

Multi-Modal 3D Object Detection in Autonomous Driving: A Survey

Yingjie Wang* · Qiuyu Mao* · Hanqi Zhu · Jiajun Deng · Yu Zhang ·
 Jianmin Ji · Houqiang Li · Yanyong Zhang

Received: date / Accepted: date

Abstract The past decade has witnessed the rapid development of autonomous driving systems. However, it remains a daunting task to achieve full autonomy, especially when it comes to understanding the ever-changing, complex driving scenes. To alleviate the difficulty of perception, self-driving vehicles are usually equipped with a suite of sensors (*e.g.*, cameras, LiDARs), hoping to capture the scenes with overlapping perspectives to minimize blind spots. Fusing these data streams and exploiting their complementary properties is thus rapidly becoming the current trend.

*equal contribution

Yingjie Wang
 University of Science and Technology of China
 E-mail: yingjiewang@mail.ustc.edu.cn

Qiuyu Mao
 University of Science and Technology of China
 E-mail: qymao@mail.ustc.edu.cn

Hanqi Zhu
 University of Science and Technology of China
 E-mail: zhuhanqi@mail.ustc.edu.cn

Jiajun Deng
 University of Science and Technology of China
 E-mail: dengjj@ustc.edu.cn

Yu Zhang
 University of Science and Technology of China
 E-mail: yuzhang@ustc.edu.cn

Jianmin Ji
 University of Science and Technology of China
 E-mail: jianmin@ustc.edu.cn

Houqiang Li
 University of Science and Technology of China
 E-mail: lihq@ustc.edu.cn

Yanyong Zhang, corresponding author
 University of Science and Technology of China
 E-mail: yanyongz@ustc.edu.cn

Nonetheless, combining data that are captured by different sensors with drastically different ranging/imaging mechanisms is not a trivial task; instead, many factors need to be considered and optimized. If not careful, data from one sensor may act as noises to data from another sensor, with even poorer results by fusing them. Thus far, there has been no in-depth guidelines to designing the multi-modal fusion based 3D perception algorithms. To fill in the void and motivate further investigation, this survey conducts a thorough study of tens of recent deep learning based multi-modal 3D detection networks (with a special emphasis on LiDAR-camera fusion), focusing on their fusion stage (*i.e.*, when to fuse), fusion inputs (*i.e.*, what to fuse), and fusion granularity (*i.e.*, how to fuse). These important design choices play a critical role in determining the performance of the fusion algorithm.

In this survey, we first introduce the background of popular sensors used for self-driving, their data properties, and the corresponding object detection algorithms. Next, we discuss existing datasets that can be used for evaluating multi-modal 3D object detection algorithms. Then we present a review of multi-modal fusion based 3D detection networks, taking a close look at their fusion stage, fusion input and fusion granularity, and how these design choices evolve with time and technology. After the review, we discuss open challenges as well as possible solutions. We hope that this survey can help researchers to get familiar with the field and embark on investigations in the area of multi-modal 3D object detection.

Keywords 3D Object Detection · Multi-modal Fusion · Sensor Fusion · Autonomous Driving

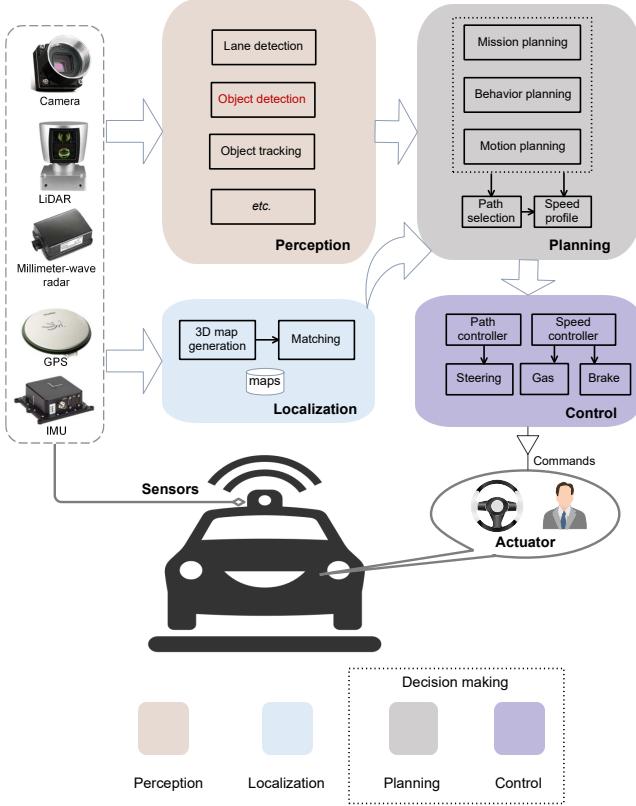


Fig. 1 The typical architecture for an autonomous driving system, consisting of three subsystems: perception, localization and decision making

1 Introduction

Recent breakthroughs in deep learning and computer vision [16, 67, 138] have fostered the rapid development of autonomous driving, which promises to free the drivers, to decrease traffic congestion, and to improve road safety. The potential of autonomous driving is, however, not yet fully unleashed, largely due to the unsatisfactory perception performance in real-world driving scenarios. As a result, even if autonomous vehicles (AVs) have seen applications in many confined and controlled environments, deploying them in urban environments still poses dire technological challenges [51, 153].

Fig. 1 illustrates a typical AV system that consists of three subsystems: perception, localization and decision making. The AV system capitalizes multiple sensors (*e.g.*, LiDAR, camera) to collect raw sensor data. Taking the raw sensor data as input, the perception and localization subsystems execute several important tasks to identify and localize objects of interest, namely, object detection, tracking, 3D map generation and mapping, *etc.* Given the objects and their locations, the decision making subsystem can navigate and make self-driving decisions. Among all the tasks, object detection,

aiming to localize and categorize objects of interest, is of great significance.

With the breakthrough of deep learning techniques, 2D object detection has drawn a great deal of attention, resulting in a plethora of algorithms [46, 47, 88, 90, 92, 128, 129]. However, localizing the objects in the 2D image plane is far from the demand of AVs to perceive the 3D real world. To this end, the task of 3D object detection is proposed with the requirement of predicting the object's three-dimensional location, shape, and rotational angles. Compared to the well-studied 2D object detection, 3D object detection is not only more important to autonomous vehicles but also more challenging. The challenges mainly stem from the fact that 3D driving scenes are also much more complex for perception [194]. For example, we need additional depth and rotation parameters to locate an object in 3D space.

In the real world, performing 3D object detection through a single type of sensor data is far from being sufficient. Firstly, each type of sensor data has its inherent limitation and shortcomings. For example, a camera-based system suffers from the lack of accurate depth information, while a LiDAR-only system is hampered by lower input data resolution, especially at long distances. As shown in Fig. 2 and Tab. 1, in average, for objects which are far from the ego-sensor ($> 60m$ in KITTI), there are usually less than 10 LiDAR points but are still with more than 400 image pixels. Secondly, the perception system must be robust against sensor malfunctioning, failure, or simply under-performing, hence mandating the necessity of having more than one type of sensor. Thirdly, data from different sensors complement each other naturally. Their combination could lead to a more comprehensive depicting of the environment and thus better detection results.

Therefore, a recent trend in 3D object detection is to *combine data streams from different sensors and develop multi-modal detection methods*. Fig. 3 shows multiple sensors in the AV system. AVs are typically equipped with cameras, LiDARs (*i.e.*, Light Detection And Ranging sensors), Radars (Radio detection and ranging sensors), GPS (Global Positioning System) and IMUs (Inertial Measurement Units) [154, 193]. In the multi-modal methods, data from multiple types of sensors that have complementary characteristics are fused to capture the scenes with overlapping perspectives, aiming at minimizing blind spots.

Though recent studies have demonstrated the benefits of fusion in various settings, conducting efficient and effective multi-modal detection in the real world still largely remains a myth and faces many challenges. Below we list some of these open challenges:

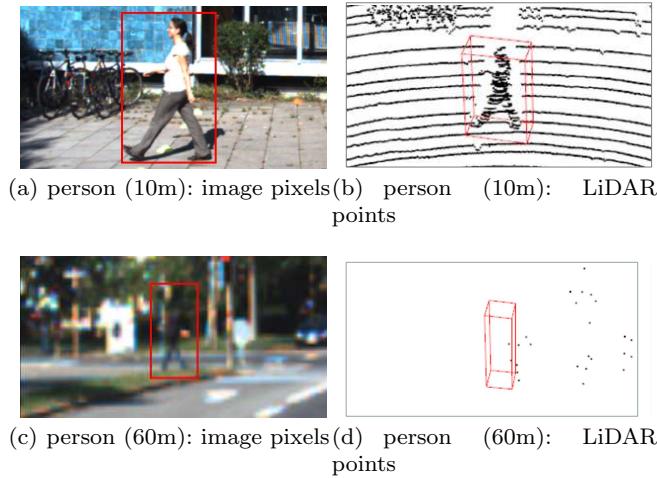


Fig. 2 The image for a person 10 meters away is shown in (a), and the corresponding point cloud data is shown in (b). The image for another person 60 meters away is shown in (c), and the corresponding point cloud data is shown in (d). It is clear that point clouds get very sparse at long distances. We modify the picture from Fig. 1 in [191].

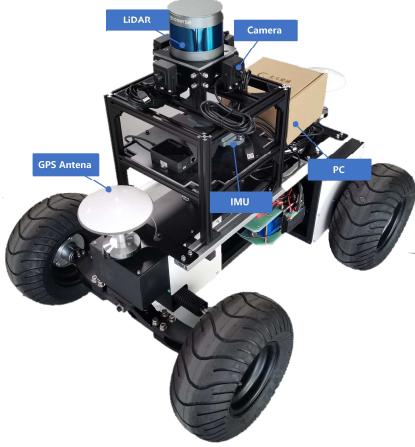


Fig. 3 The autonomous car *Sonic* is equipped with one LiDAR sensor (Velodyne VLP-16), four cameras, and one GPS. Note that the image is modified from [193].

- **Multi-Sensor Calibration:** Sensors of different types are not synchronized either temporally or spatially. In the temporal domain, it is hard to collect data at the same time due to independent acquisition cycles for each sensor. In the spatial domain, sensors have different angles of view when they are deployed. Thus, multi-sensor calibration is the first step before data fusion, which has not received much attention so far.
- **Information Loss During Fusion:** Due to the large gap between different types of sensor data (illustrated in Tab. 1), it is difficult to precisely align these data streams either in the input stage or in the feature space. To convert the sensor data into a

representation format in which they can be aligned and fused correctly, a certain amount of information loss becomes inevitable.

- **Consistent Data Augmentation Across Multiple Modalities:** Data augmentation plays a vital role in 3D object detection to enhance the size of training samples, and to ameliorate the problem of model over-fitting [186]. Augmentation strategies such as global rotation [198] and random flip [141] are widely adopted by LiDAR-based and camera-based methods but are absent in many multi-modal methods due to the concerns of leading to inconsistencies across modalities.

At present, how to address the above challenges and conduct efficient data fusion still remains an open problem. If not carefully done, data fusion may cause different data streams to act as noises to each other [5, 10], leading to even poorer results. In this paper, we set out to conduct a comprehensive review of recent fusion-based 3D object detection methods. Such a review can help pinpoint technical challenges in sensor fusion, and help us compare and contrast various models proposed to address these challenges. In particular, since cameras and LiDARs are the most common sensors for autonomous driving, our review mainly focuses on the fusion of these two types of sensor data. Specifically, when we discuss a multi-modal fusion based 3D detection algorithm, we focus on how the algorithm deals with the following three crucial design considerations:

- **Fusion Stage:** The first design consideration is concerned with at what pipeline stage the multi-modal fusion module takes place, *i.e.*, “where to fuse”. It has three options here: early fusion [156, 173], late fusion [112], and cascade fusion [123]. Early fusion usually occurs in the input stage or feature extraction stage before each branch reaches its prediction. Late fusion takes place in the prediction stage. Cascade fusion employs the hybrid mode by fusing one branch’s prediction with the other’s input. The fusion stage is the most influential design consideration as it determines the overall network architecture of the fusion based detection algorithm, and early fusion is the predominant choice.
- **Fusion Input:** The second design consideration is concerned with how the multi-modal data are input into the fusion module, *i.e.*, “what to fuse”. The fusion module can be designed to take the raw data as input, or some type of intermediate features. For example, the fusion module can take in the LiDAR data as raw point clouds [173, 65], voxel grids [187, 23, 22], and projection on the bird’s eye view (BEV) or the range view (RV) [19, 151]. Meanwhile, the fusion module can take in the camera data

as the feature maps, segmentation masks, and even pseudo-LiDAR point clouds [165].

The fusion input is a crucial design consideration because data representation plays a significant role in the overall detection performance. Among the three considerations, it has the most options. We will carefully review these options in Sec. 4, and discuss how they evolve with time and technology.

– **Fusion Granularity:** The third design consideration is concerned with at what granularity the two data streams are combined, *i.e.*, “how to fuse”. It usually has the three options: region of interest (RoI)-level, voxel-level, and point-level (with the last one at the finest granularity).

The fusion granularity plays an important role in determining the complexity and effectiveness of fusion; usually, finer fusion granularity requires more computing and leads to superior performance.

Previous surveys on deep learning based multi-modal fusion methods [3, 26, 40] cover a broad range of sensors, including radars, cameras, LiDARs, ultrasonic sensors, *etc.*, and provide a relatively brief review on a list of topics including multi-object detection, tracking, environment reconstruction, *etc.* While they are considered as a useful guide for readers to browse through the general area, our survey serves a distinctly different purpose: it targets at researchers who would like to carefully investigate the field of multi-modal 3D detection. As such, our survey intends to provide a deep and detailed review of recent research on this topic. Our contributions are summarized as below:

- We conduct an in-depth review of sensor fusion based 3D detection networks, with a special focus on LiDAR-camera fusion. We organize our discussions around the three core design considerations: fusion stage, fusion input, and fusion granularity, which answer the questions of where to fuse, what to fuse, and how to fuse, respectively.
- Most of the previous surveys have largely overlooked the fusion inputs of 3D multi-modal networks. In fact, compared to the other two design considerations, a fusion module’s input exhibits the most diversity and represents the unique idea of each design. In our survey, we discuss this design consideration thoroughly. According to their fusion input choices, we categorize the fusion based 3D detection networks into a total of five categories. We review the schemes in each category in detail, and discuss how the input combination evolve with time and technology.
- We also summarize the popular multi-modal datasets that can be employed for 3D object detection evaluation. In addition, we carefully discuss a list of open

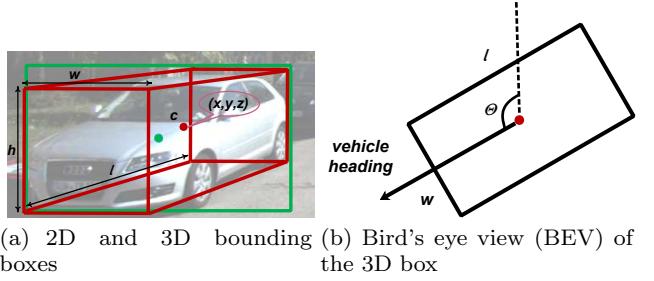


Fig. 4 (a) Example of 2D (green) and 3D (red) object detection bounding boxes, (b) Parameters of a 3D box in the BEV including its width w , length l , and vehicle heading angle θ

challenges in the field as well as possible solutions, which can hopefully inspire some future research in the area of multi-modal 3D object detection.

In this paper, we first provide a brief background of typical sensors used in autonomous driving, their data properties, and 3D object detectors through single modality respectively in Sec. 2. In Sec. 3, We present a summary of popular datasets that can be employed for evaluating multi-modal 3D object detection networks. In Sec. 4, we present a review of multi-modal fusion methods based on three crucial design choices: fusion stage, fusion input, and fusion granularity. Finally, we discuss open challenges and possible solutions in Sec. 5.

2 Background

In this section, we provide a background overview of typical sensors employed in autonomous driving, especially on 3D object detection methods that rely on each type of sensor. We mainly focus our discussions on cameras and LiDARs. Besides, we also introduce other sensors that can be employed for 3D object detection.

2.1 3D Object Detection Task

Before introducing 3D object detection methods through different camera settings, we first give an overview of 3D object detection. In the 3D object detection task, we need to provide the 3D bounding boxes of objects in the scene. As depicted in Fig. 4, we are required to predict the object center’s 3D coordinates c , length l , width w , height h as well as its deflection angle θ to obtain the red 3D bounding box.

2.2 3D Object Detection through Cameras

Cameras are the most common sensors for self-driving cars. A series of mature methods in 2D object detec-

Table 1 Using the sensors employed in the KITTI dataset as example, we compare the LiDAR and camera sensors [44]. In addition to important sensor parameters, we also qualitatively compare how external factors may affect the data quality of different sensors, with \triangle indicating minor influences, $\triangle\triangle$ moderate influences and $\triangle\triangle\triangle$ significant influences [195].

	sensor parameters and external factors in KITTI [44, 191]							external factors from [195]		
	number	cost	data format	data dimension	density/ per frame	FPS	object dis.>60m	strong light	smog vis.>2km	installation
LiDAR	1	\$8,0000	binary float matrix	3D	\approx 10,000 points	10	\approx 10 points	\triangle	$\triangle\triangle$	easy
Camera	2 grey 2 color	\$268 \$511.25	png format	2D	466,240 pixels	15	\approx 400 pixels	$\triangle\triangle\triangle$	$\triangle\triangle\triangle$	easier

tion have been developed in recent years, which can be reused in 3D object detection [46, 129]. Accordingly, image-based 3D object detection methods can achieve satisfactory performance at low expenses, often outperforming human experts [105, 146]. Several types of cameras have been widely deployed in AV, each with pros and cons. Below we talk about the 3D object detection algorithms via different camera settings.

