# RODNet: Radar Object Detection using Cross-Modal Supervision

Yizhou Wang[1], Zhongyu Jiang[1], Xiangyu Gao[1], Jenq-Neng Hwang[1],
Guanbin Xing[1], and Hui Liu [1,2]

[1]University of Washington, Seattle, WA
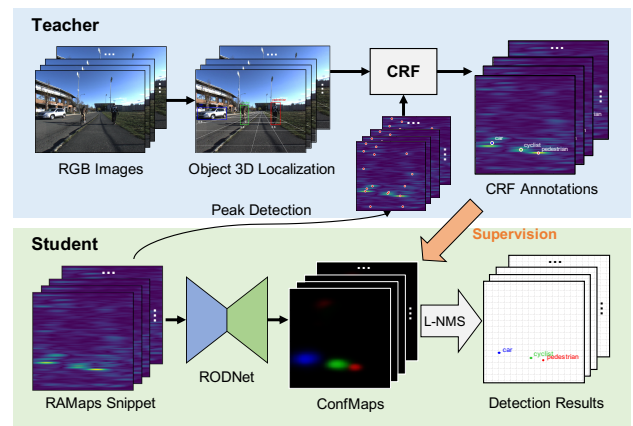[2]Silkwave Holdings Limited, Hong Kong

{ywang26, zyjiang, xygao, hwang, gxing, huiliu}@uw.edu

## Abstract

*Radar is usually more robust than the camera in severe driving scenarios, e.g., weak/strong lighting and bad weather. However, unlike RGB images captured by a camera, the semantic information from the radar signals is noticeably difficult to extract. In this paper, we propose a deep radar object detection network (RODNet), to effectively detect objects purely from the carefully processed radar frequency data in the format of range-azimuth frequency heatmaps (RAMaps). Three different 3D autoencoder based architectures are introduced to predict object confidence distribution from each snippet of the input RAMaps. The final detection results are then calculated using our post-processing method, called location-based non-maximum suppression (L-NMS). Instead of using burdensome human-labeled ground truth, we train the RODNet using the annotations generated automatically by a novel 3D localization method using a camera-radar fusion (CRF) strategy. To train and evaluate our method, we build a new dataset – CRUW, containing synchronized videos and RAMaps in various driving scenarios. After intensive experiments, our RODNet shows favorable object detection performance without the presence of the camera.*

## 1. Introduction

In autonomous or assisted driving, a camera can usually give us good semantic understandings of visual scenes. However, it is not a robust sensor under severe driving conditions, such as weak/strong lighting and bad weather, which lead to little/high exposure or blur/occluded images. Radar, on the other hand, is relatively more reliable in most harsh environments, e.g., dark, rain, fog, etc. Frequency modulated continuous wave (FMCW) radar, which operates in the millimeter-wave (MMW) band (30-300GHz) that is lower than visible light, thus, has the following proper-



**Figure 1:** The proposed cross-modal supervision pipeline. Teacher's pipeline first detects and 3D localizes the objects from the RGB images, combined with the detected peaks from the corresponding RAMaps by the proposed camera-radar fusion (CRF) algorithm. Student's pipeline learns to detect objects with radar data (RAMaps) as the input **only**.

ties: 1) great capability to penetrate through fog, smoke, and dust; 2) accurate range detection ability due to the huge bandwidth and high working frequency.

Typically, there are two kinds of data representations for the FMCW radar, i.e., radar frequency (RF) data and radar points. The RF data are generated from the raw radar signals using a series of fast Fourier transforms (FFTs), and the radar points are then derived from these RF data through peak detection [33]. Although the radar points can be directly used as the input of the LiDAR point cloud based methods, the radar points are usually much sparser, e.g., less than 5 points on a nearby car, than the point cloud from a LiDAR [9], so that it is not enough to accomplish the object detection task. Whereas, the RF data can maintain the rich Doppler information and surface texture so as to have the capability of understanding the semantic meaning of a certain object. Thus, in this work, we consider the RF data

in the range-azimuth coordinates, named RAMaps.

In this paper, we propose a radar object detection method, cross-supervised by a camera-radar fusion algorithm, that can accurately detect objects purely with the radar signal input. More specifically, we propose a novel radar object detection pipeline, which consists of two parts: teacher and student. The teacher estimates object classes and 3D locations by a reliable probabilistic-driven camera-radar fusion (CRF) strategy to automatically provide annotations for the student. The student takes radar reflection range-azimuth heatmaps (RAMaps) as the input and predicts the object confidence maps (ConfMaps). From the ConfMaps, object classes and locations are inferred using our post-processing method, called location-based non-maximum suppression (L-NMS). The aforementioned pipeline is shown in Figure 1. As for the network architecture of the RODNet, we implement 3D convolutional autoencoder networks based on [44] and [29]. Considering different temporal lengths needed for distinguishing different objects, we also propose temporal inception convolution layers, inspired by spatial inception [37], in our RODNet.

We train and evaluate the RODNet using our self-collected dataset, called CRUW, which contains about 400K camera-radar synchronized frames with various driving scenarios. Our CRUW dataset is, to the best of our knowledge, the first dataset containing synchronized stereo RGB images and RF data for autonomous driving applications. To evaluate the performance of our proposed RODNet, without the definition of bounding boxes widely used in image-based object detection on RAMaps, we introduce an evaluation method to evaluate the radar object detection performance in the radar range-azimuth coordinates. With intensive experiments, our RODNet can achieve about $83.76\%$ average precision (AP) and $85.62\%$ average recall (AR) solely based on radar input in various scenarios whether objects are visible or not in cameras.

Overall, our main contributions[1] are the following:

- A novel and accurate radar object detection network, called RODNet, for robust object detection in various driving scenarios, without camera or LiDAR.

- A camera-radar fusion (CRF) cross-modal supervision framework for training the RODNet without laborious and inconsistent human labeling.

- A new dataset, named CRUW, is collected, containing synchronized camera and radar data, which is valuable for camera-radar cross-modal research.

