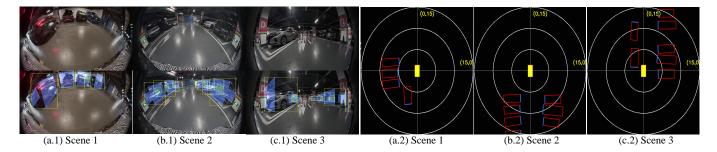


Fig. 10. More qualitative results with more objects and directions on **FPD**. (a) the input images; (b) the 2D object detection results (yellow bbox) and projected 3D object detection results (blue bbox); (c) BEV visualization from 3D results where the red bounding box denotes the predicted results and the blue segment denotes the objects' head.

# TABLE VII COMPARISON BETWEEN OUR LOOK-UP TABLE INDEXING AND SOLVEPOLY FUNCTION FROM OPENCY LIBRARY [96] WITH FPNET ON FPD. 'LTI' INDICATES THE 'LOOK-UP TABLE INDEXING'.

Method	$AP_{2D}$	$AP_{3D}$	abs rel.	time
(a)   SolvePoly [96]	53.5	66.38	0.078	120
(b)   LTI	<b>53.4</b>	<b>66.38</b>	<b>0.078</b>	<b>23</b>

TABLE VIII

COMPARISON BETWEEN INITIAL STRUCTURE DESIGN AND LIGHTWEIGHT
STRUCTURE DESIGN WITH FPNET ON FPD.

	Method	time	fps
(a)	w/o lightweight	180	5.5
(b)	w lightweight	<b>23</b>	<b>43.5</b>

Look-up table indexing. Thirdly, we discuss the look-up table indexing in the fisheye distortion module, as shown in Table VII. Before our lightweight design, we adopt the solvePoly function from Opency library [96], which costs much unaffordable time for real-time application (setting (a)). Apparently, the indexing approach (b) saves more time than (a) without disturbing the detection scores.

**Time summary**. Finally, we perform time consumption on the hardware platform, and Table VIII reveals the inference time with lightweight design and without lightweight design. From (a) and (b), lightweight design makes much sense, for we carefully develop advanced solutions to replace the elder's time-intensive operations.

4) Comparison with different bbox center settings: Different from [59], [92], we directly predict the projected 3D bbox center heatmap, instead of the 2D bbox center heatmap and offset from the projected 3D center to 2D center. Table IX exhibits that the direct 3D bbox center regression performs better. We blame that we produce 2D bbox ground truth from outer bounding rectangle of projected 3D labels, but this 2D bbox label's center may not reflect the actual object's 2D image center, such as Motorbike and Traffic\_cones demos in Figure 7, which alleviates the errors of 3D detection, especially the heading angle. While [59], [92] utilize the correct 2D bbox center from correct 2D bbox labels. Consequently, we adopt the proposal of directly regressing 3D bbox center and offset from 2D center to the projected 3D center.

# D. Qualitative Results

We provide qualitative examples on our **FPD**, as shown in Figure 11 and Figure 10. Figure 11 displays some qualitative

TABLE IX COMPARISON BETWEEN TWO BBOX CENTER SETTINGS.  $2D{\Rightarrow}3D$  INDICATES THE OFFSET FROM 2D CENTER TO THE PROJECTED 3D CENTER AND VICE VERSA.

Method	$AP_{2D}$	$AP_{3D}$	abs rel.
(a) $ $ 2D center + 2D $\Rightarrow$ 3D	48.2	60.25	0.078
(b) $ $ 3D proj. center + 3D $\Rightarrow$ 2D	<b>53.4</b>	<b>66.38</b>	<b>0.078</b>

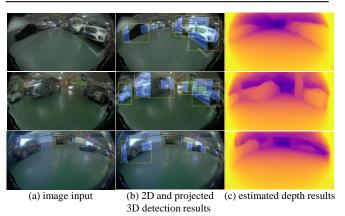


Fig. 11. Qualitative results with several valet parking scenes on **FPD**. (a) the input images; (b) the 2D object detection results (yellow bbox) and projected 3D object detection results (blue bbox); (c) estimated depth results.

results with several valet parking scenes. (b) denotes the 2D and 3D projected results and (c) denotes the estimation dense depth. Our multi-task network satisfies the real-time perception requirement in different valet parking scenes.

