

15

# Towards Deep Radar Perception for Autonomous Driving: Datasets, Methods, and Challenges

Yi Zhou <sup>1,2,4,6</sup>, Lulu Liu <sup>1,3,4,5</sup>, Haocheng Zhao<sup>1,2,4,6</sup>, Miguel López-Benítez<sup>6,7</sup>, Limin Yu<sup>2</sup>, Yutao Yue <sup>1,4,5,\*</sup>

- <sup>1</sup> Institute of Deep Perception Technology, JITRI, Wuxi 214000, China; yueyutao@idpt.org (Y.Y.)
- <sup>2</sup> Department of Electrical and Electronic Engineering, School of Advanced Technology, Xi'an Jiaotong-Liverpool University, Suzhou 215123, China; zhouyi1023@tju.edu.cn (Y.Z.); haocheng.zhao19@student.xjtlu.edu.cn (H.Z); limin.yu@xjtlu.edu.cn (L.Y.)
- <sup>3</sup> Department of Mathematical Sciences, School of Science, Xi'an Jiaotong-Liverpool University, Suzhou 215123, China; lulu.liu21@student.xjtlu.edu.cn (L.L.)
- <sup>4</sup> XJTLU-JITRI Academy of Industrial Technology, Xi'an Jiaotong-Liverpool University, Suzhou, 215123, China
- <sup>5</sup> Department of Mathematical Sciences, University of Liverpool, Liverpool L69 7ZX, U.K.
- <sup>5</sup> Department of Electrical Engineering and Electronics, University of Liverpool, Liverpool L69 3GJ, U.K.; mlopben@liverpool.ac.uk (M.LB.)
- <sup>7</sup> ARIES Research Centre, Antonio de Nebrija University, 28040 Madrid, Spain
- \* Correspondence: yueyutao@idpt.org (Y.Y.)

Abstract: With recent developments, the performance of automotive radar has improved significantly. 1 The next generation of 4D radar can achieve imaging capability in the form of high-resolution point clouds. In this context, we believe that the era of deep learning for radar perception has arrived. However, studies on radar deep learning are spread across different tasks and a holistic overview is lacking. This review paper attempts to provide a big picture of the deep radar perception stack, including signal processing, datasets, labelling, data augmentation and downstream tasks such 6 as depth and velocity estimation, object detection, and sensor fusion. For these tasks, we focus on explaining how the network structure is adapted to radar domain knowledge. In particular, we summarize three overlooked challenges in deep radar perception, including multi-path effects, 9 uncertainty problems and adverse weather effects, and present some attempts to solve them. We also 10 provide an website for browsing each reference and related codes at https://zhouyi1023.github.io/ 11 awesome-radar-perception. 12

Keywords: automotive radars; radar signal processing; object detection; multi-sensor fusion; deep13learning; autonomous driving.14

# 1. Introduction

As autonomous driving technology progresses from the demonstration stage to the 16 landing stage, it puts forward higher requirements for perception ability. Mainstream 17 autonomous driving systems rely on fusion of cameras and LiDARs for perception. Al-18 though millimeter wave radar has been widely used in mass-produced cars for active safety 19 functions such as Automatic Emergency Braking (AEB) and Forward Collision Warning 20 (FCW), it is overlooked in autonomous driving. Recently, Tesla announced the removal 21 of radar sensors from its semi-autonomous driving system Autopilot. In the CVPR 2021 22 workshop [1], Tesla's director of AI Andrej Karpathy explained their reason by presenting 23 three typical scenarios for radar malfunctions, including lost tracking due to significant 24 deceleration of the front vehicle, false slow down under bridges and missed detection of a 25 stationary vehicle parked on the side of the main road. In the first case, radar's close field 26 detection ability is related to side lobes. Conventional radars with a limited number of 27 channels are not good at side lobe compression. The second case is caused by the fact that 28 conventional radar cannot measure height information and therefore confuses the bridge 29 overhead with static objects on the road. The reason for the third case is that conventional 30

Citation: Zhou, Y.; Lu, L.; Zhao, H.; López-Benítez, M.; Yu, L.; Yue, Y. Towards Deep Radar Perception for Autonomous Driving: Datasets, Methods, and Challenges. *Sensors* 2022, 1, 0. https://doi.org/

Received: Accepted: Published:

Review

**Publisher's Note:** MDPI stays neutral with regard to jurisdictional claims in published maps and institutional affiliations.

**Copyright:** © 2022 by the authors. Submitted to *Sensors* for possible open access publication under the terms and conditions of the Creative Commons Attri- bution (CC BY) license (https:// creativecommons.org/licenses/by/ 4.0/). radar has too low angular resolution to capture the shape of a static vehicle. All these challenges can be solved with next-generation high resolution radar.

As a ranging sensor, radar is usually compared to LiDAR. A typical 77 Ghz automo-33 tive radar has a wavelength of 3.9 mm, while automotive LiDARs have a much smaller 34 wavelength of 905 nm or 1550 nm. For a small aperture radar, most of the reflected signal 35 is not received by the radar sensor because of the specular reflection. Another problem 36 with small aperture is the low angular resolution, so that two close objects cannot be sep-37 arated effectively. These two features make the radar point cloud much sparser than the 38 LiDAR point cloud. Conventional automotive radars have a low resolution in elevation 30 and therefore return a two-dimensional point cloud. The next generation of high resolution 40 radar achieves higher angular resolution in both azimuth and elevation. Because it can 41 measure 3D position and Doppler velocity, it is always referred to as a 4D radar in the 42 marketplace. In table 1, typical types of radars and LiDARs are compared. We can find that 43 conventional long-range radars have a low angular resolution in horizontal view and no resolution in vertical view. In contrast, 4D radar can achieve an angular resolution of about 45 1° in both horizontal and vertical views. Therefore, the classification of static objects is no longer a limitation for 4D radar. Although 4D radar has much higher angular resolution, as 47 shown in fig. 1, the radar point cloud is still much sparser than the 16-beam LiDAR point cloud. However, radar can measure Doppler velocity and radar cross section (RCS), which 49 is expected to better help classify road users. In addition, 4D radar has the advantages of 50 long detection range (up to 300 meters), all-weather operation, low power consumption 51 and low cost. Therefore, we believe that radar is a good complement to LiDAR and vision. 52 The fusion of these sensors enables all-weather, long-range environment perception. 53



Figure 1. Point clouds of a 4D radar and a 16-beam LiDAR from the Astyx dataset [2]

Table 1. (	Characteristics	of typical	radars an	d LiDARs
------------	-----------------	------------	-----------	----------

	Conventional Radar <sup>1</sup> (Multi-Mode)	4D Radar	16-Beam LiDAR	32-Beam LiDAR	Solid State LiDAR
Max Range	f: 250m, n:70m	300m	100m	200m	200m
FoV (H/V)	f: 20°, n: 120°/ 🗡	120°/30°	360°/30°	360°/40°	120°25°
Ang Res (H/V)	f: 1.5°, n: 4°/ 🗡	1°/1°	0.1°/ 2°	0.1°/ 0.3°	0.2°/ 0.2°
Doppler Res	0.1m/s	0.1m/s	×	×	×
Point Density	Low	Medium	High	High	High
All Weather	1	1	×	×	×
Power	5W	5W	8W	10W	15W
Expected Cost	Low	Low	Medium	High	Medium

<sup>1</sup> A typical 77GHz 4Tx-6Rx automotive radar, 2Tx-6Rx for far range and 2Tx-6Rx for near range.

In recent years, with the trend of open source, more and more datasets, models and 54 toolboxes have been released. According to our statistics, 10 radar datasets were released 55 in 2021 and 2022. Along with these datasets, some seminal papers are proposed to leverage 56 deep learning in radar perception. However, due to the limited sensing capability of 57 conventional radar, the performance of these methods is far from good enough. Since the 58 introduction of 4D radar, we believe that the era of deep radar perception has arrived. With 59 the power of deep learning, we can design a highly reliable perception system based on the 60 fusion of radar and other modalities. 61

The application of deep learning in radar perception has drawn extensive attention 62 from autonomous driving researchers. In the past two years, a number of review papers [3– 63 8] have been published in this field. Zhou *et al.* [3] categorize radar perception tasks 64 into dynamic target detection and static environment modelling. They also provide brief 65 introductions to radar-based detection, tracking and localization. Abdu et al. [4] summarize 66 the deep learning models for radar perception based on different radar representations, in-67 cluding occupancy grid maps, range-Doppler-azimuth maps, micro-Doppler spectrograms 68 and point clouds. They also introduces approaches for radar and camera fusion based on 69 the classical taxonomy of data-level, feature-level and decision-level. Scheiner et al. [5] 70 discuss the information sparsity problem and labelling challenge in learning-based radar 71 perception. Three strategies are recommended to increase radar data density, including the 72 use of pre-CFAR data, the use of high-resolution radar sensors, and the use of polarization 73 information. In this paper, we differ from other review papers in three aspects: Firstly, we 74 provide a detailed summary and description of the public available radar datasets, which 75 is very useful for developing learning-based methods. Secondly, this review does not focus 76 on specific tasks, but aims to present a big picture of the radar perception framework, as 77 illustrated in fig. 2. Thirdly, rather than simply presenting the network structure, we focus 78 on explaining why these modules work from the perspective of radar domain knowledge. 79

In this article, we systematically review the recent advancements in deep radar percep-80 tion. In section 2, we introduce radar signal processing pipeline and different radar data 81 representations. In section 3, we summarize the public-available radar datasets (section 3.1) 82 for autonomous driving, as well as the calibration (section 3.2), labelling (section 3.3) and 83 data augmentation techniques (section 3.4). In the following sections, we introduce different radar perception tasks, including radar depth completion (section 4.1), radar full 85 velocity estimation (section 4.2), radar object detection (point-cloud-based in section 5.2 and pre-CFAR-based in section 5.3). In section 6, we classify sensor fusion frameworks into 87 four categories: input fusion (section 6.1), ROI fusion (section 6.2), feature map fusion (section 6.3) and decision fusion (section 6.4). Next, we discuss three challenges toward reliable 89 radar detection, including ghost objects (section 7.1), uncertainty problems (section 7.2), 90 and adverse weather effects (section 7.3). Finally, we propose several interesting research 91 directions in section 8.

sk	Velocity Estimation				2D Detection				Segmentation				
Та	Depth	Comple	tion		3D Detection				Tracking				
e					Fusion Layer								·7
atuı	Proposal	Co	onca	tenate	Gat	ing		Attention		Evidence		Bay	esian
Fe	2D Conv	3	3D Conv		PointNet++		1	Transformer		GraphConv		LSTM	
	Cal	libration			Calibration Labelling					Augr	nentation		
ıta	Point Cloud	RA Ma	ıp	RD Map	o Grid	-map	RA	D Tensor	In	nage	Pseudo	Range	Point
Da	FFT	C	CFAI	R	DoA		Clu	Clustering		8-	LiDAR	View	Cloud
	Radar Datasets Radar Simi		ulation Radar		ndar S	Sensor		Camera		Lidar			

Figure 2. Overview of deep radar perception framework.

#### 2. Radar Signal Processing Fundamentals

Knowledge of radar signal processing is essential for the development of a deep radar 94 perception system. Different radar devices vary in their sensing capabilities. It is important to leverage radar domain knowledge to understand the performance boundary, find key 96 scenarios and solve critical problems. This section outlines the classical signal processing pipeline for automotive radar applications. 98

# 2.1. FMCW Radar Signal Processing



Figure 3. Radar Tx/Rx signals and the resulted range-Doppler map.

Off-the-shelf automotive radars operate with a sequence of linear frequency-modulated 100 continuous-wave (FMCW) signals to simultaneously measure range, angle and velocity. 101 According to regulations, automotive radar are allowed to use two frequency bands in 102 millimeter waves: 24 GHz (24-24.25 GHz) and 77 GHz (77-79 GHz). There is a trend 103 towards 77 GHz due to its larger bandwidth (76-77 GHz for long-range and 77-81 GHz for 104 short-range), higher Doppler resolution and smaller antennas [9]. As shown in fig. 3, the 105 FMCW signal is characterized by the start frequency (also known as the carrier frequency) 106  $f_c$ , the sweep bandwidth B, the chirp duration  $T_c$  and the slope  $S = B/T_c$ . During one chirp 107 duration, the frequency increases linearly from  $f_c$  to  $f_c + B$  with a slope of S. One FMCW 108 waveform is referred to as a chirp, and radar transmits a frame of  $N_c$  chirps equally spaced 109 by chirp cycle time  $T_c$ . The total time  $T_f = N_c T_c$  is called the frame time, also known as the 110 time on target (TOT). In order to avoid the need of high-speed sampling, a frequency mixer 111 combines the received signal with the transmitted signal to produce two signals with sum 112 frequency  $f_T(t) + f_R(t)$  and difference frequency  $f_T(t) - f_R(t)$ . Then a low-pass filter is 113 used to filter out the sum frequency component and obtain the intermediate frequency (IF) 114 signal. In this way, FMCW radar can achieve GHz performance with only MHz sampling. 115 In practice, a quadrature mixer is used to improve the noise figure [10], resulting in a 116 complex exponential IF signal as 117

$$x_{\rm IF}(t) = A e^{j(2\pi f_{\rm IF}t + \phi_{\rm IF})} \tag{1}$$

where A is the amplitude,  $f_{IF} = f_T(t) - f_R(t)$  is referred to as the beat frequency 118 and  $\phi_{\rm IF}$  is the phase of the IF signal. Next, the IF signal is sampled  $N_{\rm s}$  times by an ADC 119 converter, resulting in a discrete-time complex signal. Multiple frames of chirp signals are 120

93

97

assembled into a two-dimensional matrix. As shown in fig. 3, the dimension of sampling points within a chirp is referred to as fast time, and the dimension of chirp index within one frame is referred to as slow time. Assume one object moving with speed v at distance r, the frequency and phase of the IF signal are given by 122

$$f_{\rm IF} = \frac{2S(r+vT_c)}{c}, \phi_{\rm IF} = \frac{4\pi(r+vT_c)}{\lambda}$$
(2)

where  $\lambda = c/f_c$  is the wavelength of the chirp signal. From eq. (2), we can find that 125 range and Doppler velocity are coupled. Under the following assumptions: 1. the range 126 variations in slow time caused by target motion can be neglected due to the short frame 127 time. 2. the Doppler frequency in fast time can be neglected compared to the beat frequency 128 by utilizing a wideband waveform. Then, range and Doppler can be decoupled. Range 129 can be estimated from the beat frequency as  $r = c f_{\rm IF}/2S$  and Doppler velocity can be 130 estimated from the phase shift between two chirps as  $v = \Delta \phi \lambda / 4\pi T_c$ . Next, a range DFT is 131 applied in the fast-time dimension to resolve the frequency change, followed by a Doppler 132 DFT in the slow-time dimension to resolve the phase change. As a result, we obtain a 133 2D complex-valued data matrix called Range-Doppler (RD) map. In practice, a window 134 function is applied before DFT to reduce sidelobes. The range and the Doppler velocity of 135 a cell in RD map are given by 136

$$r_k = k \frac{c}{2B_{\rm IF}}, v_l = l \frac{\lambda}{2T_f} \tag{3}$$

where k and l denote the indexes of DFT,  $B_{IF}$  is the IF bandwidth, and  $T_f$  is the frame time. In practice, FFT is applied due to its computational efficiency. Accordingly, the sequence will be zero-padded to the nearest power of 2 if necessary.



**Figure 4.** MIMO radar principles.(a) Virtual array configuration of a 2Tx4Rx MIMO radar (b) In TDM mode, Tx1 and Tx2 transmit signals by turns. (c) In DDM mode, a Doppler shift is added to Tx2.

Angle information can be obtained using more than one receive or transmit channel. <sup>140</sup> Single-input multiple-output (SIMO) radars utilize a single transmit (Tx) and multiple <sup>141</sup> receive (Rx) antennas for angle estimation. Suppose one object is located at direction  $\theta$ . <sup>142</sup> Similar to Doppler processing, the induced frequency change between two adjacent receiver <sup>143</sup> antennas can be neglected, while the induced phase change can be used for calculating <sup>144</sup> the direction of angle. This phase change is given by  $\Delta \phi = 2\pi d \sin \theta / \lambda$ , where *d* is the inter-antenna spacing. To achieve maximum unambiguous angle, the spacing can be set to  $\lambda/2$ . Then, a third FFT can be applied to the received antenna dimension. For conventional radar with a small number of Rx antennas, the sequence is often padded with  $N_{\text{FFT}} - N_{Rx}$ zeros to achieve a smooth display of the spectrum. The angle at index  $\eta$  is given by

$$\theta_{\eta} = \arcsin \frac{\eta \lambda}{N_{\rm FFT}} \tag{4}$$

The angular resolution of a SIMO radar depends on the number of Rx antennas. The 150 maximum number of Rx antennas is limited by the additional cost of signal processing 151 chains on the device [11]. Multiple-input multiple-output (MIMO) radar operates with 152 multiple channels in both Tx and Rx. As illustrated in fig. 4 (a), a MIMO radar with  $N_{\text{Tx}}$ 153 Tx antennas and  $N_{\text{Rx}}$  Rx antennas can synthesize a virtual array with  $N_{\text{Tx}}N_{\text{Rx}}$  channels. 154 In order to separate the transmit signals at the receiver side, the signals from different Tx 155 antennas should be orthogonal. There are multiple ways to realize waveform orthogonal, 156 such as time-division multiplexing (TDM), frequency-division multiplexing (FDM), and 157 Doppler-division multiplexing (DDM) [12,13]. TDM is widely used for its simplicity. In 158 this mode, different Tx antennas transmit chirp signals by turns, as shown in fig. 4 (b). 159 Therefore, at the receiver side, different Tx waveforms can be easily separated in the time 160 domain. An additional phase shift compensation [14] is required to compensate for the 161 motion of detections during the Tx switching time. Another shortcoming of TDM is the 162 reduced detection range due to loss of transmitting power. DDM is also supported by many 163 radar devices. As shown in fig. 4 (c), DDM transmits all Tx waveform simultaneously and 164 separate them in Doppler domain. In order to realize waveform orthogonality, for the the 165 k-th transmitter, a Doppler shift is added to adjacent chirps as 166

$$\nu_k = \frac{2\pi(k-1)}{N} \tag{5}$$

where N is usually selected as the number of Tx antennas  $N_{\text{Tx}}$ . One drawback of DDM is its unambiguous Doppler velocity is reduced to  $\frac{1}{N}$  of the original one. Empty-168 band DDM [15] can achieve more robust velocity disambiguation by introducing several empty Doppler subbands. Some example codes are provided in RADIal dataset [16]. After 170 decoupling the received signals, we can obtain a 3D tensor by stacking RD maps with 171 respect to Tx-Rx pairs. Then, DOA can be estimated through angle FFT along the virtual 172 receiver dimension. Some super-resolution methods, such as MUSIC [17], can be applied 173 to improve angular resolution. The resulting 3D tensor is referred to as Range-Azimuth-174 Doppler (RAD) tensor or radar tensor. 175

6

In radar detection pipeline, RD maps are integrated coherently along the virtual 176 receiver dimension to increase SNR. Then, a Constant False Alarm Rate (CFAR) detector [17] 177 is applied to detect peaks in the RD map. Finally, the DOA estimation method is applied 178 for angle estimation. The output is a point cloud with measurements of range, Doppler 179 and angle. For conventional radars, only azimuth angle is resolved, while 4D radars 180 output both azimuth and elevation angles. Since radar is usually used in safety-critical 181 applications, a lower CFAR threshold ( $\leq 10$  dB) is set to achieve high recall. The accuracy 182 of detection is affected by road clutter, interferences and multi-path effects in complex 183 environments. Therefore, additional spatial-temporal filtering is required to improve accuracy. DBSCAN [18] is used to cluster radar detections into object-level targets. Clusters 185 with few detections are considered as outliers and thus be removed. Further, temporal filtering, such as Kalman filtering, is used to filter out outliers and interpolate missed 187 detections.

#### 2.2. Radar Performances

Performance of automotive radar can be evaluated in terms of maximum range, 190 maximum Doppler velocity and field of view (FoV). Equations for these attributes are 191

summarized in table 2. According to radar equation, the theoretical maximum detection <sup>192</sup> range is given by <sup>193</sup>

$$R_{\rm max} = \sqrt[4]{\frac{P_t G^2 \lambda^2 \sigma}{(4\pi)^3 P_{\rm min}}} \tag{6}$$

where  $P_t$  is the transmit power,  $P_{min}$  is the minimum detectable signal or receiver sensitivity,  $\lambda$  is the transmit wavelength,  $\sigma$  is the target RCS and G is the antenna gain. The 195 wavelength is 3.9 mm for automotive 77Ghz radar. Target RCS is a measure of ability to 196 reflect radar signals back to the radar receiver. It is a statistical quantity that varies with the 197 viewing angle and the target material. According to the test results [19], smaller objects, such as pedestrian and bike, have a average RCS value of around 2-3 dBsm, whereas 199 normal vehicles have an average RCS value of around 10 dBsm and large vehicles of around 20 dBsm. The other parameters, such as transmit power, minimum detectable 201 signal and antenna gain are design parameters aimed at meeting product requirements as 202 well as regulations. Some typical values for these parameters are summarized in table 3. 203 In practice, the maximum range is limited by the supported IF bandwidth  $B_{\rm IF}$  and ADC 204 sampling frequency. The maximum unambiguous velocity is inversely proportional to the 205 chirp duration  $T_c$ . For MIMO radar, the maximum unambiguous angle is dependent on the 206 spacing of antennas d. The theoretical maximum FoV is 180° if  $d = \lambda/2$ . In practice, FoV is 207 determined by the antenna gain pattern. Another important characteristic is resolution, *i.e.*, the ability to separate two close targets with respect to range, velocity and angle. As 209 shown in table 2, high range resolution requires a large sweep bandwidth B. High Doppler 210 resolution requires a long integration time, *i.e.*, the frame time  $N_cT_c$ . The angular resolution 211 depends on the number of virtual receivers  $N_R$ , the object angle  $\theta$  and the inter-antenna 212 spacing *d*. For the case of  $d = \lambda/2$  and  $\theta = 0^\circ$ , angular resolution is in a simple form of 213  $2/N_R$ . From the perspective of antenna theory, angular resolution can also be featured by 214 the half-power beamwidth, *i.e.*, the 3-dB beamwidth [13], which is a function of the array 215 aperture *D*. 216

Table 2. Equations for radar performance

Definition	Equation
Max Unambiguous Range	$R_m = \frac{cB_{ m F}}{2S}$
Max Unambiguous Velocity	$v_m = rac{\lambda}{4T_c}$
Max Unambiguous Angle	$ heta_{ m FoV} = \pm \arcsin(rac{\lambda}{2d})$
Range Resolution	$\Delta R = \frac{c}{2B}$
Velocity Resolution	$\Delta v = rac{\lambda}{2N_cT_c}$
Angular Resolution	$\Delta  heta_{ m res} = rac{\lambda}{N_R dcos( heta)}$
3-dB beamwidth	$\Delta \theta_{3dB} = 2 \arcsin \frac{1.4\lambda}{\pi D}$

The meaning of parameters are consistent in this section. Refer to appendix A for a quick check the meaning.

