# Radar Camera Fusion via Representation Learning in Autonomous Driving

Xu Dong       Binnan Zhuang       Yunxiang Mao       Langechuan Liu [†]

XSense.ai

11010 Roselle Street, San Diego, CA 92121

## Abstract

*Radars and cameras are mature, cost-effective, and robust sensors and have been widely used in the perception stack of mass-produced autonomous driving systems. Due to their complementary properties, outputs from radar detection (radar pins) and camera perception (2D bounding boxes) are usually fused to generate the best perception results. The key to successful radar-camera fusion is the accurate data association. The challenges in the radar-camera association can be attributed to the complexity of driving scenes, the noisy and sparse nature of radar measurements, and the depth ambiguity from 2D bounding boxes. Traditional rule-based association methods are susceptible to performance degradation in challenging scenarios and failure in corner cases. In this study, we propose to address radar-camera association via deep representation learning, to explore feature-level interaction and global reasoning. Additionally, we design a loss sampling mechanism and an innovative ordinal loss to overcome the difficulty of imperfect labeling and to enforce critical human-like reasoning. Despite being trained with noisy labels generated by a rule-based algorithm, our proposed method achieves a performance of 92.2% F1 score, which is 11.6% higher than the rule-based teacher. Moreover, this data-driven method also lends itself to continuous improvement via corner case mining.*

## 1. Introduction

LiDAR, radar, and camera are the three main sensory modalities employed by the perception system of an autonomous driving vehicle. Though LiDAR-based 3D object detection is very popular in high-level autonomy, its wide adoption is still limited by some unsolved issues. First, LiDAR is prone to adversarial conditions (e.g. rainy weather); second, current LiDAR systems still exhibit prohibitively high maintenance need and cost; third, the mass-production of LiDAR is not ready to meet the growing demand.

---

[†] indicates corresponding author `patrickl@xsense.ai`
Note: this paper was published on CVPR 2021.



Figure 1: An illustration of the associations between radar detections (radar pins) and camera detections (2D bounding boxes). The context of the scene is illustrated in the top picture, with the image captured by the camera along with the detected bounding boxes and the projected radar pins (shown as numbered blue circles). The bottom picture adds red lines to highlight the association relationships between radar pins and bounding boxes. The tiny orange line in the middle denotes *uncertain* association relationship, which will be explained later.

An automotive millimeter-wave radar can also provide a certain level of geometrical information with relatively precise range and speed estimates. Moreover, as a widely-adopted sensor in automobiles for decades, radar is relatively robust, low-cost, and low-maintenance. The fusion between radar and camera combines radar's geometrical information and camera's appearance and semantic information, which is still the mainstream perception solution
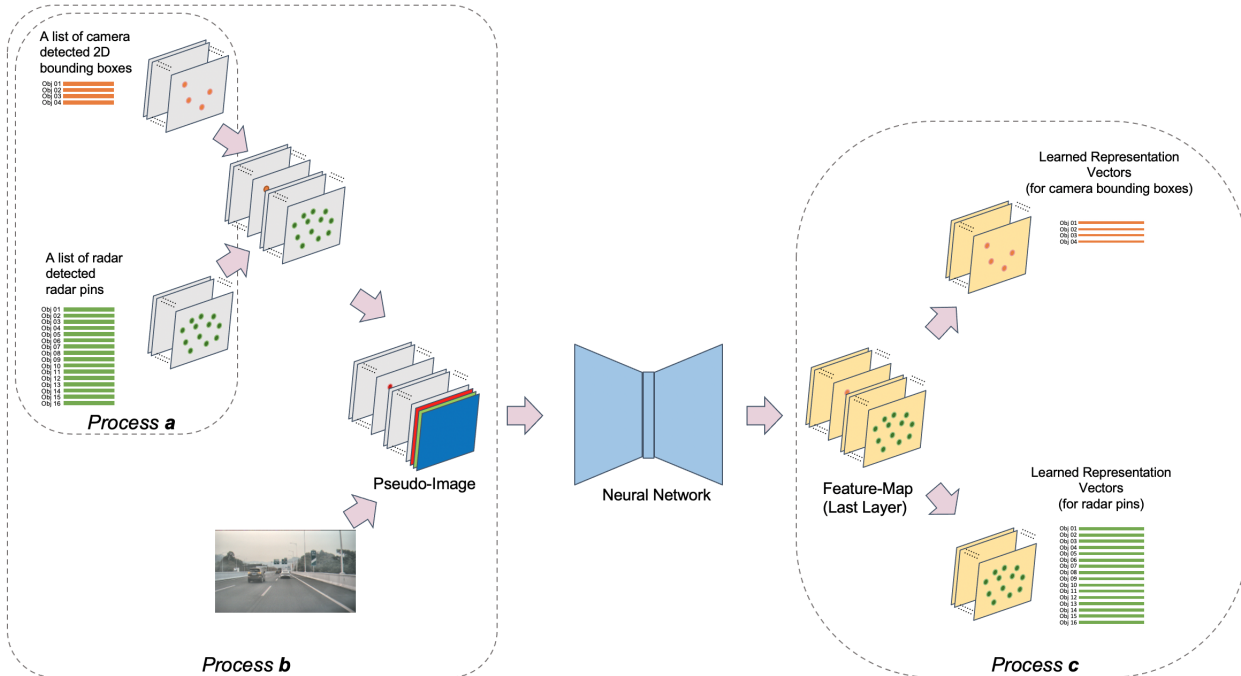
Figure 2: An overview of AssociationNet. *Process **a*** illustrates how the radar pins and 2D bounding boxes are first projected into the camera image plane and then produce a pseudo-image. *Process **b*** illustrates how the final pseudo-image is composed by concatenating all the features of radar pins, bounding boxes, and the original RGB camera image. The pseudo-image is then fed into a neural network to learn high-level semantic representations. *Process **c*** illustrates how the learned representation vectors for objects are finally extracted from the feature-map generated in the last layer of the neural network.

in many practical autonomous driving and assisted driving systems.