Figure 10 demonstrates some visualization results with more objects and other directions. Notably, (c) denotes the BEV visualization from 3D detection results. From Figure 10, our network achieves good perception performance for crowded objects in different directions.

### VI. CONCLUSION

This article presents a new large-scale fisheye dataset, the Fisheye Parking Dataset (FPD). By providing a diversity of surround-view parking scenes, the proposed dataset aims to assist the industry in constructing a more secure advanced driving assistance system in parking lots. Moreover, we provide our real-time multi-task Fisheye Perception Network (FPNet), for enhancing fisheye distortion performance and various lightweight designs. Extensive experiments on FPD validate the effectiveness of our FPNet. However, FPD has a lot of capacity for development, including how to strengthen more

data diversity, simplify our approach, and deal with sustainably increasing data and the potential for diverse vision tasks. Nevertheless, we expect **FPD** to inspire more relevant research and promote the performance of surround-view perception under parking scenes.

In the future, we will further explore the surround-view fisheye BEV perception for valet parking, in the following aspects. (1) We will explore the LiDAR-camera fused perception since the LiDAR point cloud provides more 3D information. (2) We will explore the lightweight image features extractor to achieve more robust visual features. (3) We will design specific techniques to augment the distorted fisheye images to enhance the generalization ability. (4) We will apply our surround-view fisheye BEV perception and dataset refer to other tasks, such as obstacle detection.

#### REFERENCES

- J. E. Hoffmann, H. G. Tosso, M. M. D. Santos, J. F. Justo, A. W. Malik, and A. U. Rahman, "Real-time adaptive object detection and tracking for autonomous vehicles," *IEEE Transactions on Intelligent Vehicles*, vol. 6, no. 3, pp. 450–459, 2020.
- [2] H. Wang, Y. Yu, Y. Cai, X. Chen, L. Chen, and Y. Li, "Soft-weighted-average ensemble vehicle detection method based on single-stage and two-stage deep learning models," *IEEE Transactions on Intelligent Vehicles*, vol. 6, no. 1, pp. 100–109, 2020.
- [3] T. Gao, H. Pan, and H. Gao, "Monocular 3d object detection with sequential feature association and depth hint augmentation," *IEEE Transactions on Intelligent Vehicles*, 2022.
- [4] X. Song, P. Wang, D. Zhou, R. Zhu, C. Guan, Y. Dai, H. Su, H. Li, and R. Yang, "Apollocar3D: A large 3D car instance understanding benchmark for autonomous driving," in *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2019, pp. 5452–5462.
- [5] D. Zhou, J. Fang, X. Song, L. Liu, J. Yin, Y. Dai, H. Li, and R. Yang, "Joint 3D instance segmentation and object detection for autonomous driving," in *IEEE Conf. on Computer Vision and Pattern Recognition* (CVPR), 2020, pp. 1839–1849.
- [6] P. Wu, S. Chen, and D. N. Metaxas, "MotionNet: Joint perception and motion prediction for autonomous driving based on bird's eye view maps," in *IEEE Conf. on Computer Vision and Pattern Recognition* (CVPR). IEEE, 2020, pp. 11385–11395.
- [7] K. Samal, H. Kumawat, P. Saha, M. Wolf, and S. Mukhopadhyay, "Task-driven rgb-lidar fusion for object tracking in resource-efficient autonomous system," *IEEE Transactions on Intelligent Vehicles*, vol. 7, no. 1, pp. 102–112, 2021.
- [8] Z. Chen and X. Huang, "Pedestrian detection for autonomous vehicle using multi-spectral cameras," *IEEE Transactions on Intelligent Vehicles*, vol. 4, no. 2, pp. 211–219, 2019.
- [9] D. Yang, H. Zhang, E. Yurtsever, K. A. Redmill, and Ü. Özgüner, "Predicting pedestrian crossing intention with feature fusion and spatiotemporal attention," *IEEE Transactions on Intelligent Vehicles*, vol. 7, no. 2, pp. 221–230, 2022.
- [10] Z. Wu, W. Sun, M. Wang, X. Wang, L. Ding, and F. Wang, "PSDet: Efficient and universal parking slot detection," in *IEEE Intelligent Vehicles Symposium*. IEEE, 2020, pp. 290–297.
- [11] S. Chen, N. Zhang, and H. Sun, "Collaborative localization based on traffic landmarks for autonomous driving," in *IEEE International Symposium on Circuits and Systems (ISCAS)*. IEEE, 2020, pp. 1–5.
- [12] V. R. Kumar, S. Yogamani, H. Rashed, G. Sitsu, C. Witt, I. Leang, S. Milz, and P. Mäder, "OmniDet: Surround view cameras based multitask visual perception network for autonomous driving," *IEEE Robotics* and Automation Letters, vol. 6, no. 2, pp. 2830–2837, 2021.
- [13] B. Li, T. Zhang, and T. Xia, "Vehicle detection from 3D lidar using fully convolutional network," arXiv preprint arXiv:1608.07916, 2016.
- [14] X. Chen, K. Kundu, Z. Zhang, H. Ma, S. Fidler, and R. Urtasun, "Monocular 3D object detection for autonomous driving," in *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2016, pp. 2147–2156.
- [15] Y. Su, Y. Gao, Y. Zhang, J. M. Alvarez, J. Yang, and H. Kong, "An illumination-invariant nonparametric model for urban road detection," *IEEE Transactions on Intelligent Vehicles*, vol. 4, no. 1, pp. 14–23, 2018.