- A new evaluation method for radar object detection tasks is proposed and justified for its effectiveness.

---

[1]The CRUW dataset and code are available at: `https://www.cruwdataset.org/`

The rest of this paper is organized as follows. Related works for camera and radar data learning are presented in Section 2. The proposed cross-modal supervision framework is introduced in Section 3, with training and inference of our RODNet being explained in Section 4. In Section 5, we introduce our self-collected CRUW dataset. Then, the implementation details, evaluation metrics, and experimental results are shown in Section 6. Finally, we conclude our work in Section 7.

## 2. Related Works

### 2.1. Learning of Vision Data

Image-based object detection [32, 16, 31, 24] is aimed to detect every object with its class and precise bounding box location from RGB images, which is fundamental and crucial for camera-based autonomous driving. Then, most tracking algorithms focus on exploiting the association between the detected objects in consecutive frames, the so-called tracking-by-detection framework [8, 42, 38, 40, 18, 10, 43, 19]. Among them, the TrackletNet Tracker (TNT) [40] is an effective and robust tracker to perform multiple object tracking (MOT) of the detected objects with a static or moving camera. Once the same objects among several consecutive images are associated, the missing and erroneous detections can be recovered or corrected, resulting in better subsequent 3D localization performance.

Object 3D localization has attracted many interests in autonomous and safety driving community [35, 36, 27, 28, 6]. One idea is to localize vehicles by estimating their 3D structures using a CNN, e.g., 3D bounding boxes [27] and 3D keypoints [28, 6, 22]. Then, they deform a pre-defined 3D vehicle model to fit the 2D projection, resulting in accurate vehicle locations. Another idea [35, 36], however, tries to develop a real-time monocular structure-from-motion (SfM) system, taking into account the SfM cues and object cues. Although these kinds of works achieve favorable performance in object 3D localization, they only work for vehicles since only the vehicle structure information is considered. To address this limitation, an accurate and robust object 3D localization system, based on the detected and tracked bounding boxes of objects, is proposed in [41], claiming that the system works for most common moving objects in the road scenes, such as cars, pedestrians, and cyclists. Thus, we decide to take this 3D localization system as our camera annotation method.

### 2.2. Learning of Radar Data

Significant research in radar object classification has demonstrated its feasibility as a good alternative when cameras fail to provide good results [17, 5, 13, 23, 11]. With handcrafted feature extraction, Heuel, et al. [17] classify objects using a support vector machine (SVM) to distinguish

cars and pedestrians. While, Angelov et al. [5] use neural networks to extract features from the short-time Fourier transform (STFT) heatmap. However, the above methods only focus on *classification* tasks, that assume only one object has been appropriately identified in the scene and not applicable to the complex driving scenarios. Recently, a radar object detection method is proposed in [14], which combines a statistical detection algorithm CFAR [33] with a neural network classifier VGG16 [34]. But it would easily give many *false positives*, i.e., obstacles detected as objects. Besides, the laborious human annotations required by this method are usually impossible to obtain.

Recently, the concept of cross-modal learning has been discussed in machine learning community [21, 39, 30, 20]. This concept is trying to transfer or fuse the information between two different modalities in order to help train the neural networks. Specifically, RF-Pose [44] introduces the cross-modal supervision idea into wireless signals to achieve human pose estimation based on WiFi range radio signals. As the human annotations for wireless signals are difficult to obtain, RF-Pose uses a computer vision technique, i.e., OpenPose [12], to generate annotations for training from the camera. However, radar object detection is more challenging: 1) Feature extraction for object detection (especially for classification) is more difficult than human joint detection, which could just classify different joints by their relative locations without considering object surface texture or velocity information; 2) The typical FMCW radars on the vehicles have much less resolution than the sensors used in RF-Pose. As for autonomous driving, [26] proposes a vehicle detection method using LiDAR information for cross-modal learning. However, our work is different from theirs: 1) they only consider vehicles as the target object class, while we detect pedestrians, cyclists, and cars; 2) the scenario of their dataset, mostly highway without noisy obstacles, is easier for radar object detection, while we are dealing with various traffic scenarios.

## 3. Cross-Modal Supervision

### 3.1. Radar Signal Processing and Properties

In this work, we use a common range-azimuth heatmap representation, named RAMap, to represent our radar signal reflections. RAMap can be described as a bird's-eye view (BEV) representation, where the $x$-axis shows azimuth (angle) and the $y$-axis shows range (distance). For an FMCW radar, it transmits continuous chirps and receives the reflected echoes from the obstacles. After the echoes are received and processed, we implement the fast Fourier transform (FFT) on the samples to estimate the range of the reflections. A low-pass filter (LPF) is utilized to remove the high-frequency noise across all chirps. After the LPF, we conduct a second FFT on the samples along different re-

ceiver antennas to estimate the azimuth angle of the reflections and obtain the final RAMaps. After transforming into RAMaps, the radar data become a similar format as image sequences, which can be directly processed by an image-based CNN.

Moreover, RF data also has the following special properties to be handled for object detection task.

- **Rich motion information.** According to the principle of the radio signal, rich motion information is included. The speed and its law of variation over time consist of surface texture movement details, etc. For example, the motion information of a non-rigid body, like a pedestrian, is usually random, while for a rigid body, like a car, it should be more consistent. To utilize this motion information, multiple consecutive radar frames need to be considered as the input.

- **Inconsistent resolution.** Radar usually has high-resolution in range but low-resolution in azimuth due to the limitation of radar specifications, like the number of antennas, and the distance between them.

- **Different representation.** Radar data are usually represented as complex numbers containing frequency and phase information. This kind of data is unusual to be modeled by a typical neural network.