Monocular 3D Object Detection. Monocular cameras provide dense information in the form of pixel intensity, which reveals shape and texture properties [2, 35]. The shape and texture information can also be utilized to detect lane geometry, traffic signs, and type of objects. The main disadvantage of using monocular cameras for 3D detection stems from the lack of depth information, which is necessary for accurate object size and position estimation for AVs [11]. To compensate for this, many studies have been devoted to enhancing detection accuracy through monocular cameras [18, 24, 56, 95, 99, 106, 113, 126, 165]. For example, Mousavian et al. [106] employ a designed CNN to estimate the missing depth information, which is used later to upgrade the 2D bounding box to the 3D space. Chu et al. [24] perform monocular depth estimation first and lifts the 2D pixels to pseudo 3D points. They design a novel neighbor-voting method that incorporates neighbor predictions to improve object detection from severely deformed pseudo-LiDAR point clouds. Park et al. [113] propose an end-to-end single-stage monocular-based detector. With the large-scale unlabeled data pre-training, it achieves promising detection results.

Stereo 3D Object Detection. Compared to monocular cameras, stereo cameras estimate a more accurate depth map [34, 80]. Specifically, multi-view cameras can cover different ranges of scenes through different cameras and capture depth maps more accurately [73, 74, 114]. Meanwhile, the complexity and cost involved in processing stereo images will also increase considerably. Some works exploit stereo images to generate dense point clouds to conduct 3D object detection tasks [20, 83, 120, 127, 189]. For example, Chen et al. [20] focus on generating 3D proposals by encoding

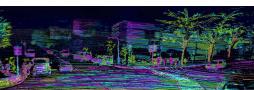
object size prior, ground-plane prior, and depth information into an energy function. Li et al. [83] add extra branches after the stereo Region Proposal Network (RPN) to predict sparse keypoints, viewpoints, and object dimensions, which are used to calculate coarse 3D object bounding boxes. Next, the accurate 3D bounding boxes are recovered by a region-based photometric alignment. Guo et al. [52] encode the depth information in the stereo cost volume, taking LiDAR features as the guidance to *distill* high-level geometry-aware representations for the stereo detection network.

Shortcomings of Camera-based 3D Object Detection. To summarize, camera-based 3D object detection has several shortcomings. Firstly, it is difficult for monocular cameras to estimate depth, which severely limits the detection accuracy [15, 64]. Secondly, camera-based detection further suffers from adverse conditions such as poor lighting, dense smoke, or heavy fog [194, 199]. So far, camera-only 3D object detection has not been able to obtain reliable performance. As far as KITTI [44] dataset is concerned, the state-of-the-art stereo-based method LIGA-Stereo [52] achieves 64.66% mAP while monocular-based DD3D [113] achieves only 16.87% mAP. To accomplish a more reliable AV system, We need to explore more powerful sensors for AVs.

2.3 3D Object Detection through LiDARs

LiDAR sensors use lasers as the light source to complete remote sensing measurements. LiDARs detect the light-wave signal between the LiDAR sensor and the detected object [158]. It continuously emits lasers and collects the information of the reflection points to obtain a full range of environmental information. When the LiDAR sensor rotates one circle, all the reflected point coordinates form a *point cloud*. As an active sensor, external illumination is not required and thus we can achieve more reliable detection under extreme lighting conditions. The typical resolution of LiDAR points ranges from 16 channels to 128 channels. As shown in Tab. 2, we conduct a detailed comparison to help readers form a clear understanding of the two popular LiDAR sensors: Velodyne HDL-64L and VLS-128. From specific

Table 2 Technical specifications for two typical Velodyne LiDARs. Note that point cloud images are from [130].

type	point cloud image	channel	range	points scanned per second	Horizontal Field of View	Vertical Field of View	price
Velodyne HDL-64E		64	120m	2.2 million	360°	26.9°	\$80,000
Velodyne VLS-128		128	220m	4.8 million	360°	40°	\$100,000

figures, it can be seen that all parameters of the 128-channel LiDAR are better than those of the 64-channel LiDAR. Obviously, LiDARs are quite costly compared to cameras. We can see the price of a Velodyne HDL-64 sensor is officially at \$80,000. The latest VLS-128 sensor has better performance but is also more expensive. Below, we briefly review existing works on 3D object detection based on the LiDAR data.

View-Based Detection. Many LiDAR-based methods project the LiDAR point clouds into the BEV, or RV, to leverage the off-the-shelf 2D Convolutional Neural Networks (CNNs). Early on, Yang et al. [176] propose an efficient, proposal-free single-stage detector. It transforms the point cloud to BEV and performs 2D CNNs to extract the point cloud features. Compact and dense RV-based methods are also proposed for 3D object detection. Recently, Liang et al. [87] employ a 2D backbone on the RV to learn the spatial features directly, and then adopt an R-CNN to get the 3D bounding boxes. H²3D R-CNN [29] first learns RV and BEV features in a sequential pattern, then fuses the two 3D representations in a multi-view fusion block.

Voxel-Based Detection. Voxel-based methods first divide points into regular 3D voxels, and then leverage the sparse convolutional neural networks [175] and transformers [103, 38] for feature extraction and bounding box prediction. VoxelNet [198] extracts discriminative voxel features to speed up the model execution. SECOND [175] reduces the computational overhead of dense 3D CNNs by applying sparse convolution. Point-Pillars [79] introduces a *pillar* representation (a particular form of the voxel) for the point cloud. Pillars are fast because all key operations can be formulated as 2D convolutions. Voxel R-CNN [28] further improves the accuracy and speed of voxel-based detectors by introducing a voxel ROI pooling operation. In addition, Mao et al. [103] introduce a Transformer-based architecture that enables long-range relationships between voxels by self-attention.

Point-Based Detection. Recent point cloud encoders such as PointNet [14], PointNet++ [122], Pointformer [110] and other point cloud backbones [68, 133] could learn

representations from raw point clouds. Point-based detectors employ them to extract the spatial geometry information for downstream tasks [141, 144, 183]. For example, Shi et al. [141] employ PointNet++ [122] as point clouds encoder and generate 3D proposals based on the extracted semantic and geometric features. Shi and Rajkumar [144] propose a graph neural network to detect objects from a LiDAR point cloud. To this end, they encode the point cloud efficiently in a fixed radius near-neighbors graph.

Point-Voxel hybrid Detection. In addition to the point and voxel representation introduced above, there are some works [53, 142, 182] that adopt a hybrid pattern, utilizing both point and voxel features for 3D object detection. For example, STD [182] first generates proposals based on the point features, then employs the voxel representation in the box refinement stage. PV-RCNN [142] integrates the multi-scale voxel features and point cloud features for accurate 3D object detection. M3DETR [49] encodes point and voxel features with multi-level scale via transformers. In general, point-voxel hybrid detectors benefit from both representations, which is superior to point or voxel-only detectors [142].

Compared with camera images, LiDAR points provide strong 3D geometric information, which is essential for 3D object detection. Furthermore, LiDAR sensors can better adapt to external factors such as strong light, which is depicted in Tab. 1. At present, LiDAR-based methods achieve better detection accuracy and higher recall than camera-based methods [19]. As far as the KITTI 3D object detection benchmark is concerned, the top monocular-images-based method DD3D [113] achieves 16.87% mAP while quite a few LiDAR-based methods [28, 29, 142, 174] achieve over 80% mAP. However, LiDAR-only algorithms are not yet ready to be widely deployed on AVs for the following reasons: 1) LiDARs are expensive and bulky, especially compared with cameras [117]. 2) The working distance of the LiDAR is rather limited, point clouds far away from the LiDAR are extremely sparse [191]. 3) LiDARs can not

work properly under extremely severe weather such as heavy rain [157].

2.4 3D Object Detection through Other Sensors.

In addition to cameras and LiDARs, AVs are often equipped with sensors such as millimeter wave (mmWave in short) radar sensors, infrared cameras, *etc.* In particular, mmWave radar has long been used on self-driving cars because it is more robust to severe weather conditions than cameras and LiDARs [195]. More importantly, radar points provide the velocity information of the corresponding object, which is crucial for avoiding dynamic objects [116]. Next, we give a brief background of mmWave radars below.

MmWave Radar Sensor. MmWave radars are active sensors that operate in the millimeter-wave bands. They can measure the reflected waves to determine the location and velocity of objects [1]. They are considerably cheaper than LiDARs, resistant to adverse weather conditions (fog, smoke, and dust), and insensitive to lighting variations [195]. However, compared with the camera and LiDAR, there are limited large-scale and public mmWave radar datasets [190]. Moreover, due to the low resolution of the mmWave radar, it is hard to directly detect the shape of an object through sparse 2D radar points. Compared with the 3D point cloud, radar points are much noisier because of multi-path reflection, rendering it hard to perform 3D detection alone [81].

The mmWave radar outputs can be organized at three levels: 1) *raw data* in the form of time-frequency spectrograms; 2) *clusters* from applying clustering algorithms [71] on raw data; and 3) *tracks* from performing object tracking on the clusters. Here, we process raw data which is collected on campus for visualization. As shown in Fig. 5, we perform two fast Fourier transforms on raw data to get the range-azimuth heatmap. The brightness in (b) represents the signal strength at that location and indicates high confidence of objects. Normally, datasets containing radar points generally utilize the representation of clusters, which are radar reflections with information containing position, velocity, and signal strength. The clusters are newly evaluated every cycle [8].

MmWave Radar based 3D Object Detection. The mmWave radar has been widely exploited in AV systems [104]. Radars usually report the detected objects as 2D points in BEV and provide the azimuth angle and radial distance to the object. For each detection, the radar also reports the instantaneous velocity of the object in the radial direction. To the best of our knowledge, Major et al. [101] propose the first radar-based



(a) the RGB image



(b) mmWave radar: the range-azimuth heatmap

Fig. 5 An RGB image (a) and a mmWave radar heatmap (b) on the same scene. The brightness in (b) indicates high confidence of objects. The data was collected by the authors at the north entrance of the West Campus of the University of Science and Technology of China.

deep neural network object detection with reliable results. However, radar-based 3D detectors face many challenges. Compared with LiDAR points, radar points are much noisier and less accurate, which brings difficulties in adapting LiDAR pipeline to the radar. Another bottleneck in radar-based detectors is the lack of publicly usable data annotated with ground-truth information [140]. In practice, mmWave radars are more often used for fusion with other sensors: *e.g.*, radar-camera, radar-LiDAR [108, 178].

2.5 Discussion

As discussed above, different sensors have different advantages, sometimes are complementary. For example, cameras are high-resolution and low-cost sensors, but lack depth information and are sensitive to light conditions. On the contrary, LiDAR points can provide 3D spatial information of the surrounding environment, but capture only sparse points at high price.

In general, camera-based methods generate less accurate 3D bounding boxes than LiDAR-based methods. Currently LiDAR-based methods lead in popularity in 3D object detection, while with some shortcomings. For example, the density of point clouds tends to decrease quickly as the distance increases while image-based de-

tectors could still detect faraway objects. To make good use of the complementary features and improve the overall performance, more methods try to design fusion networks to combine images with point clouds. These methods have achieved superior performance in 3D object detection tasks compared to methods relying on a single sensor. We will discuss this in Sec. 4 afterwards.

3 Datasets and Metrics

Datasets are an integral part of the field of deep learning. The availability of large-scale image datasets such as ImageNet, PASCAL, and COCO motivate outstanding evolution of image classification task [27, 36, 89]. Benefiting from the vigorous development of 2D images, 3D object detection eagerly requires plentiful labeled data to adapt to a changeable environment. Consequently, we discuss some widely used datasets for 3D object detection in autonomous driving.

3.1 KITTI

One of the earliest datasets for autonomous driving, KITTI [43], provides stereo color images, LiDAR point clouds, GPS coordinates, etc. The dataset supports multiple tasks: stereo matching, visual odometry, 3D tracking, 3D object detection, etc.¹ It collects data with a car equipped with a 64-channel LiDAR, 4 cameras, and a combined GPS/IMU system. There are over 20 scenes covering cities, residential and roads in the dataset. In particular, the object detection dataset contains 7,481 training and 7,518 testing frames with calibration information and annotated 2D/3D bounding boxes. KITTI annotates 8 different classes. Each class is categorized as “easy”, “moderate” and “hard” cases.

mAP (mean Average Precision) is a commonly used metrics in object detection. Some datasets containing multiple classes usually average the AP (Average Precision) score of each class, denoted as mAP. KITTI requires detection results of “car”, “pedestrian” and “cyclist” and calculates mAP of each class. It considers a predicted box as true positive (TP) if the IoU with the ground-truth box is greater than the threshold, otherwise as false positive (FP). The not detected ground-truth boxes are denoted as false negative (FN). We define *precision* and *recall* as:

$$\text{precision} = \frac{TP}{TP + FP}, \quad (1)$$

$$\text{recall} = \frac{TP}{TP + FN}. \quad (2)$$

¹ <http://www.cvlabs.net/datasets/kitti/index.php>

Based on the predicted and ground-truth boxes, we get a function of $p(r)$ with respect to recall r , the calculation of Average Precision (AP) is as below:

$$AP = \int_0^1 p(r)dr. \quad (3)$$

Remarkably, in order to facilitate the development of multi-modal detection methods in autonomous driving, the KITTI development team proposes a dataset KITTI360 [171] with richer sensor information and 360° annotations. Specifically, they annotate 3D scene elements with rough bounding primitives and then transfer this information into the image domain. As such, KITTI360 has dense semantic and instance annotations for both 3D point clouds and 2D images.

3.2 NuScenes

Developed by Motional, the nuScenes dataset is one of the largest datasets with ground-truth labels for autonomous driving [8]. It consists of 700 scenes for training, 150 scenes for validation, and 150 scenes for testing. The dataset is collected using six cameras and a 32-beam LiDAR to provide 3D annotations for 23 classes in a 360-degree field of view. NuScenes also provides 5 radar sensors for the measurement of the object velocity. The full dataset includes approximately 1.4M camera images, 390k LiDAR sweeps, 1.4M radar sweeps, and 1.4M object bounding boxes in 40k key frames. The driving scenes are collected in Boston and Singapore, which are known for their dense traffic and highly challenging driving situations. Additionally, nuScenes annotates object-level attributes such as visibility, activity, pose, etc.

As far as the object detection task² is concerned, the nuScenes requires the detection of 10 classes, including traffic cone, bicycle, pedestrian, car, bus, etc. When calculating AP for a class, instead of adopting the traditional bounding box overlap, nuScenes employs center-distance-based metrics. When matching the prediction and ground-truth, nuScenes computes their center distance and obtains AP based on a list of distance thresholds. mAP is calculated by averaging AP.

Unlike KITTI, nuScenes also considers TP’s average translation, scale, orientation, velocity, and attribute error with ground-truth, marked as ATE, ASE, AOE, AVE, and AAE, respectively. The final metric, nuScenes detection score (NDS), is derived from a weighted sum of mAP and errors, leading to a more comprehensive

² <https://www.nuscenes.org/object-detection?externalData=all&mapData=all&modalities=Any>

Table 3 Popular multi-modal dataset comparison, including year, number of LiDARs, number of LiDAR channels (we report the number of channels of the top LiDAR for Waymo dataset and the maximum number of channels among 4 LiDARs for AIODrive dataset), number of cameras, whether with radar, number of 2D boxes (we don’t distinguish between 2D boxes and 2D instance segmentation annotation), number of 3D boxes, number of annotated classes, and location (KA: Karlsruhe; SF: San Francisco; SG: Singapore; PT: Pittsburgh). Note that ApolloScape’s LiDARs scan with 1 beam at a high frequency to get dense point clouds.

dataset	year	n-LiDAR	n-chn	n-Cam	radar	n-2D	n-3D	n-cls	loc
KITTI [43]	2012	1	64	4	No	80K	80K	8	KA
ApolloScape [66]	2018	2	1	6	No	2.5M	70K	35	4x China
H3D [115]	2019	1	64	3	No	-	1.1M	8	SF
nuScenes [8]	2019	1	32	6	Yes	-	1.4M	23	Boston, SG
Argoverse [13]	2019	2	32	9	No	-	993K	15	PT, Miami
Waymo [150]	2019	5	64	5	No	9.9M	12M	4	3x USA
AIODrive [168]	2021	4	1280	10	Yes	26M	26M	23	synthetic

description of detection performance. We give the official formula below:

$$\text{NDS} = \frac{1}{10} \left[5\text{mAP} + \sum_{\text{mTP} \in \text{TP}} (1 - \min(1, \text{mTP})) \right]. \quad (4)$$

3.3 Waymo Open Dataset

The Waymo Open Dataset³ is a high-quality annotated multi-modality dataset for autonomous driving [150]. It consists of annotated data collected by Waymo’s self-driving vehicles. The dataset covers a wide variety of scenes from urban to suburban areas. There are a total of 798 scenes for training and 202 scenes for validation with 2D and 3D annotated labels, which are collected by five LiDAR sensors and five pinhole cameras. Each scene captures 20 seconds of continuous driving. The annotations provide four object categories including “car”, “pedestrian”, “cyclist” and “sign”.

Same as the KITTI dataset, Waymo Open Dataset adopts AP as the metric. Waymo Open Dataset also proposes a new metric APH which incorporates the heading accuracy of the predicted objects into the traditional AP metric. Waymo also supports the task of domain adaptation. Domain adaptation is a popular technology that learns knowledge from the source domain with sufficient annotations and transfers it to the target domain with limited or no annotations, which mitigates the lack of huge amount of labeled data [162].