Traditionally, the radar-camera fusion is achieved by the combination of rule-based association algorithms and kinematic model-based tracking. The key is data association between radar and camera detections. The noisy and sparse nature of radar detection and the depth ambiguity from a mono camera makes such association problem very challenging. Traditionally, the association process is hand-crafted based on minimizing certain distance metrics along with some heuristic rules. It not only requires a large amount of engineering and tuning but is also hard to adapt to ever-growing data.

An emerging solution is to use learning-based methods to replace the rule-based radar-camera fusion. The latest advances focus on direct 3D object detection with the combined radar and camera data as the input [16, 25, 28]. These approaches all rely on LiDAR-based ground-truth to build the link between radar and camera. This is feasible on most public datasets such as nuScenes [4], Waymo [29] etc. However, it cannot be applied to a large fleet of commercial autonomous vehicles, often equipped with only radars and cameras. In this study, we propose a scalable learning-based framework to associate radar and camera information

without the costly LiDAR-based ground-truth.

Our goal is to find representations of radar and camera detection results, such that matched pairs are close and unmatched ones are far. We convert the detection results into image channels and combine them with the original image to feed into a convolutional neural network (CNN), namely, *AssociationNet*. Training is performed based on imperfect labels obtained from a traditional rule-based association method. A loss sampling mechanism is introduced to mitigate false labels. To further boost the performance, we guide the reasoning logic of AssociationNet by adding a novel ordinal loss. The proposed AssociationNet significantly outperforms the rule-based method through scene-dependent global reasoning.

Our main contributions are summarized as follows:

- We proposed a scalable learning-based radar-camera fusion framework without using ground-truth labels from LiDAR, which is suitable for building a low-cost, production-ready perception system for autonomous driving applications.

- We designed a loss sampling mechanism to alleviate the impact of the label noise, and also invented an ordinal loss to enforce critical association logic into the

model for performance enhancement.

- We developed a robust model via representation learning, which is capable of handling various challenging scenarios, and also outperforms the traditional rule-based algorithm by 11.6% in terms of the F1 score.

## 2. Related Work

### 2.1. Sensor Fusion

Traditionally, different sensory modules process their data separately. A downstream sensor fusion module augments the sensory outputs (typically detected objects) to form a more comprehensive understanding of the surroundings. Such an object-level fusion method is the mainstream approach [9, 17, 19, 12, 31] and is still widely used on many Advanced Driver Assistance Systems (ADAS). In object-level fusion, object detection is independently performed on each sensor, and the fusion algorithm combines such object detection results to create so-called global tracks for kinematic tracking [1].

Data association is the most critical and challenging task in object-level fusion. The precise association can easily lead to 3D object detection and multiple-object tracking solutions [1, 5]. Traditional approaches tend to manually craft various distance metrics to represent the similarities between different sensory outputs. Distance minimization [9] and other heuristic rules are applied to find the associations. To handle the complexity and uncertainty, probabilistic models are also sometimes adopted in the association process [2].

### 2.2. Learning-Based Radar-Camera Fusion

The learning-based radar-camera fusion algorithms can be primarily categorized into three groups, data-level fusion, feature-level fusion, and object-level fusion. The data-level fusion and feature-level fusion combine the radar and camera information at the early stage [28, 14] and the middle stage [16, 7, 26], respectively, but both directly perform 3D object detection. Hence, they rely on LiDAR to provide ground-truth labels during training, which prohibits their usage to autonomous vehicles without LiDAR.

The learning-based object-level fusion remains underexplored due to the limited information contained in the detection results. In this study, our proposed method belongs to this category in that we focus on associating radar and camera detection results. Thus, our method is more compatible with the traditional sensor fusion pipeline. On the other hand, our method also directly takes the raw camera image for further performance enhancement.

### 2.3. CNN for Heterogeneous Data

The tremendous success of CNN on structured image data inspires its application to many other types of hetero-geneous data, such as sensor parameters, point clouds, and association relationships between two groups of data [27]. In order to get compatible with CNN, a popular approach is to adapt the heterogeneous data into a form of pseudo-images. Examples include encoding camera intrinsic into images with normalized coordinates and field of view maps [11], projecting radar data into the image plane to form new image channels [6, 28], and the various forms of projection-based LiDAR point-cloud representations [30, 23]. We adopted a similar approach in this study to handle the heterogeneous radar and camera outputs.