 Table 3. Typical automotive radar parameters[20]

Parameter	Range
Transit power (dBm)	10–13
TX/RX antenna gain (dBi)	10–25
Receiver noise figure (dB)	10–20
Target RCS (dBsm)	(-10)-20
Receiver sensitivity (dBm)	(-120) - (-115)
Minimum SNR (dB)	10-20

In practice, different types of automotive radar are designed for different scenarios. <sup>217</sup> Long Range Radar (LRR) achieves long detection range and high angular resolution at the <sup>218</sup>

229

245

254

cost of a smaller FoV. Short Range Radar (SRR) uses MIMO techniques to achieve high 219 angular resolution and large FoV. In addition, different chirp configurations [21] are used 220 for different applications. For example, Long Range Radar needs to detect fast-moving 221 vehicles at distances, and therefore utilizes small ramp slope for long distance detection, 222 long chirp integration time to increase SNR, small chirp duration to increase maximum 223 velocity and short chirp duration for high velocity resolution [22]. Short Range Radar 224 needs to detect vulnerable road users (VRUs) close to the vehicle, and therefore utilizes 225 higher sweep bandwidth for high range resolution at the cost of short range. Multi-mode 226 radar [21] can work in different modes simultaneously by sending chirps that are switched 227 sequentially with different configurations. 228

#### 2.3. Open-Source Radar Toolbox

Commercial off-the-shelf radar products can only output point clouds. It can be 230 configured to output either raw point clouds, sometimes referred to as radar detections, 231 or clustered objects wit tracked ids. The signal processing algorithm inside it is a black 232 box and cannot be modified. Alternatively, TI mm-wave radars have been widely used in 233 academic research because of their public nature and flexibility. They support configurable 234 chirps [21] and different MIMO modes [11] to adapt to different tasks. TI also provides a mmWave studio which provides GUIs for radar setup, data capturing, signal processing and 236 visualization. In addition, there are some open-source radar signal processing toolboxes 237 for TI devices, for example, RaDICaL SDK [23,24], PyRapid [25], OpenRadar [26] and 238 Pymmw [27]. These toolboxes enable researchers to build their own datasets using TI 239 devices. While there are a growing number of public radar datasets, most of them provide 240 limited information about the radar configurations they use. This makes it difficult to 241 make a fair comparison between algorithms trained on different datasets. Open radar 242 initiative [28] provides a guideline for radar configuration and encourages researchers to 243 expand this dataset by using the radar device with the same configuration. 244

# 3. Datasets, Labelling and Augmentation

Data plays a key role in the learning-based approaches. In the past, radar algorithms 246 were always evaluated on private datasets. Recently, with the trend towards open source, 247 many radar datasets have become publicly available. In this section, we summarize these 248 radar datasets with respect to their data representations, tasks, scenarios, and annotation 249 types, as shown in table 4. To motivate readers to build their own datasets, we also introduce 250 extrinsic calibration and cross-modality labelling techniques. We further investigate data 251 augmentation methods and the potential use of synthetic radar data to improve data 252 diversity. 253

#### 3.1. Radar Datasets

Different radar datasets use different types of radar. We can classify radar sensors into Low Resolution (LR) and High Resolution (HR). There are different technical routes to achieve high resolution, such as polarimetric radar [29], cooperative radars [30], multi-chip cascaded MIMO radar [13], synthetic aperture radar (SAR) [31] and spinning radar [32]. Most off-the-shelf radars can output a point cloud with range, azimuth angle, Doppler velocity and RCS. Next-generation 4D radar can also measure elevation angle. Some radar prototypes can be configured to output radar raw data, including ADC data, RA/RD maps and RAD tensors.

The role of radar in autonomous driving can be divided into localization and detection. 263 Although this paper focus on radar detection, we also introduce the localization datasets in 264 this section. Since these datasets usually provide synchronized LiDAR and image along 265 with radar data, it is possible to annotate them for detection purpose, as done in [33]. There 266 are various levels of label granularity for radar data. For radar point cloud, it is possible to 267 provide 2D bounding boxes, 3D bounding boxes or point-wise annotations. 2D bounding 268 boxes are labelled in bird's-eye view (BEV) and with orientation information, hence they 269 are sometimes referred to as pseudo-3D boxes. 3D bounding boxes further capture height 270 information and pitch angle. If properly annotated, point-wise annotation can provide 271 semantic information at a finer granularity than bounding boxes. In fact, radar detections 272 within the bounding box could also be ghost detection or clutter. Therefore, point-wise 273 annotation is a better way to capture the noisy nature of radar point cloud. Similarly, radar 274 pre-CFAR data, including RA/RD maps and RAD tensors, are also annotated point-wise. 275 Some works dilate the annotated points into a dense mask or a bounding box. However, 276 these dilated patches do not necessarily reflect the shape information. Some techniques for 277 precise dilation will be introduced later in section 4. 278

There are some large scale datasets for autonomous driving which include off-the-shelf 279 2D radars in their sensor suites. NuScenes [34] is the most popular dataset for its large-scale 280 and diverse scenarios. The capturing vehicle is equipped with a 32 beam LiDAR, 6 cameras, 281 5 long-range multi-mode radars and a GPS/IMU system. It provides 3D annotations of 282 23 classes of road users in 1,000 scenes, with a total of 1.3 million frames. However, this dataset is not a good choice for studying the role of radar in perception, because its radar 284 point clouds are too sparse. PixSet dataset [35] also aims at 3D object detection. The vehicle is equipped with a colocated sensor platform consisting of a solid-state LiDAR, a 64-beam 286 LiDAR, a TI AWR1843 radar and a GPS/IMU system. The FoVs of different modalities are largely overlapped and hence are well suited for evaluating sensor fusion algorithms. 288 RadarScenes dataset [36] is a diverse large-scale dataset for instance segmentation of radar point clouds. It uses four 77GHz radars with overlapping FoV in the front of the vehicle. 290 Each radar is in middle-range mode with maximum range of 100m and 60°FoV. Compared to NuScenes dataset, its radar point clouds are much denser. The datasets contains 100km 292 of driving in 158 different scenarios. It provides both point-wise annotations and track 293 IDs for 11 classes of moving road users. All points with zero velocity are labelled as static. 294 Pointillism [37] leverages a multi-radar setup to improve the resolution. Two TI IWR1443 295 radars were placed at the front of the car, facing forward, at a distance of 1.5 meters. The 296 aim is to study the effect of coherently integrate point clouds from two radar sensors. In 297 order provide ground truth of radar point clouds, the sensor suite also include a 16-beam 298 LiDAR and a camera with overlapping FoV. The dataset contains 54K synchronized frames 200 for five typical driving scenarios under different weather conditions. It also provides 3D box annotations of vehicles. Zendar dataset [38] is a high-resolution radar dataset that uses 301 SAR for moving vehicle detection. It provides time-synchronized image, radar ADC data, 302 2D SAR point cloud and projected LiDAR point cloud in BEV. Point-wise annotations of 303 moving vehicles are applied to the SAR point cloud. It also provides an SDK for converting raw ADC data to RD maps and visualization. 305

Several datasets utilize radar sensors in short-range (SR) or ultra-short-range (USR) 306 mode for high resolution close-field imaging. In this mode, close objects will occupy several 307 cells in both range and Doppler dimension (because of the micro-Doppler motion). To fully utilize these spatially spread range and Doppler signatures, annotations are made directly 309 on RA maps or RAD tensors. CARRADA dataset [39] uses TI AWR1843BOOST radar in 310 short-range mode, with a max distance of 50 meters. It provides real-valued RA maps, 311 RD maps and unannotated RAD tensors, as well as synchronized images, for training 312 neural networks. In both RA and RD maps, objects are annotated in point level with a 313 category from pedestrian, car, or cyclist. In addition, the dilated segmentation mask and 314 the bounding box around the cluster are also provided. The data are collected on an empty 315 test track with at most two moving objects in the FoV. RADDet dataset [40] also uses TI 316 AWR1843BOOST radar with a max distance of 50 meters, as well as a stereo camera. It 317 provides 3D bounding boxes for complex-valued RAD tensor and 2D bounding boxes for 318 RA map projected in Cartesian view. The data are captured using a tripod located on the 319 sidewalks and facing the main roads. Therefore, its scenario is much more complex than 320 the CARRADA dataset. CRUW dataset [41] uses a TI AWR1843 radar and a stereo camera for object detection. It adopts a different signal processing pipeline, which directly outputs 322 RA map using range FFT and angle FFT. Then, the object-level point-wise annotation 323

is applied to complex-valued RA maps. A probabilistic camera-radar fusion approach 324 is used to improve annotation quality. The dataset contains 3.5 hours, 400K frames of 325 camera-radar data in different driving scenarios, including parking lot, campus road, city 326 street, and highway. RaDICaL dataset [42] uses TI IWR1443BOOST radar in multiple 327 configurations for different scenarios, including indoor, parking lot, highway and single 328 human walking. It records radar ADC data together with the RGB-D images and IMU data 329 using ROS. It also provides a signal processing SDK [23] to process and annotate radar data. 330 Ghent VRU dataset [43] collects radar data specifically for VRU detection. The sensor suite 331 includes a TI AWR1243 radar, a camera and a 16-beam LiDAR. The data are recorded by a 332 vehicle driving on public roads in a crowded European city center. It provides radar RAD 333 tensors with segmentation mask annotations for VRUs. To compensate for range-dependent 334 power, many datasets apply logarithmic scaling or normalization to the pre-CFAR data 335 as default. CARRADA dataset [39] and RADDet dataset [40] apply logarithmic scaling 336 to their radar data. The normalization can be applied in different ways, including local 337 power normalizing in Ghent VRU dataset [43], min-max scaling in CRUW dataset [41], and 338 Z-score standardization in RADDet dataset [40]. Here we only summarize some operations that are explicitly mentioned. Further checks are needed when benchmarking the algorithm 340 using different datasets.

As 4D radar is just entering the market, only a few public datasets are available. 342 Astyx dataset [2] is the first public-available 4D radar dataset. The sensors include a 16-343 beam LiDAR, a camera and a Astyx 6455 HiRes 4D radar. It provides 3D bounding box 344 annotations of seven classes of road users. Each object is also featured with four levels of occlusion and three levels of uncertainty. The dataset is very small, with only 500 annotated 346 frames of short clips, each clip containing less than 10 frames. The class distribution is very 347 imbalanced, with over 90% annotated objects are car. View-of-Delft (VoD) dataset [44] is a 348 recently published 4D radar dataset especially focus on detection of VRUs. The sensor suite includes a ZF FRGen 21 4D radar, a 64-beam LiDAR and a stereo camera. It provides 8693 350 annotated frames with 3D bounding boxes and tracking ids. Each object is also annotated 351 with two levels of occlusion and four types of activity attribute (stopped, moving, parked, 352 pushed, sitting). The data is collected in campus, suburb and old-town locations, with a 353 preference on scenarios containing VRUs. It provides fine-grained annotations of vehicle, 354 trucks and 10 classes of VRUs. Different classes are equally distributed (21.6% pedestrians, 355 8.8% cyclists and 21.9% cars). RADIal dataset [45] is a 4D radar dataset for vehicle detection 356 and open space segmentation. The sensors include camera, LiDAR, 4D radar, GPS and 357 vehicle's CAN traces. The 4D radar is a 12Tx 16Rx cascaded radar. A key feature is that they also provide radar raw ADC data, which makes it possible to explore the potential of 350 neural networks in the signal processing stage. This dataset is comparable in size to the VoD dataset, with annotated 8252 frames captured in city street, highway and countryside 361 road. Two kinds of annotation are provided, including vehicle annotations and open space segmentation mask in BEV. The vehicle annotations are in the format of 2D bounding boxes 363 for image and object-level points for LiDAR and radar. Although they do not provide 364 bounding box annotations for the radar point clouds, it is possible for the researchers to 365 annotate them on their own, given the LiDAR point clouds and images. TJ4DRadSet [46] is 4D dataset for 3D detection and tracking. The sensor suite includes a 32 beam LiDAR, a 367 camera, a high performance 4D radar (Oculii Eagle) and GNSS. By utilizing Oculii's virtual 368 aperture imaging technique, this 4D radar can output a much denser point cloud than 369 others. It has a maximum detection range of 400 m and an angular resolution of less than 370 1° in both azimuth and elevation. The data are captured in a wide range of road conditions 371 in urban driving. The dataset contains a total of 40K frames of synchronized data, where 372 7757 frames of them are annotated with 3D bounding boxes and track id. 373

Radar can also be utilized for localization. Compared to camera and LiDAR, radar has advantages of long detection range and robustness to occlusions. Millimeter wave can penetrate certain non-metallic objects, such as glass, Polywood and clay bricks [32], and is less affected by dust, smoke, fog, rain, snow, and ambient lighting conditions [32].

Therefore, radar has great potential for mapping and localization in adverse weathers. The Oxford radar robotcar dataset [47] is the most popular dataset for radar SLAM. The test car 379 is equipped with a rich set of sensors, including an FMCW spinning radar, two 32 beam Li-380 DARs, a stereo camera, three monocular cameras, two 2D LiDARs and a GPS/IMU system. 381 The spinning radar can provide a 360° high-resolution intensity map of surrounding envi-382 ronments. However, it has no Doppler information and is rarely used in production cars 383 due to the high price. Mulran dataset [48] focuses on range-sensor-based place recognition. 384 It uses a spinning radar and a 64 beam LiDAR to capture the surrounding environment. 385 The recorded data are temporally (monthly revisits) and structurally (multi-city) diverse. 386 Borea dataset [49] aims at studying the effect of seasonal variation on long-term localization. 387 The sensor suite includes a spinning radar, a camera, a GPS/IMU system and a 128-beam 388 LiDAR. The data was collected by driving a repeated route over one year, thus capturing 389 seasonal variations and adverse weather conditions. It provides pose ground truth for local-390 ization task, as well as 3D bounding box annotations for object detection in sunny weather. 391 Similar to Borea dataset, EU long-term dataset [50] aims at localization in highly dynamic 392 environments and long-term autonomy. Its sensor suite includes two stereo cameras, two 393 32 beam LiDARs, two fisheye cameras, a four beam LiDAR and a 77Ghz long-range radar 394 and a 2D LiDAR facing the road. Endeavour dataset [51] adopts 5 multi-mode radar to cover the 360° surrounding environment. It is also equipped with LiDARs and RTK-GPS 396 to provide ground truth for radar odometry. ColoRadar dataset [52] utilizes a compact moving sensor rig, which consists of a 64 beam LiDAR, a TI AWR2243 cascaded 4D radar, 398 a TI AWR1843 radar and an IMU. Three levels of radar data representation are provided, including raw ADC samples, Range-Azimuth-Elevation-Doppler (RAED) tensors from the 400 4D radar and point clouds from the single-chip radar. The data are gathered in a variety of 401 scenarios, including highly diverse indoor environments, outdoor environments and an 402 underground mine. 403

There are also some radar datasets designed for specific tasks. PREVENTION [53] 404 focuses on predicting inter-vehicle interactions. The data collection car is equipped with 405 one frontal long-range radar, two corner radars, one 32 beam LiDAR and two cameras. 406 It provides annotations of 2D bounding boxes, lane change behaviours, and trajectories. 407 SCORP [54] is a radar dataset for open space segmentation in parking scenarios. It provides 408 three kinds of radar data, including RAD tensor, RA map and BEV map. Radar Ghost 409 dataset [55] aims at studying the effect of multi-path propagation in autonomous driving. 410 It provides point-wise annotations of real targets and four types of ghost targets. 411

Robust perception under adverse weather is a popular research topic for safe autonomous driving. Although there are some recently published datasets for adverse 413 weather [56–59], only a few include radar in their sensor suite. Dense dataset [60] focus on evaluating multi-modal fusion algorithms under adverse weather. In addition to LiDAR 415 and stereo camera, it is also equipped with several all-weather sensors, including one frontal long-range radar, one gated camera working on NIR band, one FIR camera and 417 one weather station sensor. The data are captured in various natural weather conditions, 418 including rain, snow, light fog and heavy fog, as well as in a controlled lab environment 419 in a fog chamber. However, the dataset only provides sparse radar targets with limited 420 FoV and poor resolution. RADIATE dataset [61] particularly focuses on leveraging radar 421 in adverse weather. The data collection car is equipped with a camera, a LiDAR and a 422 spinning radar. The datasets are captured under different weathers, such as sun, night, rain, 423 fog and snow. It provides annotations for 2D object detection, object tracking and SLAM. 424

lable 4. Kadar Datasets	Table	4.	Radar	Datasets
-------------------------	-------	----	-------	----------

Name	Year	Task	Radar Type	Data	Doppler	Range	Modalities	Scenarios	Weather	Annotations	Size
				Off-the-sh	elf Radar Data	sets for Detect	tion				
nuScenes [34]	2020	DT	LR	PC	1	SV	CLO	USH	1	3D	L
Dense [60]	2020	D	LR	PC	$\checkmark$	LR	CLO	USHT	1	3D	L
PixSet [35]	2021	DT	LR	PC	$\checkmark$	MR	CLO	USP	1	3D	Μ
RadarScenes [36]	2021	DTS	HR	PC	$\checkmark$	SV	CO	USHT	1	$\mathbf{P}_w$	L
Pointillism [37]	2020	D	2LR	PC	1	MR	CL	U	1	3D	Μ
Zendar [38]	2020	D	SAR	ADC,RD,PC	1	MR	CLO	U	×	$\mathbf{P}_w$	S
				Radar Pr	e-CFAR Datas	ets for Detection	on				
CARRADA [39]	2020	DTS	LR	RAD	1	SR	С	R	×	2D,P <sub>w</sub> ,M	М
CRUW [41]	2021	D	LR	RAD	$\checkmark$	USR	С	USHP	×	$P_o$	L
RADDet [40]	2021	D	LR	RAD	$\checkmark$	SR	С	US	×	2D	Μ
RaDICaL [42]	2021	L	LR	ADC	$\checkmark$	USR,SR	$C_dO$	USHIP	×	2D	L
Ghent VRU [43]	2020	DS	LR	RAD	1	SR	CL	U	×	М	Μ
				4D Ra	adar Datasets	for Detection					
Astyx [2]	2019	D	HR	PC	1	MR	CL	SH	X	3D	S
View-of-Delft [44]	2022	DT	HR	PC	$\checkmark$	SR	CLO	U	×	3D,T	S
RADIal [45]	2021	DS	HR	ADC,RAD,PC	$\checkmark$	MR	CLO	USH	×	P <sub>o</sub> ,M	Μ
TJ4DRadSet [46]	2022	DT	HR	PC	1	LR	CLO	U	×	3D, T	Μ
				Rada	r Datasets for	Localization					
ColoRadar [52]	2021	L	HR,LR	ADC,PC	✓	2USR	LO	SIT	✓	$P_s$	М
Oxford [47]	2020	L	SP	RA	×	SV	CLO	U	1	$P_s$	L
RADIATE [61]	2020	LDT	SP	RA	×	SV	CLO	USHP	1	2D,T, P <sub>s</sub>	Μ
Mulran [48]	2020	L	SP	RA	×	SV	LO	US	×	$P_s$	Μ
Boreas [49]	2022	LD	SP	RA	×	SV	CLO	S	1	$P_s$ ,3D	L
Endeavour [51]	2021	L	LR	PC	$\checkmark$	5LR	LO	S	×	$P_s$	Μ
EU Long-term [50]	2020	L	LR	PC	$\checkmark$	LR	CLO	U	1	$P_s$	Μ
				Rada	r Datasets for	Other Tasks					
Ghost [55]	2021	DS	LR	PC	×	LR	CLO	S	×	$P_w$	М
SCORP [54]	2020	S	LR	ADC,RAD	$\checkmark$	USR	С	Р	×	Μ	S
PREVENTION [53]	2019	DT	LR	PC	✓	1LR,2SR	CLO	UH	1	2D	L

Task: 'D', 'T', 'L', 'S' stand for detection, tracking, localization, and segmentation; Type: 'LR', 'HR', 'SP', 'SAR' stand for low-resolution, high-resolution, spinning, and SAR; Range: 'SV', 'LR', 'MR', 'SR', 'USR' stand for surrounding view, long-range (<250m), middle-range (<180m), short-range (<50m), and ultra-short-range (<25m); Modalities: 'C', '2', 'C', 'L', 'O' stand for camera, RGBD camera, LiDAR, and odometry; Scenarios: 'U', 'S', 'H', 'P', 'T', 'R', 'I' stand for urban (city), suburban, highway, parking lot, tunnel, race track and indoor; Size: 'L', 'M', 'S' stand for large, medium, and small; Weather stands for adverse weather; Label: '2D','3D','T', 'Pw','P\_0','Ps', 'M' stand for 2D bounding box, 3D bounding box, track ID, point-wise detection, object-level point, pose and segmentation mask.