3.4 Other Datasets

In addition to the three widely used datasets introduced above, there are a few recent datasets that are gaining

rapid popularity [66, 115, 13, 169, 168, 118, 41, 72]. We select some of them for detailed introduction as follows.

- ApolloScape [66] consists of data from 4 regions in China under various weather conditions. ApolloScape dataset is collected with 2 LiDAR sensors, 6 video cameras, and a combined IMU/GNSS system. It supports a variety of autonomous driving tasks such as scene parsing, lane segmentation, trajectory prediction, object detection, tracking. The dataset contains 140K+ annotated images with annotation of lane lines. For 3D object detection, ApolloScape provides 6K+ point cloud frames with annotated 3D bounding boxes. ApolloScape’s evaluation metrics are the same as KITTI. It requires the detection of vehicles, pedestrians, and bicyclists.
- H3D [115] is a large-scale full-surround 3D object detection and tracking dataset, with a special focus on crowded traffic scenes in the urban areas. The dataset is collected with 3 cameras with 260-degree field of view (FoV), and a 64-beam Velodyne LiDAR sensor. It contains over 27K frames in 160 scenes with over 1 million objects. For evaluation, H3D employs a similar protocol as KITTI with a 0.5 IoU threshold for car and a 0.25 IoU threshold for pedestrian.
- Argoverse [13, 169] supports advancements in 3D tracking, motion forecasting, and other perception tasks for self-driving vehicles. It provides rich semantic annotation for maps. For sensor setup, it is equipped with two 32-channel LiDARs, seven surround-view cameras, and two stereo cameras. It provides rich semantic information about road infrastructure and traffic rules. Argoverse dataset also provides HD maps for automatic map creation.
- Cityscapes 3D [41] extends the original Cityscapes dataset [25] with 3D bounding box annotations to

³ https://waymo.com/intl/en_us/open

support the task of 3D vehicle detection. It also provides benchmarks of pixel-level or instance-level semantic labeling and panoptic semantic labeling tasks. It annotates 3D bounding boxes and corresponding 2D instance segmentation masks for each vehicle. The 3D bounding boxes are annotated with stereo RGB images and with nine degrees of freedom. It also proposes a new metric for monocular 3D objection detection.

- AIODrive [168] is a large-scale synthetic dataset generated by the urban driving simulator, namely CARLA [33]. It synthesizes data from multiple sensors including 3D LiDARs, RGB cameras, depth cameras, radars, and IMU/GPS. All sensors collect data at a frequency of 10Hz. With the help of the simulator, it provides pretty detailed annotation with the object’s 2D/3D bounding boxes, trajectories, velocities, and accelerations. The dataset also synthesizes some adverse scenes such as terrible weather and car accidents.

3.5 Discussion

Datasets for autonomous driving are developing rapidly. From Fig. 6, we observe that the size of the three popular datasets ranges from only 15,000 frames to over 230,000 frames. However, compared to the image datasets in 2D computer vision, 3D datasets are still relatively small. For example, ImageNet [27] provides image frames of over 1.4 million. Besides, the object classes are limited and unbalanced. Fig. 6 compares the percentages of car, person, and cyclist classes. We also make a comprehensive comparison for all discussed datasets in Tab. 3.

Fig. 7 shows top-ranked methods on the three datasets. Interestingly, we observe that top-ranked methods on the nuScenes leaderboard are mainly fusion-based methods [184, 96, 21, 69]. For example, the top 8 methods on the nuScenes leaderboard are all fusion-based methods. In contrast, out of the top 10 methods on the KITTI / Waymo leaderboard, only 4 / 2 of them are multi-modal based [170, 179, 200, 100, 96, 84]. The main reason is that the LiDAR sensors employed in these datasets have different resolutions. KITTI and Waymo employ a spinning LiDAR sensor of 64 beams, while nuScenes uses a rotating 32-beam LiDAR. We may infer that multi-modal methods are much more necessary for high-performance perception when point clouds are relatively sparse.

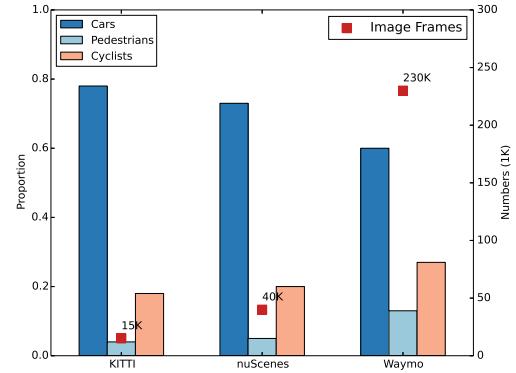


Fig. 6 Comparison of KITTI, nuScenes, and Waymo Open Dataset. From left Y-axis, we find the proportions of objects belonging to “car”, “person”, and “cyclist” classes are imbalanced clearly. From right Y-axis, we mark the total image frame number of the three datasets, ranging from 15K to 230K.

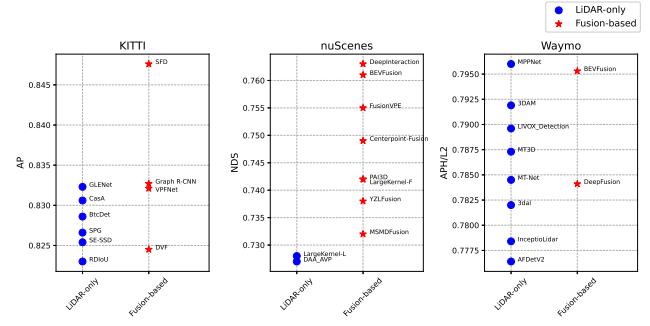


Fig. 7 Top-10 methods of the popular datasets. Among the top-10 methods, 8 methods are fusion-based on the nuScenes leaderboard, 4 methods are fusion-based on KITTI, and 2 methods are fusion-based on Waymo, respectively. Note that we only report methods with paper link on the KITTI leaderboard, while for the nuScenes and Waymo leaderboards, we report all the listed methods except repeated entries.

4 Deep Learning Based Multi-Modal 3D Detection Networks

In this section, we present our review of deep learning based multi-modal 3D detection networks, with a special focus on LiDAR and camera data. We organize our review by the following three important design considerations for fusion:

- fusion stage, *i.e.*, at what pipeline stage the fusion occurs, answering the question “where to fuse”;
- fusion input, *i.e.*, what representations are used for fusion data, answering the question “what to fuse”;
- fusion granularity, *i.e.*, at what granularity the fusion data are combined, answering the question “how to fuse”.

Fig. 8 lists the possible options for each design consideration. In the rest of this section, we discuss how the

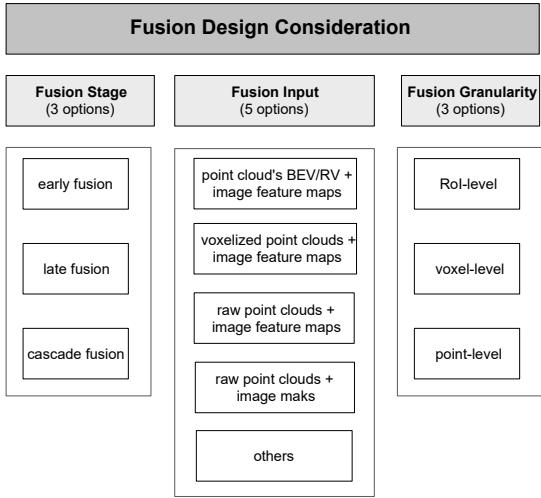


Fig. 8 Overview of the three crucial design considerations for a fusion network. We briefly summarize all the options for these considerations here.

recent deep multi-modal 3D detection networks address these three design questions.

More importantly, compared to the other two design considerations, a fusion module’s input exhibits the most diversity and represents the characteristic network design. Hence, we categorize the fusion based 3D detection networks into a total of five categories. In each category, we review the relative fusion schemes in detail.

4.1 Fusion Stage: where to fuse?

This design issue is concerned with which pipeline stage performs the fusion operation. Here, we broadly partition a detection network pipeline into the following three stages: the input stage, the feature extraction stage, and the prediction stage (illustrated in Fig. 9). Depending upon which of the three stages perform fusion, we have three options: early fusion, late fusion, and cascade fusion.

4.1.1 Early Fusion

Early fusion usually occurs in the input stage or feature extraction stage before each branch reaches its prediction [19, 156, 173]. It enables more direct interactions among multi-modal features of the intermediate layers, as shown in Fig. 9. The fused feature is then utilized to perform classification and regression tasks in the prediction stage. Early fusion can better leverage rich intermediate information from modalities, and is currently the most widely used fusion stage.

4.1.2 Late Fusion

In contrast with early fusion, late fusion employs separate branches for each modality, and then combine individual decision-level outputs through a fusion network in the prediction stage [39]. Fig. 10 outlines such a framework. Late fusion can better leverage existing networks for each modality. It also does not require to deal with issues such as how to align the data from different modalities.

Notably, Pang et al. [112] employ late fusion and outperforms single modality detectors on the KITTI leaderboard. It exploits the geometric and semantic consistencies between 2D and 3D predictions and learns the probabilistic dependencies between the two from the training data. Specifically, it obtains 2D and 3D proposals and then encodes all proposals into a sparse tensor. As summarized in Tab. 4, its shortcoming lies in the inability to exploit rich intermediate features [4, 137].

4.1.3 Cascade Fusion

Cascade fusion employs the hybrid mode by fusing one branch’s prediction with the other’s input, which builds a cascade relationship between multiple modalities. As illustrated in Fig. 11, we first obtain 2D proposals from the prediction stage of the image stream. Next, with a known camera projection matrix, a 2D proposal can be lifted to a frustum which defines a 3D search space [123]. We collect all points within the frustum to form a 3D frustum proposal that is used to classify and locate the object. Consequently, one modality provides prior information that can greatly reduce the other’s search space in cascade fusion.

The first fusion approach using the cascading structure was F-PointNet [123]. Nevertheless, its performance is greatly limited by the accuracy of the 2D detector. Subsequently, several following methods have been proposed [145, 167, 181] to further improve the accuracy.

4.1.4 Discussion

Tab. 4 summarizes the advantages and disadvantages of the three fusion stages, and Tab. 5 gives typical multi-modal methods for each fusion stage. From Tab. 5, we observe that most fusion-based algorithms employ early fusion, which is also the focus of our survey.

4.2 Fusion Input: what to fuse?

The second design choice is concerned with what form or representation the multi-modal data are input into the fusion module. A fusion module’s input exhibits

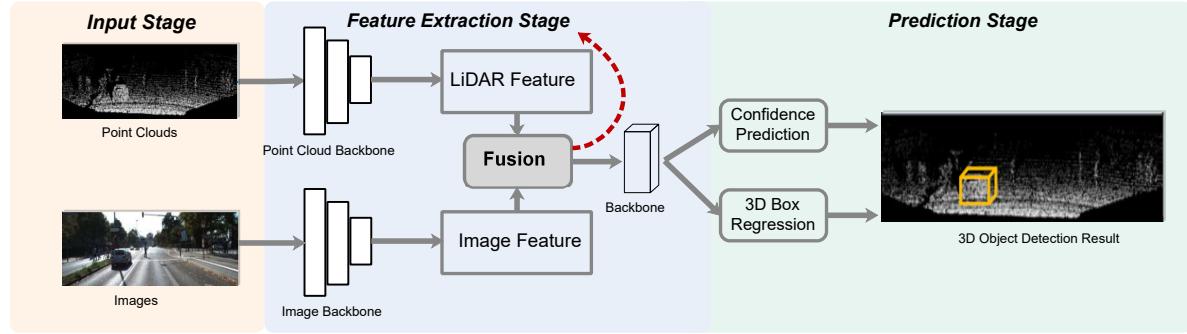


Fig. 9 An early fusion pipeline. We first extract the image and point cloud features respectively, and then conduct fusion on these features in the feature extraction stage.

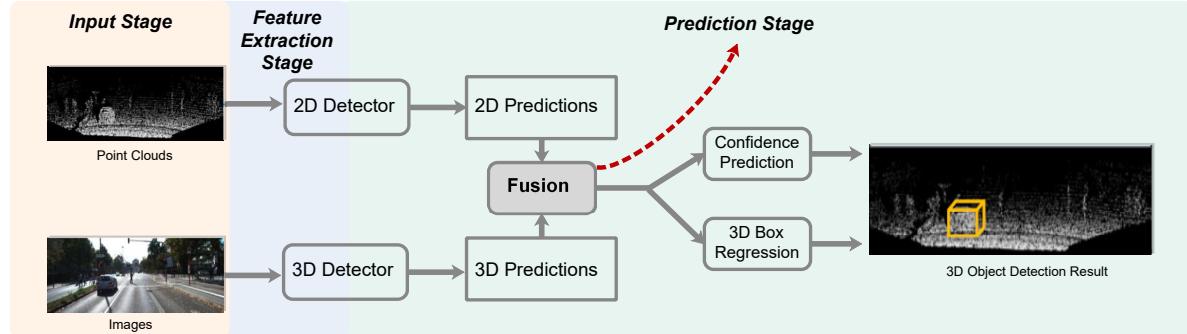


Fig. 10 A late fusion pipeline. We first get the predictions from each modality, and then take these predictions as fusion inputs. As such, the fusion network is in the prediction stage.

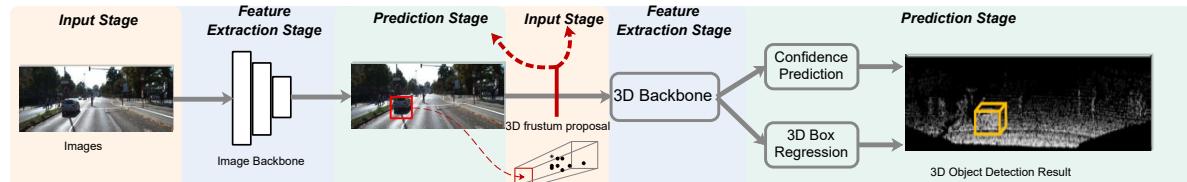


Fig. 11 A cascade fusion network. We first obtain the predictions from the image branch, and then fuse the image predictions with point cloud data for further 3D object detection. Therefore, cascade fusion usually occupies the hybrid mode by fusing one branch's prediction with the others input.

the most diversity and represents the unique idea of each design. For LiDAR-camera fusion, we are allowed to take raw sensor data, various intermediate features, or even result-level output from the image/point cloud branch as fusion input. Specifically, the fusion module can take in the LiDAR data in the form of voxel grids, raw point clouds, or point cloud's projection on BEV or RV, camera data in the form of image feature maps, segmentation masks, or even the corresponding pseudo-LiDAR point clouds.

In this section, we first present typical inputs that can be utilized by the image and point cloud branches respectively, and then we categorize the fusion based 3D detection networks into a total of five categories according to input combinations. Here, we focus our review on early fusion methods.

4.2.1 Typical Fusion Inputs for LiDAR-Camera Fusion

We first introduce the typical fusion inputs employed for the image branch and the point cloud branch, respectively, in fusion-based detection pipelines, as illustrated in Fig. 12. To be more precise, a modality's input to the fusion module is the output of a certain middle layer, be it a simple data preprocessing function or a neural network block.

Typical Fusion Input for the Image Branch. Most of the LiDAR-camera fusion methods take one of the following three fusion inputs from the image branch, namely feature maps, segmentation masks and pseudo-LiDAR point clouds (depicted in Fig. 13).

feature maps: Deep neural networks are capable of extracting appearance and geometry feature maps

Table 4 Advantages and disadvantages for different fusion stages

Categories	Advantages	Disadvantages
Early Fusion	+ Can leverage rich intermediate features from multiple modalities.	- Sensitive to inherent data misalignment between modalities.
	+ Large feature vectors can lead to better detection results with suitable learning methods.	- Large feature vectors lead to longer training/inference time.
Cascade Fusion	+ Can reduce the search space with prior information.	- Rely heavily on initial proposal generation.
Late Fusion	+ Can utilize off-the-shelf networks for each modality.	- Unable to take advantage of useful intermediate features.