### 2.4. Representation Learning

Representation learning has been considered as the key to understanding complex environments and problems [3, 20, 18]. Representation learning has been widely used in many natural language processing tasks such as word embedding [24], and many computer vision tasks, such as image classification [8], object detection [13], and keypoint matching [10]. In this study, we aim at learning a vector in the high-dimensional feature space as the representation for each object in the scene, in order to establish the interactions between objects as well as enable global reasoning about the scene.

## 3. Problem Formulation

We use a front-facing camera and a front-facing millimeter-wave mid-range radar for the proposed radar-camera fusion, yet the approach can be easily generalized to 360 perception with proper hardware setups. The camera intrinsic and the extrinsics of both sensors are obtained through offline calibration. The radar and camera operate asynchronously at 20Hz and 10Hz, respectively. The field-of-views (FOVs) of the radar and camera are 120 degrees and 52 degrees, respectively. The camera is mounted under the windshield at 1.33 meters above the ground. The output of the camera sensor at each frame is an RGB image with a size of 1828 pixels (width) by 948 pixels (height), whereas the output of the radar sensor at each frame is a list of processed points with many attributes (conventionally referred to as radar pins). Since the radar used here performs internal clustering, each output radar pin is on the object level (yet the proposed fusion technique also applies to lower level detection, e.g., radar locations). There are several tens of radar pins per frame depending on the actual scene and traffic. The attributes of each radar pin are listed in Table 1. There are two noteworthy characteristics of the radar pins. First, we only consume the 2D position information in the Bird's-Eye View (BEV) without the elevation angle, due to poor resolution and large measurement noise in the elevation dimension. Second, each radar pin either corresponds to a movable object (cars, cyclists, pedestrians, etc.) or an interfering static structure such as a traffic sign,

a street light, or a bridge.

In this study, we focus on associating 2D bounding boxes detected from a camera image to radar pins detected in the corresponding radar frame. With precise associations, many subsequent tasks like 3D object detection and tracking become much easier if not trivial.

Table 1: The Features of Each Radar Pin

| Feature | Explanation |
| --- | --- |
| object id | the id of the radar pin |
| obstacle prob | the probability of the existence of an obstacle being detected by the radar pin |
| position x | the x coordinate of the position of the detected obstacle in radar frame |
| position y | the y coordinate of the position of the detected obstacle in radar frame |
| velocity x | the velocity of the detected obstacle along the x coordinate in radar frame |
| velocity y | the velocity of the detected obstacle along the y coordinate in radar frame |

Table 2: The Features of Each 2D Bounding Box

| Feature | Explanation |
| --- | --- |
| center x | the x coordinate in the image plane of the center of the bounding box |
| center y | the y coordinate in the image plane of the center of the bounding box |
| height | the height of the bounding box in the image plane |
| width | the width of the bounding box in the image plane |
| category | the category of the detected moving object, including *sedan*, *suv*, *truck*, *bus*, *bicycle*, *tricycle*, *motorcycle*, *person*, and *unknown* |

## 4. Methods

Our proposed method mainly consists of a preprocessing step to align radar and camera data, a CNN-based deep representation learning network, AssociationNet, and a post-processing step to extract representations and make associations. An overview of the method is shown in Fig. 2 and details are explained in the following sections.

### 4.1. Radar and Camera Data Preprocessing

Temporal and spatial alignment is performed in the preprocessing stage. For each camera frame, we look for the nearest radar frame to perform data alignment. We align the nearest radar frame to the time instant of the camera frame, by moving the radar pin locations forward/backward along the time axis under a constant velocity assumption. After the temporal alignment, the radar pins are further transformed from the radar coordinate to the camera coordinate using the known extrinsics. All the attributes of the aligned radar pins will be used in AssociationNet.

Each camera frame is first fed into a 2D object detection network to produce a list of 2D bounding boxes corresponding to the movable objects in the scene. The output attributes for each detected 2D bounding box are displayed in Table 2. Though the network used in this study is an anchor-based RetinaNet [22] network, any 2D object detector will serve the purpose. After preprocessing, a list of temporally and spatially aligned radar pins and bounding boxes will be ready for association.

### 4.2. Deep Association by Representation Learning

We employ AssociationNet to learn a semantic representation (or a descriptor) of each radar pin and each bounding box. Under such representation, a pair of matched radar pin and bounding box will "look" similar, in the sense that the distance between the learned representations is small. An overview of the general process is shown in Fig. 2.

To leverage the powerful CNN architecture, we project each radar pin and 2D bounding box into the image plane to generate a pseudo-image, with each attribute occupying an independent channel. Specifically, each bounding box is assigned to the pixel location of its center. Each radar pin is assigned to the pixel location which is obtained through projecting its 3D location into the image plane using the camera intrinsic. The process is illustrated in *Process a* of the Fig. 2. Next, we concatenate the raw RGB camera image with the corresponding pseudo-image to incorporate the rich pixel-level information. AssociationNet is then applied to perform representation learning.