# 3.2. Extrinsic Calibration

Multi-sensor extrinsic calibration requires calibration targets to be observed simul-431 taneously by different modalities. The trihedral corner reflector is widely used for radar 432 calibration because of its high RCS. Multiple reflectors are usually placed outdoor to avoid 433 multi-path propagation. The difficulty lies in making the calibration target visible to both 434 radar and other sensors. El Natour *et al.* [62] builds a calibration facility by placing one 435 Luneburg lens and seven trihedral corner reflectors with known inter-distances. To make 436 the reflectors visually detectable, they paint each surface with different colors. Peršić et al. 437 [63] design a compact calibration target which can be simultaneously detected by camera, 438 radar and LiDAR. As shown in fig. 5 (a), they place a triangle-shaped checkerboard pattern 439 in front of a trihedral corner reflector. The checkerboard is made of styrofoam and is 440 transparent over a large radio frequency range. So that the millimeter wave can penetrate it 441 and detect the corner reflector behind it. In fig. 5 (b), Domhof et al. [64] design a styrofoam 442 board with four circular holes and place a corner reflector at the back. These circular holes 443 are more easily detected by the sparse LiDAR beam since they have no horizontal lines. A A A



(a) Triangle target

(b) Target with holes

**Figure 5.** Two types of radar calibration targets [63,64]. The front board is made of styrofoam. The red triangle is a radar corner reflector.

The extrinsic calibration of a 4D radar and other sensors can be easily done by modi-445 fying the classical LiDAR to camera calibration methods [65,66]. However, the calibration 446 of conventional radar is a very difficult task, since it returns a 2D point cloud with no 447 elevation resolution. This leads to the problem of vertical misalignment [63], which is 448 defined as the angular deviation between radar plane and ground plane. Sugimoto et al. 449 [67] move a corner reflector up and down to cross radar plane in multiple times. Then, 450 the plane is determined by connecting the peaks with highest intensity in the sequence. 451 Peršić et al. [63] propose a two-step optimization method to mitigate the uncertainty caused 452 by the missing elevation angle. They model radar detections as arcs by extending their 453 elevation angle. Similarly, they also convert 3D detections from other sensors to arcs by 454 neglecting the elevation angle. In the first step, they optimize the reprojection error, which 455 is the Euclidean distance of these projected arcs on the ground plane. In the second step, the 456 parameters related to the elevation measurement are refined according to the RCS error. A 457 second-order RCS model is built by fitting RCS measurements with elevation angles. Then, 458 the L2 distance between the expected and measured RCS is minimized. Experiments show 450 their method enables smaller vertical misalignment than Sugimoto's method. In order to improve efficiency, some targetless online calibration approaches [68,69] are proposed to 461 leverage target trajectories for extrinsic estimation. 462

# 3.3. Data Labelling

Before introducing the labelling process, we first discuss the time synchronization problem. Different sensors can be synchronised using pulse per second (PPS) triggering signals from the GNSS receiver [70]. However, in most of radar datasets, sensors differ in their triggering time and sampling frequency. Some of them select one sensor as the lead, and choose the closest frames from other modalities for synchronization. Assuming a tolerable time offset of 50 ms, a vehicle with a relative speed of 20 m/s will lead to

430

an offset of 1 m. Therefore, it is necessary to compensate for synchronization errors in high-speed scenarios. Kaul *et al.* [71] design a pose chain method to interpolate interframe measurements. The translational and rotational transformations are determined by a constant velocity model and Spherical Linear Interpolation (SLERP) [72] respectively.

Labelling radar data is a difficult task. Both radar point clouds and pre-CFAR data are hard to interpret by human labellers. To reduce the labelling efforts, most of datasets adopt a semi-automatic labelling framework, which includes two steps: cross-modality pre-labelling and fine-tuning.

In the first step, a well-trained detector on other modalities is leveraged for radar 478 labelling. For 3D tasks, radar point clouds can be annotated by 3D boxes predicted by 479 the detector trained with images and LiDAR point clouds [2]. If we want to get point-480 wise annotations for radar point cloud, we can first predict a dense semantic map for the 481 corresponding image using a visual segmentation network, such as mask R-CNN [73] or 482 DeepLab V3 [74]. To avoid scale ambiguity, it is better to project the masked image to radar frames using depth measured by LiDAR [71] or stereo camera [40,41]. Then, each radar 484 point can be associated with the corresponding semantic labels. Point-wise annotation of RAD tensors or RA maps is in a similar process. We can firstly use CFAR to detect 486 peaks as detections, then annotate these point detections with the aligned visual semantic map. CRUW dataset [41] propose a post-processing method to obtain point-wise object 488 annotations for RA maps. The authors define an object location similarity (OLS) metric which jointly considers the similarities in distance, scale and class. Then, they propose a 490 location-based Non-Maximum Suppression (NMS) method that selects one object point 491 out of the adjacent points based on the OLS metric. Compared to RA map, RD map alone 492 is much more difficult to label. It needs both depth and radial velocity to associate an RD 493 cell to a pixel or a LiDAR point. Radial velocity can be estimated by visual scene flow [75] 494 or by tracking [39]. 495

In the second step, manual inspection is required to correct pre-labelling errors. Identi-496 fying radar errors involves domain knowledge and therefore requires hiring of radar ex-497 perts. As a result, building a high-quality, large-scale radar dataset is both time-consuming 498 and financially expensive. To improve labelling efficiency, one way is to reduce the amount 499 of data to be labelled. Dimitrievski *et al.* [76] leverage a tracking algorithm to interpolate 500 annotations between key frames. The intermediate position is estimated by a Kalman filter 501 with optical flow as measurements. Meyer *et al.* [2] adopt an active learning [77] framework 502 to reduce labelling efforts in building the Astyx dataset. The core idea is to only label the 503 most informative data. Specifically, they first label a small number of frames and train a detector with this data subset. The trained detector is then used to make predictions on the 505 remaining unlabelled data. Next, the top N uncertain data are again manually labelled and added to the training subset. This process is repeated until convergence of the validation 507 performance.

#### 3.4. Data Augmentation

Data augmentation plays an essential role in improving generalization of deep learn-510 ing models. It is well studied for image [78], LiDAR point cloud [79] and audio spec-511 trogram [80], but overlooked in radar perception. According to the summary report of 512 Radar Object Detection 2021 (ROD2021) Challenge [81], data augmentation techniques 513 significantly improve the performance of RA-map-based radar detection. Considering the 514 radar representation, we can divide the augmentation techniques into spectral and point 515 cloud-based. Augmentation methods can also be featured by local or global depending on 516 whether the entity being augmented is a single object or the entire scene. 517

Spectral augmentation techniques are used for radar pre-CFAR data. DANet [82] adopts several global augmentation techniques borrowed from computer vision to radar RA maps. The methods include mirroring, resizing, random combination, adding Gaussian noise, and temporal reversing. Although physical fidelity is not explicitly considered, the performance gain proves the effectiveness of these augmentation methods. RADIO [83]

implements four types of spectral augmentations, including attenuation, resolution change, 523 adding speckle noise, and background shift. The first two methods are applied to a local 524 patch around each detected object. The attenuation effect is approximated by dampening 525 the cells according to an empirical relationship between the received power and range. 526 The resolution change is modelled by nearest-neighbor interpolation according to the 527 object size. The speckle noise can be approximated as multiplicative truncated exponential 528 distribution [84] or multiplicative Gaussian noise [83]. Background shift is done by adding 529 or subtracting a constant value to background cells. RAMP-CNN [85] applies global 530 geometric augmentations to RA maps. They translate and rotate RA maps in Cartesian 531 coordinate, then project back to the original polar coordinate. The out-of-boundaries 532 areas are cropped off, and the blank areas are filled with background noises. Energy loss 533 and antenna gain loss due to the transformation are compensated according to the radar 534 equation. 535



**Figure 6.** Radar data augmentation techniques. The Doppler velocity measured by radar is a scalar, so local rotation of the radar detection will cause a misalignment between the Doppler velocity and the true velocity. Global translation and rotation are free from such misalignment. When augmenting radar RA map, it is necessary to interpolate the background area and compensate the intensity of detection.

Point cloud augmentation aims to introduce invariance to geometric transformations 536 and improve the signal-to-clutter ratio. Compared with spectral augmentation, point cloud 537 augmentation methods can be easily extended to multi-modality by properly handling 538 occlusion issues [86,87]. Geometric augmentation can be applied locally or globally, de-539 pending on whether the transformation is applied to a single target or the entire scene. For 540 radar point cloud, the Doppler velocity and RCS need further consideration. As illustrated 541 in fig. 6, rotating objects locally will affect the radial velocity, and rotating the radar point cloud globally will affect the ego-motion compensated radial velocity if the ego-motion is 543 not rotated accordingly. Therefore, Palffy et al. [44] advise only use mirroring and scaling along the longitudinal axis as augmentation. Another applicable technique is the copy-545 paste augmentation, which copies the detected object from other frames and pastes it into 546 the same location in the current frame as done in [88]. A limitation of these two methods 547 is that they do not change the distribution of detections, while radar points are actually 548 randomly distributed over the object in different frames. According to experiments [89], 549 most of the radar detections are located in the proximity of the vehicle contour and wheel 550 rims. The number of detections per object is inversely proportional to the distance, and 551 the probability of detection on the contour depends heavily on the orientation. Simulation-552 based methods, which will be introduced in the next section, is more suitable to capture 553 such randomness. 554

To handle with the sparsity issue, many works utilize augmentation to increase point cloud density. One simple method is accumulating radar points from multiple frames into the current frame. However, accumulation without motion compensation will lead to point cloud aliasing. Long *et al.* [90] compensates the accumulated radar point cloud with the estimated full velocity, achieving better performance in bounding box regression. Plaffy et 550 al. [44] augment the accumulated radar point cloud by appending a temporal index to each 560 point as an additional channel. Along with the increased density, this index augmentation 561 is expected to effectively retain the temporal information. Alternatively, Bansal et al. [37] 562 leverages space coherent of two radar sensors to increase the point cloud density. They 563 fuse point clouds from two radars with overlapping FoVs in a probabilistic manner. They 564 firstly cluster the raw point clouds and then associate clusters from two radars by defining 565 a distance-dependent potential function. Points with low confidence are filtered out as 566 outliers, and the remaining points within the same cluster are coherently accumulated. 567

#### 3.5. Synthetic Data

Synthetic datasets are widely used in computer vision [91,92] and LiDAR percep-569 tion [93,94] for autonomous driving. Experiments [95] show the network trained with 570 synthetic data can generalize well in the real-world. By using synthetic radar data, the 571 labelling cost can be completely avoided. Moreover, simulation can be used to generate the 572 safety-critical long-tail scenarios [96]. Physics-based simulation methods, such as ray trac-573 ing [97,98], are widely applied to generate synthetic radar point clouds. Experiments [98] 574 show that ray tracing can successfully model the multi-path propagation and separability issue of close objects. However, it is difficult to capture the RCS variation in azimuth with 576 current methods. Another type of simulation is to build a probabilistic model of radar 577 detections, also known as model-based augmentation. The spatial distribution of radar 578 detections over the vehicle can be approximated by the surface-volume model, including 579 volcanormal measurement model [99], variational Gaussian mixture model (GMM) [99], 580 and hierarchical truncated Gaussian (HTG) [100]. Model parameters can be learned from 581 data. Using this model, we can augment new synthetic radar detections to real point 582 clouds. It is arguable that what level of fidelity is necessary for downstream tasks. In 583 [101], model-based and ray tracing methods are compared with respect to multiple tar-584 get tracking. Experiments indicate that the ray-tracing-based model achieves the lowest 585 simulation-to-reality gap. 586

There are some seminal works utilizing learning-based generative models for radar 587 simulation. For example, deep stochastic radar model [102] adopts a conditional-VAE 588 architecture. The encoder consists of two heads, one for RAD tensor and one for object list. 589 The extracted features are concatenated and further processed with an MLP. The decoder 590 generates a radar intensity map in polar grid conditioned on the encoded feature and 591 random noise. Generative models can also be used in cross-modality data generation, 592 for example, GAN-based LiDAR-to-radar generation [103], GAN-based radar-to-image 593 generation [104] and VAE-based radar-to-image generation [105]. 594

#### 4. Radar Depth and Velocity Estimation

Radar can measure range and Doppler velocity, but both of them cannot be directly used for downstream tasks. The range measurements are sparse and therefore difficult to associate with their visual correspondences. The Doppler velocity is measured in radial axis and therefore cannot be directly used for tracking. In this section, we summarize depth completion and velocity estimation methods using radar point cloud.

#### 4.1. Depth Estimation

Recently, pseudo-LiDAR-based visual object detection [106–108] has became a popular research topic. The core idea is to project pixels into a pseudo point cloud to avoid distortions induced by inverse projective mapping (IPM). The pseudo LiDAR detection is built on depth estimation. Visual depth estimation is a ill-posed problem because of the scale ambiguity. However, learning-based methods, either in supervised [109] or self-supervised [110], can successfully predict dense depth maps with camera only. Roughly speaking, these methods learn a priori knowledge of object size from the data and are therefore vulnerable to some data-related problems, such as sensitivity to input image quality [110] and learning

568

601

non-causal correlations, such as object and shadow correlations [111]. These limitations can 610 be mitigated with the help of range sensors, such as LiDAR and radar. Depth completion is 611 a sub-problem of depth estimation. It aims to recover a dense depth map for image using 612 the sparse depth measured by range sensors. Compared to LiDAR, radar has advantages of 613 low price, long range and robustness to adverse weather. Meanwhile, it faces the problems 614 of noisy detections, no height measurements and sparsity. As shown in fig. 7, due to 615 multi-path propagation, radar can see the non-line-of-sight highly reflective objects, such 616 as wheel rims and occluded vehicles. In [112], the authors refer this phenomenon as the see 617 through effect. It is beneficial in 3D coordinate, but brings difficulty in associating radar 618 detections with visual objects in image view. 619



**Figure 7.** Radar range measurements. Off-the-shelf radars return detections on a 2D radar plane. The detections are sparsely spread on objects due to specular reflection. Due to multi-path propagation, radar can see through occlusions, and meanwhile this can cause some noisy detections.

The two-stage architecture is widely applied for image guided radar depth completion 620 task. Lin et al. [113] adopt a two-stage coarse-to-fine architecture with LiDAR supervision. 621 In the first stage, a coarse radar depth is estimated by an encoder-decoder network. Radar 622 and image are processed independently by two encoders and fuse together in feature-level. 623 Then, the decoder output a coarse dense depth map in image view. The predicted depth with large errors are filtered out according to a range-dependent threshold. Next, the 625 original sensor inputs and the filtered depth map are sent to a second encoder-decoder 626 to output a fine-grained dense map. In the first stage, the quality of association can be 627 improved by expanding radar detections to better match visual objects. As shown in fig. 8 (b), Lo et al. [114] apply height extension to radar detections to compensate for the 629 missed height information. A fixed height is assumed for each detection and is projected 630 onto the image view according to the range. Then, the extended detections are sent to a 631 two-stage achitecture to output a denoised radar depth map. Long et al. [115] propose a 632 probabilistic association method to model the uncertainties of radar detections. As shown 633 in fig. 8 (c), radar points are transformed into a multi-channel enhanced radar (MER) image, 634 with each channel representing the expanded radar depth at a specific confidence level of 635 association. In this way, the occluded detections and imprecise detections at the boundary 636 are preserved but with a low confidence. Gasperini et al. [112] use radar as supervision 637 to train a monocular depth estimation model. Therefore, they apply a strict filtering to 638 only retain detections with high confidence. In the pre-processing, they remove clutters 639 inside the bounding box that exceeded the range threshold, and discard points in the upper 640 50% and outer 20% of the box as well as the overlapping regions to avoid the see-through 641 effect. All the background detections are also discarded. For association, they first apply a 642 bilateral filtering, *i.e.*, an edge preserving filtering, to constrain the expansion to be within 643 the object boundary. They further clip the association map close to the edge to get rid of 644 imprecise boundary estimations. To compensate for height information, they directly use the height of the bounding box as reference. Considering the complexity of the vehicle 646 shape, they extend the detections to the lower third of its bounding box to capture the flat 647 front surface of the vehicle. 648

675



(a) Radar detections (b) Height extension (c) Multi-channel (d) Box-based filtering

**Figure 8.** Radar detection expansion techniques. (b) Extend radar detections in height. (c) Build a probabilistic map, where the dark/light blue indicates channel with high/low confidence threshold. (d) Apply a strict filtering according to the bounding box, where only detections corresponding to the frontal surface are retained.

As ground truth, LiDAR has some inherent defects, such as sparsity, limited range and holes with no reflections. Long et al. [115] suggest to pre-process LiDAR points for better 650 supervision. They accumulate multiple frames of LiDAR point clouds to improve density. 651 Pixels with no LiDAR reaches are assigned zero values. Since LiDAR and camera do not 652 share the same FoV, the LiDAR points projected to image view also have occlusion problem. 653 Therefore, the occluded points are filtered out by two criterions, one is the difference 654 between visual optical flow and LiDAR scene flow, and the other is the difference between 655 segmentation mask and bounding boxes. Lee *et al.* [116] suggest to use both visual semantic 656 mask and LiDAR as supervision signals. Visual semantic segmentation can detect smaller 657 objects at a distance, thus compensating for the limited range of LiDAR. To extract better 658 representations, they leverage a shared decoder to learn depth estimation and semantic 659 segmentation concurrently. Both the LiDAR measurement and the visual semantic mask annotations are used as supervision. Accordingly, the loss function consists of three parts: 661 a depth loss with LiDAR points as ground truth, a visual semantic segmentation loss and a 662 semantic guided regularization term for smoothness. 663

Projecting radar to image view will lose the advantages of the see through effect. Alternatively, Niesen *et al.* [117] leverage radar RA maps for depth prediction. They 665 use a short range radar with maximum range of 40 meters. Because of the low angular resolution, the azimuth smearing effect is obvious, *i.e.*, the detections are smeared as a blurry 667 horizontal line in RA maps. It is expected that fusion of image and RA map can mitigate this effect. Therefore, they use a two branch encoder-decoder network with radar RA map 669 and image as inputs. A dense LiDAR depth map is used as ground truth. Different from the 670 above methods that align LiDAR to image, they crop, downsample and quantize LiDAR 671 detections to match radar's FoV and resolution. The proposed method is tested with their 672 self-collected data. Although the effectiveness of RA map and point cloud is not compared, 673 it provides a new direction to explore radar in depth estimation task. 674

#### 4.2. Velocity Estimation

For autonomous driving, velocity estimation is helpful for trajectory prediction and 676 path planning. Radar can accurately measure the Doppler velocity, *i.e.*, radial velocity in 677 polar coordinate. If a vehicle moves parallel to ego-vehicle at a distance, its actual velocity 678 can be approximated by the measured Doppler velocity. But this only applies in highway 679 scenarios. On urban road, it is possible for an object to move tangentially while crossing 680 the road, then its Doppler velocity will be close to zero. Therefore, Doppler velocity cannot 681 replace full velocity. Recovering full velocity from the Doppler velocity needs two steps: 682 first compensate the ego-motion, then estimate the tangential velocity. In the first step, the 683 ego-motion can be estimated by visual-inertial odometry (VIO) and GPS. Radar-inertial odometry [118,119] can also be used in visual-degraded or GPS-denied environments. 685 Then, the Doppler velocity is compensated by subtracting the ego-velocity. In the second 686 step, the full velocity is estimated according to the geometric constraints. Suppose that 687 radar observes several detections of an object and that the object is in linear motion. As

shown in fig. 9 (a), the relationship between the predicted linear velocity  $(v_x, v_y)$  and the measured Doppler velocity  $v_{r,i}$  is given by

$$v_{r,i} = v_x \cos(\theta_i) + v_y \sin(\theta_i) \tag{7}$$

where the subscript *i* denotes the *i*-th detection,  $\theta_i$  is the measured azimuth angle. By observing *N* detections per object, we can solve the linear velocity using the least square method. However, the L2 loss is not robust to outliers, such as clutter and mirco-Doppler motion of wheels. Kellner *et al.* [120] apply RANSAC to remove outliers then use Orthogonal Distance Regression (ODR) to find the optimal velocity.



#### (a) Linear motion model

(b) Curvilinear motion model

**Figure 9.** Radar motion model. (a) Linear motion model needs multiple detections for the object. (b) Curvilinear motion model requires either two radars observe the same objects, or the determination of vehicle boundary and rear axle.

Although the linear motion model is widely used for its simplicity, it will generate large position errors for motion with high curvature [121]. Alternatively, as shown in fig. 9 (b), the curvilinear motion model is given by

$$v_{r,i} = \omega(y_c - y^S)\cos(\theta_i) - \omega(x_c - x^S)\sin(\theta_i)$$
(8)

where  $\omega$  is the angular velocity,  $\theta$  is the angle of the detected point,  $(x_c, y_c)$  represents the position of the instantaneous center of rotation (ICR) and  $(x^S, y^S)$  represents the known radar position. In order to decouple angular velocity and position of the ICR, we need at least two radar sensors that observe the same object. Then, we can transform eq. (8) into a linear form as

$$y_j^{\rm S}\cos(\theta_{ji}) - x_j^{\rm S}\sin(\theta_{ji}) = y_c\cos(\theta_{ji}) - x_c\sin(\theta_{ji}) - v_{ji}^{\rm D}\omega^{-1}$$
(9)

where the subscript *j* denotes the *j*-th radar. Similarly, RANSAC and ODR can be used to find the unbiased solution of both angular velocity and position of the ICR [122]. For single radar setting, it is also possible to derive a unique solution of eq. (8) if we can correctly estimate the vehicle shape. According to the Ackermann steering geometry, the position of the ICR should be located on a line extending from the rear axle. By adding this constraints to eq. (8), the full velocity can be determined in closed form [123].