Table 5 Summary of typical multi-modal 3D detection methods: stage (fusion stage), PC-Input (point cloud Input), RGB-Input (image input), gran (fusion granularity), HW (hardware), lat (latency), DS (dataset used for evaluation), and mAP (mean average precision)

	stage	PC-Input	RGB-Input	gran	HW	lat	DS	mAP
MV3D [19]	early	view	feature map	RoI	Titan X	0.36s	KITTI	63.63%
AVOD [78]				RoI	Titan XP	0.08s	KITTI	71.76%
Confuse [85]				voxel	GTX1080	0.06s	KITTI	68.78%
SCANet [98]				RoI	GTX1080	0.09s	KITTI	66.30%
FuseSeg [151]				point	-	-	KITTI	-
MVX-Net [148]		voxel	feature map	voxel	GTX1080	-	KITTI	72.70%
3D-CVF [187]				voxel & RoI	-	0.06s	KITTI	80.45%
VPF-Net [200]				point	2080Ti	0.06s	KITTI	83.21%
PointAugmenting [159]				point	-	-	nuScenes	66.80%
PointFusion [173]	early	point	feature map	RoI	GTX1080	1.3s	KITTI	63.00%
EPNet [65]				point	Titan XP	0.1s	KITTI	81.23%
PI-RCNN [172]				point & RoI	-	0.06s	KITTI	78.53%
PointPainting [156]		point	mask	point	GTX1080	0.4s	KITTI	75.80%
CenterPointV2 [185]				point	-	-	nuScenes	67.10%
HorizonLiDAR3D [32]				point	-	-	Waymo	78.49%
MMF [86]	view	view	feature map & pseudo LiDAR	point	GTX1080	0.08s	KITTI	77.43%
MVAF [160]				voxel	Titan X	0.06s	KITTI	78.71%
F-PointNet [123]	cascade	-	-	RoI	GTX1080	0.17s	KITTI	69.79%
IPOD [181]					-	0.1s	KITTI	72.57%
F-ConvNet [167]					-	0.1s	KITTI	75.50%
RoarNet [145]				Titan X	-	-	KITTI	73.04%
SIFRNet [196]				-	-	-	KITTI	-
CLOCs [112]	late	-	-	-	-	-	KITTI	82.25%

from raw images [75, 77, 188], which are the most commonly used input for fusion between cameras and other sensors [19, 78, 85]. Compared with raw images (Fig. 13 (a)), the utilization of feature maps explores richer appearance cues and larger receptive fields, which enables more in-depth and thorough interactions between modalities. For example, as illustrated in Fig. 13 (b), we observe that the edges and textures of feature maps are more distinct than other areas. We refer the readers to Sec. 4.2.2 for more in-depth review of fusion algorithms that use image feature maps as input. Here, we list some pop-

ular backbones that can be used to obtain feature maps, which can be fed to a fusion module: *e.g.*, VGG-16 [147], ResNet [54], DenseNet [61].

masks: Images are passed through a semantic segmentation network to obtain pixel-wise segmentation masks [45, 97]. Image masks are often utilized for fusion with other sensor data, as a stand-alone product from the image processing branch. Compared with feature maps, using masks as camera data fusion input has the following advantages. Firstly, image masks can serve as more compact summary features of the image. Secondly, pixel-level image masks

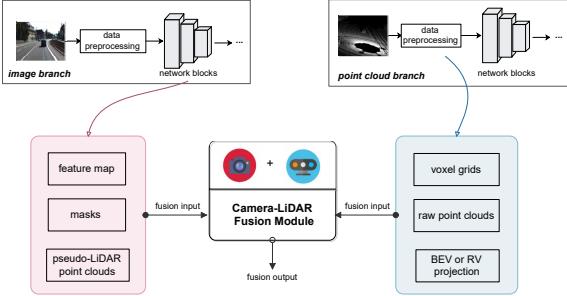


Fig. 12 Illustration of typical fusion inputs from the image branch and the point cloud branch, respectively. For the image branch, its fusion input is usually the output of a neural network block; for the point cloud branch, its fusion input is usually from simple data preprocessing such as voxelization.

can easily be used to “paint” or “decorate” the LiDAR points by conducting a point-to-pixel mapping using a known calibration matrix [156]. We refer the readers to Sec. 4.2.2 for more in-depth review of fusion algorithms that use image masks as input. Here, we list some popular image segmentation networks for the fusion-based algorithms: *e.g.* DeepLabV3 [17], Mask-RCNN [55], and lightweight network Unet [135].

pseudo-LiDAR point clouds: The camera data can also be converted to pseudo point clouds as fusion input [86]. As pointed out in [165], the pseudo point cloud representation raises image pixels to the 3D space, whose signal is much denser than actual LiDAR point cloud. On the downside, it often has a *long tail* problem since the estimated depth may not be accurate around the boundaries of the object [189], as depicted in Fig. 13 (d) with yellow circles. According to [165], the pseudo-LiDAR points are obtained by back-projecting image pixels into pseudo 3D points according to the estimated depth map. In the context of 3D multi-modal detection, this representation contributes to multi-task learning [86]. Using the pseudo point clouds as fusion input, we can readily facilitate dense feature map fusion between images and point clouds.

Typical Fusion Input for the Point Cloud Branch.

A LiDAR point cloud is often synthesized from depth measurements collected from different viewpoints. It is basically a set of points in a 3D coordinate system, commonly defined by x, y, z, and the reflection intensity. Below, we discuss the typical LiDAR inputs that are commonly used for LiDAR-camera fusion, *i.e.*, voxel grids, raw point clouds, and point cloud’s projection on BEV or RV, whose visualization results are shown in Fig. 16 respectively.

voxelized point clouds or voxel grids: A voxelized point cloud is widely utilized as the fusion input due to the efficient parallel processing potentials on a regular voxel grid [7, 60, 79, 143, 176, 177, 198]. We first discretize the 3D space into 3D voxel grids, and then obtain the voxel features through the voxel feature encoding (VFE) layer as shown in Fig. 14. We can thus utilize 3D CNNs to extract deeper point cloud features. We refer the readers to Sec. 4.2.2 for more in-depth review of fusion algorithms that employ voxel grids as input for point cloud branch. However, fusion with voxelized points also has several disadvantages. Firstly, it suffers from information loss. The voxel size is highly correlated with how much spatial information is lost. To illustrate this point, let us look at Fig. 16 (b) where voxels (in blue color) are much more sparse than the original points in (a). Secondly, during voxelization, a large number of empty voxels will be produced as the LiDAR points are only on the surface of the objects [109], which may adversely affect the fusion performance. Thirdly, processing 3D voxels require time-consuming 3D convolution operations. Accordingly, the training time of the fusion network will inevitably increase. In practice, the point cloud data is usually voxelized into an evenly spaced grid only in the x-y plane (which we call *pillars*) to meet the computation and effectiveness demands [79].

raw point clouds: Thanks to efficient 3D point cloud processing networks, the raw 3D point cloud can be directly processed to obtain suitable point features without voxelization loss [173]. Specifically, in Fig. 15, we employ the point cloud encoder [14, 122] to process raw points and obtain point feature vectors. Sec. 4.2.2 provides more detailed review of the fusion algorithms that use raw points as input from the point cloud branch. Directly taking raw points as input can retain more information, compared with voxel-based methods [124, 142, 183]. However, point-based methods are generally computationally expensive, especially when dealing with large scenes. For example, for a widely used Velodyne LiDAR HDL-64E, it collects more than 100,000 points per frame (in Tab. 1). Therefore, considering the efficiency and performance, down-sampling point cloud data appropriately is necessary for data preprocessing.

BEV or RV projection: Another typical point cloud input for the fusion module is point cloud’s BEV or RV projection. The resulting pseudo image could be thereby processed efficiently by 2D CNNs. BEV is commonly adopted to fuse with image features be-

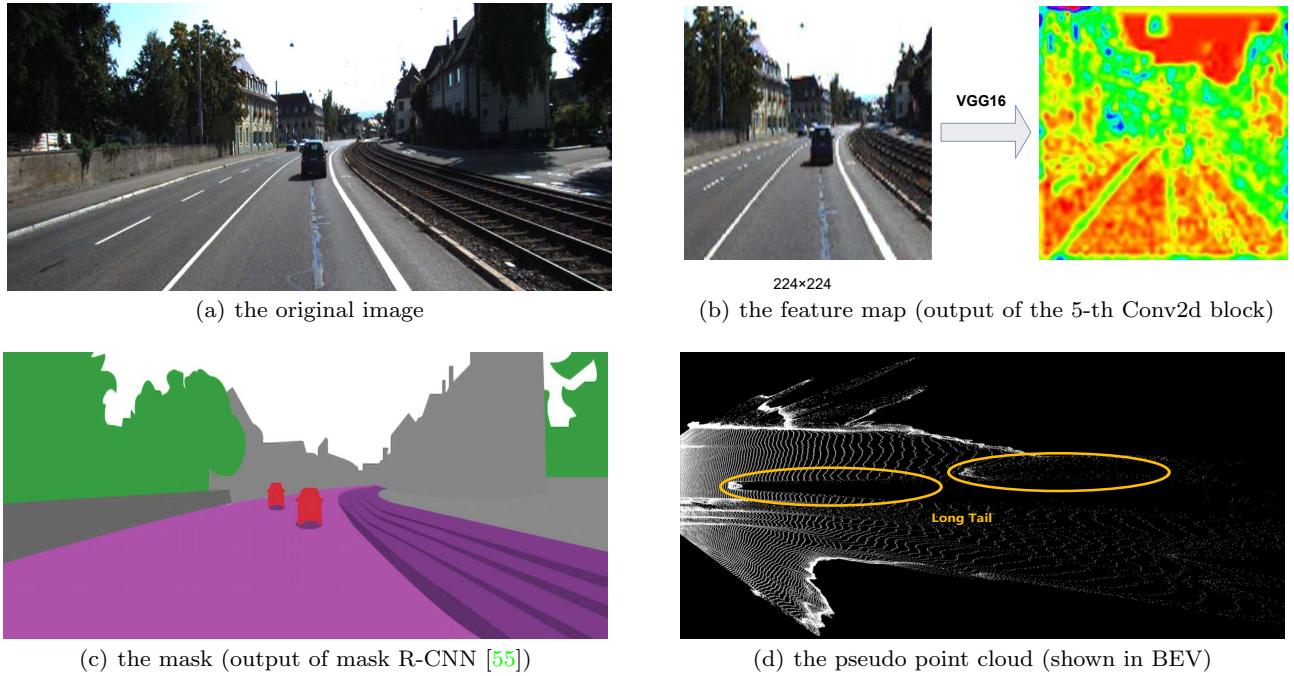


Fig. 13 Different inputs for the image branch in fusion-based 3D object detection. An RGB image (a), one of its feature maps (b), its segmentation mask (c), and its pseudo-LiDAR point cloud’s projection on BEV (d). The raw image is taken from the KITTI training set. We use a pretrained VGG16 [147] to obtain the feature map on the resized image (224×224).

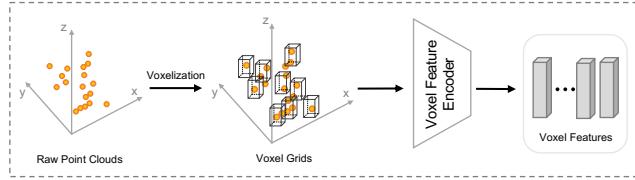


Fig. 14 A typical voxelized point cloud processing network [198]

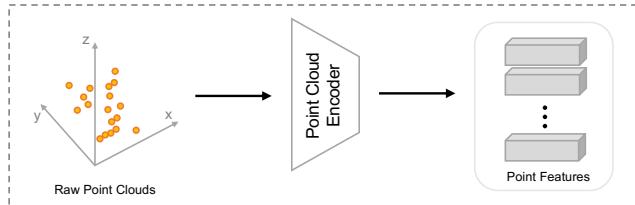


Fig. 15 A typical raw point cloud processing network [14]

cause there is much less overlapping between objects on the BEV plane. Another popular view-based input used for fusion is RV, which is also a native representation of the rotating LiDAR sensor [161]. Essentially, it forms a compact 2.5D scene [59] instead of a sparse 3D point cloud. Projecting the point cloud on RV preserves the full resolution of the LiDAR sensor data and avoiding the spatial

loss. However, RV suffers from the problem of the scale variation between nearby and far away objects [37]. In the fusion pipeline, these mentioned point cloud’s projections are usually first processed with 2D CNNs to get view-based features, and then pooled to the same size as the image features.

4.2.2 Typical Input Combinations for LiDAR-Camera Fusion

In our survey, we find the following combinations of fusion inputs are the most popular for the LiDAR-camera fusion module: (1) point clouds’ BEV/RV + image feature maps, (2) voxelized point clouds + image feature maps, (3) raw point clouds + image feature maps and (4) raw point clouds + image masks. In addition, exploiting more than one type of inputs for images/point clouds to form a more comprehensive fusion has become a recent trend. We also review these methods here. Below we discuss the fusion-based methods that fall into one of these 5 fusion input categories, with a special focus on how these fusion input combinations evolve with time and technology.

point cloud’s BEV/RV + image feature maps. Before 3D object detection became popular, 2D object detection based on images had drawn a great deal of attention [129]. Therefore, as soon as LiDAR was considered for 3D object detection, several LiDAR-camera

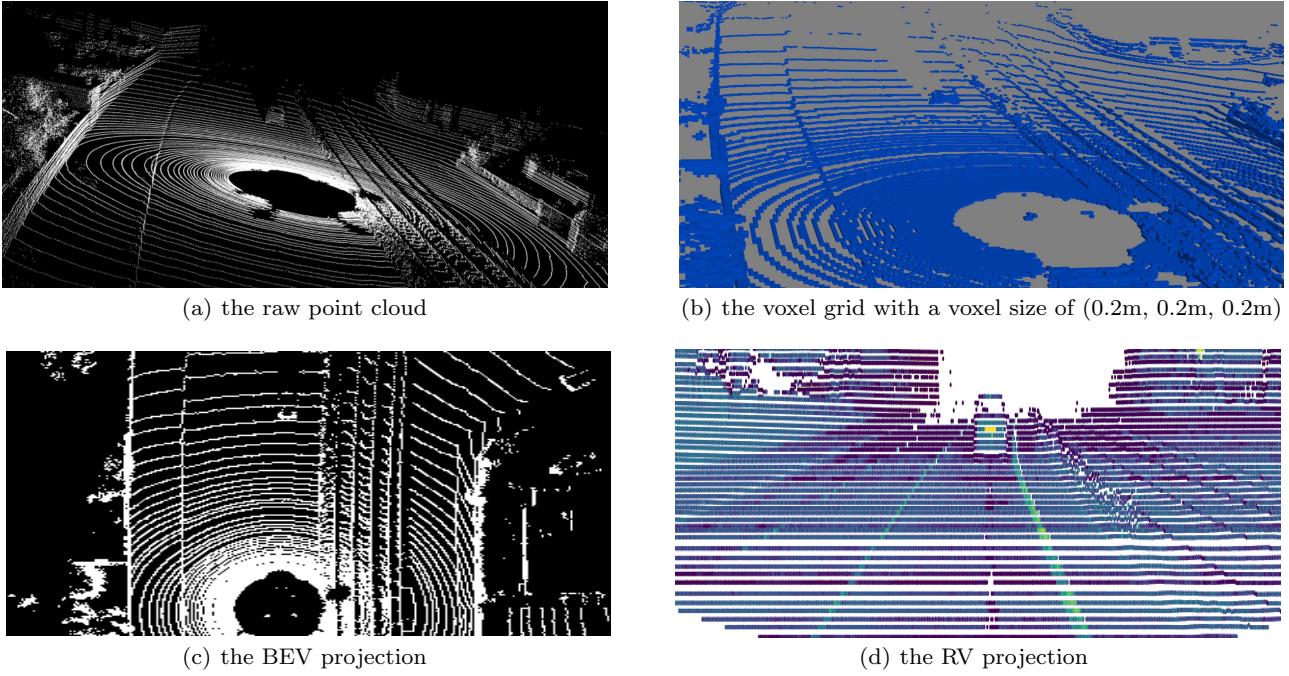


Fig. 16 Different inputs for the point cloud branch in fusion-based 3D object detection. A raw point cloud (a), its voxel grids with the voxel size of [0.2m, 0.2m, 0.2m] (b), its projection on BEV (c), and its projection on RV (d). The raw point cloud data is taken from KITTI training set.

fusion algorithms were proposed to project 3D point clouds on a 2D plane, and combine the resulting 2D view of the point cloud with image feature maps. We discuss typical algorithms in this category below.

MV3D [19] is a pioneering work in this category. As shown in Fig. 17, it takes the FV (front view) and BEV of a point cloud as input and exploits a 3D Region Proposal Network (RPN) to generate 3D proposals. Next, MV3D integrates multi-view features vectors from multi-proposals into the same length and puts them through a region-based fusion network. AVOD [78] achieves better performance than MV3D, especially in the small object category by designing a more advanced RPN that employs high-resolution feature maps. It also merges features from multiple views in the RPN phase to generate more accurate positive proposals. AVOD only takes point cloud’s BEV and image as input, which effectively decreases the computation cost. Based on AVOD, SCANet [98] utilizes an encoder-decoder based proposal network with a Spatial-Channel Attention (SCA) module to capture multi-scale contextual information and an Extension Spatial Upsample (ESU) module to recover the spatial information.

Nevertheless, these methods have limitations, especially when detecting small objects such as pedestrians and cyclists. To overcome these limitations, Contfuse [85] performs continuous convolutions [163] to extract multi-scale convolutional feature maps from point

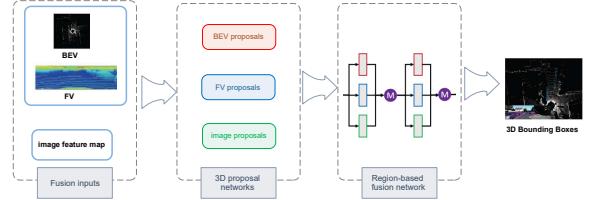


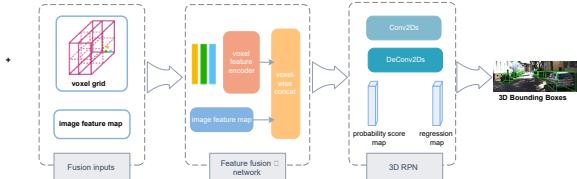
Fig. 17 The MV3D pipeline that fuses point cloud’s projections and image feature map [19]

cloud’s BEV and fuse them with image features. The engagement of continuous convolution captures local information from neighboring observations and leads to less geometric information loss. In addition, another downside of using point cloud’s BEV or FV as fusion input lies in the inevitable 3D spatial information loss when projecting the 3D point cloud to the 2D plane.