### 3.2. Camera-Only (CO) Supervision

The annotations for the radar are in the radar range-azimuth coordinates (similar to those of the BEV of a camera). To recover the 3D information from 2D images, we take advantage of a recent work on an effective and robust system for visual object 3D localization based on a monocular camera [41]. Even though stereo cameras can also be used for object 3D localization, however, high computational cost and sensitivity to camera setup configurations (e.g., baseline) result in the limitation of the stereo camera localization system. The proposed system takes a CNN inferred depth map as the input, incorporating adaptive ground plane estimation and multi-object tracking results, to effectively estimate object classes and 3D locations relative to the camera.

However, the above camera-only system may not be accurate enough after transforming to the radar range-azimuth coordinates because: 1) The systematic bias in the camera-radar sensor system that the peaks in the RF images may not be consistent with the 3D geometric center of the object; 2) Cameras' performance can be easily affected by lighting or weather conditions. Since we do have the radar information available, camera-radar cross calibration and supervision should be used. Therefore, an even more accurate self-annotation method, based on camera-radar fusion, is required for training the RODNet.

## 3.3. Camera-Radar Fusion (CRF) Supervision

The camera-only annotation can be improved by radar, which has a plausible capability of distance estimation. The CFAR algorithm [33] is commonly used in signal processing community to detect peaks from RAMaps. As shown in Fig. 1, a number of peaks are detected from the input RAMaps. However, these peaks cannot be directly used as the supervision because 1) CFAR cannot provide the object classes for the peaks; 2) CFAR usually gives a large number of false positives. Thus, these radar peaks are fused with the above CO supervision using an effective CRF strategy.

First, the CO annotations are projected from 3D camera coordinates to radar range-azimuth coordinates by the sensor system calibration. After the coordinates between camera and radar are aligned, a probabilistic CRF algorithm is developed. The basic idea of this algorithm is to generate two probability maps for camera and radar locations separately, and then fuse them by element-wise product. The probability map for camera locations with object class $cls$ is generated by

$$\mathcal{P}^c_{(cls)}(\mathbf{x}) = \max_i \left\{ \mathcal{N} \left( \frac{1}{2\pi\sqrt{|\boldsymbol{\Sigma}^c_{i(cls)}|}} \exp \left\{ -\frac{1}{2}(\mathbf{x}-\mu^c_i)^\top (\boldsymbol{\Sigma}^c_{i(cls)})^{-1}(\mathbf{x}-\mu^c_i) \right\} \right) \right\},$$

$$\mu^c_i = \begin{bmatrix} \rho^c_i \\ \theta^c_i \end{bmatrix}, \boldsymbol{\Sigma}^c_{i(cls)} = \begin{bmatrix} \left(d_i s_{(cls)}/c_i\right)^2 & 0 \\ 0 & \delta_{(cls)} \end{bmatrix}. \quad (1)$$

Here, $d_i$ is the object depth, $s_{(cls)}$ is the scale constant, $c_i$ is the depth confidence, and $\delta_{(cls)}$ is the typical azimuth error for camera localization. $\mathcal{N}(\cdot)$ represents the normalization operation for each object's probability map. Similarly, the probability map for radar locations is generated by

$$\mathcal{P}^r(\mathbf{x}) = \max_j \left\{ \mathcal{N} \left( \frac{1}{2\pi\sqrt{|\boldsymbol{\Sigma}^r_j|}} \exp \left\{ -\frac{1}{2}(\mathbf{x}-\mu^r_j)^\top (\boldsymbol{\Sigma}^r_j)^{-1}(\mathbf{x}-\mu^r_j) \right\} \right) \right\},$$

$$\mu^r_j = \begin{bmatrix} \rho^r_j \\ \theta^r_j \end{bmatrix}, \boldsymbol{\Sigma}^r_j = \begin{bmatrix} \delta^r_j & 0 \\ 0 & \epsilon(\theta^r_j) \end{bmatrix}. \quad (2)$$

Here, $\delta^r_j$ is the radar's range resolution, and $\epsilon(\cdot)$ is the radar's azimuth resolution. Then, an element-wise product is used to obtain the fused probability map for each class,

$$\mathcal{P}^{CRF}_{(cls)}(\mathbf{x}) = \mathcal{P}^c_{(cls)}(\mathbf{x}) * \mathcal{P}^r(\mathbf{x}). \quad (3)$$

Finally, the fused annotations are derived from the fused probability maps $\mathcal{P}^{CRF}$ by peak detection.

## 3.4. ConfMap Generation

After objects are accurately localized in the radar range-azimuth coordinates, the annotations need to be transformed into a proper representation that is compatible with our RODNet. Considering the idea in [12] that defines the human joint heatmap to represent joint locations, we define the confidence map (ConfMap) in range-azimuth coordinates to represent object locations. One set of ConfMaps has multiple channels, where each channel represents one specific
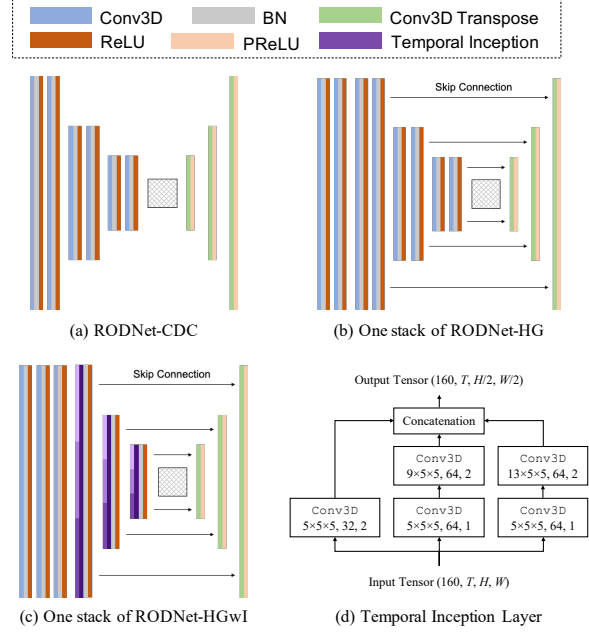


**Figure 2:** The architectures of our three RODNet models.

class label. The value at the pixel in the $cls$-th channel represents the probability of an object with class $cls$ existing at that range-azimuth location. Thus, we use Gaussian distributions to set the ConfMap values around the object locations, whose mean is the object location, and the variance is related to the object class and scale information.