Above methods predict velocity in object level under the assumption of rigid motion. 710 However, the micro-motion of object parts, such as the swinging arms of pedestrians, are 711 also useful for classification. Capturing these non-rigid motions requires velocity estimation 712 at the point level. It can be achieved by fusing with other modalities or by using temporal 713 consistency between adjacent radar frames. Long et al. [90] estimate point-wise velocity 714 by fusion of radar and camera. They first estimate the dense global optical flow and the 715 association between radar points and image pixels through neural network models. Next, 716 they derive the closed-form full velocity based on the geometric relationship between 717 optical flow and Doppler velocity. Ding et al. [124] estimate the scene flow for 4D radar 718 point cloud in a self-supervised learning framework. Scene flow is a 3D motion field 719 and can be roughly considered as the linear velocity field. Their model consists of two 720 steps: flow estimation and static flow refinement. In the flow estimation step, they adopt a 721 similar structure with PointPWCNet [125]. To compensate for the positional randomness 722 of detections between frames, a cost-volume layer is utilized for patch-to-patch correlation. 723 The features and correlation maps are then sent to a decoder network for flow regression. 724 In the static flow refinement step, they assume that most radar detections are static and 725 therefore use the Kabsch algorithm [126] to robustly estimate the ego-motion. They then 726 filter out moving objects based on the coarse ego-motion, and apply the Kabsch algorithm 727 again to all static points for fine-grained ego-motion estimation. The self-supervised loss 728 consists of three parts: a radial displacement loss, which penalize errors between the 729 estimated velocity projected along radial axis and the measured Doppler velocity, a soft 730 Chamfer distance loss, which encourage temporal consistence between two consecutive 731 point clouds, and a soft spatial smoothness loss, which encourage the spatial consistence 732 for the estimated velocities with their neighbours. The soft version of loss is used to model 733 spatial sparsity and temporal randomness of radar point cloud. 734

#### 5. Radar Object Detection

Due to low resolution, classical radar detection algorithm has limited classification 736 capability. In recent years, the performance of automotive radar has been greatly improved. 737 At hardware level, next generation imaging radars can output high-resolution point clouds. At algorithmic level, neural networks show their potentials to learn better features from 739 the dataset. In this section, we consider a broader definition of radar detection, including 740 point-wise detection, 2D/3D bounding box detection and instance segmentation. We first 741 introduce the classical detection pipeline and recent improvements on clustering and feature 742 selection. As shown in fig. 10, neural networks can be applied to different stages in the 743 classical pipeline. According to input data structure, we classify the deep radar detection 744 into point-cloud-based and pre-CFAR-based. Radar point cloud and pre-CFAR data are 745 similar to LiDAR point cloud and visual image respectively. Accordingly, architectures for 746 LiDAR and vision tasks can be adapted for radar detection. We focus on how knowledge 747 from the radar domain can be incorporated into these networks to address the low SNR 748 problem. 749



**Figure 10.** Overview of radar detection frameworks: Blue boxes indicate classical radar detection modules. Orange boxes represent AI-based substitutions.

# 5.1. Classical Detection Pipeline

As shown in fig. 10, conventional radar detection pipeline consists of four steps: CFAR detection, clustering, feature extraction, and classification. Firstly, a CFAR detector is applied to detect peaks in RD heatmap as a list of targets. Then, the moving targets are projected to Cartesian coordinate and clustered by DBSCAN [18]. Static targets are usually filtered out before clustering because they are indistinguishable from environmental clutter. Within each cluster, hand-crafted features, such as statistics of measurements and shape descriptors, are extracted and sent to a machine learning classifier. Improvements can be

735

799

made upon each of these four steps. CFAR is usually executed in an on-chip DSP, so the choice of method is restricted by hardware support. Cell-Averaging (CA) CFAR [17] is 759 widely used due to its efficiency. It estimates the noise as the average power of neighbouring 760 cells around the cell under test (CUT) within a CFAR window. A threshold is set to achieve a 761 constant false alarm rate for Rayleigh distributed noise. The next generation high-resolution 762 radar chips also support Order-Statistics (OS) CFAR [17]. It sorts neighbouring cells around 763 the CUT according to the received power, and selects the k-th cell to represent the noise 764 value. OS-CFAR has advantages in distinguishing close targets, but introduces a slightly 765 increased false alarm rate and additional computational costs. More sophisticated CFAR 766 variants are summarized in [127], but are rarely used in automotive applications. Deep 767 learning methods can be used to improve noise estimation [128] and peak classification [127] 768 in CFAR. Clustering is the most important stage in radar detection pipeline, especially for 769 the next generation high resolution radar [129]. DBSCAN is favored for several reasons: 770 It does not require a pre-specified number of clusters, it fits arbitrary shapes, and it runs fast [130]. Some works improve DBSCAN by explicitly considering characteristics of radar 772 point cloud. Grid-based DBSCAN [131] suggests clustering radar points in a RA grid map 773 to avoid the range-dependent resolution variations in Cartesian coordinate. Multi-stage 774 clustering [132] proposes a coarse-to-fine two-stage framework to alleviate the negative impact of clutter. It applies a second cluster merging based on the velocity and spatial 776 trajectory of clusters estimated from the first stage. 777

With the improvement of automotive radar resolution, radar target classification has 778 become a hot research topic. For moving objects, the micro-Doppler velocity of moving 779 components such as wheels and arms, can be useful for classification. To better observe 780 these micro-motion, short-time Fourier transform (STFT) is applied to extract Doppler 781 spectrograms. Different types of VRUs can be classified according to their micro-Doppler 782 signatures[133,134]. For static objects, Cai et al. [135] suggest the use of statistical RCS and 783 time-domain RCS as useful features for classification of vehicles and pedestrians. Some 784 researchers work on exploiting a large number of features for better classification. Scheiner 785 et al. [136] consider a large set of 98 features and use the heuristic-guided backward 786 elimination for feature selection. They find range and Doppler features are most important 787 for classification, while angle and shape features are usually discarded, probably because of the low angular resolution. Schumann et al. [137] compares the performance of random 789 forest and LSTM for radar classification. Experiments show that LSTM with an input of 8-790 frame sequences performs slightly better than random forests, especially in the classification 791 of classes with similar shape, such as pedestrians and pedestrian groups, and for false alarms. But LSTM is more sensitive to the amount of training examples. To cope with 793 class imbalance in radar datasets, Scheiner et al. [138] suggest using classifier binarization techniques, which can be divided into two variants: one-vs-all (OVA) and one-vs-one 795 (OVO). OVA trains N classifiers to separate one class from the other N - 1 classes, and OVO trains  $\binom{N}{2}$  classifiers for every class pair. During inference, the results are decided by 797 max-voting. 792

# 5.2. Point Cloud Detector

End-to-end object detectors are expected to replace the conventional pipelines based 800 on hand-crafted features. However, convolutional neural network is not well designed for 801 sparse data structure [139]. It is necessary to increase the input density of radar point cloud 802 for better performance. Dreher et al. [140] accumulate radar points into an occupancy grid 803 mapping (OGM), then apply YOLOv3 [141] for object detection. Some works [142–144] 804 utilize point cloud segmentation networks, such as PointNet [145] and PointNet++ [146], 805 followed by a bounding box regression module for 2D radar detection. The original 3D 806 point cloud input is replaced by a 4D radar point cloud with two spatial coordinates in x-y 807 plane, Doppler velocity and RCS. Scheiner *et al.* [144] compare performances of two-stage 808 clustering method, OGM-based method and PointNet-based method with respect to 2D 809 detection. Experiments show that OGM-based method performs best, while PointNet-810 based method performs far worse than others probably due to sparsity. Liu *et al.* [147] suggest that incorporating global information can help with the sparsity issue of radar point cloud. Therefore, they add a gMLP [148] block to each set abstraction layer in PointNet++. The gMLP block is expected to extract global features at an affordable computational cost.

Most radar detection methods only apply to moving targets, since static objects are 815 difficult to classify due to low angular resolution. Schumann et al. [149] propose a scene 816 understanding framework to detect both static and dynamic objects simultaneously. For 817 static objects, they first build a RCS histogram grid map through temporal integration of 818 multiple frames, and send it to an fully convolutional network (FCN) [150] for semantics 819 segmentation. For dynamic objects, they adopt a two-branch recurrent architecture: one is 820 the point feature generation module, which uses PointNet++ to extract features from the 821 input point cloud. The other is the memory abstraction module, which learns temporal 822 features from the temporal neighbors in the memorized point cloud. The resulted features 823 are concatenated together and sent to a instance segmentation head. In addition, a memory update module is proposed to integrate targets into the memorized point cloud. Finally, 825 static and dynamic points are combined into a single semantic point cloud. The proposed framework can successfully detect moving targets such as cars and pedestrians, as well as 827 static targets like parked cars, infrastructures, poles and vegetation.

As 4D radars have gradually come to the market, radar point cloud density has 829 increased considerably. A major advantage of 4D radar is that static objects can be classified 830 based on elevation measurements without the need to build an occupancy grid map. 831 Therefore, it is possible to train a single detector for both static and dynamic objects. Plaffy 832 et al. [44] applies PointPillars [151] to 4D radar point clouds for 3D detection of multi-class 833 road users. They find the performance can be improved by temporal integration and by 834 introducing of additional features, such as elevation, Doppler velocity and RCS. Among 835 them, the Doppler velocity is essential for detecting pedestrians and bicyclists. However, 836 the performance of the proposed 4D radar detector (mAP 47.0) is still far inferior than 837 their LiDAR detector on 64-beam LiDAR (mAP 62.1). They argue this performance gap 838 comes from radar's poor ability in determining the exact 3D position of objects. RPFA-Net 839 [152] improves PointPillars by introducing a Radar Pillar Features Attention (PFA) module. 840 It leverages self-attention to extract the global context feature from pillars. The global 841 features are then residually connected to the original feature map and sent to a CNN-based 842 detection network. The idea behind is to explore the global relationship between objects 843 for a better heading angle estimation. In fact, self-attention is basically an set operator, so 844 it is well suited for sparse point clouds. Radar transformer *et al.* [153] is a classification network constructed entirely by self-attention modules. The 4D radar point cloud is first 846 sent to an MLP network for input embedding. The following feature extraction network consists of two branches. In the local feature branch, it uses three stacked set abstraction 848 modules [146] and vector attention modules [154] to extract hierarchical local features. In the global feature branch, the extracted local features at each hierarchy are concatenated 850 with global feature map at the previous hierarchy and fed into a vector attention module 851 for feature extraction. In the last hierarchy, a scalar-attention, *i.e.*, the conventional self-852 attention, is used for feature integration. Finally, the feature map is sent to a classification 853 head. Experiments show the proposed radar transformer outperforms other point cloud 854 networks in terms of classification. The above two attention-based approaches show their 855 potential in modeling the global context and extracting semantic information. Further 856 works should focus on combine these two advantages into a fully attention-based detection 857 network.

# 5.3. Pre-CFAR Detector

There are some attempts to explore the potential of pre-CFAR data for detection. Radar pre-CFAR data encode rich information of both targets and backgrounds, but is hard to interpret by human. Neural network is expected to better utilize these information. One option is to use neural network to replace CFAR [155] or DOA estimation[75,156]. Readers

can refer to [157] for a detailed survey of learning-based DOA estimation. Alternatively, there are also some efforts to perform end-to-end detection through neural networks. Deep 865 radar detector [158] jointly trains two cascaded networks for CFAR and DOA estimation 866 respectively. Zhang et al. [159] use stacked complex RD maps as input to a FCN for 867 3D detection. In order to remove the DC component in phase, they perform a phase normalization by using RD cells in the first receiver as normalizers. They argue that phase 869 normalization is crucial for successful training. Rebut et al. [45] design a DDM-MIMO 870 encoder with complex RD map as input. In DDM configuration, as illustrated in fig. 3, all 871 Tx antennas transmit signals at the same time. Instead of performing waveform separation, 872 they directly apply range FFT and Doppler FFT to ADC signals received by Rx antennas. 873 In this way, targets detected from different Tx antennas should be located separately with a 874 fixed Doppler shifts in RA map. To extract these features, they design a two-layer MIMO 875 encoder, consisting of a dilated convolutional layer to separate Tx channels, followed by a 876 convolutional layer to mix the information. This MIMO encoder are jointly trained with the following RA encoder, detection head and segmentation head. 878

In close-field applications that require large bandwidth and high resolution, RD maps are not suitable because the extended Doppler profile can lead to false alarms. RA map, on 880 the other hand, does not suffer from the same problem. For each detection point on RA map, the micro-Doppler information in the slow time can be utilized for better classification. 882 RODNet [41] uses complex RA maps as input for object detection. It performs range FFT 883 followed by angle FFT to get a complex RA map for each sampled chirp. It is difficult to 884 separate static clutter and moving objects using RA map alone without Doppler dimension. 885 To utilize the motion information, it samples a few chirps within a frame. Then, the 886 sequences of RA maps corresponding to these chirps are sent to a temporal convolution 887 layer. Specifically, it first uses 1x1 convolutions along the chirp dimension to aggregate 888 temporal information. Then, a 3D convolution layer is used to extract temporal features. 880 Finally, the features are merged along the chirp dimension by max-pooling. Experiments 890 indicate sampling 8 chirps out of 255 can achieve a comparable performance with using the 891 full chirp sequences. 892

Training neural network to utilize phase information in complex RA or RD map is 803 a difficult task. Alternatively, some works attempt to use the real-valued RAD tensor as 894 input. A key issue in using the 3D RAD tensor as input is the curse of dimensionality. 895 Therefore, many techniques are proposed to reduce the computational cost of 3D tensor processing. RADDet [40] normalizes and reshapes the RAD tensor to an image-like data 897 structure. The Doppler dimension is treated as channel of 2D RA maps. Then, YOLO is applied to the RA map for object detection. One disadvantage is that this method fails 800 to utilize the spatial distribution of Doppler velocities. Alternatively, 3D convolution can be used to extract features from all three dimensions in a 3D tensor, but requires huge 901 computation and memory overheads [160]. RODNet [41] samples chirp sequences, as described above, to reduce input dimensionality. RTCNet [161] reduces tensor size by 903 cropping a small cube around each point detected by CFAR, and then uses 3D CNN to 904 classify these small cubes. However, its detection performance is limited by the CFAR 905 detector. To fully exploit the information encoded in RAD tensors, some works [85,162,163] 906 adopt the multi-view encoder-decoder architecture. Major et al. [162] and Ouaknine et al. 907 [163] both utilize a similar multi-view structure. The RAD tensor is projected into three 2D 908 views. Then, three decoders extract features from these views respectively. To fuse these 909 features, Ouaknine et al. directly concatenate three feature maps. Major et al. recover the 910 tensor shape by duplicating these 2D feature maps along the missing dimension, then use a 911 3D convolution layer to fuse them. Next, the Doppler dimension is suppressed by pooling 912 to recover the shape of RA feature map. Finally the fused feature maps are sent to a decoder 913 for downstream segmentation tasks. Another difference is Major et al. use skip-connection 914 while Ouaknine et al. adopt a ASPP [74] pathway to encode information from different resolutions. RAMP-CNN et al. [85] is also built in a multi-view architecture, but it uses three 916 encoder-decoders for feature map extraction. Their fusion method is similar to Major's but in 2D. 918

Radar pre-CFAR data are captured in a polar coordinate. For object detection, polar-010 to-Cartesian transformation is necessary to obtain the correct bounding box. Major et al. 920 [162] compare three configurations for coordinate transformation: pre-processed input 921 transformation, learning from neural networks, and transformation on middle-layer fea-922 ture map. Experiments show applying explicit polar-to-Cartesian transformation to the 923 last-layer feature map achieves the best performance, the implicit learning-based transfor-924 mation is slightly worse and the pre-processed transformation is far inferior than other 925 two. They attribute this poor performance to distorted azimuth side-lobes in the input. 926 In fact, conventional 2D convolution is not the best choice for radar pre-CFAR data, since 927 range, Doppler and azimuth dimension vary in their dynamic ranges and resolutions. 928 Instead of 2D convolution, PolarNet [164] uses cascade of two 1D convolutions, including 929 a column-wise convolution to extract range-dependent features, followed by a row-wise convolution to mix information from spatial neighbours. A similar idea is used in google's 931 RadarNet [165] for gesture recognition. They first extract range-wise features then summa-932 rized them together in the later stage. Meyer *et al.* [166] use a isotropic graph convolution 033 network (GCN) [167] to encode the RAD tensor and achieves more than 10% improvement in AP for 3D detection. They argue that the performance gain comes from the ability of 935 GCN to aggregate information from neighboring nodes.

Incorporating temporal information is an effective way to improve the performance of 937 pre-CFAR detectors. There are multiple ways to add temporal information to the network. 938 Major *et al.* [162] use a convolutional LSTM layer to process a sequence of feature maps 939 from the encoder network. Experiments indicate the temporal layer enables more accurate 940 detection and significantly better velocity estimation. Ouaknine *et al.* [163] compare the 941 performance between the static model with accumulated inputs and the temporal model 942 with stacked inputs. For the static model, RAD tensors within 3 frames are accumulated 943 into one single tensor, and fed to a multi-view encoder-decoder for segmentation. For the 944 temporal model, RAD tensors within 5 frames are stacked to form a 4D tensor and then 945 sent to a multi-view encoder-decoder. In each branch, multiple 3D convolution layers are 946 used to leverage spatial-temporal information. The results show that the introduction of the temporal dimension can significantly improve detection performance. Pervsic *et al.* 948 [68] discuss the effect of the number of stacked radar frames. They find too long frames will introduce many background clutter, which in turn makes the model difficult to learn 950 target correspondences. According to their experiments, stacking of 5 frames is the most suitable choice. RODNet [41] investigate on stacking multiple frames on feature-level. It 952 concatenates the extracted per-frame features and sends them to a 3D CNN layer. For 953 motion compensation, they apply deformable convolution [168] on the chirp dimension 954 in the first few layers. In addition, an inception module with different temporal length are used in the later layers. Despite the introduction of additional computational costs, 956 these two temporal modules significantly improve the average precision. Li et al. [169] 957 explicitly model the temporal relationship between features extracted from two consecutive 958 frames using a attention module. Firstly, they stack RA maps in two orders, *i.e.*current 050 frame on top and previous frame on top. Then, they use two encoders to extract features from these two inputs and concatenate the features together. A positional encoding is 961 further added to compensate the positional imprecision. Next, the features are sent to a 962 masked attention module. The mask is used to disable cross-object attention in the same 963 frame. Finally, the temporally enhanced features are sent to a encoder for object detection. This attention-based approach is more semantically interpretable and avoids the locality 965 constraint induced by convolution.

#### 6. Sensor Fusion for Detection

Different sensors observe and represent an object with different features. Sensor fusion 968 can be considered as the mapping of different modalities into a common latent space where different features of the same object can be associated together. In this section, we 970 focus on sensor fusion for detection. We argue that the conventional taxonomy of fusion 971 architectures into early (input), middle (feature) and late (decision) fusion is ambiguous for 972 neural network based detection. For example, in the definition of late fusion, we cannot 973 distinguish between ROI-level (without category information) fusion and object-level (with 974 category information) fusion. Therefore, we explicitly classify fusion methods according to 975 the fusion stage. This is beneficial because different fusion stages correspond to different 976 levels of semantics, *i.e.*, the classification capabilities. As shown in fig. 11, we classify fusion 977 architectures into four categories: input fusion, ROI fusion, feature map fusion and decision 978 fusion.



**Figure 11.** Overview of radar and camera fusion frameworks. We classify the fusion frameworks into input fusion, ROI fusion, feature map fusion and decision fusion. For ROI fusion, we further investigate two architectures: cascade fusion which projects radar proposals to image view, and parallel fusion which fuses radar ROIs and visual ROIs.

# 6.1. Input Fusion

Input fusion is applied to radar point cloud. It projects radar points into a pseudoimage with range, velocity, and RCS as channels [170,171]. Then, similar to a RGB-Depth image, the radar pseudo-image and the visual image are concatenated as a whole. Finally, a visual detector can be applied to this multi-channel image for detection. Input fusion does not make independent use of the detection capability of radar. In other words, the radar and vision modalities are tightly coupled. Assuming good alignment between modalities, it makes the network easier to learn joint feature embeddings. However, an obvious disadvantage is that the architecture is not robust to sensor failures.

The fusion performance depends on the alignment of radar detections with visual pixels. As mentioned in section 4.1, the difficulties lie in three aspects: Firstly, the radar point cloud is highly sparse. Many reflections from the surface are bounced away due to specular reflections. As a result, the detected points are sparsely distributed over the object. In addition to the sparsity, the lateral imprecision of radar measurements leads to further difficulties. The radar points can be out of the visual bounding box. The imprecision of comes from different aspects, *e.g.*, imprecise extrinsic calibration, multi-path effects and provide the sparse of the

979

- 00

low angular resolution. The third limitation is that low-resolution radar does not provide height information. To address these difficulties, some association techniques are required. 997 Relying on the network to implicitly learn association is a hard task, because the network 998 tends to simply ignore the weak modality such as radar. The expansion methods described 999 in section 4.1 can be applied as a pre-processing stage for input fusion. However, object detection does not require such a strict association as depth completion, so some of the 1001 expansion methods are too costly for real-time processing. Nobis et al. [170] utilize the light- 1002 weight height extension as pre-processing. Both Chadwick *et al.* [171] and Yadav *et al.* [172] 1003 add a one-layer convolution to radar input before concatenation. This convolutional layer can be considered as a lightweight version of the association network. Radar detections at 1005 different ranges requires different size of receptive field for association. Therefore, Nobis et 1006 *al.* [170] concatenate the radar pseudo-image with image feature maps at multi-scales. 1007

#### 6.2. ROI Fusion

ROI fusion is adapted from the classical two-stage detection framework [173]. Region of Interests (ROIs) can be considered as a set of object candidates without category information. The fusion architecture can be further divided into cascade fusion and parallel fusion. In cascade fusion, radar detections are directly used for region proposal. Radar points are projected into image view as the candidate locations for anchors. Then, the ROI is determined with the help of visual semantics. In the second stage, each ROI is classified and its position is refined. Nabati *et al.* [174] adopt two techniques to improve the anchor quality. They add offsets to anchors to model the positional imprecision of radar detections. To mitigate the scale ambiguity in the image view, they rescale anchor size according to the range measurements. In their following work [175], they directly propose 3D bounding boxes and then map these boxes to the image view. In this way, the rescaling step can be avoided. It is also possible to propose region on radar point cloud using visual ROIs. For example, CenterFusion [176] propose a frustum-based association to generate radar ROI frustums using visual bounding boxes.