The point cloud’s range view (RV) can avoid the mentioned information loss problem. Compared to the BEV and FV projections, a RV is a compact, and more importantly, an intrinsic representation from LiDAR. As such, a very recent trend is to combine the point cloud’s RV with RGB image feature maps directly without incurring the projection loss. With the RV representation as input, FuseSeg [151] establishes the point-pixel mapping and maximizes the multi-modal information. In Tab. 6, we summarize the contributions and limitations of fusion methods in this category.

Table 6 Summary of methods that fuse point cloud’s projections and image feature maps

Methods	Year	Venue	image backbone	point cloud’s projections	Contributions
MV3D [19]	2017	CVPR	VGG-16	BEV, FV	<ul style="list-style-type: none"> Pioneer in exploiting BEV and FV LiDAR projections and monocular camera frames to detect vehicles. Design a deep fusion architecture which allows interaction between LiDAR and camera data.
AVOD [78]	2018	IROS	VGG-16	BEV	<ul style="list-style-type: none"> Improve the detection of small targets via a feature extractor that produces high-resolution feature maps.
Confuse [85]	2018	ECCV	ResNet	BEV	<ul style="list-style-type: none"> Exploit continuous convolutions to fuse at different levels of resolution.
SCANet [98]	2019	ICASSP	VGG-16	BEV	<ul style="list-style-type: none"> Propose a spatial-channel attention module that is capable of encoding multi-scale and global context information.
FuseSeg [151]	2020	WACV	MobileNetV2	RV	<ul style="list-style-type: none"> Pioneer to utilize dense range views of point clouds. Establish point-wise correspondences between the range view and image features.

**Fig. 18** The MVX-Net pipeline that fuses voxelized point clouds and image feature maps [148]

voxelized point clouds + image feature maps. Voxelization turns irregular point clouds into regular 3D voxels. With voxelization becomes popular with point cloud processing [28, 142, 175, 198], voxel grids have been commonly used as fusion input. The methods that fall in this category are summarized in Tab. 7.

As shown in Fig. 18, Sindagi et al. [148] use voxels and image feature maps as input, projecting non-empty voxel features to the image plane through calibration. The image features are then concatenated to the voxel features through a designed fusion network. At the last stage, the 3D RPN processes the aggregated data and produces the 3D detection results. Similarly, in another work 3D-CVF [187], the spatial attention maps [155] are applied to weigh each modality depending on their contributions to the detection task. 3D-CVF [187] employs auto-calibrated projection to construct smooth joint LiDAR-camera features.

Nonetheless, methods in this category face the *feature blurring* problem when only the center point of every voxel grid is projected onto the image feature. This results in the loss of detailed spatial information within each voxel. Recently, to overcome this bottleneck, VPF-Net[200] cleverly aligns and aggregates the point cloud and image features at the “virtual” points. Particularly, with the density lying between 3D vox-

els and 2D pixels, the virtual points can nicely bridge the resolution gap between the two sensors and preserve more information for processing. Later, PointAugmenting [159] solves the blurring problem by first “decorating” raw points with corresponding features extracted by pre-trained 2D detection models. Then decorated points are voxelized and further processed. PointAugmenting also benefits from an occlusion-aware point filtering algorithm, which consistently pastes virtual objects into images and point clouds during training.

raw point clouds + image feature maps. As mentioned before, voxelized point clouds could cause certain degree of information loss due to voxelization process. Later, the advent of PointNet [14] makes it possible to process the raw point cloud directly without any projection or voxelization. Consequently, PointNet inspires a series of studies to combine points directly with feature maps as fusion input (summarized in Tab. 8).

Different from previous fusion-based methods, PointFusion [173] combines global image features from ResNet-50 and point cloud features from PointNet in a concatenation fashion, as demonstrated in Fig. 19. Such concatenation operations, though simple, cannot align the multi-modal features nicely. Therefore, Huang et al. [65] propose LI fusion layer that explicitly establishes the mapping between point features and camera image features, thus providing finer and more discriminative representations. It also exploits point cloud features to estimate the importance of corresponding image features, which reduces the influence of occlusion and depth uncertainty.

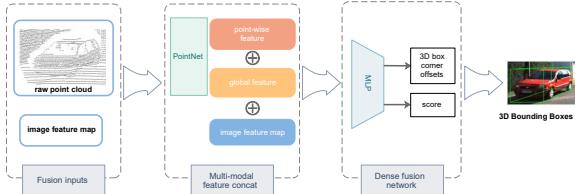
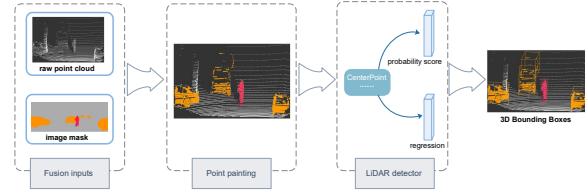
raw point clouds + image masks. Xie et al. [172] conduct continuous convolution directly on 3D points, and meanwhile, it retrieves deeper semantic features instead of image features as image input. The main ratio-

Table 7 Summary of methods that fuse voxelized point clouds and image feature maps

Methods	Year	Venue	image backbone	Contributions
MVX-Net [148]	2019	ICRA	VGG-16	<ul style="list-style-type: none"> Propose two fusion schemes to fuse multi-modal information. <i>PointFusion</i>: to aggregate 3D point is aggregated by an image feature to capture a dense context; <i>VoxelFusion</i>: a relatively later fusion strategy where image features are appended at the voxel level.
3D-CVF [187]	2020	ECCV	ResNet	<ul style="list-style-type: none"> Combine the camera and LiDAR features using the cross-view spatial feature fusion strategy. Employ an attention map to weigh the information from each modality depending on their contributions.
VPF-Net [200]	2021	TMM	2D CNNs	<ul style="list-style-type: none"> Effectively alleviate the resolution mismatch problem in fusing LiDAR and camera data. Explore further optimization through cut-n-paste based data augmentation.
PointAugmenting [159]	2021	CVPR	ResNet	<ul style="list-style-type: none"> Decorate point clouds with the corresponding CNN features. Design a novel cross-modal data augmentation algorithm considering the modality consistency.

Table 8 Summary of methods that fuse raw point clouds and image feature maps

Methods	Year	Venue	Image backbone	Contributions
PointFusion [173]	2018	CVPR	ResNet	<ul style="list-style-type: none"> Can directly unitize ResNet and PointNet. Integrate the global and local features to predict the bounding box.
PI-RCNN [172]	2020	AAAI	U-Net	<ul style="list-style-type: none"> Directly apply continuous convolution on raw points to preclude the quantization loss. Employ representation of deeper semantic features.
EPNet [65]	2020	ECCV	Four light-weighted convolutional blocks	<ul style="list-style-type: none"> Enhance the point features with semantic image features at a point-wise level. Exploit a consistency loss to encourage both localization and classification.

**Fig. 19** The PointFusion pipeline that fuses raw point clouds and image feature maps [173]**Fig. 20** The PointPainting pipeline that fuses raw point clouds and image masks [156]

nales lie in the two aspects. 1) Features learned under the supervision of semantic segmentation are generally more expressive and compact when representing image. 2) It is feasible to obtain the homogeneous transformation matrix, which can build relationship between 2D masks and 3D points [156]. As such, quite a few recent

studies use result-level features such as segmentation masks to fuse with raw points, which is shown in Tab. 9.

In order to fuse the point cloud data and image masks, the LiDAR points are projected by a homogeneous transformation into the image plane, which establishes the 3D-2D mapping between the two. This transformation on the KITTI dataset [44] is $T_{\text{camera} \leftarrow \text{LiDAR}}$, while it requires extra care for the nuScenes [8] trans-

formation since the LiDAR and camera sensors operate at different frequencies. Let $T_{\text{car} \leftarrow \text{LiDAR}}$ be the transformation from the LiDAR sensor to the reference frame of cars, and let $T_{\text{camera} \leftarrow \text{car}}$ be the transformation from the reference frame of cars to the camera sensor. The complete matrix calculation is as below:

$$T_{\text{camera} \leftarrow \text{LiDAR}} = T_{\text{camera} \leftarrow \text{car}} \times T_{\text{car} \leftarrow \text{LiDAR}}. \quad (5)$$

Consequently, we can append the 2D mask as an additional channel of the corresponding 3D point.

In Fig. 20, we present PointPainting [156] as a typical example fusion network. It takes raw points and segmentation results as fusion input. Next, in the fusion module (PointPainting dotted box of Fig. 20), we first project the points onto the image, and then append segmentation scores to the raw LiDAR point. More importantly, PointPainting could be freely applied to both point-based and voxel-based LiDAR detectors and further improves the overall performance.

Inspired by the successful PointPainting, Center-PointV2 [185] gets almost the state-of-the-art result on nuScenes, and HorizonLiDAR3D [32] ranks the top on Waymo Open Dataset Challenge.

Multi-Inputs for a Modality. Meanwhile, as deep learning networks that are designed to process point clouds and images become increasingly diverse, it is also common for a single modality to adopt multiple inputs for fusion.

MMF [86] is a pioneer in this category. It presents an end-to-end architecture that performs multiple tasks including 2D and 3D object detection, depth completion, etc. Specifically, the fusion module takes the image feature map as well as the pseudo-LiDAR point clouds from the image branch, and BEV from the point cloud branch. These inputs are then fused jointly for 3D object detection. Recently, Wang et al. [160] propose a multi-representation fusion framework that takes voxel grids, point cloud's RV projection, and image feature maps as input. They further estimate the importance of these three sources with attention modules to achieve adaptive fusion.

4.2.3 Discussion

To summarize, point cloud inputs evolve from point cloud's projections, voxel grids, to raw points, which strives to minimize the information loss incurred in point cloud projection or voxelization. Meanwhile, RGB image inputs evolve from lower-level feature maps to higher-level semantic segmentation results, which strives to exploit the richness of image data.

Furthermore, to leverage the different perspectives from the same data stream, recent fusion networks start

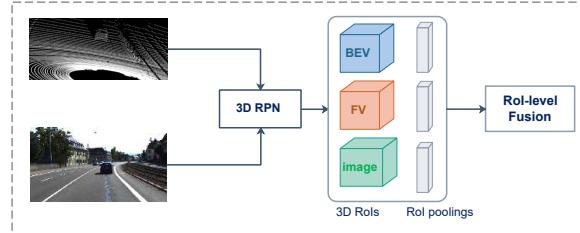


Fig. 21 Illustration of ROI-level fusion. To perform fusion at ROI-level, we first obtain 3D RoIs from a shared set of 3D proposals and we employ ROI pooling [46] to get feature vectors of the same length.

to take advantage of multiple inputs from the same modality. These trends are enabled by the rapidly growing computing capabilities as well as the fast development of powerful deep learning networks. These factors combined, more accurate LiDAR-camera fusion results are delivered.

4.3 Fusion Granularity: how to fuse?

In this section, we discuss the third design choice for fusion-based algorithms. It is defined by at what granularity the two data streams are combined, which also addresses the question of “how to fuse”.

There are usually three options: ROI-level, voxel-level, and point-level (with the last one at the finest granularity). In general, fusion granularity is crucial to the complexity and effectiveness of the fusion framework. Finer fusion granularity requires more computation but often leads to superior performance. Below we discuss the three granularity levels in detail (summarized in Tab. 5).

4.3.1 ROI-level

In essence, ROI-level fusion only fuses features at selected object regions instead of dense locations on the feature maps. Hence ROI-level fusion is normally performed at a relatively late stage (*i.e.*, after the 3D region proposal generation stage). This fusion granularity happens when applying ROI-pooling for each view to obtain feature vectors of the same length [19, 78], as illustrated in Fig. 21. Also, it usually happens at the object proposal level in order to get 3D frustums from 2D RoIs through geometrical relationships [166, 173].

As a result, ROI-level fusion limits the ability of the neural network to capture the cross-modality interactions at earlier stages. To overcome this drawback, ROI-level fusion is often combined with other fusion granularity for further refinement of proposals [86, 187].

Table 9 Summary of methods that fuse raw point clouds and image masks

Methods	Year	Venue	Contributions
PointPainting [156]	2020	CVPR	<ul style="list-style-type: none"> Pioneer in painting LiDAR point clouds with image-based semantic mask. Achieve fine-grained point-wise correspondence.
HorizonLiDAR3D [32]	2020	Arxiv	<ul style="list-style-type: none"> Introduce a one-stage, anchor-free, and NMS-free detector. Effectively enhance the point cloud using point painting and test time augmentation.
CenterPointV2 [185]	2021	CVPR	<ul style="list-style-type: none"> Represent, detect, and track 3D objects as points. The representation is compatible with off-the-shelf 3D encoders.

4.3.2 Voxel-level

Voxel-level fusion exploits a relatively earlier fusion stage compared with ROI-level. Voxelized point cloud data is usually projected onto the image plane so we can append the image feature to each voxel, which is described in Fig. 22. Here, we establish a relatively approximate correspondence between the voxel features and image features. Specifically, we project each voxel feature center to the image plane through camera projection matrix. After obtaining a reference point in the image domain, the corresponding image feature is appended to the LiDAR voxel feature branch. Voxel-level fusion leads to a certain degree of information loss, resulting from both the spatial information loss in voxelization and the non-smooth camera feature maps. To address this issue, one can combine neighboring image feature pixels by the interpolated projection to correct the spatial offsets, which can achieve more accurate correspondence between voxels and the image feature. [85, 187]. Furthermore, instead of adopting a one-to-one matching pattern, we could explore cross attention mechanism that enables each voxel to perceive the whole image domain and adaptively attend corresponding 2D features.

In contrast to ROI-level fusion, this voxel-level granularity is finer and more precise. Besides, to deal with empty voxels derived from LiDAR sparsity, voxel-level fusion could aggregate dense image information to compensate for sparse LiDAR features [148].

4.3.3 Point-level

Point-level fusion is usually early fusion, where every 3D point is aggregated by an image feature or mask in order to capture a dense context. By “lifting” the corresponding image features or masks to the coordinates of the 3D points, point-level fusion provides an additional channel for each 3D point. Specifically, we use the known transformation matrix [149] to project 3D points to 2D image pixels and thereby establish a 3D-2D mapping. Next, we can decorate the point or voxel features with

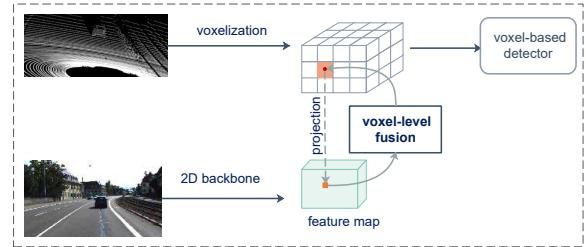


Fig. 22 Illustration of voxel-level fusion. We can gain the relationship between voxel centers and image features through: 1) use camera projection matrix to obtain reference point; 2) fetch the corresponding feature in the image domain.

the corresponding image masks through the mapping index. Fig. 23 outlines this process. The outstanding advantage of point-level fusion is the capability of summarizing useful information from both modalities since the image features are concatenated at a very early stage. Compared with the above two fusion granularity levels, we can simply build corresponding relations between dense images and sparse point clouds without the blurring problem (refer to Sec. 4.2.2) [156, 172].

Although experimental results show that point-level fusion effectively improves the overall performance [156, 185], there are still limitations. Firstly, due to the inherent occlusion problem in the image domain, 3D points which are mapped to the occluded image region may get the invalid image information [156]. Secondly, point-level fusion is less efficient in terms of memory consumption as compared to voxel-level fusion, as pointed out in [148].

4.3.4 Discussion

Fig. 24 clearly shows the years in which the deep learning based multi-modal 3D detection methods appeared. We also mark the fusion granularity of each method. With the passage of time, we observe that the granularity was relatively coarse at first and becomes finer. Meanwhile, some fusion methods adopt more than one fusion granularity level for further refinement.

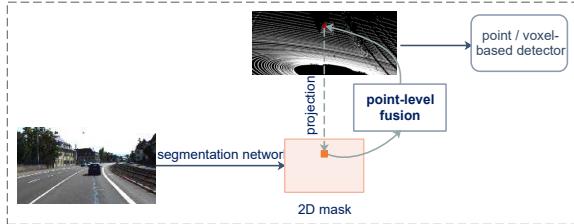


Fig. 23 Illustration of point-level fusion. We first perform a 3D-2D projection between points and pixels, and then the corresponding image mask can be added as the additional channel to decorate the 3D point cloud. For each decorated point, we flexibly select the voxel-based or point-based detector.

4.4 LiDAR-Camera Fusion: summary and development

In summary, ROI-level fusion is rather limited as this fusion lacks deep feature interaction. The later voxel-level and point fusion methods allow deep feature exchange and has its own merits. However, some researches recently reveal that such methods are easily affected by the sensor misalignment due to the hard association between points and pixels established by calibration matrices.

Most recently, the success of BEV-based methods in BEV map segmentation encourages us to extend it to the fusion-based 3D object detection task [134, 180, 197, 93, 94, 63, 62]. Follow-up works [6, 96] have proved that fusing LiDAR features with camera features in BEV is robust against degenerated image quality and sensor misalignment. As such, a new **BEV-level** paradigm for LiDAR-camera fusion has emerged. Instead of collecting 2D masks or features by the 3D-2D hard association, these methods directly lift image features to the 3D world, and these lifted features can be processed to the BEV level to fuse with the LiDAR BEV feature at a certain stage of the detection pipeline. For example, BEVfusion [96] lifts every image feature to the BEV space with an off-the-shelf depth estimator LSS [119] in a learnable fashion, then these lifted points are processed by a separate 3D encoder to produce a BEV map, the LiDAR-camera fusion happens at the BEV level by merging the two BEV maps from both modalities.