# 4. Radar Object Detection

## 4.1. RODNet Architecture

The three different network architectures for the RODNet are shown in Figure 2, named 3D Convolution-Deconvolution (RODNet-CDC), 3D stacked hourglass (RODNet-HG), and 3D stacked hourglass with temporal inception (RODNet-HGwI), respectively. RODNet-CDC is a shallow 3D CNN network that squeeze the features in both spatial and temporal domains to better extract temporal information. While the RODNet-HG is adopted from [29], but we replace 2D convolution layers with 3D convolution layers and adjust the parameters for our task. As for the RODNet-HGwI, we replace the 3D convolution layers in each hourglass by the temporal inception layers [37] with different temporal kernel scales $(5, 9, 13)$ to extract different lengths of temporal features from the RAMaps.

Overall, our RODNet is fed with a snippet of RAMaps $\mathbf{R}$ with dimension $(C_{RF}, \tau, w, h)$ and predicts ConfMaps $\hat{\mathbf{D}}$ with dimension $(C_{cls}, \tau, w, h)$, where $C_{RF}$ is the number of channels in each RAMap, referring [44], where real and imaginary values are treated as two different channels, i.e., $C_{RF} = 2$; $\tau$ represents the snippet length; $C_{cls}$ is the

number of object classes; $w$ and $h$ are width and height of RAMaps or ConfMaps respectively. Thus, RODNet predicts separate ConfMaps for individual radar frames. With systematically derived CRF annotations, we train our RODNet using binary cross entropy loss,

$$\ell = -\sum_{cls}\sum_{i,j} \boldsymbol{D}_{i,j}^{cls}\log \hat{\boldsymbol{D}}_{i,j}^{cls}+\left(1-\boldsymbol{D}_{i,j}^{cls}\right)\log\left(1-\hat{\boldsymbol{D}}_{i,j}^{cls}\right).$$
(4)

Here, $\boldsymbol{D}$ represents the ConfMaps generated from camera annotations, $\hat{\boldsymbol{D}}$ represents the ConfMaps prediction, $(i,j)$ represents the pixel indices, and $cls$ is the class label.

## 4.2. L-NMS: Identify Detections from ConfMaps

To obtain the final detections from ConfMaps, a post-processing step is still required. Here, we adopt the idea of non-maximum suppression (NMS), which is frequently used in image-based object detection to remove the redundant bounding boxes from the detectors. Here, NMS uses intersection over union (IoU) as the criterion to determine if a bounding box should be removed. However, there is no bounding box definition in our problem. Thus, inspired by object keypoint similarity (OKS) defined for human pose evaluation in COCO dataset [25], we define a new metric, called object location similarity (OLS), to take the role of IoU, which describes the relationship between two detections considering their distance, classes and scale information on ConfMaps. More specifically,

$$\text{OLS} = \exp\left\{\frac{-d^2}{2(s\kappa_{cls})^2}\right\},$$
(5)

where $d$ is the distance (in meters) between the two points on RAMap, $s$ is the object distance from the sensors, representing object scale information, and $\kappa_{cls}$ is a per-class constant which represents the error tolerance for class $cls$, which can be determined by the object average size of the corresponding class. Moreover, we empirically tune $\kappa_c$ to make OLS distributed reasonably between 0 and 1. Here, we try to interpret OLS as a definition of Gaussian probability, where distance $d$ acts as bias and $(s\kappa_{cls})^2$ acts as variance. Therefore, OLS is a distance metric in a similarity manner, which also considers object sizes and distances, so that more reasonable than other traditional distance metrics, such as Euclidean distance, Mahalanobis distance, etc. This OLS metric is also used to *match detections and ground truth* for evaluation purpose, mentioned in Section 6.1.

After OLS is defined, we propose a location-based NMS (L-NMS), whose procedure can be summarized as follows:

1) Get all the peaks in all $C$ channels in ConfMaps within a $3 \times 3$ window as a peak set $P = \{p_n\}_{n=1}^N$.

2) Pick the peak $p^* \in P$ with the highest confidence and remove it from the peak set. Calculate OLS with each of the rest peaks $p_i$, where $p_i \neq p^*$.
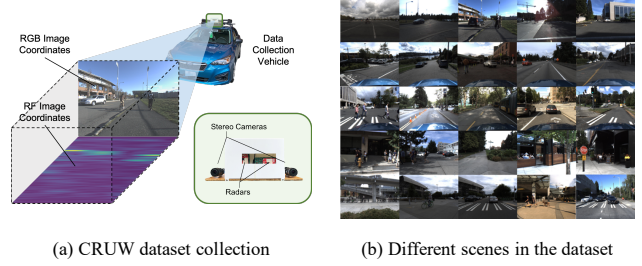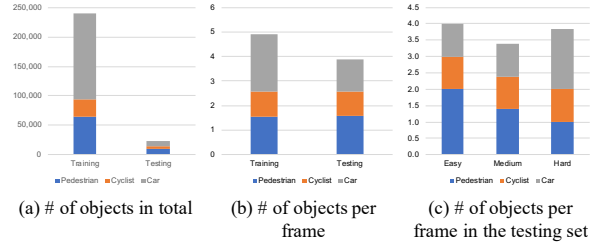


(a) CRUW dataset collection　　　(b) Different scenes in the dataset

**Figure 3:** Sensor platform and driving scenes in CRUW dataset.



(a) # of objects in total　　(b) # of objects per frame　　(c) # of objects per frame in the testing set

| Scenarios | # of Seqs | # of Frames | Vision-Fail % |
|---|---|---|---|
| Parking Lot | 124 | 106,057 | 15% |
| Campus Road | 112 | 94,416 | 11% |
| City Street | 216 | 175,392 | 6% |
| Highway | 12 | 20,376 | 0% |
| Overall | 464 | 396,241 | 9% |