As shown in Fig. 3, the network consists of a ResNet-50 [15] as the backbone, a Feature Pyramid Network [21] for feature-map decoding, and two extra layers to restore the output feature-map size to the original input size. The output feature-map contains the high-level semantic representations of radar pins and bounding boxes. As each radar pin or bounding box has a unique pixel location in the feature-map, we extract the representation vector of each of those
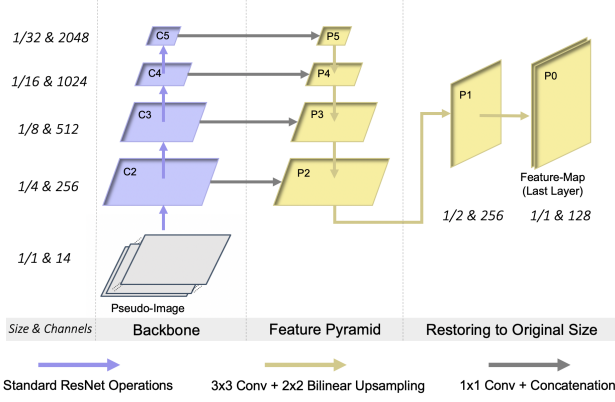
Figure 3: The architecture of the neural network. It consists of a ResNet-50 as the backbone, a feature pyramid for feature decoding, and two extra layers to restore the feature-map size. The feature-map in the last layer will be used to extract the representation vectors as shown in the *Process c* of the Fig. 2

on the output feature-map at its corresponding pixel location. The process is illustrated in *Process c* of the Fig. 2.

The input pseudo-image contains seven radar pin channels, four bounding box channels, and three raw camera image *RGB* channels. The radar pin channels include *object-id*, *obstacle-prob*, *position-x*, *position-y*, *velocity-x*, *velocity-y*[1], and a *heatmap* to indicate the projected pixel location. The bounding box channels include *height*, *width*, *category*, and also a *heatmap* to indicate the pixel location. The output feature-map contains 128 channels, resulting in the dimension of the representation vector to be 64 for each radar pin and bounding box.

The obtained representation vector captures the semantic meaning of each radar pin and each bounding box in a high dimension space. If a radar pin and a bounding box come from the same object in the real world, we treat the pair of radar pin and bounding box as a positive sample, otherwise, it is considered as a negative sample. We try to minimize the distance between the representation vectors of any positive sample and maximize the distance between the representation vectors of any negative sample. Based on such logic, we design loss functions according to the association ground-truth labels. We pull together the representation vectors of positive samples with the following pull loss:

$$L_{pull} = \frac{1}{n_{pos}} \sum_{(i_1,i_2)\in\mathbb{POS}} \max(0, \|h_{i_1} - h_{i_2}\| - m_1) \quad (1)$$

---

[1]Positions and velocities used here are under camera coordinate, as it is after the spatial alignment step in the preprocessing.
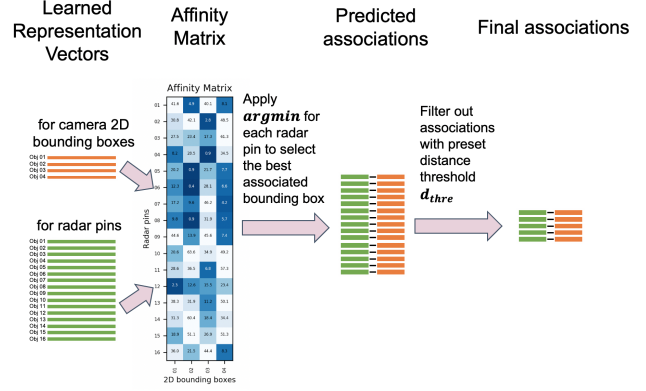


Figure 4: An overview of the process of obtaining final associations from the learned representation vectors.

And we push apart the representation vectors of negative samples with the following push loss:

$$L_{push} = \frac{1}{n_{neg}} \sum_{(i_1,i_2)\in\mathbb{NEG}} \max(0, m_2 - \|h_{i_1} - h_{i_2}\|) \quad (2)$$

Here $\mathbb{POS}$ and $\mathbb{NGE}$ are the set of positive samples and the set of negative samples, respectively in each frame; $n_{pos}$ and $n_{neg}$ are the total number of associations in $\mathbb{POS}$ and $\mathbb{NGE}$ respectively; $(i_1, i_2)$ denotes the $i^{th}$ association pair consisting of radar pin $i_1$ and bounding box $i_2$; $h_{i_1}$ and $h_{i_2}$ denotes the learned representation vectors; and $m_1$ and $m_2$ are the thresholds for the desired distances of representations among positive associations and negative associations, which were preset to be 2.0 and 8.0 in our experiments.

During inference, we calculate the Euclidean distance between the representation vectors of all possible radar-pin-bounding-box pairs. If the distance falls below a certain threshold, the radar pin and the bounding box will be considered as a successful association. More details of the inference process will be explained later.