4.5 Fusion with Other Sensors

So far, we have discussed LiDAR-camera fusion methods in depth. We next briefly summarize methods that involve the fusion with millimeter wave radar (which we refer to as mmWave radar in this paper for brevity) sensors. The employment of mmWave radar is getting popular recently due to its long ranges, low cost, and

sensitivity to motions [76]. Accordingly, we briefly discuss Radar-Camera fusion and LiDAR-Radar fusion.

For Radar-Camera fusion, Chadwick et al. [12] project radar detection results to the image plane to boost the object detection accuracy for distant objects. Similarly, Nabati and Qi [107] use radar detection results to first generate 3D object proposals, then project them to the image plane to perform joint 2D object detection and depth estimation. CenterFusion [108] proposes to exploit both radar and camera data for 3D object detection. It first utilizes a center point detection network to detect objects by identifying their center points on the image. Next, it solves the key data association problem using a novel frustum-based method to associate radar detections with the corresponding 2D proposals. The above methods all directly use radar detection results without exploring features of radar points. Instead, Kim et al. [76] propose a low-level sensor fusion 3D object detector that combines two ROIs from radar and camera feature maps by a Gated ROI Fusion (GRIF), which provides more robust vehicle detection performance.

For LiDAR-Radar fusion, RadarNet [178] fuses radar and LiDAR data for 3D object detection. It employs an early fusion approach to learn joint representations from the two sensors and a decision fusion mechanism to exploit the radars radial velocity evidence. Disappointingly, RadarNet faces significant performance degradation in rare but critical adverse weather conditions. To remedy this, Qian et al. [125] exploit complementary radar which is less impacted by adverse weather and becomes prevalent on vehicles. They present a two-stage deep fusion detector to enhance the overall detection results. Specifically, this method first generates 3D proposals from LiDAR and complementary radar and then fuse region-wise features between multi-modal sensor streams.

Finally, we would like to point out that it is also of use to fuse multiple sensors of the same kind. HorizonLiDAR3D [32] combines all point clouds generated by five LiDAR sensors to augment the information of the point cloud data. In this work, a simple concatenation of point clouds from all LiDAR sensors is performed.

5 Open Challenges and Possible Solutions

Sensor modalities hold different properties and capture the same scene from various perspectives, rendering it a challenging task to combine data from multiple modalities into a coherent data stream. In this section, we discuss open challenges and possible solutions for multi-modal 3D object detection, which we hope to provide helpful guidelines on how to improve the performance of the multi-sensor perception systems.

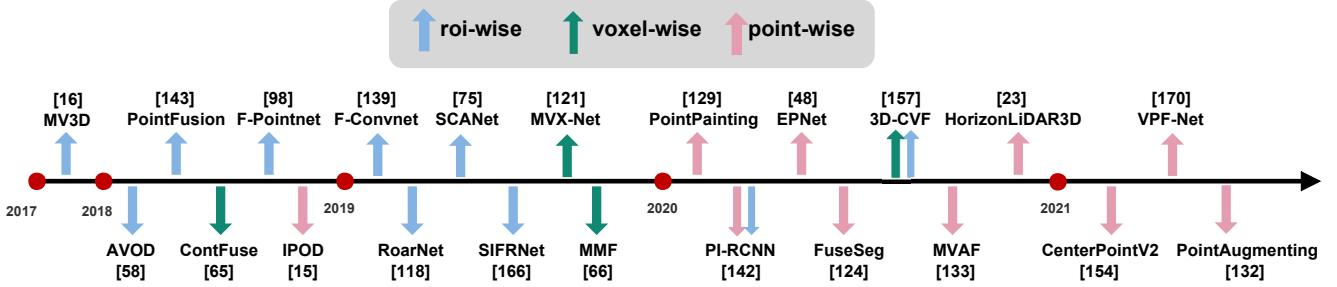


Fig. 24 Timeline of the fusion-based 3D object detection methods. We use different colors to mark their fusion granularity.

5.1 Open Challenge I: Multi-Sensor Calibration

As shown in Fig. 25, multiple sensors mounted on the autonomous vehicle are from different sensor coordinates. Fusion based methods are required the alignment of these sensor data. Here, we use LiDAR-camera fusion as example to explain this challenge. Point clouds are a set of points indicating 3D coordinates of the objects. RGB images are matrices of pixels, with each pixel's coordinate represented as (x, y) , where x, y is the pixel's row and column index, respectively. To build the map from 3D LiDAR coordinates to the 2D image plane, we must perform calibration between the two.

Traditional calibration methods use a calibration target to derive the intrinsic and extrinsic camera parameters. This cumbersome process requires lots of manual efforts. A common practice is to develop a target-less, automatic calibration method that can continuously calibrate the LiDAR sensor and camera on the fly. Target-less calibration is currently an active topic of research in this field [82, 111, 139]. These methods automatically calibrate among multiple sensors without human experts. However, inevitable bumps and jitters when driving AVs lead to the variation of the extrinsic parameters for the well calibrated LiDAR-camera system. Much worse, the error will gradually accumulate if not corrected in time, which may eventually affect the perception results. A possible solution to prevent this problem is to integrate the LiDAR and camera in a suite[66, 150], preventing their relative displacement to the greatest extent.

5.2 Open Challenge II: Information Loss during Fusion

When fusing data from multiple modalities, a certain amount of information will be lost inevitably due to projection, quantization, feature burring, etc. When devising a multi-modal fusion network, we need to pay attention to the stage, input, and granularity of the fusion operation, in order to minimize the information loss.

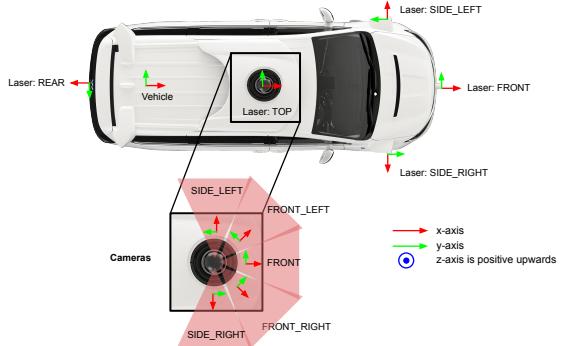


Fig. 25 Cameras and LiDAR sensors deployed on the Waymo autonomous vehicle [150]

The choice of fusion stage results in a different level of information loss. A later fusion stage is easy to implement, but cannot enjoy the rich information embedded in the raw data or earlier feature maps. Considering the complexity of the problem, it is very challenging, if at all possible, to pinpoint the optimal fusion stage that balances information loss and ease of implementation. To this end, a possible solution is to consider utilizing Neural Architecture Search (NAS) technique [91, 152] to find the appropriate fusion stage within a pipeline. It defines search space and then devises a search algorithm to propose near-optimal neural architectures.

The choice of fusion inputs has the greatest bearing on the amount of information loss, as a result of data projection or voxelization. For example, converting a point cloud to its BEV compresses the point cloud in the vertical direction and thus leads to the loss of height information. Converting a point cloud to its RV suffers from the problem of scale variation. Accordingly, it's important to find suitable input representations that reserve rich geometric and semantic information as much as possible. Moving forward, we could investigate several possible solutions. Specifically, we could exploit the attention mechanism [155, 164] to enhance certain features for each modality. Or, we could employ multiple representations to retain important information. For

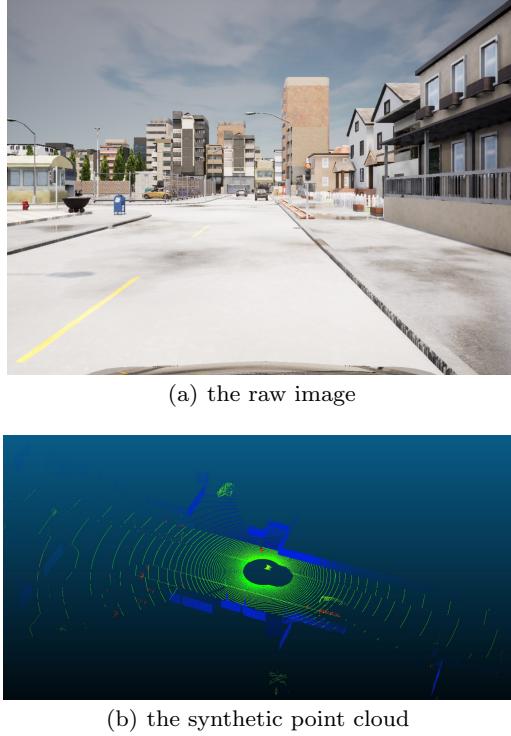


Fig. 26 An example RGB image (a) and the corresponding synthetic point cloud with semantic segmentation annotations (b). Both images are obtained from KITTI-CARLA synthetic dataset [31].

example, we can utilize both the point cloud and the corresponding voxel grid as fusion input for the point cloud branch [28, 29]. However, existing approaches do not take full advantage of temporal fusion input, which potentially limits the performance of multi-modal 3D object detection. In the further, we believe it is of significance to learn the 4D spatio-temporal information fusion across sensor and time.

The choice of fusion granularity can also affect the amount of information loss, *e.g.*, aligning the multi-modal data in a coarse granularity leads to the problem of feature blurring. A possible solution is to employ learnable calibration offsets to aggregate neighbor spatial information [187]. In this way, we can maximize the effect of data fusion.

5.3 Open Challenge III: Efficient Multi-Modal Data Augmentation

Due to the limited number of objects in the dataset, data augmentation is usually adopted to ensure efficient learning and avoid overfitting. Existing data augmentation techniques for each single data stream can be applied to deep fusion methods, such as object cut-and-paste, random flipping, scaling, rotation, and so

on [175, 198]. However, to keep data augmentation consistent across multiple modalities, we need to build the fine-grained mapping between data elements (such as points or pixels). Unfortunately, the augmentation operations usually choose to work on randomly selected objects and are thus inconsistent across the modalities.

Recently, several methods are proposed [159, 192] to address this problem. Zhang et al. [192] present a new multi-modality augmentation approach by cutting point cloud and imagery patches of ground-truth objects and pasting them into different scenes in a consistent manner, which prevents misalignment between multi-modal data. When projecting 3D points to 2D pixels, it first performs the reverse operation of translation, rotation, flip, *etc.* to restore the original point cloud, then gets point-pixel mapping based on the calibration information. In the future, more efficient multi-modal augmentation techniques need to be investigated.

5.4 Open Challenge IV: Low-Cost Multi-Modal 3D Object Detection

Monocular or stereo cameras are the most common low-cost sensors that can meet the requirements of mass production. However, without accurate 3D geometry information, relying on cameras alone cannot yield 3D detection results comparable to LiDAR-based methods. In fact, the state-of-the-art monocular method DD3D [113] achieves only 16.87% mAP on the KITTI 3D object detection leader board; the best stereo method LIGA-Stereo [52] can achieve 64.66% mAP. Nevertheless, the best LiDAR-only method BtcDet [174] has obtained 82.86% mAP.

Moving forward, with the development of *knowledge distillation* [57], one could exploit LiDAR data to distill 3D geometric information for camera-based detectors using large-scale and well-calibrated multi-modal data. Such a method can potentially achieve accurate detection as well as low system cost.

5.5 Open Challenge V: Shortage of Large Datasets

Another bottleneck in multi-modal 3D detection is the availability of high-quality, publicly usable datasets annotated with ground-truth information. Currently, popular datasets in 3D detection have the following issues: small size, class imbalance, and labeling errors, as discussed in Sec. 3.

Unsupervised and weakly-supervised fusion networks could allow the networks to be trained on large, unlabeled or partially labeled datasets [9].

There are also emerging works on generating synthetic datasets for RGB images and point clouds [31, 33, 42, 70, 102, 121, 131, 136, 168], which provide large-scale data with rich annotations. Fig. 26 shows an example of the KITTI-CARLA [31] synthetic dataset. However, there may be a domain gap between synthetic datasets and real-world datasets. Some recent works [58, 30, 50, 132] try to utilize technologies such as photorealistic rendering, unsupervised domain adaptation, and generative adversarial networks (GANs) [48] to bridge the gap between synthetic and real-world data. Still, how to use the models trained on the synthetic data to deal with real-world scenarios remains to be further investigated.

6 Conclusion

Due to the increasing importance of 3D vision in applications such as autonomous driving, this paper reviews the recent multi-modal 3D object detection networks, especially those that fuse camera images and LiDAR point clouds. We first carefully compare popular sensors and discuss their advantages and disadvantages and summarize the common problems of single-modal methods. We then provide an in-depth summary of several popular datasets that are commonly used for autonomous driving. In order to provide a systematic review, we discuss the multi-modal fusion methods based upon their choices for the following three design considerations: (1) fusion stage, *i.e.*, where does the fusion take place in the pipeline, (2) fusion input, *i.e.*, what data inputs are used for fusion, and (3) fusion granularity, *i.e.*, at what granularity level the two data streams are combined. Finally, we discuss open challenges and potential solutions in multi-modal 3D object detection.

References

- Ahmad WA, Wessel J, Ng HJ, Kissinger D (2020) IoT-ready millimeter-wave radar sensors. In: IEEE Global Conference on Artificial Intelligence and Internet of Things (GCAIoT), pp 1–5
- Andriluka M, Roth S, Schiele B (2010) Monocular 3d pose estimation and tracking by detection. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp 623–630
- Arnold E, Al-Jarrah OY, Dianati M, Fallah S, Oxtoby D, Mouzakitis A (2019) A survey on 3d object detection methods for autonomous driving applications. *IEEE Transactions on Intelligent Transportation Systems (TITS)* 20(10):3782–3795
- Asvadi A, Garrote L, Premebida C, Peixoto P, Nunes U (2017) Multimodal vehicle detection: Fusing 3d-lidar and color camera data. *Pattern Recognition Letters* 115:20–29
- Bai X, Hu Z, Zhu X, Huang Q, Chen Y, Fu H, Tai CL (2022) Transfusion: Robust lidar-camera fusion for 3d object detection with transformers. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp 1090–1099
- Beltrn J, Guindel C, Moreno FM, Cruzado D, Garca F, De La Escalera A (2018) Birdnet: A 3d object detection framework from lidar information. In: 2018 21st International Conference on Intelligent Transportation Systems (ITSC), pp 3517–3523
- Caesar H, Bankiti V, Lang AH, Vora S, Liang VE, Xu Q, Krishnan A, Pan Y, Baldan G, Beijbom O (2020) nuscenes: A multimodal dataset for autonomous driving. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp 11618–11628
- Caine B, Roelofs R, Vasudevan V, Ngiam J, Chai Y, Chen Z, Shlens J (2021) Pseudo-labeling for scalable 3d object detection. CoRR abs/2103.02093
- Caltagirone L, Bellone M, Svensson L, Wahde M (2019) Lidar-camera fusion for road detection using fully convolutional neural networks. *Robotics and Autonomous Systems* 111:125–131
- Carr P, Sheikh Y, Matthews I (2012) Monocular object detection using 3d geometric primitives. In: Fitzgibbon A, Lazebnik S, Perona P, Sato Y, Schmid C (eds) European Conference on Computer Vision (ECCV), pp 864–878
- Chadwick S, Maddern W, Newman P (2019) Distant vehicle detection using radar and vision. In: IEEE International Conference on Robotics and Automation (ICRA), pp 8311–8317
- Chang MF, Lambert J, Sangkloy P, Singh J, Bak S, Hartnett A, Wang D, Carr P, Lucey S, Ramanan D, Hays J (2019) Argoverse: 3d tracking and forecasting with rich maps. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp 8740–8749
- Charles RQ, Su H, Kaichun M, Guibas LJ (2017) Pointnet: Deep learning on point sets for 3d classification and segmentation. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp 77–85
- Chen L, Zou Q, Pan Z, Lai D, Cao D (2019) Surrounding vehicle detection using an fpga panoramic camera and deep cnns. *IEEE Transactions on Intelligent Transportation Systems PP(99):1–13*
- Chen L, Lin S, Lu X, Cao D, Wu H, Guo C, Liu C, Wang F (2021) Deep neural network based vehicle and pedestrian detection for autonomous driving: A survey. *IEEE Transactions on Intelligent Transportation Systems (TITS)* 22(6):3234–3246
- Chen LC, Papandreou G, Kokkinos I, Murphy K, Yuille AL (2018) Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *IEEE Transactions on Pattern Recognition and Machine Intelligence (TPAMI)* 40(4):834–848
- Chen X, Kundu K, Zhang Z, Ma H, Fidler S, Urtasun R (2016) Monocular 3d object detection for autonomous driving. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp 2147–2156
- Chen X, Ma H, Wan J, Li B, Xia T (2017) Multi-view 3d object detection network for autonomous driving. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp 1907–1915