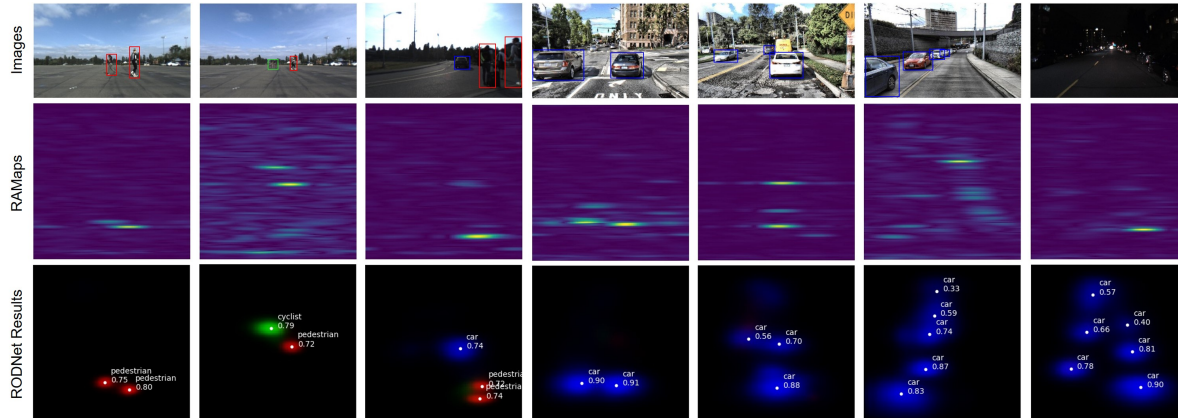(d) Driving scenarios statistics for CRUW dataset

**Figure 4:** Illustration for our CRUW dataset distribution.

3) If OLS between $p^*$ and $p_i$ is greater than a threshold, remove $p_i$ from the peak set.

4) Repeat Steps 2 and 3 until the peak set becomes empty.

Moreover, during the inference stage, we can send overlapped RAMap snippets into the RODNet, which provides different ConfMaps predictions for a single radar frame. Then, we merge these different ConfMaps together to obtain the final ConfMaps results. This scheme can improve the system's robustness and can be considered as a performance-speed trade-off, discussed in Section 6.2.

## 5. CRUW Dataset

Going through some existing autonomous driving datasets [15, 3, 4, 9, 7], only nuScenes [9] and Oxford RobotCar [7] consider radar. However, the format is 3D radar points, which are usually sparse and without motion and surface texture information that needed for our task. In order to efficiently train and evaluate our RODNet using RF data, we collect a new dataset – CRUW dataset. Our sensor platform contains a pair of stereo cameras [1] and two 77GHz FMCW radar antenna arrays [2]. The sensors, assembled and mounted together as shown in Figure 3 (a), are

**Figure 5:** Example results from our RODNet. The first row shows the images and the second row is the corresponding Radar frames in RAMap format. The ConfMaps predicted by the RODNet is shown in the third row, where the white dots represent the final detections after post-processing. Different colors in the ConfMaps represent different detected object classes.

well-calibrated and synchronized. Even though our cross-modal supervision requires just one monocular camera, the stereo cameras are setup to provide some ground truth of depth for object 3D localization performance validation.

The CRUW dataset contains more than 3 hours with 30 FPS (about 400K frames) of camera/radar data under different driving scenarios, including campus road, city street, highway, parking lot, etc. Some sample visual data are shown in Figure 3 (b). Besides, we also collect several vision-fail scenarios where the image qualities are pretty bad, i.e., dark, strong light, blur, etc. These data are only used for testing to illustrate that our method can still be reliable when vision techniques fail.

The object distribution in CRUW is shown in Figure 4. The statistics only consider the objects within the radar field of view. There are about 260K objects in CRUW dataset in total, including $92\%$ training and $8\%$ testing. The average number of objects in each frame is similar between training and testing data. From each scenario, we randomly select several complete sequences as testing sequences, which are not used for training. Thus, the training and testing sequences are captured at different locations and different time to show the generalization capability of the proposed system. For the ground truth needed only for evaluation purposes, we annotate $10\%$ of the visible data and $100\%$ vision-fail data. The annotation is operated on RAMaps by labeling the object classes and locations according to the corresponding images and RAMap reflection magnitude.

## 6. Experiments

### 6.1. Evaluation Metrics

To evaluate the performance, we utilize our proposed OLS (Eq. 5), replacing the role of IoU in image-level ob-

ject detection, to determine whether the detection result can be matched with a ground truth. During the evaluation, we first calculate OLS between each detection result and ground truth in every frame. Then, we use different thresholds from $0.5$ to $0.9$ with a step of $0.05$, for OLS and calculate the average precision (AP) and average recall (AR) for all different OLS thresholds. Here, we use AP and AR to represent the average values among different OLS thresholds, and use $AP^{OLS}$ and $AR^{OLS}$ to represent the values at a certain OLS threshold. Overall, we use AP and AR as our main evaluation metrics for the radar object detection task.
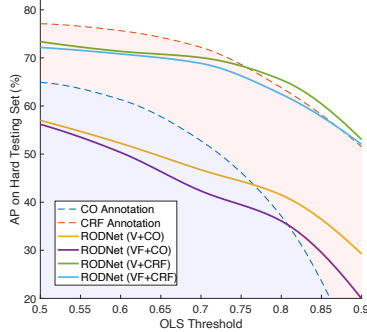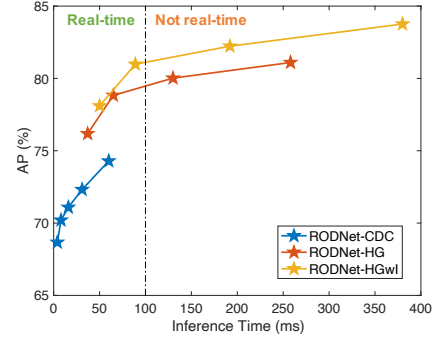
### 6.2. Radar Object Detection Results