Cascade fusion is particularly well suited for low-resolution radars, where the radar 1023 point cloud has a high detection recall but is very sparse. However, there are two potential 1024 problems with the cascade structure. Firstly, the performance is limited by the completeness 1025 of proposed ROIs in the first stage. In other words, if an object is missed, we cannot recover 1026 it in the second stage. The second problem is that the cascade structure cannot take 1027 advantage of modality redundancy. If the radar sensor is nonfunctional, the whole sensing 1028 system will fail. Therefore, it is necessary to introduce a parallel structure to ROI fusion. 1029 Nabati et al. [175] adopt a two-branch structure for ROI fusion. The radar and visual 1030 ROIs are generated independently. Then, the fusion module merges radar ROIs and visual ROIs by taking an set union, while the redundant ROIs are removed through NMS. To 1032 enable adaptive fusion of modalities, Kim *et al.* [177] propose a Gated Region of Interest Fusion (GRIF) module for ROI fusion. It first predicts a weight for each ROI through a 1034 convolutional-sigmoid layer. Then, the ROIs from radar and vision are multiplied by their corresponding weights and element-wise added together. 1036

#### 6.3. Feature Map Fusion

Feature-map fusion leverage the semantics from both radar and image. From section 5.3, we find that high resolution radars can provide sufficient semantic cues for classification. Therefore, feature-map fusion utilizes two encoders to map radar and image into the same latent space with high-level semantics. The detection frameworks are flexible, including one stage methods [178,179] and two-stage methods [33,180,181]. The one-stage method leverage two branches of neural networks to extract feature maps from radar and image respectively, and then concatenate the feature maps together. The two-stage fusion methods are adapted from the classical fusion architecture AVOD [182]. It firstly fuses the ROIs proposed from radar and image in the first stage. In the second stage, the fused ROIs are projected to radar and visual feature maps respectively. The feature maps inside 1047

1037

the ROIs are cropped and resized to an equal-sized feature crop. The feature crop pairs 1048 from radar and image are then fused by element-wise mean and sent to a detection head. 1049 Generally speaking, the two-stage method has better performance, but it is much slower 1050 than the one-stage method. Anchor free methods [183,184] further avoid the complicated 1051 computation related to anchor boxes, such as calculating IOU score during training. 1052

Feature-map fusion allows the network to flexibly combine radar and visual semantics. 1053 However, the fusion network may face the problem of overlooking weak modalities and 1054 modality synergies [185]. Some training techniques are needed to force the network to 1055 learn from radar input. Nobis et al. [170] adopt a modality-wise dropout approach that 1056 randomly deactivates image branch during training. Lim *et al.* [178] use a weight freezing strategy to fix the weights of the pre-trained feature extractors when training the fusion 1058 branch. Experiments show that freezing only the image branch works best. However, 1059 fusion of multiple modalities is not guaranteed to always be better than using single 1000 modality. Sometimes we want the network to lower the weight of radar branch if it gives noisy inputs. To achieve adaptive fusion, Cheng et al. [186] adopt self-attention and global 1002 channel attention [187] in their network. The self-attention is used to enhance real target 1063 points and weaken clutter points. Then, the global attention module is applied to estimate 1004 modality-wise weights. Bijelic et al. [60] estimate the sensor entropy as the modality weight. 1065 For each modality, the entropy is evaluated pixel-wise as a weight mask. Then these weight masks are multiplied with the corresponding feature maps at each fusion layer. 1067

# 6.4. Decision Fusion

Decision fusion assumes that objects are detected independently by different modalities and fuses them according to their spatial-temporal relationships. This structure realizes sensing redundancy at the system level and is therefore robust to modality-wise error. Due to the low resolution of radar, most existing studies do not explicitly consider the category information estimated by radar. In other words, they only fuse the location information from radar and vision branches, while retaining the category information estimated by vision. Since the next generation 4D radar can provide classification capabilities, it is expected that future fusion frameworks should consider both location and category information.

The location can be optimal fused in a tracking framework. Different objects are first 1077 associated and then sent to a Bayesian tracking module for fusion. Due to the low resolution 1078 of radar, association is difficult to achieve in some scenarios, e.g., a truck splitting into 1079 two vehicles or two close objects merging into one. Such association ambiguity can be 1000 mitigated using a track-to-track fusion architecture [188]. By estimating tracks, temporal information can be leveraged to filter out false alarms and interpolate missed detections. 1002 Some researchers exploit deep learning to make a better association between radar and 1083 other modalities. RadarNet [183] propose an attention-based late fusion to optimize the 1004 estimated velocity. Firstly, they train a fiver-layer MLP with softmax to estimate the normalized association scores between each bounding box and its nearby radar detections. 1086 Then, they predict the velocity by weighted averaging the radar-measured velocities using 1087 the association scores. AssociationNet [189] attempts to map the radar detections to a better 1088 representation space in contrastive learning framework. It first projects radar objects and 1089 visual bounding boxes to the image plane as pseudo images. To utilize the visual semantics, 1000 they concatenate these pseudo images with the original image. Next, the concatenated 1001 images are sent to an encoder-decoder network to output a feature map. Representation 1002 vectors are extracted from the feature map according to the locations of radar detections. A 1003 contrastive loss is designed to pull together the representation vectors of positive samples 1004 and push away the representation vectors of negative examples. During inference, they 1005 compute the Euclidean distance between the representation vectors of all possible radarvisual pairs. The pairs with distance below the threshold are considered associative. 1097

Category information, especially the conflict in category predictions, is difficult to 1098 handle in sensor fusion. BayesOD [190] proposes a probabilistic framework for fusing 1099 bounding boxes with category. The locations of bounding boxes are modelled by Gaussian 1100

distributions. The category prior is modelled as a Dirichlet distribution, thereby allowing a Dirichlet posterior to be computed in closed form. Then, the bounding box with the 1102 highest categorical score is considered as the cluster center, while the other bounding boxes 1103 are treated as measurements. Finally, Bayesian inference is used to optimally fuse the 1104 location and category information of these bounding boxes. Probabilistic methods have 1105 their inherent shortage in modelling the lack of knowledge [191]. For example, a uniform 1106 distribution brings confusion of either the network has no confidence in its prediction or the 1107 input is indeed ambiguous for classification. In contrast, set-based methods have no such 1108 problem. Chavez et al. [192] leverage the evidential theory to fuse LiDAR, camera and radar. 1109 They consider the frame of discernment, *i.e.*, the set of mutually exclusive hypotheses, as 1110  $\Omega = \{pedestrains(p), bikes(b), cars(c), truck(t)\}, and assign each possible hypothesis,$ *i.e.*, 1111a subset of  $\Omega$ , with a belief. In the case of object detection, possible hypotheses are selected 1112 according to sensor characteristics. For example, a car is sometimes confused as part of a 1113 truck. Thus, if a car is detected, evidence should be also put into the set  $\{c, t\}$  and the set of 1114 ignorance  $\Omega$ . Accordingly, we can assign the belief *m* to a car detection as 1115

$$m(\lbrace c \rbrace) = \gamma_c \alpha_c, \quad m(\lbrace c, t \rbrace) = \gamma_c (1 - \alpha_c), \quad m(\Omega) = 1 - \gamma_c$$
(10)

where  $\gamma_c$  is a discounting factors to model the uncertainty of misdetection, and  $\alpha_c$  is 1110 the accurateness, *i.e.*, the rate of correct predictions in car detecting. Suppose there are two 1117 sources of evidence  $S_1$  and  $S_2$  from different modalities. Each of these sources provides 1118 a list of detections as  $A = \{a_1, a_2, ..., a_m\}$  and  $B = \{b_1, b_2, ..., b_n\}$ . Then, three propositions 1119 can be defined regarding the possible association of two detections  $a_i$  and  $b_j$  as 1120

- {1} if  $a_i$  and  $b_j$  are the same object;
- {0} if  $a_i$  and  $b_j$  are not the same object 1122
- $\{0,1\}$  for the ignorance of association.

The belief of association can be determined according to both location and category similarities. The evidence for location similarity is defined according to the Mahalanobis distance as

$$m_{a_i,b_j}^p(\{0\}) = \alpha(1 - f(d_{a_i,b_j}))m_{a_i,b_j}^p$$

$$m_{a_i,b_j}^p(\{1\}) = \alpha f(d_{a_i,b_j}) \quad m_{a_i,b_j}^p(\{1,0\}) = 1 - \alpha$$
(11)

where  $f(d_{a_i,b_j}) = \exp(-\lambda d_{a_i,b_j}) \in [0,1]$  measure the similarity with respect to the 1127 Mahalanobis distance  $d_{a_i,b_j}$  and a scaling factor  $\lambda$ , and  $\alpha$  is an evidence discounting factor. 1128 For the category similarity, two detections belong to the same category is too weak to 1129 provide evidence that they are the same object. However, if two detections are of different 1130 categories, it is reasonable to assign evidence to the proposition that they are not the same object. Accordingly, the evidence for category similarity is given by 1132

$$m_{a_{i},b_{j}}^{c}(\{0\}) = \sum_{A \cap B = \emptyset} m_{a_{i}}^{c}(A)m_{b_{j}}^{c}(B) \quad \forall A, B \subset \Omega$$
  
$$m_{a_{i},b_{j}}^{c}(\{1\}) = 0, \quad m_{a_{i},b_{j}}(\{0,1\}) = 1 - m_{a_{i},b_{j}}^{c}(\{0\})$$
(12)

where the mass evidences are fused if no common category hypothesis is shared, 1133 *i.e.*,  $A \cap B = \emptyset$ . And the rest of the evidence are placed in the ignorance hypothesis. 1134 Finally, for each detection pairs, the category similarity and the location similarity are fused according to Yager's combination rule [193]. Evidential fusion provides an reliable framework for information fusion. However, it cannot be directly applied to neuralnetwork-based detectors that make predictions on a single hypothesis. To address this problem, conformal prediction [194] can be used to generate confidence sets from a trained network using a small amount of calibration data.

# 7. Challenges

Although deep radar perception shows good performance on datasets, there are few studies investigating the generalization of these methods. In fact, some challenging situations are overlooked, but may prohibit the use of these methods in real-world scenarios. For example, the ghost objects caused by multi-path propagation are common in complex scenarios. Over-confidence is a general problem with neural networks. Since radar is always used for safety-critical applications, it is important to calibrate the detection network and output the predictive uncertainty. Even though we always refer to radar as an all-weather sensor, robustness in adverse weather is not well tested in many radar fusion methods. In this section, we present these three challenges and summarize some recent works that attempts to solve them.

# 7.1. Ghost Object Detection

Multi-path is a phenomenon in the physics of waves where a wave from a target 1153 travels to a detector through two or more paths. Because of the multi-path propagation, 1154 radar receives both direct reflections and indirect time-shifted reflections of targets. If the target reflections and the multi-path reflections occupy the same RD cell, the performance 1156 of DOA estimation is affected. Otherwise, if they occupy different cells, it can produce 1157 ghost targets in multi-path directions. In the latter case, since ghost detection has similar 1158 dynamics to the real target, it is difficult to eliminate them in the traditional detection pipeline. The multi-path effect can be classified into three types [195]. The first type is the reflection between ego-vehicle and targets. Therefore, the distance and velocity of clutter should be multiple times of the true measurement. The second type is the underbody 1102 reflection. It usually happens under the truck, resulting in points with longer distances. 1103 This see-through effect is sometimes beneficial, since occluded vehicles can be detected. 1104 The third type is mirrored ghost detections caused by the reflective surface. Because of the 1165 large wavelength of automotive 77GHz radar, many flat facilities, such as concrete walls, 1166 guardrails and noise cancellation walls, can be regarded as reflective surfaces. As shown 1167 in fig. 12, this kind of multi-path effect can be further categorized into type 1 and type 2 1168 depending on whether the final reflection occurs on the target or the surface [196]. The 1109 number of reflections is referred to as the order of the multi path. Usually, only orders 1170 below 3 need to be considered, since higher order reflections return little energy due to 1171 signal diffusion. 1172



**Figure 12.** Multi-path effect: The solid orange box is the real object. Dotted boxes are ghost objects caused by multi-path propagation.

A high quality dataset is necessary for the performance evaluation of ghost detection. 1173 However, the labelling of ghost objects is a difficult task and requires expert knowledge. 1174 Chamseddine *et al.* [88] propose a method to automatically identify radar ghost objects by 1175 comparing with the LiDAR point cloud. However, LiDAR measurements are not perfect. It 1176 has its own inherent defects, such as sparsity, limited range and holes with no reflections. 1177

Therefore, using LiDAR as ground truth could be sometimes problematic. In the Radar Ghost dataset [55], ghost objects are manually annotated with the help of a helper tool. This tool can automatically calculate the locations of potential ghosts based on real objects and reflective surfaces. As a result, four types of multi-path effects are annotated, including type-1 second-order bounces, type-2 second-order bounces, type-2 third-order bounces, and the other higher-order bounces. In addition, they also provide a synthetic dataset by overlaying objects from different frames within the same scene.

Unlike clutter, ghost objects cannot be filtered by temporal tracking because they 1185 have the same kinematic properties as real targets. Instead, they can be detected by 1186 geometric methods [195,197]. With a radar ghost dataset, it is also possible to train a neural 1187 network for ghost detection, such as PointNet-based methods [88] and PointNet++-based 1188 methods [196,198]. Because of the signal diffusion, the higher order reflections can be 1189 safely ignored. Thus, ghost objects usually occur in a ring-shaped region with the similar 1100 distance as the real target. Accordingly, Griebel et al. [198] design a ring grouping to replace 1191 the multi-scale grouping in PointNet++. The scene structure and relationship between 1192 detections are important cues to identify ghost objects. Garcia et al. [199] suggest the 1193 occupancy grid map can provide information of scene structure. Therefore, they use the 1194 occupancy grid map and the list of moving objects as inputs to FCN, to predict a heatmap 1195 of moving ghost detections. Wang et al. [200] propose to use multimodal transformers to 1106 capture the semantic affinity between ghost objects and real objects with LiDAR as reference. 1197 They design a multimodal attention block, which consists of two modules. The first one is a self-attention module for radar point cloud. It is expected to model the similarities of 1199 real objects and mirrored ones. The feature maps from the radar and LiDAR branches are 1200 then fused by a second multimodal attention module. This fusion module can be seen as 1201 calculating the correlation between LiDAR detections and real radar detections. 1202

#### 7.2. Uncertainty in Radar Detection

Learning-based radar detection shows its potential in classifying different road users. 1204 However, the performance evaluated on research datasets could be biased due to class 1205 imbalance and simple scenarios. Palffy et al. [44] summarize some failure cases for radar 1206 detection in VoD dataset: Two close pedestrians can be detected as one bicyclist. One large 1207 object, for example a truck or a bus, can be split into two smaller ones. Distant objects 1208 with few reflections may be missed by the detector. Strong reflections from metal poles 1209 and high curbs can mask real objects. Most of these failures come from the imperfec- 1210 tion of radar sensors with respect to angular resolution and dynamic range. To make it 1211 worse, neural networks tend to be overconfident in their incorrect predictions [201]. For autonomous driving, the misspecified confidence in perception can leak to downstream 1213 tasks like sensor fusion and decision making, potentially leading to catastrophic failure. 1214 Patel *et al.* [202] investigate the class uncertainty of a learning-based radar classifier under different perturbations, including domain shift, signal corruptions and out-of-distribution 1216 data. Experiments indicate their baseline network are severely over-confident under these perturbations. 1218

There are two kinds of uncertainty: Data uncertainty, also known as aleatoric uncertainty, is caused by noisy input. Model uncertainty, also known as epistemic uncertainty, is caused by insufficient or inappropriate training of the network. Sources of model uncertainty include three cases: covariate shift (p(x) changes), label shift (p(y) changes) and open set recognition (unseen y) [203]. The sum of data uncertainty and model uncertainty is referred to as predictive uncertainty. For the task of probabilistic object detection [204], uncertainties of two parameters are of interest: class uncertainty, which encodes the confidence in the classification, and spatial uncertainty, which represents the reliability of the bounding boxes. Class uncertainty can be seen as model uncertainty, while spatial uncertainty is more relevant to data uncertainty introduced by noisy input.

For classification tasks, the simplest way is to learn a function that maps the pseudo probability output by softmax layer into true probability. The true probability is defined 1230

as the class-wise accuracy on the training set. This process is usually called network calibration. Since it is a post-processing method, both the model size and inference time are not affected. The calibration method is mainly concerned with the aleatoric part of the overall uncertainty [191]. To calibrate the radar classifier, Patel *et al.* [205] compares different post-processing techniques, including temperature scaling [201], latent Gaussian process (GP) [206] and mutual information maximization [207]. Mutual information maximization achieves the best balance between performance and inference time. Some recent researches indicate that soft-label augmentation techniques, such as label smoothing [208] and mixup [209,210], can effectively mitigate the over-confidence problem, thus helps network calibration. Patel *et al.* [205] suggests the use of label smoothing regularization in radar classification. The core idea is that the classifier should give lower confidence to distant objects with low received power. Therefore, they propose two label smoothing techniques to generate soft labels according to the range and the received power respectively. Experiments show that both of them can significantly improve the calibration performance. 1240

In addition to calibrating the class uncertainty, we are also interested in estimating the spatial uncertainty in bounding box regression. Monte-Carlo Dropout [211] and Deep Ensembles [212] are popular in estimating predictive uncertainty. However, experiments [213] show that these methods only provide marginal improvements in object detection, but at a high cost. Direct modelling [214] is widely used to estimate the aleatoric uncertainty in bounding box regression. The idea is to let the network estimate both mean and variance of a prediction. The loss is constructed as

$$\mathcal{L}(\theta) = \frac{1}{N} \sum_{i=1}^{N} \frac{1}{2\sigma(x_i)^2} \|y_i - f(x_i)\|^2 + \frac{1}{2} \log \sigma(x_i)^2$$
(13)

where  $\sigma(x_i)$  is the estimated variance which reduces the penalty with high variance and penalises high variance at the same time. Dong *et al.* [215] estimate the spatial uncertainty in radar detection using direct modelling. Experiments indicate that adding variance prediction for bounding box parameters can improve detection performance, especially under high IoU threshold.

#### 7.3. Fusion in Adverse Weather

Adverse weather conditions, such as heavy rain, snow and fog, can be a significant 1258 threat to safe driving. Different sensors operate in different electromagnetic wavebands, 1259 thus having different robustness to environments. A comparison of the weather effects on 1260 different sensors can be found in [216]. Visual perception is susceptible to blur, noise, and 1201 brightness distortions [217,218]. In adverse weather, LiDAR suffers from reduced detection 1202 range and blocked view in powder snow [219], heavy rain [220], and strong fog [220]. In 1263 contrast, radar is more robust under adverse weather. The effect of weather on radar can be divided into attenuation and backscattering [221]. The attenuation effect decreases the 1265 received power of the signal, and the backscattering effect increases the interference at 1206 the receiver. Experiments [60,222–224] reveal attenuation and backscattering under dust, 1207 fog and light rain are negligible for radar, while the performance of radar degrades under 1208 heavy rainfall. Zang et al. [221] summarize the mathematical models for attenuation and 1209 backscattering effects of rain. They suggest the detection range of radar can be reduced by 1270 up to 45% under severe rainfall conditions (150 mm/h). For close targets with small RCS, 1271 the backscattering effect is more severe and can cause additional performance degradation. 1272

Driving in adverse weather can be considered as a corner case [225] for autonomous <sup>1273</sup> driving. The concept of operational design domain(ODD) [226] is proposed to define <sup>1274</sup> the conditions under which autonomous vehicles are designed to operate safely. If the <sup>1275</sup> monitoring system [227] detects a violation of OOD requirements, control will be handed <sup>1276</sup> over to the driver. However, changes in the operational environment are usually rapid <sup>1277</sup> and unpredictable. Therefore, the hand-over mechanism is controversial in terms of safety. <sup>1278</sup> In the future, a fully autonomous vehicle (SAE Level 5) is expected to work under all <sup>1279</sup> environmental and weather conditions. However, most fusion methods are not designed <sup>1280</sup>

to explicitly consider weather effects. Networks trained in good weather may experience 1281 performance degradation in adverse weather. 1282

There are some possible ways to adapt the network to different weather conditions. 1283 One way is to add a scene switching module [228], then use different networks for different weathers. This method is straightforward but introduces additional computational and 1285 memory costs. The other option is to add some dynamic mechanisms into the network. 1286 Qian et al. [33] add a two-stage attention block in the fusion module. They first apply self- 1287 attention to each modality and then mix them through cross attentions. Experiments show 1288 that the fusion mechanism performs robustly in foggy weather. They further investigate 1289 the domain generalization problem, *i.e.*, training with a good-weather dataset and inference in foggy weather. The result shows a significant accuracy drop compared with training 1201 with data in both good and foggy weather, indicating that their model relies on data to 1202 generalize. Malawade et al. [229] propose a gating strategy to rank each modality and 1203 pick the top 3 reliable modalities for fusion. They compares three types of gating methods: 1294 knowledge-based, CNN-based and attention-based gating. Knowledge-based gating use a 1295 set of pre-defined modality-wise weights for each weather condition, while CNN-based 1206 and attention-based learn the weights from data. Experiments on RADIATE dataset [61] 1207 indicate gating methods outperform fusion methods under adverse weather, and attentionbased gating can achieve the best performance. Alternatively, Bijelic et al. [60] proposes an 1299 entropy-steered fusion network which uses the sensor entropy as modality-wise weights. 1300 Specifically, they use a deep fusion architecture that continuously fuses feature maps from 1301 different modalities. The pixel-wise entropy is used as the attention map for each sensor 1302 branch.Since the entropy map is conditioned only on sensor inputs, the fusion network 1303 can perform robustly in unseen adverse weather. According to uncertainty theory, sensor 1304 entropy can be considered as a measure of data uncertainty. To utilize both data and model 1305 uncertainties, Ahuja et al. [230] propose an uncertainty-aware fusion framework. They 1306 leverage a decision-level fusion architecture and expect each branch to output both data 1307 uncertainty and model uncertainty. A gating function is used to apply a weighted average 1308 to each modality according to the predicted uncertainty. Then, they design two modules 1309 to handle the data with high uncertainty. One for failure detection. A sensor with data 1310 uncertainty consistently above a threshold is considered to be malfunctioning. The other 1311 is used for continuous learning. Data with model uncertainty above a threshold will be 1312 added to the training set for continuous learning. 1313