20. Chen X, Kundu K, Zhu Y, Ma H, Fidler S, Urtasun R (2018) 3d object proposals using stereo imagery for accurate object class detection. *IEEE Transactions on Pattern Recognition and Machine Intelligence (TPAMI)* 40(5):1259–1272
21. Chen Y, Liu J, Qi X, Zhang X, Sun J, Jia J (2022) Scaling up kernels in 3d cnns. arXiv preprint arXiv:220610555
22. Chen Z, Li Z, Zhang S, Fang L, Jiang Q, Zhao F (2022) Autoalignv2: Deformable feature aggregation for dynamic multi-modal 3d object detection. In: European Conference on Computer Vision (ECCV)
23. Chen Z, Li Z, Zhang S, Fang L, Jiang Q, Zhao F, Zhou B, Zhao H (2022) Autoalign: Pixel-instance feature aggregation for multi-modal 3d object detection. In: IJCAI
24. Chu X, Deng J, Li Y, Yuan Z, Zhang Y, Ji J, Zhang Y (2021) Neighbor-vote: Improving monocular 3d object detection through neighbor distance voting. In: ACM International Conference on Multimedia (ACM MM), ACM, pp 5239–5247
25. Cordts M, Omran M, Ramos S, Rehfeld T, Enzweiler M, Benenson R, Franke U, Roth S, Schiele B (2016) The cityscapes dataset for semantic urban scene understanding. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp 3213–3223
26. Cui Y, Chen R, Chu W, Chen L, Tian D, Li Y, Cao D (2021) Deep learning for image and point cloud fusion in autonomous driving: A review. *IEEE Transactions on Intelligent Transportation Systems (TITS)* pp 1–18
27. Deng J, Dong W, Socher R, Li LJ, Li K, Fei-Fei L (2009) Imagenet: A large-scale hierarchical image database. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp 248–255
28. Deng J, Shi S, Li P, Zhou W, Zhang Y, Li H (2020) Voxel r-cnn: Towards high performance voxel-based 3d object detection. arXiv:201215712
29. Deng J, Zhou W, Zhang Y, Li H (2021) From multi-view to hollow-3d: Hallucinated hollow-3d R-CNN for 3d object detection. *IEEE Transactions on Circuits and Systems for Video Technology (TCSVT)* 31(12):4722–4734
30. Denninger M, Sundermeyer M, Winkelbauer D, Olefir D, Hodan T, Zidan Y, Elbadrawy M, Knauer M, Katam H, Lodhi A (2020) Blenderproc: Reducing the reality gap with photorealistic rendering. In: International Conference on Robotics: Science and Systems, RSS 2020
31. Deschaud JE (2021) Kitti-carla: a kitti-like dataset generated by carla simulator. arXiv preprint arXiv:210900892
32. Ding Z, Hu Y, Ge R, Huang L, Chen S, Wang Y, Liao J (2020) 1st place solution for waymo open dataset challenge - 3d detection and domain adaptation. CoRR abs/2006.15505
33. Dosovitskiy A, Ros G, Codevilla F, Lopez A, Koltun V (2017) CARLA: An open urban driving simulator. In: Proceedings of the Annual Conference on Robot Learning, pp 1–16
34. Engelberg T, Niem W (2009) Method for classifying an object using a stereo camera. US Patent App. 10/589,641
35. Enzweiler M, Gavrila DM (2009) Monocular pedestrian detection: Survey and experiments. *IEEE Transactions on Pattern Recognition and Machine Intelligence (TPAMI)* 31:2179–2195
36. Everingham M, Van Gool L, Williams CK, Winn J, Zisserman A (2010) The pascal visual object classes (voc) challenge. *International journal of computer vision* 88(2):303–338
37. Fan L, Xiong X, Wang F, Wang N, Zhang Z (2021) Rangedet: In defense of range view for lidar-based 3d object detection. CoRR abs/2103.10039
38. Fan L, Pang Z, Zhang T, Wang YX, Zhao H, Wang F, Wang N, Zhang Z (2022) Embracing single stride 3d object detector with sparse transformer. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp 8458–8468
39. Fayyad J, Jaradat M, Gruyer D, Najjaran H (2020) Deep learning sensor fusion for autonomous vehicle perception and localization: A review. *Sensors* 20:4220
40. Feng D, Haase-Schütz C, Rosenbaum L, Hertlein H, Gläser C, Timm F, Wiesbeck W, Dietmayer K (2021) Deep multi-modal object detection and semantic segmentation for autonomous driving: Datasets, methods, and challenges. *IEEE Transactions on Intelligent Transportation Systems (TITS)* 22(3):1341–1360
41. Gähler N, Jourdan N, Cordts M, Franke U, Denzler J (2020) Cityscapes 3d: Dataset and benchmark for 9 dof vehicle detection. In: IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)
42. Gaidon A, Wang Q, Cabon Y, Vig E (2016) Virtual worlds as proxy for multi-object tracking analysis. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp 4340–4349
43. Geiger A, Lenz P, Urtasun R (2012) Are we ready for autonomous driving? the kitti vision benchmark suite. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp 3354–3361
44. Geiger A, Lenz P, Stiller C, Urtasun R (2013) Vision meets robotics: The kitti dataset. *The International Journal of Robotics Research (IJRR)* 32(11):1231–1237
45. Geiger D, Yuille AL (1991) A common framework for image segmentation. *International Journal on Computer Vision (IJCV)* 6(3):227–243
46. Girshick R (2015) Fast r-cnn. In: IEEE International Conference on Computer Vision (ICCV), pp 1440–1448
47. Girshick R, Donahue J, Darrell T, Malik J (2014) Rich feature hierarchies for accurate object detection and semantic segmentation. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp 580–587
48. Goodfellow I, Pouget-Abadie J, Mirza M, Xu B, Warde-Farley D, Ozair S, Courville A, Bengio Y (2020) Generative adversarial networks. *Communications of the ACM* 63(11):139–144
49. Guan T, Wang J, Lan S, Chandra R, Wu Z, Davis L, Manocha D (2022) M3detr: Multi-representation, multi-scale, mutual-relation 3d object detection with transformers. In: Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision, pp 772–782
50. Guizilini V, Li J, Ambru R, Gaidon A (2021) Geometric unsupervised domain adaptation for semantic segmentation. In: Proceedings of the IEEE/CVF International Conference on Computer Vision, pp 8537–8547
51. Guo J, Kurup U, Shah M (2019) Is it safe to drive? an overview of factors, metrics, and datasets for driveability assessment in autonomous driving. *IEEE Transactions on Intelligent Transportation Systems PP(99):1–17*
52. Guo X, Shi S, Wang X, Li H (2021) Liga-stereo: Learning lidar geometry aware representations for stereo-based 3d detector. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp 3153–3163

53. He C, Zeng H, Huang J, Hua XS, Zhang L (2020) Structure aware single-stage 3d object detection from point cloud. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR)
54. He K, Zhang X, Ren S, Sun J (2016) Deep residual learning for image recognition. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp 770–778
55. He K, Gkioxari G, Dollár P, Girshick RB (2017) Mask R-CNN. In: IEEE International Conference on Computer Vision (ICCV), pp 2980–2988
56. He T, Soatto S (2019) Mono3d++: Monocular 3d vehicle detection with two-scale 3d hypotheses and task priors. In: Association for the Advancement of Artificial Intelligence (AAAI), vol 33, pp 8409–8416
57. Hinton G, Vinyals O, Dean J (2015) Distilling the knowledge in a neural network. arXiv preprint arXiv:150302531
58. Hodan T, Vineet V, Gal R, Shalev E, Hanzelka J, Connell T, Urbina P, Sinha SN, Guenter B (2019) Photorealistic image synthesis for object instance detection. In: 2019 IEEE international conference on image processing (ICIP), IEEE, pp 66–70
59. Hu P, Ziglar J, Held D, Ramanan D (2020) What you see is what you get: Exploiting visibility for 3d object detection. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR), Computer Vision Foundation / IEEE, pp 10998–11006
60. Hu Y, Ding Z, Ge R, Shao W, Huang L, Li K, Liu Q (2021) Afdetv2: Rethinking the necessity of the second stage for object detection from point clouds. arXiv preprint arXiv:211209205
61. Huang G, Liu Z, Van Der Maaten L, Weinberger KQ (2017) Densely connected convolutional networks. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp 2261–2269
62. Huang J, Huang G (2022) Bevdet4d: Exploit temporal cues in multi-camera 3d object detection. arXiv preprint arXiv:220317054
63. Huang J, Huang G, Zhu Z, Du D (2021) Bevdet: High-performance multi-camera 3d object detection in bird-eye-view. arXiv preprint arXiv:211211790
64. Huang P, Cheng M, Chen Y, Luo H, Wang C, Li J (2017) Traffic sign occlusion detection using mobile laser scanning point clouds. IEEE Transactions on Intelligent Transportation Systems 18(9):2364–2376
65. Huang T, Liu Z, Chen X, Bai X (2020) Epnet: Enhancing point features with image semantics for 3d object detection. In: European Conference on Computer Vision (ECCV), vol 12360, pp 35–52
66. Huang X, Wang P, Cheng X, Zhou D, Geng Q, Yang R (2019) The apolloscape open dataset for autonomous driving and its application. IEEE Transactions on Pattern Recognition and Machine Intelligence (TPAMI) 42(10):2702–2719
67. Ioannidou A, Chatzilari E, Nikolopoulos S, Kompati-siaris I (2017) Deep learning advances in computer vision with 3d data: A survey. ACM Computing Survey 50(2):20:1–20:38
68. Jiang M, Wu Y, Lu C (2018) Pointsift: A sift-like network module for 3d point cloud semantic segmentation. CoRR abs/1807.00652
69. Jiao Y, Jie Z, Chen S, Chen J, Wei X, Ma L, Jiang YG (2022) Msmdfusion: Fusing lidar and camera at multiple scales with multi-depth seeds for 3d object detection. arXiv preprint arXiv:220903102
70. Kar A, Prakash A, Liu MY, Cameracci E, Yuan J, Rusiniak M, Acuna D, Torralba A, Fidler S (2019) Meta-sim: Learning to generate synthetic datasets. In: IEEE International Conference on Computer Vision (ICCV), pp 4550–4559
71. Kellner D, Klappstein J, Dietmayer K (2012) Grid-based DBSCAN for clustering extended objects in radar data. In: IEEE Intelligent Vehicles Symposium (IV), pp 365–370
72. Kesten R, Usman M, Houston J, Pandya T, Nadhamuni K, Ferreira A, Yuan M, Low B, Jain A, Ondruska P, Omari S, Shah S, Kulkarni A, Kazakova A, Tao C, Platinsky L, Jiang W, Shet V (2019) Level 5 perception dataset 2020. <https://level-5.global/level5/data/>
73. Kim K, Woo W (2005) A Multi-view Camera Tracking for Modeling of Indoor Environment. Springer Berlin Heidelberg
74. Kim K, Woo W (2005) A multi-view camera tracking for modeling of indoor environment. In: Aizawa K, Nakamura Y, Satoh S (eds) Advances in Multimedia Information Processing - PCM 2004, pp 288–297
75. Kim Y (2014) Convolutional neural networks for sentence classification. Eprint Arxiv
76. Kim Y, Choi JW, Kum D (2020) Grif net: Gated region of interest fusion network for robust 3d object detection from radar point cloud and monocular image. In: IROS, pp 10857–10864
77. Krizhevsky A, Sutskever I, Hinton GE (2012) Imagenet classification with deep convolutional neural networks. In: Advances in Neural Information Processing Systems (NeurIPS), vol 25
78. Ku J, Mozifian M, Lee J, Harakeh A, Waslander SL (2018) Joint 3d proposal generation and object detection from view aggregation. In: IEEE International Conference on Intelligent Robots and Systems (IROS), pp 1–8
79. Lang AH, Vora S, Caesar H, Zhou L, Yang J, Beijbom O (2019) Pointpillars: Fast encoders for object detection from point clouds. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp 12697–12705
80. Lee CH, Lim YC, Kwon S, Lee JH (2011) Stereo vision-based vehicle detection using a road feature and disparity histogram. Optical Engineering 50(2):027004–027004–23
81. Lee S (2020) Deep learning on radar centric 3d object detection. CoRR abs/2003.00851
82. Levinson J, Thrun S (2013) Automatic online calibration of cameras and lasers. In: Robotics: Science and Systems, vol 2, p 7
83. Li P, Chen X, Shen S (2019) Stereo r-cnn based 3d object detection for autonomous driving. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp 7644–7652
84. Li Y, Yu AW, Meng T, Caine B, Ngiam J, Peng D, Shen J, Lu Y, Zhou D, Le QV, et al. (2022) Deepfusion: Lidar-camera deep fusion for multi-modal 3d object detection. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp 17182–17191
85. Liang M, Yang B, Wang S, Urtasun R (2018) Deep continuous fusion for multi-sensor 3d object detection. In: European Conference on Computer Vision (ECCV), pp 663–678
86. Liang M, Yang B, Chen Y, Hu R, Urtasun R (2019) Multi-task multi-sensor fusion for 3d object detection. In: IEEE Conference on Computer Vision and Pattern

- Recognition (CVPR), pp 7337–7345
87. Liang Z, Zhang M, Zhang Z, Zhao X, Pu S (2020) Rangercnn: Towards fast and accurate 3d object detection with range image representation. CoRR abs/2009.00206
 88. Lin T, Dollr P, Girshick R, He K, Hariharan B, Belongie S (2017) Feature pyramid networks for object detection. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp 936–944
 89. Lin TY, Maire M, Belongie S, Hays J, Perona P, Ramanan D, Dollár P, Zitnick CL (2014) Microsoft coco: Common objects in context. In: European Conference on Computer Vision (ECCV), pp 740–755
 90. Lin TY, Goyal P, Girshick R, He K, Dollr P (2017) Focal loss for dense object detection. IEEE Transactions on Pattern Recognition and Machine Intelligence (TPAMI) PP(99):2999–3007
 91. Liu H, Simonyan K, Yang Y (2019) DARTS: differentiable architecture search. In: International Conference on Learning Representations (ICLR)
 92. Liu W, Anguelov D, Erhan D, Szegedy C, Reed S, Fu CY, Berg AC (2016) Ssd: Single shot multibox detector. In: Leibe B, Matas J, Sebe N, Welling M (eds) European Conference on Computer Vision (ECCV), pp 21–37
 93. Liu Y, Wang T, Zhang X, Sun J (2022) Petr: Position embedding transformation for multi-view 3d object detection. arXiv preprint arXiv:220305625
 94. Liu Y, Yan J, Jia F, Li S, Gao Q, Wang T, Zhang X, Sun J (2022) Petrv2: A unified framework for 3d perception from multi-camera images. arXiv preprint arXiv:220601256
 95. Liu Z, Wu Z, Tth R (2020) Smoke: Single-stage monocular 3d object detection via keypoint estimation. In: IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW), pp 4289–4298
 96. Liu Z, Tang H, Amini A, Yang X, Mao H, Rus D, Han S (2022) Bevfusion: Multi-task multi-sensor fusion with unified bird’s-eye view representation. arXiv preprint arXiv:220513542
 97. Long J, Shelhamer E, Darrell T (2015) Fully convolutional networks for semantic segmentation. IEEE Transactions on Pattern Recognition and Machine Intelligence (TPAMI) 39(4):640–651
 98. Lu H, Chen X, Zhang G, Zhou Q, Ma Y, Zhao Y (2019) Scanet: Spatial-channel attention network for 3d object detection. In: IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp 1992–1996
 99. Ma X, Wang Z, Li H, Zhang P, Ouyang W, Fan X (2019) Accurate monocular 3d object detection via color-embedded 3d reconstruction for autonomous driving. In: IEEE International Conference on Computer Vision (ICCV), pp 6851–6860
 100. Mahmoud A, Hu JS, Waslander SL (2022) Dense voxel fusion for 3d object detection. arXiv preprint arXiv:220300871
 101. Major B, Fontijne D, Ansari A, Sukhavasi RT, Gowaiker R, Hamilton M, Lee S, Grzechnik SK, Subramanian S (2019) Vehicle detection with automotive radar using deep learning on range-azimuth-doppler tensors. In: IEEE International Conference on Computer Vision Workshop (ICCVW), pp 924–932
 102. Manivasagam S, Wang S, Wong K, Zeng W, Sazanovich M, Tan S, Yang B, Ma WC, Urtasun R (2020) Lidarsim: Realistic lidar simulation by leveraging the real world. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp 11167–11176
 103. Mao J, Xue Y, Niu M, Bai H, Feng J, Liang X, Xu H, Xu C (2021) Voxel transformer for 3d object detection. In: IEEE International Conference on Computer Vision (ICCV)
 104. Marchand , Chaumette F (1999) An autonomous active vision system for complete and accurate 3d scene reconstruction. International Journal on Computer Vision (IJCV) 32(3):171–194
 105. Mnih V, Kavukcuoglu K, Silver D, Rusu A, Veness J, Bellemare M, Graves A, Riedmiller M, Fidjeland A, Ostrovski G, Petersen S, Beattie C, Sadik A, Antonoglou I, King H, Kumaran D, Wierstra D, Legg S, Hassabis D (2015) Human-level control through deep reinforcement learning. Nature 518:529–33
 106. Mousavian A, Anguelov D, Flynn J, Koeck J (2017) 3d bounding box estimation using deep learning and geometry. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp 5632–5640
 107. Nabati R, Qi H (2019) RRPN: radar region proposal network for object detection in autonomous vehicles. In: IEEE International Conference on Image Processing (ICIP), pp 3093–3097
 108. Nabati R, Qi H (2021) Centerfusion: Center-based radar and camera fusion for 3d object detection. In: IEEE Winter Conference on Applications of Computer Vision (WACV), pp 1527–1536
 109. Nießner M, Zollhöfer M, Izadi S, Stamminger M (2013) Real-time 3d reconstruction at scale using voxel hashing. ACM Transactions on Graphics (TOG) 32(6):169:1–169:11
 110. Pan X, Xia Z, Song S, Li LE, Huang G (2021) 3d object detection with pointformer. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp 7463–7472
 111. Pandey G, McBride JR, Savarese S, Eustice RM (2012) Automatic targetless extrinsic calibration of a 3d lidar and camera by maximizing mutual information. In: Association for the Advancement of Artificial Intelligence (AAAI), p 20532059
 112. Pang S, Morris D, Radha H (2020) Clocs: Camera-lidar object candidates fusion for 3d object detection. In: IEEE International Conference on Intelligent Robots and Systems (IROS), pp 10386–10393
 113. Park D, Ambrus R, Guizilini V, Li J, Gaidon A (2021) Is pseudo-lidar needed for monocular 3d object detection? In: IEEE International Conference on Computer Vision (ICCV), pp 3142–3152
 114. Park JY, Chu CW, Kim HW, Lim SJ, Park JC, Koo BK (2009) Multi-view camera color calibration method using color checker chart
 115. Patil A, Malla S, Gang H, Chen YT (2019) The h3d dataset for full-surround 3d multi-object detection and tracking in crowded urban scenes. In: IEEE International Conference on Robotics and Automation (ICRA), pp 9552–9557
 116. Patole SM, Torlak M, Wang D, Ali M (2017) Automotive radars: A review of signal processing techniques. IEEE Signal Processing Magazine 34(2):22–35
 117. de Paula Veronese L, Auat-Cheein F, Mutz F, Oliveira-Santos T, Guivant JE, de Aguiar E, Badue C, De Souza AF (2020) Evaluating the limits of a lidar for an autonomous driving localization. IEEE Transactions on Intelligent Transportation Systems (TITS) 22(3):1449–1458
 118. Pham QH, Sevestre P, Pahwa RS, Zhan H, Pang CH, Chen Y, Mustafa A, Chandrasekhar V, Lin J (2020) A*