- [16] T. Hehn, J. Kooij, and D. Gavrila, "Fast and compact image segmentation using instance stixels," *IEEE Transactions on Intelligent Vehicles*, vol. 7, no. 1, pp. 45–56, 2021.
- [17] I. Baek, A. Davies, G. Yan, and R. R. Rajkumar, "Real-time detection, tracking, and classification of moving and stationary objects using multiple fisheye images," in *IEEE Intelligent Vehicles Symposium*. IEEE, 2018, pp. 447–452.
- [18] V. R. Kumar, C. Eising, C. Witt, and S. Yogamani, "Surround-view fisheye camera perception for automated driving: Overview, survey and challenges," arXiv preprint arXiv:2205.13281, 2022.
- [19] Z. Wu, W. Zhang, J. Wang, M. Wang, Y. Gan, X. Gou, M. Fang, and J. Song, "Disentangling and Vectorization: A 3D visual perception approach for autonomous driving based on surround-view fisheye cameras," in *IEEE Int. Conf. on Intelligent Robots and Systems (IROS)*. IEEE, 2021, pp. 5576–5582.
- [20] H. Rashed, E. Mohamed, G. Sistu, V. R. Kumar, C. Eising, A. El-Sallab, and S. Yogamani, "Generalized object detection on fisheye cameras for autonomous driving: Dataset, representations and baseline," in *IEEE Winter Conf. on Applications of Computer Vision (WACV)*, 2021, pp. 2272–2280.
- [21] A. Geiger, P. Lenz, and R. Urtasun, "Are we ready for autonomous driving? the kitti vision benchmark suite," in *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2012, pp. 3354–3361.
- [22] S. Zhang, R. Benenson, and B. Schiele, "Citypersons: A diverse dataset for pedestrian detection," in *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2017, pp. 3213–3221.
- [23] S. Yogamani, C. Hughes, J. Horgan, G. Sistu, P. Varley, D. O'Dea, M. Uricár, S. Milz, M. Simon, K. Amende et al., "WoodScape: A multi-task, multi-camera fisheye dataset for autonomous driving," in *IEEE Int. Conf. on Computer Vision (ICCV)*. IEEE, 2019, pp. 9308–9318.
- [24] Y. Liao, J. Xie, and A. Geiger, "KITTI-360: A novel dataset and benchmarks for urban scene understanding in 2D and 3D," arXiv preprint arXiv:2109.13410, 2021.
- [25] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, and B. Schiele, "The Cityscapes dataset for semantic urban scene understanding," in *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 3213–3223.