#### 8. Future Research Directions

In this paper, we summarize the recent developments on deep radar perception. As <sup>1315</sup> we can see, many research efforts have focused on developing models for detection tasks. <sup>1316</sup> However, there are also some unexplored research topics or fundamental questions to <sup>1317</sup> be addressed. In this section, we propose some interesting research directions to the <sup>1318</sup> automotive radar community. <sup>1319</sup>

# 8.1. High Quality Dataset

Deep learning revolution started with the introduction of ImageNet dataset [231]. 1321 However, radar perception has not yet seen its ImageNet moment. Although many datasets 1322 exist, they differ in scale, resolution, data representation, scenario, and labelling granularity. 1323 The granularity and quality of labelling is also a key issue for radar datasets. Therefore, it is 1324 hard to fairly compare different models trained on different datasets. Since the introduction 1325 of 4D imaging radar to the market, we anticipate an urgent need for datasets with high quality annotations and diverse scenes. 1327

#### 8.2. Radar Domain Knowledge

In the absence of high-quality datasets, we need to avoid treating AI in radar as a data fitting game. It is essential to exploit domain knowledge to develop a generalizable perception model. Radar domain knowledge need to be considered at many stages, such as 1331

1320

1314

33 of 43

labelling, data augmentation, model structure, training techniques and evaluation metrics. 1332 Take ghost detection as an example. From a data perspective, we need to use our expert 1333 knowledge to label ghost objects [55]. From a network perspective, we can design an 1334 attention module [200] or utilize graph convolution [166] to model the relationship between 1335 real and ghost objects. We hope that researchers put more focus on solving these critical 1336 problems in radar perception. 1337

# 8.3. Uncertainty Quantification

As introduced in section 7.2, uncertainty quantification is important for applying AI 1330 in safety-critical applications. Due to the low SNR of radar data and the small size of 1340 radar datasets, both high data and model uncertainties are expected for CNN-based radar 1341 detectors. However, there is still very little work touching on this topic. Although many 1342 uncertainty quantification methods have been proposed, they are not necessarily helpful 1343 for a specific task. For example, Feng *et al.* [232] find that sampling-based methods are not 1344 useful for visual object detection. Similarly, we need empirical experiments and theoretical 1345 explanations to demonstrate the necessity and effectiveness of uncertainty quantization 1346 methods for radar perception. 1347

# 8.4. Motion Forecasting

An overlooked feature of radar is Doppler velocity. In addition to being a feature of 1349 moving road users, Doppler velocity is valuable for motion forecasting. Motion forecasting 1350 is a popular research topic in autonomous driving [233]. By accurately estimating the 1351 motion of road users, the down-stream path planning module can better react to future 1352 interactions. Lin et al. [234] predict trajectories by building a constant velocity model 1353 with binarized RA maps as input. However, experiments show that the constant velocity 1354 model performs poorly in predicting vehicle trajectories [233]. As mentioned in section 4, 1355 a second-order nonlinear motion model can be developed using the measured Doppler 1356 velocity. We believe that radar has great potential to play an important role in motion 1357 forecasting. 1358

# 8.5. Interference Mitigation

For FMCW radar, mutual interference is a challenging task to solve. It occurs when 1300 multiple radars operate simultaneously in direct line of sight [235]. Depending on if the 1361 chirp configuration, *i.e.* slope and chirp duration, is same between interferer and victim 1302 radar, interference can be classified as coherent and incoherent [236]. Coherent interference 1363 occurs when the same chirp configuration is used and leads to ghost detections. Incoherent 1364 interference is caused by different types of chirps, resulting in significantly increased noise 1305 floor, masked weak target and thus reduced probability of detection. In reality, partially 1366 coherent interference is more widely seen where interferer has a slightly different chirp 1367 configuration. Oyedare et al. [237] summarize deep learning methods for interference 1368 mitigation. Although these methods achieve better performance than classical zeroing 1369 methods, they are generally designed for specific types of disturbances and require sig- 1370 nificant computational costs. Future research should consider interference mitigation and 1371 downstream tasks (e.g., detection) as a whole, and build an end-to-end learning framework 1372 to optimize them together. 1373

# 9. Conclusions

The purpose of this review article is to provide a big picture of deep radar perception. 1375 We first summarize the principles of radar signal processing. Then, we present a detailed 1376 summary of radar datasets for autonomous driving. To encourage researchers to build 1377 their own datasets, we also present methods for calibration and labelling. We further 1378 investigate data augmentation and synthetic radar data to improve data diversity. Radar 1370 can be used for depth completion and velocity estimation. For ease of depth completion, 1380 several expansion methods are introduced to better associate the radar detections with the 1381

1348

1338

1359

image pixels. The full velocity can be recovered as a geometric optimization problem or in 1382 a self-supervised learning way. 1383

Radar detection is the main focus of this paper. We classify deep radar detectors 1384 into point-cloud-based and pre-CFAR-based. PointNet variant networks and multi-view 1385 encoedr-decoders are popular choices for radar point clouds and radar tensors respectively. 1386 By increasing spatial density and exploiting temporal information, significant performance 1387 improvements can be achieved. Some new operators, such as depth-wise convolution, 1388 attention and graph convolution, are leveraged for larger receptive field. In practical applications, radar is often fused with cameras and LIDAR. We classify fusion frameworks 1300 into four categories. Input fusion requires a lightweight pre-processing to explicitly handle 1391 radar position imprecision. Cascaded ROI fusion is not robust to sensor failures, while 1392 parallel ROI fusion improves it. Feature map fusion provides the network with greater 1303 flexibility to combine radar and visual semantics, but requires specific training techniques 1394 for effective learning. Decision fusion takes advantage of modal redundancy and is there- 1395 fore popular in real-world applications. Location information can be robustly fused in a 1306 track-to-track architecture or with the help of network semantics. Category information 1307 can be fused with Bayesian inference or evidence theory. 1398

We summarize three challenges for deep radar perception. Firstly, multi-path effects 1399 need to be explicitly considered in object detection. Secondly, we need to alleviate the 1400 problem of overconfidence in radar classification and estimate the uncertainty in bounding 1401 box regression. Thirdly, the fusion architecture should have adaptive mechanisms to take 1402 full advantage of radar's all-weather capabilities. Finally, some future research directions 1403 are proposed. There is an urgent need for high quality radar datasets. Radar domain 1404 knowledge and uncertainty quantification can help us to develop a generalizable AI model. 1405 Considering the perceptual system as a whole, we can extend the end-to-end learning 1406 framework forward, *i.e.*, joint learning with interference mitigation, or backward, *i.e.*, 1407 predicting motion. 1408

Author Contributions: Conceptualization, Y.Z., L.L., H.Z.; investigation, Y.Z.; writing—original draft preparation, Y.Z.; writing-review and editing, Y.Z.; visualization, Z.Y.; supervision, M.LB., L.Y., 1410 Y.Y.; project administration, Y.Y. All authors have read and agreed to the published version of the 1411 manuscript. 1412

Funding: This research was partially funded by XJTLU-JITRI Academy of Industrial Technology, 1413 Institute of Deep Perception Technology (IDPT) and the Research Enhancement Fund of XJTLU 1414 (REF-19-01-04). 1415

Conflicts of Interest: The authors declare no conflict of interest.

Abbreviations	1417
The following abbreviations are used in this manuscript:	1418

1416

RA	Range-Azimuth
RD	Range-Doppler
RAD	Range-Azimuth-Doppler
RCS	Radar Cross Section
FoV	Field of View
BEV	Bird's Eye View
FMCW	Frequency-Modulated Continuous-Wave
IF	Intermediate Frequency
TDM	Time-Division Multiplexing
DDM	Doppler-Division Multiplexing
SNR	Signal-to-Noise Ratio
VRU	Vulnerable Road User
ICR	Instantaneous Center of Rotation
OGM	Occupancy Grid Mapping
FCN	Fully Convolutional Network
ROI	Region of Interest
NMS	Non-Maximum Suppression
ODD	Operational Design Domain
RFS	Random Finite Set

# Appendix A

**Table A1.** Parameters used in the radar signal processing section.

Parameter	Meaning	Parameter	Meaning
С	Light Speed $(m/s^2)$	NTr	Number of Tx
-	8 of con (, c )	- • 1 λ	Antennas
λ	Wavelength (m)	Np	Number of Rx
74	wavelengur (III)	INKX	Antennas
£	Carrier Frequency	đ	Inter Antenna
fc	(Hz)	u	Spacing (m)
В	Sweep Bandwidth	D	Array Aperture (m)
D	(dB)	2	Thirdy Tip ercene (iii)
S	Chirp Slope	Р,	Transmit Power
5	Chilp Slope	1 t	(dBW)
$N_c$	Number of Chirps	G	Antenna Gain (dB)
T	China Dunatain (a)	מ	Minimum Detectable
1 <sub>C</sub>	Chirp Duratoin (s)	P <sub>min</sub>	Power (dBw)
T	Encode Diversities (a)	_	Radar Cross Section
$I_f$	Frame Duration (s)	σ	(dBm <sup>2</sup> )
P	IE Ponduridth (dP)	CNIP	Signal-to-Noise
DIF	ir bandwidth (db)	SINK	Ration

# References

- 1. Karpathy, A. Keynotes at CVPR Workshop on Autonomous Driving. https://cvpr2021.wad.vision/, 2021.
- Meyer, M.; Kuschk, G. Automotive radar dataset for deep learning based 3d object detection. 2019 16th European Radar Conference (EuRAD). IEEE, 2019, pp. 129–132.
- Zhou, T.; Yang, M.; Jiang, K.; Wong, H.; Yang, D. MMW Radar-Based Technologies in Autonomous Driving: A Review. Sensors 1420 2020, 20, 7283.
- Abdu, F.J.; Zhang, Y.; Fu, M.; Li, Y.; Deng, Z. Application of Deep Learning on Millimeter-Wave Radar Signals: A Review. Sensors 1428 2021, 21, 1951.
- Scheiner, N.; Weishaupt, F.; Tilly, J.F.; Dickmann, J. New Challenges for Deep Neural Networks in Automotive Radar Perception. 1430 In Automatisiertes Fahren 2020; Springer, 2021; pp. 165–182.
- Wei, Z.; Zhang, F.; Chang, S.; Liu, Y.; Wu, H.; Feng, Z. MmWave Radar and Vision Fusion for Object Detection in Autonomous
   Driving: A Review. arXiv preprint arXiv:2108.03004 2021.
- Tang, X.; Zhang, Z.; Qin, Y. On-road object detection and tracking based on radar and vision fusion: A review. IEEE Intelligent Transportation Systems Magazine 2021.
- Ravindran, R.; Santora, M.J.; Jamali, M.M. Multi-Object Detection and Tracking, Based on DNN, for Autonomous Vehicles: A Review. IEEE Sensors Journal 2021, 21, 5668–5677.

1420

1421

9.	Hakobyan, G.; Yang, B. High-performance automotive radar: A review of signal processing algorithms and modulation schemes. <i>IEEE Signal Processing Magazine</i> <b>2019</b> , <i>36</i> , 32–44.	1438 1439
10.	Ramasubramanian, K.; Instruments, T. Using a complex-baseband architecture in FMCW radar systems. <i>Texas Instruments</i> 2017,	1440
11	13. Pag S MIMO Padar Amplication Parart SWP 15511 Targe Instruments 2017	1441
11. 12	Kao, S. WIINO Radai. Application Report SWRASS4A, lexus instruments 2017.	1442
12.	IEEE, 2014, pp. 1–6.	1443 1444
13.	Sun, S.; Petropulu, A.P.; Poor, H.V. MIMO radar for advanced driver-assistance systems and autonomous driving: Advantages and challenges. <i>IEEE Signal Processing Magazine</i> <b>2020</b> , <i>37</i> , 98–117.	1445 1446
14.	Bechter, J.; Roos, F.; Waldschmidt, C. Compensation of motion-induced phase errors in TDM MIMO radars. <i>IEEE Microwave and Wireless Components Letters</i> <b>2017</b> , 27, 1164–1166.	1447 1448
15.	Gupta, J. High-End Corner Radar Reference Design. Design Guide TIDEP-01027, Texas Instruments 2022.	1449
16.	Rebut, J.; Ouaknine, A.; Malik, W.; Pérez, P. RADIal Dataset. https://github.com/valeoai/RADIal, 2022.	1450
17.	Richards, M.A. Fundamentals of radar signal processing; Tata McGraw-Hill Education, 2005.	1451
18.	Schubert, E.; Sander, J.; Ester, M.; Kriegel, H.P.; Xu, X. DBSCAN revisited, revisited: why and how you should (still) use DBSCAN. <i>ACM Transactions on Database Systems (TODS)</i> <b>2017</b> , <i>42</i> , 1–21.	1452 1453
19.	Muckenhuber, S.: Museliic, E.: Stettinger, G. Performance evaluation of a state-of-the-art automotive radar and corresponding	1454
	modeling approaches based on a large labeled dataset. <i>Journal of Intelligent Transportation Systems</i> <b>2021</b> , pp. 1–20.	1455
20.	Gamba, J. Radar signal processing for autonomous driving; Springer, 2020.	1456
21.	Dham, V. Programming chirp parameters in TI radar devices. <i>Application Report SWRA553, Texas Instruments</i> <b>2017</b> .	1457
22.	Hasch, J.; Topak, E.; Schnabel, R.; Zwick, T.; Weigel, R.; Waldschmidt, C. Millimeter-wave technology for automotive radar	1458
	sensors in the 77 GHz frequency band. IEEE Transactions on Microwave Theory and Techniques 2012, 60, 845–860.	1459
23.	Lim, T.Y.; Markowitz, S.; Do, M.N. RaDICaL Dataset SDK. https://github.com/moodoki/radical_sdk, 2021.	1460
24.	Lim, T.Y.; Markowitz, S.; Do, M.N. IWR Raw ROS Node. https://github.com/moodoki/iwr_raw_rosnode, 2021.	1461
25.	Mostafa, A. pyRAPID. http://radar.alizadeh.ca, 2020.	1462
26.	Pan, E.; Tang, J.; Kosaka, D.; Yao, R.; Gupta, A. OpenRadar. https://github.com/presenseradar/openradar, 2019.	1463
27.	Constapel, M.; Cimdins, M.; Hellbrück, H. A Practical Toolbox for Getting Started with mmWave FMCW Radar Sensors.	1464
	Proceedings of the 4th KuVS/GI Expert Talk on Localization, 2019.	1465
28.	Gusland, D.; Christiansen, J.M.; Torvik, B.; Fioranelli, F.; Gurbuz, S.Z.; Ritchie, M. Open Radar Initiative: Large Scale Dataset for	1466
29.	Benchmarking of micro-Doppler Recognition Algorithms. 2021 IEEE Radar Conference (RadarConf21). IEEE, 2021, pp. 1–6. Visentin, T. <i>Polarimetric radar for automotive applications</i> ; Vol. 90, KIT Scientific Publishing, 2019.	1467 1468
30.	Gottinger, M.; Hoffmann, M.; Christmann, M.; Schütz, M.; Kirsch, F.; Gulden, P.; Vossiek, M. Coherent automotive radar networks:	1469
	The next generation of radar-based imaging and mapping. IEEE Journal of Microwaves 2021, 1, 149–163.	1470
31.	Laribi, A.; Hahn, M.; Dickmann, J.; Waldschmidt, C. Performance investigation of automotive SAR imaging. 2018 IEEE MTT-S	1471
	International Conference on Microwaves for Intelligent Mobility (ICMIM). IEEE, 2018, pp. 1-4.	1472
32.	Adams, M.; Adams, M.D.; Jose, E. Robotic navigation and mapping with radar; Artech House, 2012.	1473
33.	Qian, K.; Zhu, S.; Zhang, X.; Li, L.E. Robust Multimodal Vehicle Detection in Foggy Weather Using Complementary Lidar and	1474
	Radar Signals. Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2021, pp. 444–453.	1475
34.	Caesar, H.; Bankiti, V.; Lang, A.H.; Vora, S.; Liong, V.E.; Xu, Q.; Krishnan, A.; Pan, Y.; Baldan, G.; Beijbom, O. nuscenes:	1476
	A multimodal dataset for autonomous driving. Proceedings of the IEEE/CVF conference on computer vision and pattern	1477
	recognition, 2020, pp. 11621–11631.	1478
35.	Déziel, J.L.; Merriaux, P.; Tremblay, F.; Lessard, D.; Plourde, D.; Stanguennec, J.; Goulet, P.; Olivier, P. PixSet: An Opportunity for	1479
	3D Computer Vision to Go Beyond Point Clouds With a Full-Waveform LiDAR Dataset. <i>arXiv preprint arXiv:2102.12010</i> <b>2021</b> .	1480
36.	Schumann, O.; Hahn, M.; Scheiner, N.; Weishaupt, F.; Tilly, J.F.; Dickmann, J.; Wöhler, C. RadarScenes: A real-world radar point cloud data set for automotive applications. <i>arXiv preprint arXiv:</i> 2104.02493 <b>2021</b> .	1481 1482
37.	Bansal, K.; Rungta, K.; Zhu, S.; Bharadia, D. Pointillism: Accurate 3d bounding box estimation with multi-radars. Proceedings of the 18th Conference on Embedded Networked Sensor Systems, 2020, pp. 340–353.	1483 1484
38.	Mostajabi, M.; Wang, C.M.; Ranjan, D.; Hsyu, G. High-Resolution Radar Dataset for Semi-Supervised Learning of Dynamic	1485
	Objects. Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops, 2020, pp. 100–101.	1486
39.	Ouaknine, A.; Newson, A.; Rebut, J.; Tupin, F.; Pérez, P. CARRADA Dataset: Camera and Automotive Radar with Range-Angle-	1487
	Doppler Annotations. arXiv preprint arXiv:2005.01456 2020.	1488
40.	Zhang, A.; Nowruzi, F.E.; Laganiere, R. RADDet: Range-Azimuth-Doppler based radar object detection for dynamic road users.	1489
	2021 18th Conference on Robots and Vision (CRV). IEEE, 2021, pp. 95–102.	1490
41.	Wang, Y.; Jiang, Z.; Li, Y.; Hwang, J.N.; Xing, G.; Liu, H. RODNet: A Real-Time Radar Object Detection Network Cross-Supervised	1491
	by Camera-Radar Fused Object 3D Localization. IEEE Journal of Selected Topics in Signal Processing 2021, 15, 954–967.	1492
42.	Lim, T.Y.; Markowitz, S.; Do, M.N. RaDICaL: A Synchronized FMCW Radar, Depth, IMU and RGB Camera Data Dataset with	1493
	Low-Level FMCW Radar Signals. IEEE Journal of Selected Topics in Signal Processing 2021.	1494