We train our RODNet using the training data with CRF annotations in CRUW dataset. For testing, we perform inference and evaluation on the human-annotated visible data. The quantitative results are shown in Table 1. We compare our results with the following radar-only baselines: 1) a decision tree using handcrafted features from radar data [14]; 2) a radar object classification network [5] appended after the CFAR detection algorithm; 3) radar object detection method reported in [14]. To evaluate the performance under different scenarios, we split the test set into three levels, i.e., easy, medium, and hard. Among all the three competing methods, the AR performance for [14], [5] is relatively stable in the three different test sets, but their APs vary a lot. Especially, the APs drop from around $80\%$ to $10\%$ for easy to hard testing sets. This is caused by a large number of false positives detected by the traditional CFAR algorithm, which would significantly decrease the precision. Comparing with the baseline and competing methods, our RODNet outperforms significantly on both AP and AR metrics, achieving the best performance of $83.76\%$ AP and $85.62\%$ AR with the RODNet-HGwI architecture and CRF supervi-

**Table 1:** Radar object detection performance evaluated on CRUW dataset.

| Methods | Overall | | Easy | | Medium | | Hard | |
|---|---|---|---|---|---|---|---|---|
| | AP | AR | AP | AR | AP | AR | AP | AR |
| Decision Tree [14] | 4.70 | 44.26 | 6.21 | 47.81 | 4.63 | 43.92 | 3.21 | 37.02 |
| CFAR+ResNet [5] | 40.49 | 60.56 | 78.92 | 85.26 | 11.00 | 33.02 | 6.84 | 36.65 |
| CFAR+VGG-16 [14] | 40.73 | 72.88 | 85.24 | 88.97 | 47.21 | 62.09 | 10.97 | 45.03 |
| **RODNet (Ours)** | 83.76 | 85.62 | 94.52 | 95.94 | 72.49 | 75.59 | 66.77 | 71.24 |

**Table 2:** Ablation studies on the performance improvement with different architectures and annotations.

| Architectures | Supervision | AP | $AP^{0.5}$ | $AP^{0.7}$ | $AP^{0.9}$ | AR | $AR^{0.5}$ | $AR^{0.7}$ | $AR^{0.9}$ |
|---|---|---|---|---|---|---|---|---|---|
| RODNet-CDC | CO | 52.62 | 78.21 | 54.66 | 18.92 | 63.95 | 84.13 | 68.76 | 30.71 |
| | CRF | 74.29 | 78.42 | 76.06 | 64.58 | 77.85 | 80.05 | 78.93 | 71.72 |
| RODNet-HG | CO | 73.86 | 80.34 | 74.94 | 61.16 | 79.87 | 83.94 | 80.73 | 71.39 |
| | CRF | 81.10 | 84.71 | 83.08 | 70.21 | 84.26 | 86.54 | 85.42 | 77.44 |
| RODNet-HGwI | CO | 77.75 | 82.88 | 79.93 | 61.88 | 81.11 | 85.13 | 82.78 | 68.63 |
| | CRF | 83.76 | 87.99 | 86.00 | 70.88 | 85.62 | 88.79 | 87.37 | 76.26 |



**Figure 6:** Performance of vision-based and our RODNet on "Hard" testing set with different localization error tolerance. (V: visible data; VF: vision-fail data)



**Figure 7:** Performance-speed trade-off for the RODNet real-time implementation.

**Table 3:** The mean localization error (standard deviation) of CO/CRF annotations on CRUW dataset (in meters).

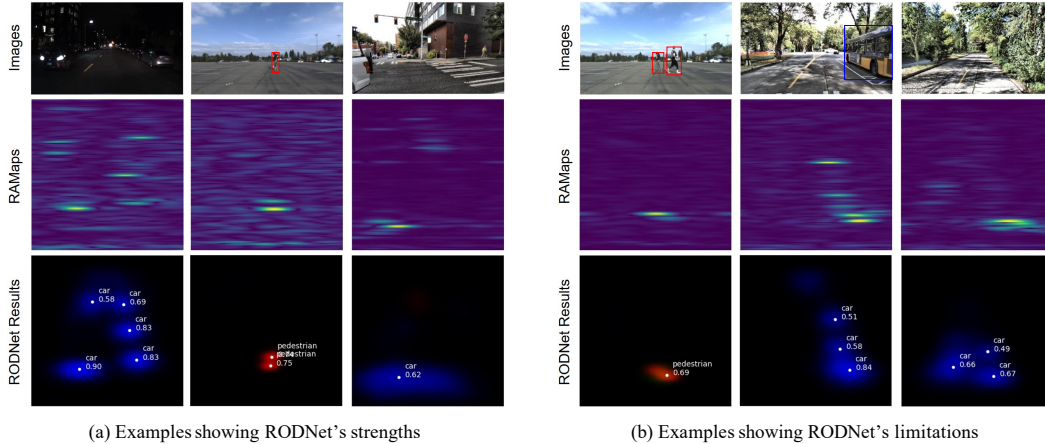| Supervision | Pedestrian | Cyclist | Car |
|---|---|---|---|
| CO | 0.69 ($\pm$0.77) | 0.87 ($\pm$0.89) | 1.57 ($\pm$1.12) |
| CRF | 0.67 ($\pm$0.55) | 0.82 ($\pm$0.59) | 1.26 ($\pm$0.64) |

sion. From now on, the RODNet discussed is referring to RODNet-HGwI, unless specified.