- [26] H. Caesar, V. Bankiti, A. H. Lang, S. Vora, V. E. Liong, Q. Xu, A. Krishnan, Y. Pan, G. Baldan, and O. Beijbom, "nuScenes: A multimodal dataset for autonomous driving," in *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 11621–11631.
- [27] P. Sun, H. Kretzschmar, X. Dotiwalla, A. Chouard, V. Patnaik, P. Tsui, J. Guo, Y. Zhou, Y. Chai, B. Caine et al., "Scalability in perception for autonomous driving: Waymo open dataset," in *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 2446–2454.
- [28] N. Silberman, D. Hoiem, P. Kohli, and R. Fergus, "Indoor segmentation and support inference from RGBD images," in *European Conf. on Computer Vision (ECCV)*. Springer, 2012, pp. 746–760.
- [29] X. Pan, J. Shi, P. Luo, X. Wang, and X. Tang, "Spatial as deep: Spatial cnn for traffic scene understanding," in AAAI Conf. on Artificial Intell. (AAAI), vol. 32, no. 1, 2018.
- [30] M. Toromanoff, E. Wirbel, F. Wilhelm, C. Vejarano, X. Perrotton, and F. Moutarde, "End to end vehicle lateral control using a single fisheye camera," in *IEEE Int. Conf. on Intelligent Robots and Systems (IROS)*. IEEE, 2018, pp. 3613–3619.
- [31] V. R. Kumar, S. A. Hiremath, M. Bach, S. Milz, C. Witt, C. Pinard, S. Yogamani, and P. Mäder, "FisheyeDistanceNet: Self-supervised scaleaware distance estimation using monocular fisheye camera for autonomous driving," in *IEEE Int. Conf. on Robotics and Automation* (ICRA). IEEE, 2020, pp. 574–581.
- [32] K. Yang, X. Hu, Y. Fang, K. Wang, and R. Stiefelhagen, "Omnisupervised omnidirectional semantic segmentation," *IEEE Transactions on Intelligent Transportation Systems*, 2020.
- [33] I. Armeni, S. Sax, A. R. Zamir, and S. Savarese, "Joint 2d-3d-semantic data for indoor scene understanding," arXiv preprint arXiv:1702.01105, 2017.
- [34] A. Chang, A. Dai, T. Funkhouser, M. Halber, M. Niessner, M. Savva, S. Song, A. Zeng, and Y. Zhang, "Matterport3D: Learning from rgb-d data in indoor environments," arXiv preprint arXiv:1709.06158, 2017.
- [35] N. Zioulis, A. Karakottas, D. Zarpalas, and P. Daras, "Omnidepth: Dense depth estimation for indoors spherical panoramas," in *European Conf.* on Computer Vision (ECCV), 2018, pp. 448–465.
- [36] F.-E. Wang, H.-N. Hu, H.-T. Cheng, J.-T. Lin, S.-T. Yang, M.-L. Shih, H.-K. Chu, and M. Sun, "Self-supervised learning of depth and camera motion from 360° videos," in *Asian Conf. on Computer Vision (ACCV)*. Springer, 2018, pp. 53–68.