### 4.2.1 Loss Sampling

The association labels used for supervising the learning process are ultimately from the traditional rule-based method, and hence are far from 100% accurate. To mitigate the impact of the inaccurate labels, we first purify the labels by applying some simple filters to remove low-confidence associations, which increases the precision in the remaining association labels at the cost of the undermined recall. During the push loss calculation in the training of AssociationNet, instead of exhausting all negative pairs (a pair of a radar pin and a bounding box that is not present in the association labels), we only sample a fraction of those to be used for push
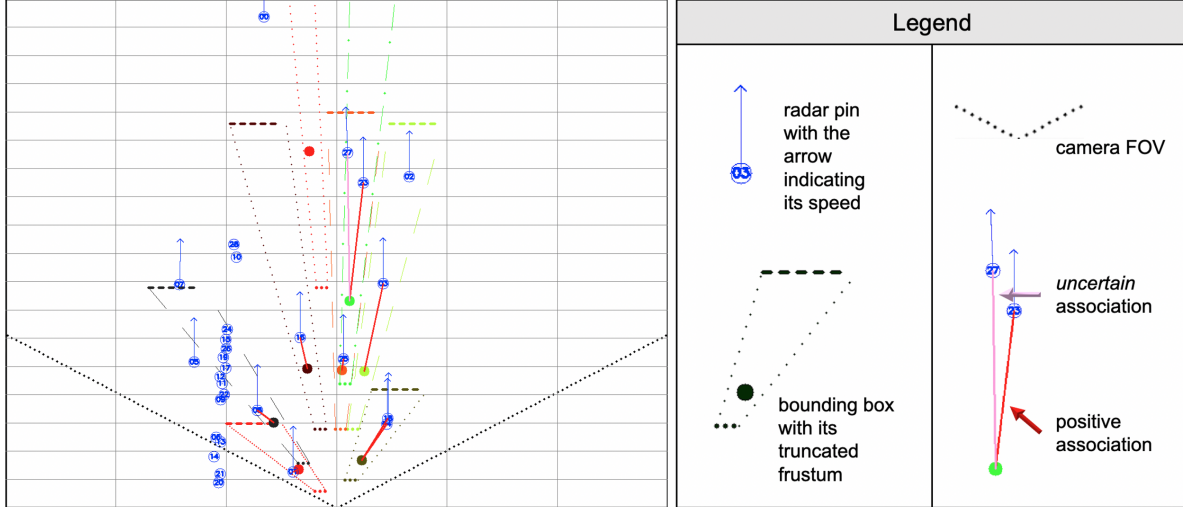
Figure 5: An illustration of radar pins, bounding boxes, and their association relationships under BEV perspective. This BEV image corresponds to the same scene as displayed in Fig. 1. Each grid in the image represents a 10-meter-by-10-meter square in physical space. The bounding boxes are represented as solid cycles in this image. The location of a bounding box is estimated by the Inverse Projective Mapping (IPM) method from the bounding box's center, to provide a rough reference for its real 3D location. A truncated frustum accompanying each bounding box is also plotted, for better assisting human curators to determine the association relationships[2].

loss calculation to alleviate pushing apart positive pairs by mistake. The number of sampled negative pairs is set to be equal to the number of positive ones at each frame.

#### 4.2.2 Ordinal Loss

One particular kind of error made by AssociationNet is that it could violate the simple ordinal rule, i.e., given two pairs of associated radar pins and bounding boxes, the farther radar pin associates to the closer bounding box. To solve this issue, an ordinal loss is introduced.

Denote the $y$ coordinate of bounding box $i$'s bottom edge as $y^i_{\max}$ and the bounding box's depth in 3D world as $d^i_b$ (which is supposedly the same as the depth of the associated radar pin $d^i_r$). For any two random bounding boxes on the same image we have the property:

$$y^i_{\max} > y^j_{\max} \iff d^i_b > d^j_b \qquad (3)$$

The ordering of the objects in the 3D world can be interpreted as the relative vertical ordering of the bottom edges of the corresponding bounding boxes.

---

[2]The frustum is calculated also by the IPM method, from the two side edges of each bounding box. According to projective geometry, the real object detected by a bounding box has to be within the bounding box's frustum, and hence the possibly matched radar pins as well. We truncated the frustums for the ease of visualizing. The widths of each frustum at the truncated positions are set to be one meter and five meters, respectively. As the physical width of a vehicle is most likely to be within the range, the possibly matched radar pins also tend to lie within the truncated frustum.

Hence, we design an additional ordinal loss to enforce the self-consistency within any two associations according to the ordinal rule, which is written as:

$$L_{ord} = \frac{2}{\widehat{n_{pos}} \cdot (\widehat{n_{pos}} - 1)} \cdot \\ \sum_{\substack{i \in \widehat{\mathbb{POS}} \\ j \in \widehat{\mathbb{POS}}}} \sigma(-(d^i_r - d^j_r) \cdot (y^i_{\max} - y^j_{\max})), \qquad (4)$$

where $\widehat{\mathbb{POS}}$ denotes the set of predicted positive associations and $\widehat{n_{pos}}$ is the size of the set; $i$ and $j$ are two random associations in $\widehat{\mathbb{POS}}$; $d^*_r$ represents the depth of the radar pin in an association , and $y^*_{\max}$ represents the $y$ coordinate of the bounding box's bottom edge in an association; and $\sigma$ is the sigmoid function to smooth the loss values.

Finally, the total loss is calculated as:

$$L_{tot} = L_{pull} + L_{push} + w_{ord} \cdot L_{ord}, \qquad (5)$$

where the $w_{ord}$ is the adjustable weight to balance losses.

### 4.3. Training and Inference

The AssociationNet was trained with a batch size of 48 frames at four NVIDIA GeForce RTX 2080 Ti GPUs. The SGD optimizer was used for training at a total of 10K iterations. The learning rate was set to be $10^{-4}$ initially, and then

was decreased by a factor of 10 at the end of 8K iterations and 9K iterations, respectively.