- Dimitrievski, M.; Shopovska, I.; Van Hamme, D.; Veelaert, P.; Philips, W. Weakly supervised deep learning method for vulnerable road user detection in FMCW radar. 2020 IEEE 23rd International Conference on Intelligent Transportation Systems (ITSC). IEEE, 2020, pp. 1–8.
- Palffy, A.; Pool, E.; Baratam, S.; Kooij, J.; Gavrila, D. Multi-class Road User Detection with 3+ 1D Radar in the View-of-Delft Dataset. IEEE Robotics and Automation Letters 2022.
- Rebut, J.; Ouaknine, A.; Malik, W.; Pérez, P. Raw High-Definition Radar for Multi-Task Learning. arXiv preprint arXiv:2112.10646 1500
   2021.
- Zheng, L.; Ma, Z.; Zhu, X.; Tan, B.; Li, S.; Long, K.; Sun, W.; Chen, S.; Zhang, L.; Wan, M.; et al. TJ4DRadSet: A 4D Radar Dataset for Autonomous Driving. arXiv preprint arXiv:2204.13483 2022.
- 47. Barnes, D.; Gadd, M.; Murcutt, P.; Newman, P.; Posner, I. The oxford radar robotcar dataset: A radar extension to the oxford 1504 robotcar dataset. 2020 IEEE International Conference on Robotics and Automation (ICRA). IEEE, 2020, pp. 6433–6438.
- 48. Kim, G.; Park, Y.S.; Cho, Y.; Jeong, J.; Kim, A. Mulran: Multimodal range dataset for urban place recognition. 2020 IEEE 1506 International Conference on Robotics and Automation (ICRA). IEEE, 2020, pp. 6246–6253.
- Burnett, K.; Yoon, D.J.; Wu, Y.; Li, A.Z.; Zhang, H.; Lu, S.; Qian, J.; Tseng, W.K.; Lambert, A.; Leung, K.Y.; et al. Boreas: A Multi-Season Autonomous Driving Dataset. arXiv preprint arXiv:2203.10168 2022.
- Yan, Z.; Sun, L.; Krajník, T.; Ruichek, Y. EU long-term dataset with multiple sensors for autonomous driving. 2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS). IEEE, 2020, pp. 10697–10704.
- Huang, R.; Zhu, K.; Chen, S.; Xiao, T.; Yang, M.; Zheng, N. A High-precision and Robust Odometry Based on Sparse MMW
   Radar Data and A Large-range and Long-distance Radar Positioning Data Set. 2021 IEEE International Intelligent Transportation
   Systems Conference (ITSC). IEEE, 2021, pp. 98–105.
- 52. Kramer, A.; Harlow, K.; Williams, C.; Heckman, C. ColoRadar: The Direct 3D Millimeter Wave Radar Dataset. *arXiv preprint* 1515 *arXiv:*2103.04510 2021. 1516
- Izquierdo, R.; Quintanar, A.; Parra, I.; Fernández-Llorca, D.; Sotelo, M. The prevention dataset: a novel benchmark for prediction of vehicles intentions. 2019 IEEE Intelligent Transportation Systems Conference (ITSC). IEEE, 2019, pp. 3114–3121.
- Nowruzi, F.E.; Kolhatkar, D.; Kapoor, P.; Al Hassanat, F.; Heravi, E.J.; Laganiere, R.; Rebut, J.; Malik, W. Deep open space segmentation using automotive radar. 2020 IEEE MTT-S International Conference on Microwaves for Intelligent Mobility (ICMIM). IEEE, 2020, pp. 1–4.
- 55. Kraus, F.; Scheiner, N.; Ritter, W.; Dietmayer, K. The Radar Ghost Dataset–An Evaluation of Ghost Objects in Automotive Radar
   Data. 2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS). IEEE, 2021, pp. 8570–8577.
- Sakaridis, C.; Dai, D.; Van Gool, L. ACDC: The adverse conditions dataset with correspondences for semantic driving scene understanding. Proceedings of the IEEE/CVF International Conference on Computer Vision, 2021, pp. 10765–10775.
- Kenk, M.A.; Hassaballah, M. DAWN: vehicle detection in adverse weather nature dataset. *arXiv preprint arXiv:2008.05402* 2020.
   Jin, J.; Fatemi, A.; Lira, W.M.P.; Yu, F.; Leng, B.; Ma, R.; Mahdavi-Amiri, A.; Zhang, H. Raidar: A rich annotated image dataset of 1527
- rainy street scenes. Proceedings of the IEEE/CVF International Conference on Computer Vision, 2021, pp. 2951–2961.
   Pitropov, M.; Garcia, D.E.; Rebello, J.; Smart, M.; Wang, C.; Czarnecki, K.; Waslander, S. Canadian adverse driving conditions
   dataset. *The International Journal of Robotics Research* 2021, 40, 681–690.
- Bijelic, M.; Gruber, T.; Mannan, F.; Kraus, F.; Ritter, W.; Dietmayer, K.; Heide, F. Seeing through fog without seeing fog: Deep 1531 multimodal sensor fusion in unseen adverse weather. Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2020, pp. 11682–11692.
- Sheeny, M.; De Pellegrin, E.; Mukherjee, S.; Ahrabian, A.; Wang, S.; Wallace, A. RADIATE: A Radar Dataset for Automotive Perception. arXiv preprint arXiv:2010.09076 2020.
- El Natour, G.; Aider, O.A.; Rouveure, R.; Berry, F.; Faure, P. Radar and vision sensors calibration for outdoor 3D reconstruction. 1536 2015 IEEE International Conference on Robotics and Automation (ICRA). IEEE, 2015, pp. 2084–2089.
- Peršić, J.; Marković, I.; Petrović, I. Extrinsic 6dof calibration of a radar–lidar–camera system enhanced by radar cross section <sup>1538</sup> estimates evaluation. *Robotics and Autonomous Systems* 2019, 114, 217–230.
- 64. Domhof, J.; Kooij, J.F.; Gavrila, D.M. An extrinsic calibration tool for radar, camera and lidar. 2019 International Conference on Robotics and Automation (ICRA). IEEE, 2019, pp. 8107–8113.
- Geiger, A.; Moosmann, F.; Car, Ö.; Schuster, B. Automatic camera and range sensor calibration using a single shot. 2012 IEEE 1542 international conference on robotics and automation. IEEE, 2012, pp. 3936–3943.
- Dhall, A.; Chelani, K.; Radhakrishnan, V.; Krishna, K.M. LiDAR-camera calibration using 3D-3D point correspondences. arXiv 1544 preprint arXiv:1705.09785 2017.
- Sugimoto, S.; Tateda, H.; Takahashi, H.; Okutomi, M. Obstacle detection using millimeter-wave radar and its visualization on 1546 image sequence. Proceedings of the 17th International Conference on Pattern Recognition, 2004. ICPR 2004. IEEE, 2004, Vol. 3, 1547 pp. 342–345.
- Peršić, J.; Petrović, L.; Marković, I.; Petrović, I. Spatio-temporal multisensor calibration based on gaussian processes moving <sup>1549</sup> object tracking. arXiv preprint arXiv:1904.04187 2019.
- Peršić, J.; Petrović, L.; Marković, I.; Petrović, I. Online multi-sensor calibration based on moving object tracking. Advanced Robotics 1551 2021, 35, 130–140.

- Faizullin, M.; Kornilova, A.; Ferrer, G. Open-Source LiDAR Time Synchronization System by Mimicking GPS-clock. arXiv preprint arXiv:2107.02625 2021.
- Kaul, P.; De Martini, D.; Gadd, M.; Newman, P. Rss-net: Weakly-supervised multi-class semantic segmentation with FMCW radar. 2020 IEEE Intelligent Vehicles Symposium (IV). IEEE, 2020, pp. 431–436.
- Shoemake, K. Animating rotation with quaternion curves. Proceedings of the 12th annual conference on Computer graphics and interactive techniques, 1985, pp. 245–254.
- He, K.; Gkioxari, G.; Dollár, P.; Girshick, R. Mask r-cnn. Proceedings of the IEEE international conference on computer vision, 1550 2017, pp. 2961–2969.
- 74. Chen, L.C.; Papandreou, G.; Schroff, F.; Adam, H. Rethinking atrous convolution for semantic image segmentation. *arXiv preprint* 1561 *arXiv:*1706.05587 2017.
- Grimm, C.; Fei, T.; Warsitz, E.; Farhoud, R.; Breddermann, T.; Haeb-Umbach, R. Warping of Radar Data into Camera Image for Cross-Modal Supervision in Automotive Applications. arXiv preprint arXiv:2012.12809 2020.
- Dimitrievski, M.; Shopovska, I.; Van Hamme, D.; Veelaert, P.; Philips, W. Automatic labeling of vulnerable road users in multi-sensor data. 2021 IEEE International Intelligent Transportation Systems Conference (ITSC). IEEE, 2021, pp. 2623–2630.
- 77. Settles, B. Active learning literature survey 2009.
- 78. Shorten, C.; Khoshgoftaar, T.M. A survey on image data augmentation for deep learning. *Journal of Big Data* 2019, *6*, 1–48.
- Hahner, M.; Dai, D.; Liniger, A.; Van Gool, L. Quantifying data augmentation for lidar based 3d object detection. arXiv preprint arXiv:2004.01643 2020.
- Park, D.S.; Chan, W.; Zhang, Y.; Chiu, C.C.; Zoph, B.; Cubuk, E.D.; Le, Q.V. SpecAugment: A Simple Data Augmentation Method for Automatic Speech Recognition. *Proc. Interspeech* 2019 2019, pp. 2613–2617.
- Wang, Y.; Hwang, J.N.; Wang, G.; Liu, H.; Kim, K.J.; Hsu, H.M.; Cai, J.; Zhang, H.; Jiang, Z.; Gu, R. ROD2021 Challenge: A 1573 Summary for Radar Object Detection Challenge for Autonomous Driving Applications. Proceedings of the 2021 International 1574 Conference on Multimedia Retrieval, 2021, pp. 553–559.
- Ju, B.; Yang, W.; Jia, J.; Ye, X.; Chen, Q.; Tan, X.; Sun, H.; Shi, Y.; Ding, E. DANet: Dimension Apart Network for Radar Object
   Detection. Proceedings of the 2021 International Conference on Multimedia Retrieval, 2021, pp. 533–539.
- Sheeny, M.; Wallace, A.; Wang, S. Radio: parameterized generative radar data augmentation for small datasets. *Applied Sciences* 1578 2020, 10, 3861.
- Bing, J.; Chen, B.; Liu, H.; Huang, M. Convolutional neural network with data augmentation for SAR target recognition. *IEEE* 1580 *Geoscience and remote sensing letters* 2016, 13, 364–368.
- Gao, X.; Xing, G.; Roy, S.; Liu, H. Ramp-cnn: A novel neural network for enhanced automotive radar object recognition. *IEEE* 1582 Sensors Journal 2020, 21, 5119–5132.
- Wang, C.; Ma, C.; Zhu, M.; Yang, X. Pointaugmenting: Cross-modal augmentation for 3d object detection. Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2021, pp. 11794–11803.
- 87. Zhang, W.; Wang, Z.; Change Loy, C. Multi-modality cut and paste for 3d object detection. arXiv e-prints 2020, pp. arXiv-2012. 1596
- Chamseddine, M.; Rambach, J.; Stricker, D.; Wasenmuller, O. Ghost Target Detection in 3D Radar Data using Point Cloud based Deep Neural Network. 2020 25th International Conference on Pattern Recognition (ICPR). IEEE, 2021, pp. 10398–10403.
- Berthold, P.; Michaelis, M.; Luettel, T.; Meissner, D.; Wuensche, H.J. Radar reflection characteristics of vehicles for contour and feature estimation. 2017 Sensor Data Fusion: Trends, Solutions, Applications (SDF). IEEE, 2017, pp. 1–6.
- 90. Long, Y.; Morris, D.; Liu, X.; Castro, M.; Chakravarty, P.; Narayanan, P. Full-Velocity Radar Returns by Radar-Camera Fusion. 1591 Proceedings of the IEEE/CVF International Conference on Computer Vision, 2021, pp. 16198–16207.
- 91. Cabon, Y.; Murray, N.; Humenberger, M. Virtual kitti 2. arXiv preprint arXiv:2001.10773 2020.
- 92. Tremblay, J.; Prakash, A.; Acuna, D.; Brophy, M.; Jampani, V.; Anil, C.; To, T.; Cameracci, E.; Boochoon, S.; Birchfield, S. Training deep networks with synthetic data: Bridging the reality gap by domain randomization. Proceedings of the IEEE conference on computer vision and pattern recognition workshops, 2018, pp. 969–977.
- Hurl, B.; Czarnecki, K.; Waslander, S. Precise synthetic image and lidar (presil) dataset for autonomous vehicle perception. 2019 IEEE Intelligent Vehicles Symposium (IV). IEEE, 2019, pp. 2522–2529.
- P4. Rong, G.; Shin, B.H.; Tabatabaee, H.; Lu, Q.; Lemke, S.; Možeiko, M.; Boise, E.; Uhm, G.; Gerow, M.; Mehta, S.; et al. Lgsvl 1509 simulator: A high fidelity simulator for autonomous driving. 2020 IEEE 23rd International conference on intelligent transportation 1600 systems (ITSC). IEEE, 2020, pp. 1–6.
- 95. Johnson-Roberson, M.; Barto, C.; Mehta, R.; Sridhar, S.N.; Rosaen, K.; Vasudevan, R. Driving in the matrix: Can virtual worlds replace human-generated annotations for real world tasks? *arXiv preprint arXiv:1610.01983* **2016**.
- Wang, J.; Pun, A.; Tu, J.; Manivasagam, S.; Sadat, A.; Casas, S.; Ren, M.; Urtasun, R. Advsim: Generating safety-critical scenarios for self-driving vehicles. Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2021, pp. 9909–9918.
- Schüßler, C.; Hoffmann, M.; Bräunig, J.; Ullmann, I.; Ebelt, R.; Vossiek, M. A Realistic Radar Ray Tracing Simulator for Large MIMO-Arrays in Automotive Environments. *IEEE Journal of Microwaves* 2021, 1, 962–974.
- Holder, M.; Rosenberger, P.; Winner, H.; D'hondt, T.; Makkapati, V.P.; Maier, M.; Schreiber, H.; Magosi, Z.; Slavik, Z.; Bringmann, 1609
   O.; et al. Measurements revealing challenges in radar sensor modeling for virtual validation of autonomous driving. 2018 21st 1610
   International Conference on Intelligent Transportation Systems (ITSC). IEEE, 2018, pp. 2616–2622.

1567

- Scheel, A.; Dietmayer, K. Tracking multiple vehicles using a variational radar model. *IEEE Transactions on Intelligent Transportation* Systems 2018, 20, 3721–3736.
- Xia, Y.; Wang, P.; Berntorp, K.; Svensson, L.; Granström, K.; Mansour, H.; Boufounos, P.; Orlik, P.V. Learning-Based Extended Object Tracking Using Hierarchical Truncation Measurement Model With Automotive Radar. *IEEE Journal of Selected Topics in Signal Processing* 2021, 15, 1013–1029.
- 101. Ngo, A.; Bauer, M.P.; Resch, M. A Multi-Layered Approach for Measuring the Simulation-to-Reality Gap of Radar Perception for 1617 Autonomous Driving. arXiv preprint arXiv:2106.08372 2021.
- Wheeler, T.A.; Holder, M.; Winner, H.; Kochenderfer, M.J. Deep stochastic radar models. 2017 IEEE Intelligent Vehicles 1619 Symposium (IV). IEEE, 2017, pp. 47–53.
- Wang, L.; Goldluecke, B.; Anklam, C. L2R GAN: LiDAR-to-radar translation. Proceedings of the Asian Conference on Computer Vision, 2020.
- 104. Lekic, V.; Babic, Z. Automotive radar and camera fusion using generative adversarial networks. Computer Vision and Image 1623 Understanding 2019, 184, 1–8.
- Ditzel, C.; Dietmayer, K. GenRadar: Self-Supervised Probabilistic Camera Synthesis Based on Radar Frequencies. *IEEE Access* 1025
   2021, 9, 148994–149042.
- 106. Wang, Y.; Chao, W.L.; Garg, D.; Hariharan, B.; Campbell, M.; Weinberger, K.Q. Pseudo-lidar from visual depth estimation:
   Bridging the gap in 3d object detection for autonomous driving. Proceedings of the IEEE/CVF Conference on Computer Vision
   and Pattern Recognition, 2019, pp. 8445–8453.
- 107. Weng, X.; Kitani, K. Monocular 3d object detection with pseudo-lidar point cloud. Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops, 2019, pp. 0–0.
- Qian, R.; Garg, D.; Wang, Y.; You, Y.; Belongie, S.; Hariharan, B.; Campbell, M.; Weinberger, K.Q.; Chao, W.L. End-to-end pseudo-lidar for image-based 3d object detection. Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2020, pp. 5881–5890.
- Fu, H.; Gong, M.; Wang, C.; Batmanghelich, K.; Tao, D. Deep ordinal regression network for monocular depth estimation.
   Proceedings of the IEEE conference on computer vision and pattern recognition, 2018, pp. 2002–2011.
- Godard, C.; Mac Aodha, O.; Firman, M.; Brostow, G.J. Digging into self-supervised monocular depth estimation. Proceedings of the IEEE/CVF International Conference on Computer Vision, 2019, pp. 3828–3838.
- 111. Dijk, T.v.; Croon, G.d. How do neural networks see depth in single images? Proceedings of the IEEE/CVF International Conference on Computer Vision, 2019, pp. 2183–2191.
- Gasperini, S.; Koch, P.; Dallabetta, V.; Navab, N.; Busam, B.; Tombari, F. R4Dyn: Exploring radar for self-supervised monocular depth estimation of dynamic scenes. 2021 International Conference on 3D Vision (3DV). IEEE, 2021, pp. 751–760.
- Lin, J.T.; Dai, D.; Van Gool, L. Depth estimation from monocular images and sparse radar data. 2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS). IEEE, 2020, pp. 10233–10240.
- Lo, C.C.; Vandewalle, P. Depth Estimation From Monocular Images And Sparse Radar Using Deep Ordinal Regression Network. 2021 IEEE International Conference on Image Processing (ICIP). IEEE, 2021, pp. 3343–3347.
- Long, Y.; Morris, D.; Liu, X.; Castro, M.; Chakravarty, P.; Narayanan, P. Radar-camera pixel depth association for depth completion.
   Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2021, pp. 12507–12516.
- Lee, W.Y.; Jovanov, L.; Philips, W. Semantic-guided radar-vision fusion for depth estimation and object detection. BMVC, the 32nd British Machine Vision Conference, Proceedings. BMVA Press, 2021, p. 13.
- Niesen, U.; Unnikrishnan, J. Camera-Radar Fusion for 3-D Depth Reconstruction. 2020 IEEE Intelligent Vehicles Symposium (IV). 1651 IEEE, 2020, pp. 265–271.
- Kramer, A.; Stahoviak, C.; Santamaria-Navarro, A.; Agha-Mohammadi, A.A.; Heckman, C. Radar-inertial ego-velocity estimation for visually degraded environments. 2020 IEEE International Conference on Robotics and Automation (ICRA). IEEE, 2020, pp. 1654 5739–5746.
- Cen, S.H.; Newman, P. Precise ego-motion estimation with millimeter-wave radar under diverse and challenging conditions.
   2018 IEEE International Conference on Robotics and Automation (ICRA). IEEE, 2018, pp. 6045–6052.
- Kellner, D.; Barjenbruch, M.; Dietmayer, K.; Klappstein, J.; Dickmann, J. Instantaneous lateral velocity estimation of a vehicle using Doppler radar. Proceedings of the 16th International Conference on Information Fusion. IEEE, 2013, pp. 877–884.
- Schubert, R.; Richter, E.; Wanielik, G. Comparison and evaluation of advanced motion models for vehicle tracking. 2008 11th international conference on information fusion. IEEE, 2008, pp. 1–6.
- Kellner, D.; Barjenbruch, M.; Klappstein, J.; Dickmann, J.; Dietmayer, K. Instantaneous full-motion estimation of arbitrary objects using dual Doppler radar. 2014 IEEE Intelligent Vehicles Symposium Proceedings. IEEE, 2014, pp. 324–329.
- Schlichenmaier, J.; Yan, L.; Stolz, M.; Waldschmidt, C. Instantaneous actual motion estimation with a single high-resolution radar sensor. 2018 IEEE MTT-S International Conference on Microwaves for Intelligent Mobility (ICMIM). IEEE, 2018, pp. 1–4.
- 124. Ding, F.; Pan, Z.; Deng, Y.; Deng, J.; Lu, C.X. Self-Supervised Scene Flow Estimation with 4D Automotive Radar. *arXiv preprint arXiv:2203.01137* **2022**.
- 125. Sun, D.; Yang, X.; Liu, M.Y.; Kautz, J. Pwc-net: Cnns for optical flow using pyramid, warping, and cost volume. Proceedings of the IEEE conference on computer vision and pattern recognition, 2018, pp. 8934–8943.

- 126. Kabsch, W. A solution for the best rotation to relate two sets of vectors. Acta Crystallographica Section A: Crystal Physics, Diffraction, 1670 Theoretical and General Crystallography 1976, 32, 922–923.
- 127. Cao, Z.; Fang, W.; Song, Y.; He, L.; Song, C.; Xu, Z. DNN-Based Peak Sequence Classification CFAR Detection Algorithm for
   High-Resolution FMCW Radar. *IEEE Transactions on Geoscience and Remote Sensing* 2021.
- Lin, C.H.; Lin, Y.C.; Bai, Y.; Chung, W.H.; Lee, T.S.; Huttunen, H. DL-CFAR: A Novel CFAR target detection method based on deep learning. 2019 IEEE 90th Vehicular Technology Conference (VTC2019-Fall). IEEE, 2019, pp. 1–6.
- Scheiner, N.; Schumann, O.; Kraus, F.; Appenrodt, N.; Dickmann, J.; Sick, B. Off-the-shelf sensor vs. experimental radar-How much resolution is necessary in automotive radar classification? 2020 IEEE 23rd International Conference on Information Fusion (FUSION). IEEE, 2020, pp. 1–8.
- Pedregosa, F.; Varoquaux, G.; Gramfort, A.; Michel, V.; Thirion, B.; Grisel, O.; Blondel, M.; Prettenhofer, P.; Weiss, R.; Dubourg, V.; 1679 et al. Comparing different clustering algorithms on toy datasets. https://scikit-learn.org/0.15/auto\_examples/cluster/plot\_ cluster\_comparison.html#example-cluster-plot-cluster-comparison-py.
- Kellner, D.; Klappstein, J.; Dietmayer, K. Grid-based DBSCAN for clustering extended objects in radar data. 2012 IEEE Intelligent Vehicles Symposium. IEEE, 2012, pp. 365–370.
- Scheiner, N.; Appenrodt, N.; Dickmann, J.; Sick, B. A multi-stage clustering framework for automotive radar data. 2019 IEEE 1684 Intelligent Transportation Systems Conference (ITSC). IEEE, 2019, pp. 2060–2067.
- Angelov, A.; Robertson, A.; Murray-Smith, R.; Fioranelli, F. Practical classification of different moving targets using automotive radar and deep neural networks. *IET Radar, Sonar & Navigation* 2018, 12, 1082–1089.
- 134. Gao, X.; Xing, G.; Roy, S.; Liu, H. Experiments with mmwave automotive radar test-bed. 2019 53rd Asilomar Conference on Signals, Systems, and Computers. IEEE, 2019, pp. 1–6.
- Cai, X.; Giallorenzo, M.; Sarabandi, K. Machine Learning-Based Target Classification for MMW Radar in Autonomous Driving. *IEEE Transactions on Intelligent Vehicles* 2021, 6, 678–689.
- Scheiner, N.; Appenrodt, N.; Dickmann, J.; Sick, B. Radar-based road user classification and novelty detection with recurrent neural network ensembles. 2019 IEEE Intelligent Vehicles Symposium (IV). IEEE, 2019, pp. 722–729.
- Schumann, O.; Wöhler, C.; Hahn, M.; Dickmann, J. Comparison of random forest and long short-term memory network 1694 performances in classification tasks using radar. 2017 Sensor Data Fusion: Trends, Solutions, Applications (SDF). IEEE, 2017, pp. 1695 1–6.
- Scheiner, N.; Appenrodt, N.; Dickmann, J.; Sick, B. Radar-based feature design and multiclass classification for road user recognition. 2018 IEEE Intelligent Vehicles Symposium (IV). IEEE, 2018, pp. 779–786.
- 139. Graham, B.; van der Maaten, L. Submanifold sparse convolutional networks. arXiv preprint arXiv:1706.01307 2017.
- Dreher, M.; Erçelik, E.; Bänziger, T.; Knol, A. Radar-based 2D Car Detection Using Deep Neural Networks. 2020 IEEE 23rd International Conference on Intelligent Transportation Systems (ITSC). IEEE, 2020, pp. 1–8.
- 141. Redmon, J.; Farhadi, A. Yolov3: An incremental improvement. arXiv preprint arXiv:1804.02767 2018.
- Schumann, O.; Hahn, M.; Dickmann, J.; Wöhler, C. Semantic segmentation on radar point clouds. 2018 21st International Conference on Information Fusion (FUSION). IEEE, 2018, pp. 2179–2186.
- Danzer, A.; Griebel, T.; Bach, M.; Dietmayer, K. 2d car detection in radar data with pointnets. 2019 IEEE Intelligent Transportation Systems Conference (ITSC). IEEE, 2019, pp. 61–66.
- Scheiner, N.; Kraus, F.; Appenrodt, N.; Dickmann, J.; Sick, B. Object detection for automotive radar point clouds-a comparison. 1707 AI Perspectives 2021, 3, 1–23.
- 145. Qi, C.R.; Su, H.; Mo, K.; Guibas, L.J. Pointnet: Deep learning on point sets for 3d classification and segmentation. Proceedings of the IEEE conference on computer vision and pattern recognition, 2017, pp. 652–660.
- 146. Qi, C.R.; Yi, L.; Su, H.; Guibas, L.J. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. Advances in neural information processing systems 2017, 30.
- 147. Liu, J.; Xiong, W.; Bai, L.; Xia, Y.; Huang, T.; Ouyang, W.; et al. Deep Instance Segmentation with Automotive Radar Detection Points 2022.
- 148. Liu, H.; Dai, Z.; So, D.; Le, Q. Pay attention to MLPs. Advances in Neural Information Processing Systems 2021, 34.
- Schumann, O.; Lombacher, J.; Hahn, M.; Wöhler, C.; Dickmann, J. Scene understanding with automotive radar. *IEEE Transactions* 1716 on Intelligent Vehicles 2019, 5, 188–203.
- Long, J.; Shelhamer, E.; Darrell, T. Fully convolutional networks for semantic segmentation. Proceedings of the IEEE conference on computer vision and pattern recognition, 2015, pp. 3431–3440.
- Lang, A.H.; Vora, S.; Caesar, H.; Zhou, L.; Yang, J.; Beijbom, O. Pointpillars: Fast encoders for object detection from point clouds.
   Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2019, pp. 12697–12705.
- Xu, B.; Zhang, X.; Wang, L.; Hu, X.; Li, Z.; Pan, S.; Li, J.; Deng, Y. RPFA-Net: a 4D RaDAR Pillar Feature Attention Network for 3D
   Object Detection. 2021 IEEE International Intelligent Transportation Systems Conference (ITSC). IEEE, 2021, pp. 3061–3066.
- 153. Bai, J.; Zheng, L.; Li, S.; Tan, B.; Chen, S.; Huang, L. Radar transformer: An object classification network based on 4d mmw imaging radar. Sensors 2021, 21, 3854.
- Zhao, H.; Jia, J.; Koltun, V. Exploring self-attention for image recognition. Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2020, pp. 10076–10085.