- 3d dataset: Towards autonomous driving in challenging environments. In: IEEE International Conference on Robotics and Automation (ICRA), pp 2267–2273
119. Philion J, Fidler S (2020) Lift, splat, shoot: Encoding images from arbitrary camera rigs by implicitly unprojecting to 3d. In: European Conference on Computer Vision, Springer, pp 194–210
 120. Pon AD, Ku J, Li C, Waslander SL (2020) Object-centric stereo matching for 3d object detection. In: IEEE International Conference on Robotics and Automation (ICRA), pp 8383–8389
 121. Prakash A, Boochoon S, Brophy M, Acuna D, Cameracci E, State G, Shapira O, Birchfield S (2019) Structured domain randomization: Bridging the reality gap by context-aware synthetic data. In: IEEE International Conference on Robotics and Automation (ICRA), pp 7249–7255
 122. Qi CR, Yi L, Su H, Guibas LJ (2017) Pointnet++: Deep hierarchical feature learning on point sets in a metric space. In: Advances in Neural Information Processing Systems (NeurIPS), vol 30
 123. Qi CR, Liu W, Wu C, Su H, Guibas LJ (2018) Frustum pointnets for 3d object detection from rgbd data. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp 918–927
 124. Qi CR, Litany O, He K, Guibas L (2019) Deep hough voting for 3d object detection in point clouds. In: IEEE International Conference on Computer Vision (ICCV), pp 9276–9285
 125. Qian K, Zhu S, Zhang X, Li LE (2021) Robust multi-modal vehicle detection in foggy weather using complementary lidar and radar signals. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp 444–453
 126. Qin Z, Wang J, Lu Y (2019) Monogrnet: A geometric reasoning network for monocular 3d object localization. In: Association for the Advancement of Artificial Intelligence (AAAI), vol 33, pp 8851–8858
 127. Qin Z, Wang J, Lu Y (2019) Triangulation learning network: From monocular to stereo 3d object detection. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp 7615–7623
 128. Redmon J, Divvala S, Girshick R, Farhadi A (2016) You only look once: Unified, real-time object detection. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp 779–788
 129. Ren S, He K, Girshick R, Sun J (2017) Faster r-cnn: Towards real-time object detection with region proposal networks. *IEEE Transactions on Pattern Recognition and Machine Intelligence (TPAMI)* 39(6):1137–1149
 130. Repairer Driven News (2018) Velodyne: Leading LIDAR price halved, new high-res product to improve self-driving cars. <https://www.repairerdrivennews.com/2018/01/02/velodyne-leading-lidar-price-halved-new-high-res-product-to-improve-self-driving-cars/>
 131. Richter SR, Vineet V, Roth S, Koltun V (2016) Playing for data: Ground truth from computer games. In: Leibe B, Matas J, Sebe N, Welling M (eds) European Conference on Computer Vision (ECCV), pp 102–118
 132. Richter SR, Al Hajja HA, Koltun V (2022) Enhancing photorealism enhancement. *IEEE Transactions on Pattern Analysis and Machine Intelligence*
 133. Riegler G, Ulusoy AO, Geiger A (2017) Octnet: Learning deep 3d representations at high resolutions. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR), IEEE Computer Society, pp 6620–6629
 134. Roddick T, Cipolla R (2020) Predicting semantic map representations from images using pyramid occupancy networks. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp 11138–11147
 135. Ronneberger O, Pfisterer, Brox T (2015) U-net: Convolutional networks for biomedical image segmentation. In: Medical Image Computing and Computer-Assisted Intervention (MICCAI), vol 9351, pp 234–241
 136. Ros G, Sellart L, Materzynska J, Vazquez D, Lopez AM (2016) The synthia dataset: A large collection of synthetic images for semantic segmentation of urban scenes. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp 3234–3243
 137. Schlosser J, Chow CK, Kira Z (2016) Fusing lidar and images for pedestrian detection using convolutional neural networks. In: IEEE International Conference on Robotics and Automation (ICRA), pp 2198–2205
 138. Schmidhuber J (2015) Deep learning in neural networks: An overview. *Neural Networks* 61:85–117
 139. Schneider N, Piewak F, Stiller C, Franke U (2017) Regnet: Multimodal sensor registration using deep neural networks. In: IEEE Intelligent Vehicles Symposium (IV), pp 1803–1810
 140. Sheeny M, Pellegrin ED, Mukherjee S, Ahrabian A, Wang S, Wallace AM (2021) RADIATE: A radar dataset for automotive perception. In: IEEE International Conference on Robotics and Automation (ICRA)
 141. Shi S, Wang X, Li H (2019) Pointrcnn: 3d object proposal generation and detection from point cloud. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp 770–779
 142. Shi S, Guo C, Jiang L, Wang Z, Shi J, Wang X, Li H (2020) Pv-rcnn: Point-voxel feature set abstraction for 3d object detection. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp 10526–10535
 143. Shi S, Wang Z, Shi J, Wang X, Li H (2020) From points to parts: 3d object detection from point cloud with part-aware and part-aggregation network. *IEEE Transactions on Pattern Recognition and Machine Intelligence (TPAMI)* pp 1–1
 144. Shi W, Rajkumar R (2020) Point-gnn: Graph neural network for 3d object detection in a point cloud. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp 1711–1719
 145. Shin K, Kwon YP, Tomizuka M (2019) Roarnet: A robust 3d object detection based on region approximation refinement. In: IEEE Intelligent Vehicles Symposium (IV), pp 2510–2515
 146. Silver D, Huang A, Maddison C, Guez A, Sifre L, Driessche G, Schrittwieser J, Antonoglou I, Panneershelvam V, Lanctot M, Dieleman S, Grewe D, Nham J, Kalchbrenner N, Sutskever I, Lillicrap T, Leach M, Kavukcuoglu K, Graepel T, Hassabis D (2016) Mastering the game of go with deep neural networks and tree search. *Nature* 529:484–489
 147. Simonyan K, Zisserman A (2014) Very deep convolutional networks for large-scale image recognition. *Computer Science*
 148. Sindagi VA, Zhou Y, Tuzel O (2019) Mvx-net: Multi-modal voxelnet for 3d object detection. In: IEEE International Conference on Robotics and Automation (ICRA), pp 7276–7282

149. Strecha C, von Hansen W, Van Gool L, Fua P, Thoennessen U (2008) On benchmarking camera calibration and multi-view stereo for high resolution imagery. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp 1–8
150. Sun P, Kretzschmar H, Dotiwala X, Chouard A, Patnaik V, Tsui P, Guo J, Zhou Y, Chai Y, Caine B, Vasudevan V, Han W, Ngiam J, Zhao H, Timofeev A, Ettinger S, Krivokon M, Gao A, Joshi A, Zhang Y, Shlens J, Chen Z, Anguelov D (2020) Scalability in perception for autonomous driving: Waymo open dataset. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR)
151. Sun Y, Zuo W, Yun P, Wang H, Liu M (2020) Fuseseg: Semantic segmentation of urban scenes based on rgb and thermal data fusion. *IEEE Transactions on Automation Science and Engineering* PP(99):1–12
152. Tang H, Liu Z, Zhao S, Lin Y, Lin J, Wang H, Han S (2020) Searching efficient 3d architectures with sparse point-voxel convolution. In: European Conference on Computer Vision (ECCV), pp 685–702
153. Urmson C, Anhalt J, Bagnell D, Baker C, Bittner R, Clark M, Dolan J, Duggins D, Galatali T, Geyer C, et al. (2008) Autonomous driving in urban environments: Boss and the urban challenge. *Journal of Field Robotics* 25(8):425–466
154. Urmson C, Baker C, Dolan J, Rybski P, Salesky B, Whittaker WR, Ferguson D, Darmas M (2009) Autonomous driving in traffic: Boss and the urban challenge. *Ai Magazine* 30(2):17–28
155. Vaswani A, Shazeer N, Parmar N, Uszkoreit J, Jones L, Gomez AN, Kaiser Lu, Polosukhin I (2017) Attention is all you need. In: Advances in Neural Information Processing Systems (NeurIPS), vol 30, p 60006010
156. Vora S, Lang AH, Helou B, Beijbom O (2020) Point-painting: Sequential fusion for 3d object detection. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp 4603–4611
157. Wallace AM, Halimi A, Buller GS (2020) Full waveform lidar for adverse weather conditions. *IEEE Transactions on Vehicular Technology (TVT)* 69(7):7064–7077
158. Wandinger U (2005) Introduction to Lidar. Brooks/Cole Pub. Co.,
159. Wang C, Ma C, Zhu M, Yang X (2021) Pointaugmenting: Cross-modal augmentation for 3d object detection. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp 11794–11803
160. Wang G, Tian B, Zhang Y, Chen L, Cao D, Wu J (2020) Multi-View Adaptive Fusion Network for 3D Object Detection. arXiv e-prints p arXiv:2011.00652
161. Wang J, Zhou L (2019) Traffic light recognition with high dynamic range imaging and deep learning. *IEEE Transactions on Intelligent Transportation Systems* 20(4):1341–1352
162. Wang M, Deng W (2018) Deep visual domain adaptation: A survey. *Neurocomputing* 312:135–153
163. Wang S, Suo S, Ma W, Pokrovsky A, Urtasun R (2018) Deep parametric continuous convolutional neural networks. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp 2589–2597
164. Wang X, Girshick RB, Gupta A, He K (2018) Non-local neural networks. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR), Computer Vision Foundation / IEEE Computer Society, pp 7794–7803
165. Wang Y, Chao WL, Garg D, Hariharan B, Campbell M, Weinberger KQ (2019) Pseudo-lidar from visual depth estimation: Bridging the gap in 3d object detection for autonomous driving. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp 8437–8445
166. Wang Z, Jia K (2019) Frustum convnet: Sliding frustums to aggregate local point-wise features for amodal. In: IEEE International Conference on Intelligent Robots and Systems (IROS), pp 1742–1749
167. Wang Z, Jia K (2019) Frustum convnet: Sliding frustums to aggregate local point-wise features for amodal 3d object detection. In: IEEE International Conference on Intelligent Robots and Systems (IROS), pp 1742–1749
168. Weng X, Man Y, Cheng D, Park J, O'Toole M, Kitani K (2020) All-In-One Drive: A Large-Scale Comprehensive Perception Dataset with High-Density Long-Range Point Clouds. arXiv
169. Wilson B, Qi W, Agarwal T, Lambert J, Singh J, Khan-delwal S, Pan B, Kumar R, Hartnett A, Pontes JK, Ramanan D, Carr P, Hays J (2021) Argoverse 2: Next generation datasets for self-driving perception and forecasting. In: Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks (NeurIPS Datasets and Benchmarks 2021)
170. Wu X, Peng L, Yang H, Xie L, Huang C, Deng C, Liu H, Cai D (2022) Sparse fuse dense: Towards high quality 3d detection with depth completion. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp 5418–5427
171. Xie J, Kiefel M, Sun MT, Geiger A (2016) Semantic instance annotation of street scenes by 3d to 2d label transfer. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp 3688–3697
172. Xie L, Xiang C, Yu Z, Xu G, Yang Z, Cai D, He X (2020) Pi-rcnn: An efficient multi-sensor 3d object detector with point-based attentive cont-conv fusion module. In: Association for the Advancement of Artificial Intelligence (AAAI), vol 34, pp 12460–12467
173. Xu D, Anguelov D, Jain A (2018) Pointfusion: Deep sensor fusion for 3d bounding box estimation. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp 244–253
174. Xu Q, Zhong Y, Neumann U (2021) Behind the curtain: Learning occluded shapes for 3d object detection. arXiv preprint
175. Yan Y, Mao Y, Li B (2018) SECOND: sparsely embedded convolutional detection. *Sensors* 18(10):3337
176. Yang B, Luo W, Urtasun R (2018) PIXOR: real-time 3d object detection from point clouds. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp 7652–7660
177. Yang B, Luo W, Urtasun R (2018) Pixor: Real-time 3d object detection from point clouds. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp 7652–7660
178. Yang B, Guo R, Liang M, Casas S, Urtasun R (2020) Radarnet: Exploiting radar for robust perception of dynamic objects. In: European Conference on Computer Vision (ECCV), vol 12363, pp 496–512
179. Yang H, Liu Z, Wu X, Wang W, Qian W, He X, Cai D (2022) Graph r-cnn: Towards accurate 3d object detection with semantic-decorated local graph. In: European Conference on Computer Vision (ECCV)
180. Yang W, Li Q, Liu W, Yu Y, Ma Y, He S, Pan J (2021) Projecting your view attentively: Monocular road

- scene layout estimation via cross-view transformation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp 15536–15545
181. Yang Z, Sun Y, Liu S, Shen X, Jia J (2018) IPOD: intensive point-based object detector for point cloud. CoRR
 182. Yang Z, Sun Y, Liu S, Shen X, Jia J (2019) Std: Sparse-to-dense 3d object detector for point cloud. In: IEEE International Conference on Computer Vision (ICCV), pp 1951–1960
 183. Yang Z, Sun Y, Liu S, Jia J (2020) 3dssd: Point-based 3d single stage object detector. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp 11037–11045
 184. Yang Z, Chen J, Miao Z, Li W, Zhu X, Zhang L (2022) Deepinteraction: 3d object detection via modality interaction. arXiv preprint arXiv:220811112
 185. Yin T, Zhou X, Krähenbühl P (2021) Center-based 3d object detection and tracking. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp 11784–11793
 186. Yoo J, Ahn N, Sohn K (2020) Rethinking data augmentation for image super-resolution: A comprehensive analysis and a new strategy. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp 8372–8381
 187. Yoo JH, Kim Y, Kim J, Choi JW (2020) 3d-cvf: Generating joint camera and lidar features using cross-view spatial feature fusion for 3d object detection. In: European Conference on Computer Vision (ECCV), pp 720–736
 188. Yosinski J, Clune J, Bengio Y, Lipson H (2014) How transferable are features in deep neural networks? In: Advances in Neural Information Processing Systems (NeurIPS), vol 27
 189. You Y, Wang Y, Chao W, Garg D, Pleiss G, Hariharan B, Campbell ME, Weinberger KQ (2020) Pseudo-lidar++: Accurate depth for 3d object detection in autonomous driving. In: International Conference on Learning Representations (ICLR)
 190. Zewge NS, Kim Y, Kim J, Kim JH (2019) Millimeter-wave radar and rgb-d camera sensor fusion for real-time people detection and tracking. In: 2019 7th International Conference on Robot Intelligence Technology and Applications (RiTA), pp 93–98
 191. Zhang H, Yang D, Yurtsever E, Redmill KA, mit zgner (2021) Faraway-frustum: Dealing with lidar sparsity for 3d object detection using fusion
 192. Zhang W, Wang Z, Loy CC (2020) Multi-modality cut and paste for 3d object detection
 193. Zhang Y, Zhang S, Zhang Y, Ji J, Duan Y, Huang Y, Peng J, Zhang Y (2020) Multi-modality fusion perception and computing in autonomous driving. Journal of Computer Research and Development 57(9):1781
 194. Zhang Y, Carballo A, Yang H, Takeda K (2021) Autonomous driving in adverse weather conditions: A survey. arXiv preprint arXiv:211208936
 195. Zhang Y, Carballo A, Yang H, Takeda K (2021) Autonomous driving in adverse weather conditions: A survey. CoRR abs/2112.08936
 196. Zhao X, Liu Z, Hu R, Huang K (2019) 3d object detection using scale invariant and feature reweighting networks. In: Association for the Advancement of Artificial Intelligence (AAAI), pp 9267–9274
 197. Zhou B, Krähenbühl P (2022) Cross-view transformers for real-time map-view semantic segmentation. In: Pro-ceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp 13760–13769
 198. Zhou Y, Tuzel O (2018) Voxelnet: End-to-end learning for point cloud based 3d object detection. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp 4490–4499
 199. Zhou Y, Wan G, Hou S, Yu L, Wang G, Rui X, Song S (2020) Da4ad: End-to-end deep attention-based visual localization for autonomous driving. In: European Conference on Computer Vision (ECCV), pp 271–289
 200. Zhu H, Deng J, Zhang Y, Ji J, Mao Q, Li H, Zhang Y (2022) Vpfnet: Improving 3d object detection with virtual point based lidar and stereo data fusion. IEEE Transactions on Multimedia (TMM)