At the inference time, the representation vectors for all radar pins and bounding boxes are first predicted using the trained model. An affinity matrix is then calculated, where each matrix element corresponds to the distance between the representations of a radar pin and a bounding box. In reality, each bounding box may be associated with multiple radar pins (this is usually the case where the corresponding vehicles are of large sizes, such as trailer trucks and buses.), while each radar pin can only match to at most one bounding box. As a result, we associate each radar pin to the bounding box with the smallest distance in the affinity matrix. Lastly, the improbable associations with a distance larger than a threshold are filtered out, which usually consists of radar pins from interfering static objects. The whole inference process is summarized in Fig. 4.

### 4.4. Evaluation

The predicted associations are compared against human-annotated ground-truth associations in the test dataset. We use precision, recall, and F1 score as the metrics for evaluating the performance.

In some very complicated scenes, correctly associating all radar pins and bounding boxes is very challenging even for human annotators. Therefore, we mark those plausible but dubious associations as "*uncertain*" in the evaluation process. An example is shown in Fig. 5. For those "*uncertain*" associations, they are counted as neither positive nor negative associations, which will be excluded from both true and false positive predictions.

## 5. Experiments and Discussion

### 5.1. Dataset

The AssociationNet was trained and evaluated on an in-house dataset with 12 driving sequences collected by a testing fleet, which consists of 14.8 hours of driving in various driving scenarios, including highway, urban, and city roads. The radar and camera were synchronized at 10 Hz initially and further downsampled to 2 Hz, in order to reduce the temporal correlation among adjacent frames. Eleven sequences out of the twelve were used for training with the other one left for the test. Therefore, there are 104,314 synchronized radar and camera frames in the training dataset, and 2,714 in the test dataset. For the training data, the association labels were generated by a traditional rule-based algorithm with additional filtering to increase the precision. For the test data, we manually curated the labels with human annotators to obtain high-quality ground-truth labels.

Table 3: The Effect of Loss Sampling

| Sample Ratio | Performance |
| --- | --- |
| | Precision / Recall / F1 |
| *no sampling* | 0.896 / 0.925 / 0.911 |
| 1:2 | 0.901 / 0.931 / 0.916 |
| 1:1 | 0.906 / 0.939 / 0.922 |
| 2:1 | 0.899 / 0.933 / 0.915 |
| 3:1 | 0.899 / 0.929 / 0.914 |

Table 4: The Effect of Ordinal Loss

| Loss Weight $w_{ord}$ | Performance |
| --- | --- |
| | Precision / Recall / F1 |
| 0.0 | 0.897 / 0.912 / 0.904 |
| 0.5 | 0.897 / 0.923 / 0.910 |
| 1.0 | 0.899 / 0.931 / 0.915 |
| 2.0 | 0.906 / 0.939 / 0.922 |
| 5.0 | 0.889 / 0.918 / 0.903 |

### 5.2. Effect of Loss Sampling

We studied the effect of loss sampling on the AssociationNet's performance. Experiments were conducted with *no sampling* (meaning that all the negative pairs present in the label are used for push loss calculation), and loss sampling with different sampling ratios. The sampling ratio is defined as the ratio between the number of positive pairs and the number of negative pairs at each frame. The result is shown in Table 3. We can see that the best sampling ratio is 1:1 with the loss sampling mechanism, which boosts the performance by 1.1% in terms of the F1 score.

### 5.3. Effect of Ordinal Loss

The effect of the ordinal loss is shown in Table 4. The ordinal loss can facilitate both precision and recall to some degree. With the optimal loss weight, the performance is boosted by 1.8% in terms of the F1 score.