Some qualitative results are shown in Figure 5, where we can see that the RODNet can accurately localize and classify multiple objects in different scenarios. The examples on the left of Figure 5 are the scenarios that are relatively clean with fewer noises on the RAMaps, while the right ones are more complex with different kinds of obstacles, like trees, traffic sign, walls, etc. Especially, in the second to the last example, we can see high reflections on the right of the RAMap, which comes from the walls. The resulting ConfMap shows that the RODNet does not recognize them as any object, which is quite promising.
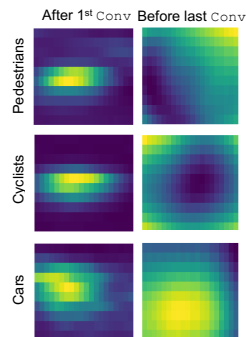
## 6.3. Ablation Studies

First, AP and AR under different OLS thresholds are analyzed in Table 2. Besides, we compare the teacher's performance on object 3D localization for both CO and CRF annotations, shown in Table 3. The CRF annotations are more accurate than CO annotations especially for the cars. From Table 2 and 3, we can find that, with more robust CRF annotations, the performance of our RODNet can increase significantly for all the three architectures. In Figure 6, the performance of teacher and student are compared on "Hard" testing set. Our RODNet shows its superiority and robustness on its localization performance.

Second, real-time implementation is important for autonomous driving applications. As mentioned in Section 4.2, we use different overlapping lengths during the inference, running on an NVIDIA TITAN XP, and report the time consumed in Figure 7. Here, we show the AP of three building architectures for the RODNet, and use 100 ms as a reasonable real-time threshold.

(a) Examples showing RODNet's strengths

(b) Examples showing RODNet's limitations

**Figure 8:** Examples illustrate strengths and limitations of our RODNet.

After the RODNet is well-trained, we analyze the features learned from the radar data. In Figure 9, we show two different kinds of feature maps, i.e., the features after the first convolution layer and the features before the last layer. These feature maps are generated by cropping some randomly chosen objects from the original feature maps and average them into one. From the visualization, we notice that the feature maps are similar in the beginning, but they become more discriminative at the end of the RODNet. Note that the visualized features are pixel-wise averaged within each object class to better represent the general class-level features.



**Figure 9:** RODNet feature visualization.

### 6.4. Strengths and Limitations

**RODNet Strengths.** Some examples to illustrate the RODNet's advantages are shown in Figure 8 (a). First, the RODNet has similar performance in some severe conditions, like during the night, shown in the first example. Moreover, the RODNet can handle some occlusion cases when the camera usually fails. In the second example, two pedestrians are nearly fully occluded in the image, but our RODNet can still detect both of them. This is because they are separate in the radar point of view. Last but not least, the RODNet has a wider field of view (FoV) than vision so that it can see more information. As shown in the third example, there is only a small part of the car visible in the camera view, which can hardly be detected from the camera side, but the RODNet can successfully detect it.

**RODNet Limitations.** Some failure cases are shown in Figure 8 (b). When two objects are very near, the RODNet often fails to distinguish them due to the limited resolution of radar. In the first example, the RAMap patterns of the two pedestrians are intersected, so that our result only shows one pedestrian detected. Another problem is, for huge objects like bus and train, the RODNet often detects it as multiple objects as shown in the second example. Lastly, the RODNet is sometimes affected by noisy surroundings. In the third example, there is no object in the view, but the RODNet detects the obstacles as several cars. The last two problems should be solved with a larger training dataset.

## 7. Conclusion

Object detection is crucial in autonomous driving and many other areas. Computer vision society has been focusing on this topic for decades and come up with many good solutions. However, vision-based detection is still suffering from many severe conditions. This paper proposed a brand-new and novel object detection method purely from radar information, which is more robust than vision. The proposed RODNet can accurately and robustly detect objects in various autonomous driving scenarios even during the night or bad weather. Moreover, this paper presented a new way to learn radar data using cross-modal supervision, which can potentially improve the role of radar in autonomous driving applications.

## Acknowledgement

# References

[1] Flir systems. `https://www.flir.com/`.

[2] Texas instruments. `http://www.ti.com/`.

[3] Apollo scape dataset. `http://apolloscape.auto/`, 2018.

[4] Waymo open dataset: An autonomous driving dataset. `https://www.waymo.com/open`, 2019.

[5] A. Angelov, A. Robertson, R. Murray-Smith, and F. Fioranelli. Practical classification of different moving targets using automotive radar and deep neural networks. *IET Radar, Sonar Navigation*, 12(10):1082–1089, 2018.

[6] Junaid Ahmed Ansari, Sarthak Sharma, Anshuman Majumdar, J Krishna Murthy, and K Madhava Krishna. The earth ain't flat: Monocular reconstruction of vehicles on steep and graded roads from a moving camera. In *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 8404–8410. IEEE, 2018.

[7] Dan Barnes, Matthew Gadd, Paul Murcutt, Paul Newman, and Ingmar Posner. The oxford radar robotcar dataset: A radar extension to the oxford robotcar dataset. In *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*, Paris, 2020.

[8] Philipp Bergmann, Tim Meinhardt, and Laura Leal-Taixe. Tracking without bells and whistles. *arXiv preprint arXiv:1903.05625*, 2019.

[9] Holger Caesar, Varun Bankiti, Alex H. Lang, Sourabh Vora, Venice Erin Liong, Qiang Xu, Anush Krishnan, Yu Pan, Giancarlo Baldan, and Oscar Beijbom. nuscenes: A multimodal dataset for autonomous driving. *arXiv preprint arXiv:1903.11027*, 2019.

[10] Jiarui Cai, Yizhou Wang, Haotian Zhang, Hung-Min Hsu, Chengqian Ma, and Jenq-Neng Hwang. Ia-mot: Instance-aware multi-object tracking with motion consistency. *arXiv preprint arXiv:2006.13458*, 2020.

[11] Peibei Cao, Weijie Xia, Ming Ye, Jutong Zhang, and Jianjiang Zhou. Radar-id: human identification based on radar micro-doppler signatures using deep convolutional neural networks. *IET Radar, Sonar & Navigation*, 12(7):729–734, 2018.

[12] Zhe Cao, Tomas Simon, Shih-En Wei, and Yaser Sheikh. Realtime multi-person 2d pose estimation using part affinity fields. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7291–7299, 2017.