- [37] D. L. Stone, S. Ravi, E. Benli, and Y. Motai, "Deepfusenet of omnidirectional far-infrared and visual stream for vegetation detection," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 59, no. 11, pp. 9057–9070, 2021.
- [38] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, "Focal loss for dense object detection," in *IEEE Int. Conf. on Computer Vision (ICCV)*. IEEE, 2017, pp. 2980–2988.
- [39] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You only look once: Unified, real-time object detection," in *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2016, pp. 779–788.
- [40] R. Girshick, "Fast R-CNN," in IEEE Int. Conf. on Computer Vision (ICCV). IEEE, 2015, pp. 1440–1448.
- [41] S. Ren, K. He, R. Girshick, and J. Sun, "Faster R-CNN: Towards realtime object detection with region proposal networks," *Conference and Workshop on Neural Information Processing Systems (NeurIPS)*, vol. 28, 2015.
- [42] K. He, G. Gkioxari, P. Dollár, and R. Girshick, "Mask R-CNN," in *IEEE Int. Conf. on Computer Vision (ICCV)*. IEEE, 2017, pp. 2961–2969.
- [43] C.-Y. Fu, W. Liu, A. Ranga, A. Tyagi, and A. C. Berg, "DSSD: Deconvolutional single shot detector," arXiv preprint arXiv:1701.06659, 2017.
- [44] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg, "SSD: Single shot multibox detector," in *European Conf. on Computer Vision (ECCV)*. Springer, 2016, pp. 21–37.
- [45] B. X. Zhang S, Wen L, "Single-shot refinement neural network for object detection," in *IEEE Conf. on Computer Vision and Pattern Recognition* (CVPR). IEEE, 2018.
- [46] Q. Chen, Y. Wang, T. Yang, X. Zhang, J. Cheng, and J. Sun, "You only look one-level feature," in *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2021, pp. 13039–13048.
- [47] S. Ren, K. He, R. Girshick, and J. Sun, "Faster R-CNN: Towards realtime object detection with region proposal networks," *Conference and Workshop on Neural Information Processing Systems (NeurIPS)*, vol. 28, 2015.
- [48] Y. Wu, Y. Chen, L. Yuan, Z. Liu, L. Wang, H. Li, and Y. Fu, "Rethinking classification and localization for object detection," in *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2020, pp. 10186–10195.
- [49] P. Yang, G. Zhang, L. Wang, L. Xu, Q. Deng, and M.-H. Yang, "A part-aware multi-scale fully convolutional network for pedestrian detection," IEEE Transactions on Intelligent Transportation Systems, vol. 22, no. 2, pp. 1125–1137, 2020.
- [50] G. Feng, J. Meng, L. Zhang, and H. Lu, "Encoder deep interleaved network with multi-scale aggregation for RGB-D salient object detection," *Pattern Recognition*, vol. 128, p. 108666, 2022.
- [51] J. Ku, A. D. Pon, and S. L. Waslander, "Monocular 3D object detection leveraging accurate proposals and shape reconstruction," in *IEEE Conf.* on Computer Vision and Pattern Recognition (CVPR), 2019, pp. 11867– 11876.
- [52] G. Brazil and X. Liu, "M3D-RPN: monocular 3D region proposal network for object detection," in *IEEE Int. Conf. on Computer Vision* (ICCV). IEEE, 2019, pp. 9287–9296.
- [53] A. Simonelli, S. R. Bulo, L. Porzi, M. López-Antequera, and P. Kontschieder, "Disentangling monocular 3D object detection," in IEEE Int. Conf. on Computer Vision (ICCV), 2019, pp. 1991–1999.
- [54] L. Liu, J. Lu, C. Xu, Q. Tian, and J. Zhou, "Deep fitting degree scoring network for monocular 3D object detection," in *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 1057–1066.
- [55] 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 AAAI Conf. on Artificial Intell. (AAAI), 2020, pp. 10478–10485.
- [56] Y. Zhang, J. Lu, and J. Zhou, "Objects are Different: Flexible monocular 3D object detection," in *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2021, pp. 3289–3298.
- [57] X. Zhou, D. Wang, and P. Krähenbühl, "Objects as points," arXiv preprint arXiv:1904.07850, 2019.
- [58] Z. Liu, Z. Wu, and R. Tóth, "SMOKE: Single-stage monocular 3D object detection via keypoint estimation," in *IEEE Conf. on Computer Vision* and Pattern Recognition Workshops (CVPRW). IEEE, 2020, pp. 996– 997.
- [59] P. Li and H. Zhao, "Monocular 3D detection with geometric constraint embedding and semi-supervised training," *IEEE Robotics and Automa*tion Letters, vol. 6, no. 3, pp. 5565–5572, 2021.
- [60] P. Li, H. Zhao, P. Liu, and F. Cao, "RTM3D: Real-time monocular 3D detection from object keypoints for autonomous driving," in *European Conf. on Computer Vision (ECCV)*. Springer, 2020, pp. 644–660.