### 5.4. Comparison with Rule-Based Algorithm

We compared the performance of AssociationNet with the traditional rule-based algorithm, as shown in Table 5. Notably, though the traditional rule-based algorithm was used to generate association labels to supervise the training of AssociationNet, AssociationNet significantly outperforms the rule-based alternative. This demonstrates the inherent robustness of learning-based algorithms in handling complex scenarios.
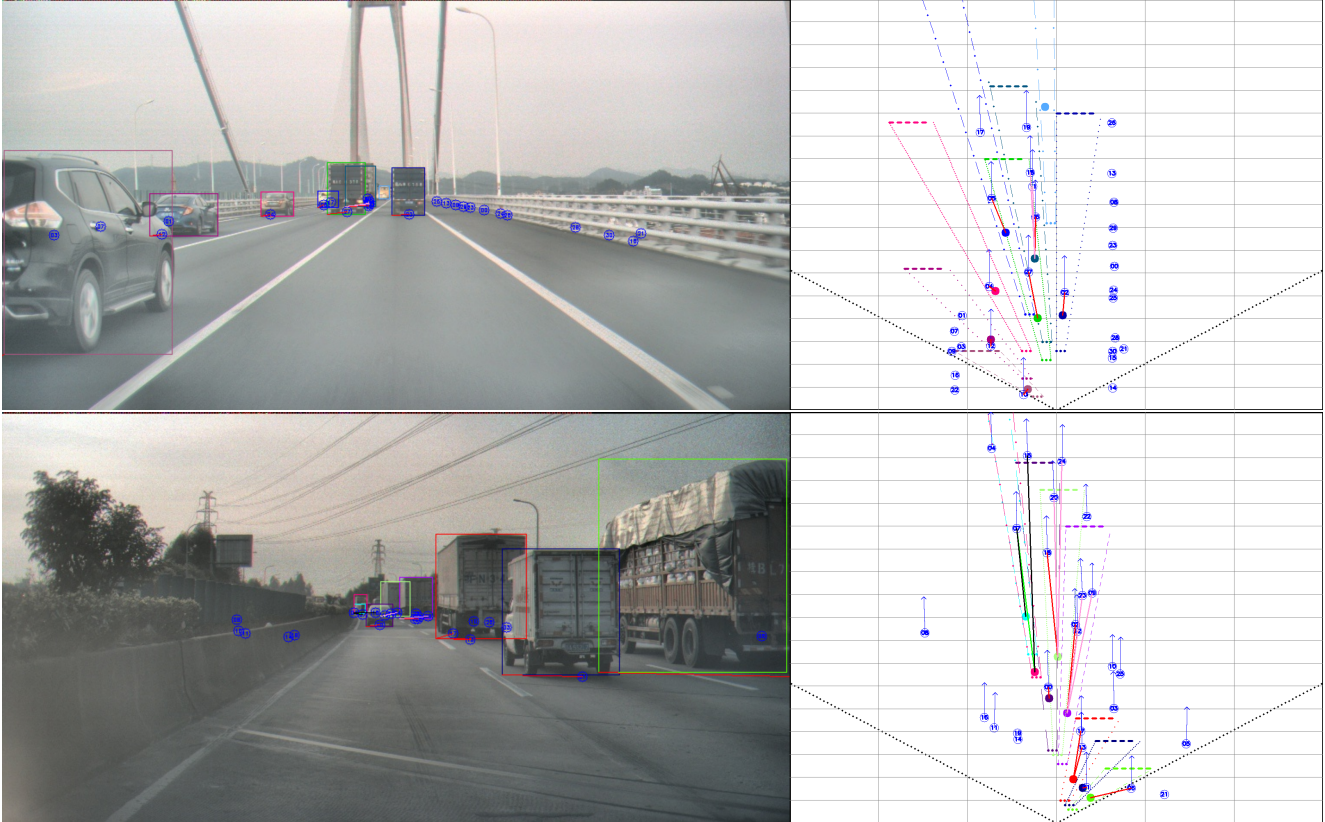
Figure 6: Examples of AssociationNet predictions. Here, the red solid lines represent the true-positive associations; and the pink solid lines represent predicted positive associations but labeled as *uncertain* in the ground-truth. In the second example, the added green lines represent the false-positive predictions; and the added black lines represent the false-negative predictions. Also, note that each bounding box on the left corresponds to a solid circle with the same color on the right.

Table 5: Comparison with Rule-based Algorithm

| Algorithm | Performance |
|---|---|
| | Precision / Recall / F1 |
| Rule-based | 0.890 / 0.736 / 0.806 |
| Learning-based | 0.906 / 0.939 / 0.922 |

## 5.5. Visualization

Examples of the predicted associations are shown in Fig. 6. Despite multiple big trucks present in both examples, AssociationNet correctly predicted their associations, which demonstrates the robustness of the algorithm. On the other hand, in the second example, there are two bounding boxes incorrectly associated, with one bounding box having no predicted associations and the other associated to a wrong radar pin. The two bounding boxes correspond to vehicles at the very far range. The mistakes are largely due to the small sizes of the objects in the camera image and also the heavy occlusions.

## 6. Conclusion

In this work, we developed a scalable learning-based radar-camera fusion algorithm, without using LiDAR for ground-truth labels generation. Such a solution has many practical merits at the current technological stage, including low cost, low maintenance, high reliability, and more importantly, readiness for mass production. We employed deep representation learning to tackle the challenging association problem, with the benefits of enabled feature-level interaction and global reasoning. We also designed a loss sampling mechanism and a novel ordinal loss to mitigate the impact of label noise and enforce critical human logic into the learning process. Although imperfect labels generated by a traditional rule-based algorithm were used to train the network, our proposed algorithm outperforms the rule-based teacher by 11.6% in terms of the F1 score.

# 7. Acknowledgements

# References

[1] Michael Aeberhard, Stefan Schlichtharle, Nico Kaempchen, and Torsten Bertram. Track-to-track fusion with asynchronous sensors using information matrix fusion for surround environment perception. *IEEE Transactions on Intelligent Transportation Systems*, 13(4):1717–1726, 2012. 3

[2] Yaakov Bar-Shalom, Fred Daum, and Jim Huang. The probabilistic data association filter. *IEEE Control Systems Magazine*, 29(6):82–100, 2009. 3

[3] Yoshua Bengio, Aaron Courville, and Pascal Vincent. Representation learning: A review and new perspectives. *IEEE transactions on pattern analysis and machine intelligence*, 35(8):1798–1828, 2013. 3

[4] 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. 2

[5] Josip Ćesić, Ivan Marković, Igor Cvišić, and Ivan Petrović. Radar and stereo vision fusion for multitarget tracking on the special euclidean group. *Robotics and Autonomous Systems*, 83:338–348, 2016. 3