[13] Samuele Capobianco, Luca Facheris, Fabrizio Cuccoli, and Simone Marinai. Vehicle classification based on convolutional networks applied to fmcw radar signals. In *Italian Conference for the Traffic Police*, pages 115–128. Springer, 2017.

[14] Xiangyu Gao, Guanbin Xing, Sumit Roy, and Hui Liu. Experiments with mmwave automotive radar test-bed. In *Asilomar Conference on Signals, Systems, and Computers*, 2019.

[15] Andreas Geiger, Philip Lenz, Christoph Stiller, and Raquel Urtasun. Vision meets robotics: The kitti dataset. *The International Journal of Robotics Research*, 32(11):1231–1237, 2013.

[16] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 2961–2969, 2017.

[17] S. Heuel and H. Rohling. Two-stage pedestrian classification in automotive radar systems. In *2011 12th International Radar Symposium (IRS)*, pages 477–484, Sep. 2011.

[18] Hung-Min Hsu, Tsung-Wei Huang, Gaoang Wang, Jiarui Cai, Zhichao Lei, and Jenq-Neng Hwang. Multi-camera tracking of vehicles based on deep features re-id and trajectory-based camera link models.

[19] Hung-Min Hsu, Yizhou Wang, and Jenq-Neng Hwang. Traffic-aware multi-camera tracking of vehicles based on reid and camera link model. In *Proceedings of the 28th ACM International Conference on Multimedia*, pages 964–972, 2020.

[20] Longlong Jing and Yingli Tian. Self-supervised visual feature learning with deep neural networks: A survey. *arXiv preprint arXiv:1902.06162*, 2019.

[21] Andrej Karpathy and Li Fei-Fei. Deep visual-semantic alignments for generating image descriptions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3128–3137, 2015.

[22] Hyungjin Kim, Bingbing Liu, and Hyun Myung. Road-feature extraction using point cloud and 3d lidar sensor for vehicle localization. In *2017 14th International Conference on Ubiquitous Robots and Ambient Intelligence (URAI)*, pages 891–892. IEEE, 2017.

[23] Jihoon Kwon and Nojun Kwak. Human detection by neural networks using a low-cost short-range doppler radar sensor. In *2017 IEEE Radar Conference (RadarConf)*, pages 0755–0760. IEEE, 2017.

[24] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision*, pages 2980–2988, 2017.

[25] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014.

[26] Bence Major, Daniel Fontijne, Amin Ansari, Ravi Teja Sukhavasi, Radhika Gowaikar, Michael Hamilton, Sean Lee, Slawomir Grzechnik, and Sundar Subramanian. Vehicle detection with automotive radar using deep learning on range-azimuth-doppler tensors. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*, 2019.

[27] Arsalan Mousavian, Dragomir Anguelov, John Flynn, and Jana Kosecka. 3d bounding box estimation using deep learning and geometry. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7074–7082, 2017.

[28] J Krishna Murthy, GV Sai Krishna, Falak Chhaya, and K Madhava Krishna. Reconstructing vehicles from a single image: Shape priors for road scene understanding. In *2017 IEEE International Conference on Robotics and Automation (ICRA)*, pages 724–731. IEEE, 2017.

[29] Alejandro Newell, Kaiyu Yang, and Jia Deng. Stacked hourglass networks for human pose estimation. In *European conference on computer vision*, pages 483–499. Springer, 2016.

[30] Yonggang Qi, Yi-Zhe Song, Honggang Zhang, and Jun Liu. Sketch-based image retrieval via siamese convolutional neural network. In *2016 IEEE International Conference on Image Processing (ICIP)*, pages 2460–2464. IEEE, 2016.

[31] Joseph Redmon and Ali Farhadi. Yolov3: An incremental improvement. *arXiv preprint arXiv:1804.02767*, 2018.

[32] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *Advances in neural information processing systems*, pages 91–99, 2015.

[33] Mark A Richards. *Fundamentals of radar signal processing*. Tata McGraw-Hill Education, 2005.

[34] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.

[35] Shiyu Song and Manmohan Chandraker. Robust scale estimation in real-time monocular sfm for autonomous driving. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1566–1573, 2014.

[36] Shiyu Song and Manmohan Chandraker. Joint sfm and detection cues for monocular 3d localization in road scenes. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3734–3742, 2015.

[37] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1–9, 2015.

[38] Zheng Tang and Jenq-Neng Hwang. Moana: An online learned adaptive appearance model for robust multiple object tracking in 3d. *IEEE Access*, 7:31934–31945, 2019.

[39] Subhashini Venugopalan, Marcus Rohrbach, Jeffrey Donahue, Raymond Mooney, Trevor Darrell, and Kate Saenko. Sequence to sequence-video to text. In *Proceedings of the IEEE international conference on computer vision*, pages 4534–4542, 2015.

[40] Gaoang Wang, Yizhou Wang, Haotian Zhang, Renshu Gu, and Jenq-Neng Hwang. Exploit the connectivity: Multiobject tracking with trackletnet. In *Proceedings of the 27th ACM International Conference on Multimedia*, pages 482–490. ACM, 2019.

[41] Yizhou Wang, Yen-Ting Huang, and Jenq-Neng Hwang. Monocular visual object 3d localization in road scenes. In *Proceedings of the 27th ACM International Conference on Multimedia*, pages 917–925. ACM, 2019.

[42] Linjie Yang, Yuchen Fan, and Ning Xu. Video instance segmentation. *arXiv preprint arXiv:1905.04804*, 2019.

[43] Haotian Zhang, Yizhou Wang, Jiarui Cai, Hung-Min Hsu, Haorui Ji, and Jenq-Neng Hwang. Lifts: Lidar and monocular image fusion for multi-object tracking and segmentation.

[44] Mingmin Zhao, Tianhong Li, Mohammad Abu Alsheikh, Yonglong Tian, Hang Zhao, Antonio Torralba, and Dina Katabi. Through-wall human pose estimation using radio signals. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7356–7365, 2018.