- [61] B. Arsenali, P. Viswanath, and J. Novosel, "RotInvMTL: Rotation invariant multinet on fisheye images for autonomous driving applications," in *IEEE Int. Conf. on Computer Vision Workshops (ICCVW)*, 2019, pp. 0–0.
- [62] E. Plaut, E. Ben Yaacov, and B. El Shlomo, "Monocular 3D object detection in cylindrical images from fisheye cameras," arXiv e-prints, pp. arXiv-2003, 2020.
- [63] X. Gu, W. Yuan, Z. Dai, C. Tang, S. Zhu, and P. Tan, "DRO: Deep recurrent optimizer for structure-from-motion," in AAAI Conf. on Artificial Intell. (AAAI), 2021.
- [64] S. Lee, J. Lee, B. Kim, E. Yi, and J. Kim, "Patch-wise attention network for monocular depth estimation," in AAAI Conf. on Artificial Intell. (AAAI), vol. 35, no. 3, 2021, pp. 1873–1881.
- [65] V. Patil, C. Sakaridis, A. Liniger, and L. Van Gool, "P3Depth: Monocular depth estimation with a piecewise planarity prior," in *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2022.
- [66] D. Eigen, C. Puhrsch, and R. Fergus, "Depth map prediction from a single image using a multi-scale deep network," Conference and Workshop on Neural Information Processing Systems (NeurIPS), vol. 27, 2014.
- [67] I. Laina, C. Rupprecht, V. Belagiannis, F. Tombari, and N. Navab, "Deeper depth prediction with fully convolutional residual networks," in *Int. Conf. on 3D Vision (3DV)*, 2016, pp. 239–248.
- [68] R. Ranftl, A. Bochkovskiy, and V. Koltun, "Vision transformers for dense prediction," in *IEEE Int. Conf. on Computer Vision (ICCV)*, 2021, pp. 12 179–12 188.
- [69] G. Yang, H. Tang, M. Ding, N. Sebe, and E. Ricci, "Transformers solve the limited receptive field for monocular depth prediction," arXiv preprint arXiv:2103.12091, 2021.
- [70] G. Blott, M. Takami, and C. Heipke, "Semantic segmentation of fisheye images," in *European Conf. on Computer Vision (ECCV)*, 2018, pp. 181–196.
- [71] V. R. Kumar, S. Milz, C. Witt, M. Simon, K. Amende, J. Petzold, S. Yogamani, and T. Pech, "Near-field depth estimation using monocular fisheye camera: A semi-supervised learning approach using sparse lidar data," in *IEEE Conf. on Computer Vision and Pattern Recognition* Workshops (CVPRW), vol. 7, 2018.
- [72] V. R. Kumar, M. Klingner, S. Yogamani, S. Milz, T. Fingscheidt, and P. Mader, "Syndistnet: Self-supervised monocular fisheye camera distance estimation synergized with semantic segmentation for autonomous driving," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2021, pp. 61–71.
- [73] T. Li, G. Tong, H. Tang, B. Li, and B. Chen, "Fisheyedet: A self-study and contour-based object detector in fisheye images," *IEEE Access*, vol. 8, pp. 71739–71751, 2020.
- [74] M. Yahiaoui, H. Rashed, L. Mariotti, G. Sistu, I. Clancy, L. Yahiaoui, V. R. Kumar, and S. Yogamani, "Fisheyemodnet: Moving object detection on surround-view cameras for autonomous driving," arXiv preprint arXiv:1908.11789, 2019.
- [75] Z. Wu, M. Wang, L. Yin, W. Sun, J. Wang, and H. Wu, "Vehicle reid for surround-view camera system," arXiv preprint arXiv:2006.16503, 2020.
- [76] D. Dooley, B. McGinley, C. Hughes, L. Kilmartin, E. Jones, and M. Glavin, "A blind-zone detection method using a rear-mounted fisheye camera with combination of vehicle detection methods," *IEEE Transactions on Intelligent Transportation Systems*, vol. 17, no. 1, pp. 264–278, 2015.
- [77] C.-Y. Yang and H. H. Chen, "Efficient face detection in the fisheye image domain," *IEEE Trans. Image Proc. (TIP)*, vol. 30, pp. 5641–5651, 2021.
- [78] L. Mariotti and C. Eising, "Spherical formulation of geometric motion segmentation constraints in fisheye cameras," *IEEE Transactions on Intelligent Transportation Systems*, 2020.
- [79] V. R. Kumar, M. Klingner, S. Yogamani, M. Bach, S. Milz, T. Fingscheidt, and P. Mäder, "SVDistNet: Self-supervised near-field distance estimation on surround view fisheye cameras," *IEEE Transactions on Intelligent Transportation Systems*, 2021.
- [80] I.-C. Lo, K.-T. Shih, and H. H. Chen, "Efficient and accurate stitching for 360° dual-fisheye images and videos," *IEEE Trans. Image Proc.* (TIP), vol. 31, pp. 251–262, 2021.
- [81] Y.-H. Li, I.-C. Lo, and H. H. Chen, "Deep face rectification for 3600 dual-fisheye cameras," *IEEE Trans. Image Proc. (TIP)*, vol. 30, pp. 264– 276, 2020.
- [82] M. Shere, H. Kim, and A. Hilton, "Temporally consistent 3d human pose estimation using dual 360deg cameras," in *IEEE Winter Conf. on Applications of Computer Vision (WACV)*, 2021, pp. 81–90.