1699

1702

- Cheng, Y.; Su, J.; Chen, H.; Liu, Y. A New Automotive Radar 4D Point Clouds Detector by Using Deep Learning. ICASSP 1728 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). IEEE, 2021, pp. 8398–8402.
- 156. Gall, M.; Gardill, M.; Horn, T.; Fuchs, J. Spectrum-based single-snapshot super-resolution direction-of-arrival estimation using deep learning. 2020 German Microwave Conference (GeMiC). IEEE, 2020, pp. 184–187.
- Fuchs, J.; Gardill, M.; Lübke, M.; Dubey, A.; Lurz, F. A Machine Learning Perspective on Automotive Radar Direction of Arrival Estimation. *IEEE Access* 2022.
- 158. Brodeski, D.; Bilik, I.; Giryes, R. Deep radar detector. 2019 IEEE Radar Conference (RadarConf). IEEE, 2019, pp. 1–6.
- 159. Zhang, G.; Li, H.; Wenger, F. Object detection and 3d estimation via an FMCW radar using a fully convolutional network. ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). IEEE, 2020, pp. 4487–4491.
- 160. Mittal, S.; et al. A survey of accelerator architectures for 3D convolution neural networks. *Journal of Systems Architecture* 2021, 1737 115, 102041.
- Palffy, A.; Dong, J.; Kooij, J.F.; Gavrila, D.M. CNN based road user detection using the 3D radar cube. *IEEE Robotics and Automation Letters* 2020, 5, 1263–1270.
- Major, B.; Fontijne, D.; Ansari, A.; Teja Sukhavasi, R.; Gowaikar, R.; Hamilton, M.; Lee, S.; Grzechnik, S.; Subramanian, S. 1741
   Vehicle detection with automotive radar using deep learning on range-azimuth-doppler tensors. Proceedings of the IEEE/CVF 1742
   International Conference on Computer Vision Workshops, 2019, pp. 0–0.
- 163. Ouaknine, A.; Newson, A.; Pérez, P.; Tupin, F.; Rebut, J. Multi-View Radar Semantic Segmentation. Proceedings of the IEEE/CVF International Conference on Computer Vision, 2021, pp. 15671–15680.
- 164. Nowruzi, F.E.; Kolhatkar, D.; Kapoor, P.; Heravi, E.J.; Hassanat, F.A.; Laganiere, R.; Rebut, J.; Malik, W. PolarNet: Accelerated Deep Open Space Segmentation Using Automotive Radar in Polar Domain. arXiv preprint arXiv:2103.03387 2021.
- 165. Hayashi, E.; Lien, J.; Gillian, N.; Giusti, L.; Weber, D.; Yamanaka, J.; Bedal, L.; Poupyrev, I. Radarnet: Efficient gesture recognition technique utilizing a miniature radar sensor. Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems, 2021, pp. 1–14.
- Meyer, M.; Kuschk, G.; Tomforde, S. Graph convolutional networks for 3d object detection on radar data. Proceedings of the IEEE/CVF International Conference on Computer Vision, 2021, pp. 3060–3069.

167. Kipf, T.N.; Welling, M. Semi-supervised classification with graph convolutional networks. arXiv preprint arXiv:1609.02907 2016. 1753

- Dai, J.; Qi, H.; Xiong, Y.; Li, Y.; Zhang, G.; Hu, H.; Wei, Y. Deformable convolutional networks. Proceedings of the IEEE 1754 international conference on computer vision, 2017, pp. 764–773.
- Li, P.; Wang, P.; Berntorp, K.; Liu, H. Exploiting Temporal Relations on Radar Perception for Autonomous Driving. arXiv preprint arXiv:2204.01184 2022.
- 170. Nobis, F.; Geisslinger, M.; Weber, M.; Betz, J.; Lienkamp, M. A deep learning-based radar and camera sensor fusion architecture for object detection. 2019 Sensor Data Fusion: Trends, Solutions, Applications (SDF). IEEE, 2019, pp. 1–7.
- 171. Chadwick, S.; Maddern, W.; Newman, P. Distant vehicle detection using radar and vision. 2019 International Conference on Robotics and Automation (ICRA). IEEE, 2019, pp. 8311–8317.
- Yadav, R.; Vierling, A.; Berns, K. Radar+ RGB Fusion For Robust Object Detection In Autonomous Vehicle. 2020 IEEE International Conference on Image Processing (ICIP). IEEE, 2020, pp. 1986–1990.
- 173. Girshick, R. Fast r-cnn. Proceedings of the IEEE international conference on computer vision, 2015, pp. 1440–1448.
- 174. Nabati, R.; Qi, H. Rrpn: Radar region proposal network for object detection in autonomous vehicles. 2019 IEEE International Conference on Image Processing (ICIP). IEEE, 2019, pp. 3093–3097.
- 175. Nabati, R.; Qi, H. Radar-camera sensor fusion for joint object detection and distance estimation in autonomous vehicles. *arXiv* 1767 preprint arXiv:2009.08428 2020.
- 176. Nabati, R.; Qi, H. Centerfusion: Center-based radar and camera fusion for 3d object detection. Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision, 2021, pp. 1527–1536.
- 177. Kim, Y.; Choi, J.W.; Kum, D. GRIF Net: Gated Region of Interest Fusion Network for Robust 3D Object Detection from Radar
   Point Cloud and Monocular Image. 2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS). IEEE,
   2020, pp. 10857–10864.
- 178. Lim, T.Y.; Ansari, A.; Major, B.; Fontijne, D.; Hamilton, M.; Gowaikar, R.; Subramanian, S. Radar and camera early fusion for vehicle detection in advanced driver assistance systems. Machine Learning for Autonomous Driving Workshop at the 33rd Conference on Neural Information Processing Systems, 2019, Vol. 2.
- 179. Zhang, J.; Zhang, M.; Fang, Z.; Wang, Y.; Zhao, X.; Pu, S. RVDet: Feature-level Fusion of Radar and Camera for Object Detection.
   2021 IEEE International Intelligent Transportation Systems Conference (ITSC). IEEE, 2021, pp. 2822–2828.
- Kim, J.; Kim, Y.; Kum, D. Low-level Sensor Fusion Network for 3D Vehicle Detection using Radar Range-Azimuth Heatmap and Monocular Image. Proceedings of the Asian Conference on Computer Vision, 2020.
- Meyer, M.; Kuschk, G. Deep learning based 3d object detection for automotive radar and camera. 2019 16th European Radar Conference (EuRAD). IEEE, 2019, pp. 133–136.
- 182. Ku, J.; Mozifian, M.; Lee, J.; Harakeh, A.; Waslander, S.L. Joint 3d proposal generation and object detection from view aggregation.
   2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS). IEEE, 2018, pp. 1–8.
- Yang, B.; Guo, R.; Liang, M.; Casas, S.; Urtasun, R. Radarnet: Exploiting radar for robust perception of dynamic objects. European Conference on Computer Vision. Springer, 2020, pp. 496–512.

1734

- 184. Shah, M.; Huang, Z.; Laddha, A.; Langford, M.; Barber, B.; Zhang, S.; Vallespi-Gonzalez, C.; Urtasun, R. Liranet: End-to-end trajectory prediction using spatio-temporal radar fusion. arXiv preprint arXiv:2010.00731 2020.
- 185. Liu, Y.; Fan, Q.; Zhang, S.; Dong, H.; Funkhouser, T.; Yi, L. Contrastive multimodal fusion with tupleinfonce. Proceedings of the IEEE/CVF International Conference on Computer Vision, 2021, pp. 754–763.
- Cheng, Y.; Xu, H.; Liu, Y. Robust Small Object Detection on the Water Surface Through Fusion of Camera and Millimeter Wave Radar. Proceedings of the IEEE/CVF International Conference on Computer Vision, 2021, pp. 15263–15272.
- Hu, J.; Shen, L.; Sun, G. Squeeze-and-excitation networks. Proceedings of the IEEE conference on computer vision and pattern recognition, 2018, pp. 7132–7141.
- 188. Matzka, S.; Altendorfer, R. A comparison of track-to-track fusion algorithms for automotive sensor fusion. In *Multisensor Fusion* 1795 and Integration for Intelligent Systems; Springer, 2009; pp. 69–81.
- Dong, X.; Zhuang, B.; Mao, Y.; Liu, L. Radar Camera Fusion via Representation Learning in Autonomous Driving. Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2021, pp. 1672–1681.
- Harakeh, A.; Smart, M.; Waslander, S.L. Bayesod: A bayesian approach for uncertainty estimation in deep object detectors. 2020
   IEEE International Conference on Robotics and Automation (ICRA). IEEE, 2020, pp. 87–93.
- 191. Hüllermeier, E.; Waegeman, W. Aleatoric and epistemic uncertainty in machine learning: An introduction to concepts and methods. *Machine Learning* **2021**, *110*, 457–506.
- 192. Chavez-Garcia, R.O.; Aycard, O. Multiple sensor fusion and classification for moving object detection and tracking. *IEEE* 1803 *Transactions on Intelligent Transportation Systems* 2015, 17, 525–534.
- Florea, M.C.; Jousselme, A.L.; Bossé, É.; Grenier, D. Robust combination rules for evidence theory. *Information Fusion* 2009, 1805 10, 183–197.
- 194. Angelopoulos, A.N.; Bates, S. A gentle introduction to conformal prediction and distribution-free uncertainty quantification. 1807 arXiv preprint arXiv:2107.07511 2021.
- 195. Kopp, J.; Kellner, D.; Piroli, A.; Dietmayer, K. Fast Rule-Based Clutter Detection in Automotive Radar Data. 2021 IEEE
   International Intelligent Transportation Systems Conference (ITSC). IEEE, 2021, pp. 3010–3017.
- 196. Kraus, F.; Scheiner, N.; Ritter, W.; Dietmayer, K. Using machine learning to detect ghost images in automotive radar. 2020 IEEE 1811
   23rd International Conference on Intelligent Transportation Systems (ITSC). IEEE, 2020, pp. 1–7.
- 197. Kamann, A.; Held, P.; Perras, F.; Zaumseil, P.; Brandmeier, T.; Schwarz, U.T. Automotive radar multipath propagation in uncertain environments. 2018 21st International Conference on Intelligent Transportation Systems (ITSC). IEEE, 2018, pp. 859–864.
- 198. Griebel, T.; Authaler, D.; Horn, M.; Henning, M.; Buchholz, M.; Dietmayer, K. Anomaly Detection in Radar Data Using PointNets. 2021 IEEE International Intelligent Transportation Systems Conference (ITSC). IEEE, 2021, pp. 2667–2673.
- 199. Garcia, J.M.; Prophet, R.; Michel, J.C.F.; Ebelt, R.; Vossiek, M.; Weber, I. Identification of ghost moving detections in automotive scenarios with deep learning. 2019 IEEE MTT-S International Conference on Microwaves for Intelligent Mobility (ICMIM). IEEE, 2019, pp. 1–4.
- 200. Wang, L.; Giebenhain, S.; Anklam, C.; Goldluecke, B. Radar ghost target detection via multimodal transformers. *IEEE Robotics* 1820 and Automation Letters 2021, 6, 7758–7765.
- Guo, C.; Pleiss, G.; Sun, Y.; Weinberger, K.Q. On calibration of modern neural networks. International Conference on Machine Learning. PMLR, 2017, pp. 1321–1330.
- 202. Patel, K.; Beluch, W.; Rambach, K.; Cozma, A.E.; Pfeiffer, M.; Yang, B. Investigation of Uncertainty of Deep Learning-based Object 1824 Classification on Radar Spectra. 2021 IEEE Radar Conference (RadarConf21). IEEE, 2021, pp. 1–6.
- 203. Geng, C.; Huang, S.j.; Chen, S. Recent advances in open set recognition: A survey. *IEEE transactions on pattern analysis and machine intelligence* 2020, 43, 3614–3631.
   1826
- Hall, D.; Dayoub, F.; Skinner, J.; Zhang, H.; Miller, D.; Corke, P.; Carneiro, G.; Angelova, A.; Sünderhauf, N. Probabilistic object detection: Definition and evaluation. Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision, 2020, pp. 1031–1040.
- Patel, K.; Beluch, W.; Rambach, K.; Pfeiffer, M.; Yang, B. Improving Uncertainty of Deep Learning-based Object Classification on 1831
   Radar Spectra using Label Smoothing. arXiv preprint arXiv:2109.12851 2021.
- Wenger, J.; Kjellström, H.; Triebel, R. Non-parametric calibration for classification. International Conference on Artificial Intelligence and Statistics. PMLR, 2020, pp. 178–190.
- 207. Patel, K.; Beluch, W.H.; Yang, B.; Pfeiffer, M.; Zhang, D. Multi-Class Uncertainty Calibration via Mutual Information Maximization based Binning. International Conference on Learning Representations, 2020.
- 208. Müller, R.; Kornblith, S.; Hinton, G.E. When does label smoothing help? Advances in neural information processing systems 2019, 32. 1837
- 209. Thulasidasan, S.; Chennupati, G.; Bilmes, J.A.; Bhattacharya, T.; Michalak, S. On mixup training: Improved calibration and predictive uncertainty for deep neural networks. *Advances in Neural Information Processing Systems* **2019**, 32.
- Hendrycks, D.; Mu, N.; Cubuk, E.D.; Zoph, B.; Gilmer, J.; Lakshminarayanan, B. AugMix: A Simple Data Processing Method to Improve Robustness and Uncertainty. International Conference on Learning Representations, 2019.
- 211. Gal, Y.; Ghahramani, Z. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. international conference on machine learning. PMLR, 2016, pp. 1050–1059.
- 212. Fort, S.; Hu, H.; Lakshminarayanan, B. Deep ensembles: A loss landscape perspective. arXiv preprint arXiv:1912.02757 2019.

- 213. Feng, D.; Wang, Z.; Zhou, Y.; Rosenbaum, L.; Timm, F.; Dietmayer, K.; Tomizuka, M.; Zhan, W. Labels are not perfect: Inferring spatial uncertainty in object detection. *IEEE Transactions on Intelligent Transportation Systems* 2021.
- 214. Kendall, A.; Gal, Y. What uncertainties do we need in bayesian deep learning for computer vision? Advances in neural information 1847 processing systems 2017, 30.
- Dong, X.; Wang, P.; Zhang, P.; Liu, L. Probabilistic oriented object detection in automotive radar. Proceedings of the IEEE/CVF 1849 Conference on Computer Vision and Pattern Recognition Workshops, 2020, pp. 102–103.
- Mohammed, A.S.; Amamou, A.; Ayevide, F.K.; Kelouwani, S.; Agbossou, K.; Zioui, N. The perception system of intelligent ground vehicles in all weather conditions: a systematic literature review. *Sensors* 2020, 20, 6532.
- Hendrycks, D.; Dietterich, T. Benchmarking Neural Network Robustness to Common Corruptions and Perturbations. International Conference on Learning Representations, 2019.
- 218. Secci, F.; Ceccarelli, A. On failures of RGB cameras and their effects in autonomous driving applications. 2020 IEEE 31st 1855 International Symposium on Software Reliability Engineering (ISSRE). IEEE, 2020, pp. 13–24.
- Jokela, M.; Kutila, M.; Pyykönen, P. Testing and validation of automotive point-cloud sensors in adverse weather conditions. *Applied Sciences* 2019, 9, 2341.
- Carballo, A.; Lambert, J.; Monrroy, A.; Wong, D.; Narksri, P.; Kitsukawa, Y.; Takeuchi, E.; Kato, S.; Takeda, K. LIBRE: The multiple 3D LiDAR dataset. 2020 IEEE Intelligent Vehicles Symposium (IV). IEEE, 2020, pp. 1094–1101.
- Zang, S.; Ding, M.; Smith, D.; Tyler, P.; Rakotoarivelo, T.; Kaafar, M.A. The impact of adverse weather conditions on autonomous vehicles: how rain, snow, fog, and hail affect the performance of a self-driving car. *IEEE vehicular technology magazine* 2019, 1862 14, 103–111.
- Brooker, G.; Hennessey, R.; Lobsey, C.; Bishop, M.; Widzyk-Capehart, E. Seeing through dust and water vapor: Millimeter wave radar sensors for mining applications. *Journal of Field Robotics* 2007, 24, 527–557.
- 223. Guan, J.; Madani, S.; Jog, S.; Gupta, S.; Hassanieh, H. Through fog high-resolution imaging using millimeter wave radar.
   Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2020, pp. 11464–11473.
- 224. Gourova, R.; Krasnov, O.; Yarovoy, A. Analysis of rain clutter detections in commercial 77 GHz automotive radar. 2017 European 1866
   Radar Conference (EURAD). IEEE, 2017, pp. 25–28.
- 225. Breitenstein, J.; Termöhlen, J.A.; Lipinski, D.; Fingscheidt, T. Corner Cases for Visual Perception in Automated Driving: Some 1870 Guidance on Detection Approaches. arXiv preprint arXiv:2102.05897 2021.
- 226. Koopman, P.; Fratrik, F. How many operational design domains, objects, and events? Safeai@ aaai, 2019.

 227. Antonante, P.; Spivak, D.I.; Carlone, L. Monitoring and diagnosability of perception systems. 2021 IEEE/RSJ International 1873 Conference on Intelligent Robots and Systems (IROS). IEEE, 2020, pp. 168–175.

 Zheng, Z.; Yue, X.; Keutzer, K.; Sangiovanni Vincentelli, A. Scene-aware Learning Network for Radar Object Detection. 1875 Proceedings of the 2021 International Conference on Multimedia Retrieval, 2021, pp. 573–579.

- Malawade, A.V.; Mortlock, T.; Faruque, M.A.A. HydraFusion: Context-Aware Selective Sensor Fusion for Robust and Efficient 1877 Autonomous Vehicle Perception. arXiv preprint arXiv:2201.06644 2022.
- Ahuja, N.; Alvarez, I.J.; Krishnan, R.; Ndiour, I.J.; Subedar, M.; Tickoo, O. Robust multimodal sensor fusion for autonomous driving vehicles, 2020. US Patent App. 16/911,100.
- Deng, J.; Dong, W.; Socher, R.; Li, L.J.; Li, K.; Fei-Fei, L. Imagenet: A large-scale hierarchical image database. 2009 IEEE conference 1881 on computer vision and pattern recognition. Ieee, 2009, pp. 248–255.
- 232. Feng, D.; Harakeh, A.; Waslander, S.L.; Dietmayer, K. A review and comparative study on probabilistic object detection in autonomous driving. *IEEE Transactions on Intelligent Transportation Systems* 2021.
- 233. Ettinger, S.; Cheng, S.; Caine, B.; Liu, C.; Zhao, H.; Pradhan, S.; Chai, Y.; Sapp, B.; Qi, C.R.; Zhou, Y.; et al. Large scale interactive motion forecasting for autonomous driving: The waymo open motion dataset. Proceedings of the IEEE/CVF International Conference on Computer Vision, 2021, pp. 9710–9719.
- Lin, Y.C.; Gu, M.X.; Lin, C.H.; Lee, T.S. Deep-Learning Based Decentralized Frame-to-Frame Trajectory Prediction Over Binary Range-Angle Maps for Automotive Radars. *IEEE Transactions on Vehicular Technology* 2021, 70, 6385–6398.
- Kunert, M. The EU project MOSARIM: A general overview of project objectives and conducted work. 2012 9th European Radar Conference. IEEE, 2012, pp. 1–5.
- Alland, S.; Stark, W.; Ali, M.; Hegde, M. Interference in automotive radar systems: Characteristics, mitigation techniques, and current and future research. *IEEE Signal Processing Magazine* 2019, 36, 45–59.
- 237. Oyedare, T.; Shah, V.K.; Jakubisin, D.J.; Reed, J.H. Interference Suppression Using Deep Learning: Current Approaches and Open
   1894
   Challenges. arXiv preprint arXiv:2112.08988 2021.