[6] Simon Chadwick, Will Maddern, and Paul Newman. Distant vehicle detection using radar and vision. In *2019 International Conference on Robotics and Automation (ICRA)*, pages 8311–8317. IEEE, 2019. 3

[7] Shuo Chang, Yifan Zhang, Fan Zhang, Xiaotong Zhao, Sai Huang, Zhiyong Feng, and Zhiqing Wei. Spatial attention fusion for obstacle detection using mmwave radar and vision sensor. *Sensors*, 20(4):956, 2020. 3

[8] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PMLR, 2020. 3

[9] Hyunggi Cho, Young-Woo Seo, BVK Vijaya Kumar, and Ragunathan Raj Rajkumar. A multi-sensor fusion system for moving object detection and tracking in urban driving environments. In *2014 IEEE International Conference on Robotics and Automation (ICRA)*, pages 1836–1843. IEEE, 2014. 3

[10] Daniel DeTone, Tomasz Malisiewicz, and Andrew Rabinovich. Superpoint: Self-supervised interest point detection and description. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, pages 224–236, 2018. 3

[11] Jose M. Facil, Benjamin Ummenhofer, Huizhong Zhou, Luis Montesano, Thomas Brox, and Javier Civera. Camconvs: Camera-aware multi-scale convolutions for single-view depth. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019. 3

[12] Fernando Garcia, Pietro Cerri, Alberto Broggi, Arturo de la Escalera, and José María Armingol. Data fusion for overtaking vehicle detection based on radar and optical flow. In *2012 IEEE Intelligent Vehicles Symposium*, pages 494–499. IEEE, 2012. 3

[13] Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2014. 3

[14] Xiao-peng Guo, Jin-song Du, Jie Gao, and Wei Wang. Pedestrian detection based on fusion of millimeter wave radar and vision. In *Proceedings of the 2018 International Conference on Artificial Intelligence and Pattern Recognition*, pages 38–42, 2018. 3

[15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 4

[16] Vijay John and Seiichi Mita. Rvnet: deep sensor fusion of monocular camera and radar for image-based obstacle detection in challenging environments. In *Pacific-Rim Symposium on Image and Video Technology*, pages 351–364. Springer, 2019. 2, 3

[17] Naoki Kawasaki and Uwe Kiencke. Standard platform for sensor fusion on advanced driver assistance system using bayesian network. In *IEEE Intelligent Vehicles Symposium, 2004*, pages 250–255. IEEE, 2004. 3

[18] Alexander Kolesnikov, Xiaohua Zhai, and Lucas Beyer. Revisiting self-supervised visual representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1920–1929, 2019. 3

[19] Dirk Langer and Todd Jochem. Fusing radar and vision for detecting, classifying and avoiding roadway obstacles. In *Proceedings of Conference on Intelligent Vehicles*, pages 333–338. IEEE, 1996. 3

[20] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *nature*, 521(7553):436–444, 2015. 3

[21] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature pyramid networks for object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2117–2125, 2017. 4

[22] 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. 4

[23] Gregory P. Meyer, Ankit Laddha, Eric Kee, Carlos Vallespi-Gonzalez, and Carl K. Wellington. Lasernet: An efficient probabilistic 3d object detector for autonomous driving. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019. 3

[24] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*, 2013. 3

[25] Ramin Nabati and Hairong Qi. Rrpn: Radar region proposal network for object detection in autonomous vehicles. In *2019*

*IEEE International Conference on Image Processing (ICIP)*, pages 3093–3097. IEEE, 2019. 2

[26] Ramin Nabati and Hairong Qi. Centerfusion: Center-based radar and camera fusion for 3d object detection. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 1527–1536, 2021. 3

[27] Alejandro Newell and Jia Deng. Pixels to graphs by associative embedding. *arXiv preprint arXiv:1706.07365*, 2017. 3

[28] Felix Nobis, Maximilian Geisslinger, Markus Weber, Johannes Betz, and Markus Lienkamp. A deep learning-based radar and camera sensor fusion architecture for object detection. In *2019 Sensor Data Fusion: Trends, Solutions, Applications (SDF)*, pages 1–7. IEEE, 2019. 2, 3

[29] Pei Sun, Henrik Kretzschmar, Xerxes Dotiwalla, Aurelien Chouard, Vijaysai Patnaik, Paul Tsui, James Guo, Yin Zhou, Yuning Chai, Benjamin Caine, Vijay Vasudevan, Wei Han, Jiquan Ngiam, Hang Zhao, Aleksei Timofeev, Scott Ettinger, Maxim Krivokon, Amy Gao, Aditya Joshi, Yu Zhang, Jonathon Shlens, Zhifeng Chen, and Dragomir Anguelov. Scalability in perception for autonomous driving: Waymo open dataset, 2019. 2

[30] Yue Wang, Alireza Fathi, Abhijit Kundu, David Ross, Caroline Pantofaru, Tom Funkhouser, and Justin Solomon. Pillar-based object detection for autonomous driving. *arXiv preprint arXiv:2007.10323*, 2020. 3

[31] Ziguo Zhong, Stanley Liu, Manu Mathew, and Aish Dubey. Camera radar fusion for increased reliability in adas applications. *Electronic Imaging*, 2018(17):258–1, 2018. 3