- [83] M. Xu, L. Yang, X. Tao, Y. Duan, and Z. Wang, "Saliency prediction on omnidirectional image with generative adversarial imitation learning," *IEEE Trans. Image Proc. (TIP)*, vol. 30, pp. 2087–2102, 2021.
- [84] C. Mao, A. Gupta, V. Nitin, B. Ray, S. Song, J. Yang, and C. Vondrick, "Multitask learning strengthens adversarial robustness," in *European Conf. on Computer Vision (ECCV)*. Springer, 2020, pp. 158–174.
- [85] S. Vandenhende, S. Georgoulis, W. Van Gansbeke, M. Proesmans, D. Dai, and L. Van Gool, "Multi-task learning for dense prediction tasks: A survey," *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 2021.
- [86] M. Teichmann, M. Weber, M. Zoellner, R. Cipolla, and R. Urtasun, "MultiNet: Real-time joint semantic reasoning for autonomous driving," in *IEEE Intelligent Vehicles Symposium (IV)*. IEEE, 2018, pp. 1013– 1020.
- [87] A. Kendall, Y. Gal, and R. Cipolla, "Multi-task learning using uncertainty to weigh losses for scene geometry and semantics," in *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 7482–7491
- [88] G. Sistu, I. Leang, S. Chennupati, S. Yogamani, C. Hughes, S. Milz, and S. Rawashdeh, "NeurAll: Towards a unified visual perception model for automated driving," in *International Conference on Intelligent Transportation Systems (ITSC)*. IEEE, 2019, pp. 796–803.
- [89] K. Li, H. Xiong, J. Liu, Q. Xu, and J. Wang, "Real-time monocular joint perception network for autonomous driving," *IEEE Transactions* on *Intelligent Transportation Systems*, 2022.
- [90] J. Ku, A. Harakeh, and S. L. Waslander, "In defense of classical image processing: Fast depth completion on the cpu," in *Proc. Canadian Conf.* on Computer and Robot Vision. IEEE, 2018, pp. 16–22.
- [91] F. Yu, D. Wang, E. Shelhamer, and T. Darrell, "Deep layer aggregation," in *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2018, pp. 2403–2412.
- [92] X. Liu, N. Xue, and T. Wu, "Learning auxiliary monocular contexts helps monocular 3D object detection," in AAAI Conf. on Artificial Intell. (AAAI), 2021.
- [93] A. Mousavian, D. Anguelov, J. Flynn, and J. Kosecka, "3D bounding box estimation using deep learning and geometry," in *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 7074–7082.
- [94] K. Duan, S. Bai, L. Xie, H. Qi, Q. Huang, and Q. Tian, "CenterNet: Keypoint triplets for object detection," in *IEEE Int. Conf. on Computer Vision (ICCV)*. IEEE, 2019, pp. 6569–6578.
- [95] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, and D. Ramanan, "Microsoft COCO: Common objects in context," in *European Conf. on Computer Vision (ECCV)*. Springer, 2014, pp. 740–755.
- [96] G. Bradski, "The opency library." Dr. Dobb's Journal: Software Tools for the Professional Programmer, vol. 25, no. 11, pp. 120–123, 2000.