

Radar-Lidar Fusion for Object Detection by Designing Effective Convolution Networks

Farzeen Munir^{1,2}, Shoaib Azam^{1,2}, Tomasz Kucner¹, Ville Kyrki¹, Moongu Jeon³

¹Department of Electrical Engineering and Automation, Aalto University, Finland

²Finnish Center for Artificial Intelligence, Finland

³School of Electrical Engineering and Computer Science, Gwangju Institute of Science and Technology, South Korea

Email: (farzeen.munir, shoaib.azam, tomasz.kucner, ville.kyrki)@aalto.fi, mgjeon@gist.ac.kr

Abstract—Object detection is a core component of perception systems, providing the ego vehicle with information about its surroundings to ensure safe route planning. While cameras and Lidar have significantly advanced perception systems, their performance can be limited in adverse weather conditions. In contrast, millimeter-wave technology enables radars to function effectively in such conditions. However, relying solely on radar for building a perception system doesn't fully capture the environment due to the data's sparse nature. To address this, sensor fusion strategies have been introduced. We propose a dual-branch framework to integrate radar and Lidar data for enhanced object detection. The primary branch focuses on extracting radar features, while the auxiliary branch extracts Lidar features. These are then combined using additive attention. Subsequently, the integrated features are processed through a novel Parallel Forked Structure (PFS) to manage scale variations. A region proposal head is then utilized for object detection. We evaluated the effectiveness of our proposed method on the Radiate dataset using COCO metrics. The results show that it surpasses state-of-the-art methods by 1.89% and 2.61% in favorable and adverse weather conditions, respectively. This underscores the value of radar-Lidar fusion in achieving precise object detection and localization, especially in challenging weather conditions.

Index Terms—radar object detection, sensor fusion, attention

I. INTRODUCTION

Environmental perception is crucial for crafting dependable decision-making algorithms in autonomous vehicles [1]. Modern autonomous vehicles are outfitted with a variety of advanced sensors, including Lidars, cameras, and radars [2] [3]. While leveraging these diverse sensor modalities enhances accuracy and robustness in environmental perception, it also poses challenges in creating cohesive perception systems for such vehicles. Moreover, many of these sensors face operational constraints [4]. Fusing data from multiple modalities offers a potential solution, providing a richer environmental representation and addressing the individual limitations of each sensor [5] [6] [7].

In the literature, most existing sensor fusion approaches discuss the fusion of Lidar and camera. Although these two sensors give rich information about the environment, yet there are some limitations in using them. For instance, the point

cloud from the Lidar gives dense 3D measurement at a close range but becomes sparse at long ranges, thus limiting its ability to detect objects at a far length [8]. Similarly, cameras provide a rich spatial representation of the environment but do not provide depth information compared to Lidar data. However, the fusion of Lidar and camera complement each other and give promising results in 3D object detection for autonomous driving, both of these sensors are prone to adverse weather conditions, which can significantly reduce their performance [5] [6]. Conversely, because radars operate at millimetre wavelength, they can penetrate or diffract around the tiny particles in adverse weather conditions such as fog, snow, rain and dust, thus can be utilized for the long-range perception of the environment [9] [10] [11].

Recently, some research efforts have been made to use radar data as a single modality for object detection. Some common approaches that utilize radar data either represent the radar data in image format or as point cloud data. For instance, [12] have adopted the radar image representation and designed a channel boosting approach in an ensemble feature extractor network for feature representation. Later, they utilized the transformer network to learn the contextual representation between the embedded features and used them for object detection. Similarly, [10] have transformed the radar data into an ego-centric bird-eye-view radar image to exploit the temporal relations and used them for object recognition. However, majority of works have utilized radar data for object detection only as a point cloud representation. That said, these approaches provide promising improvement in radar-based object detection, yet usage of standalone radar for the environment perception is limited due to its low spatial resolution. To address this, radar data can be fused with Lidar and camera to complement each other [13]. The camera-radar fusion improves the spatial resolution of the perception system but results in the sparse representation of the environment. [14]. Some works exploit this and fuse the radar and Lidar data for object detection [9]. However, sensor fusion between radar and Lidar provides a better representation of the environment, but how to fuse the data between these two modalities is an open research problem.

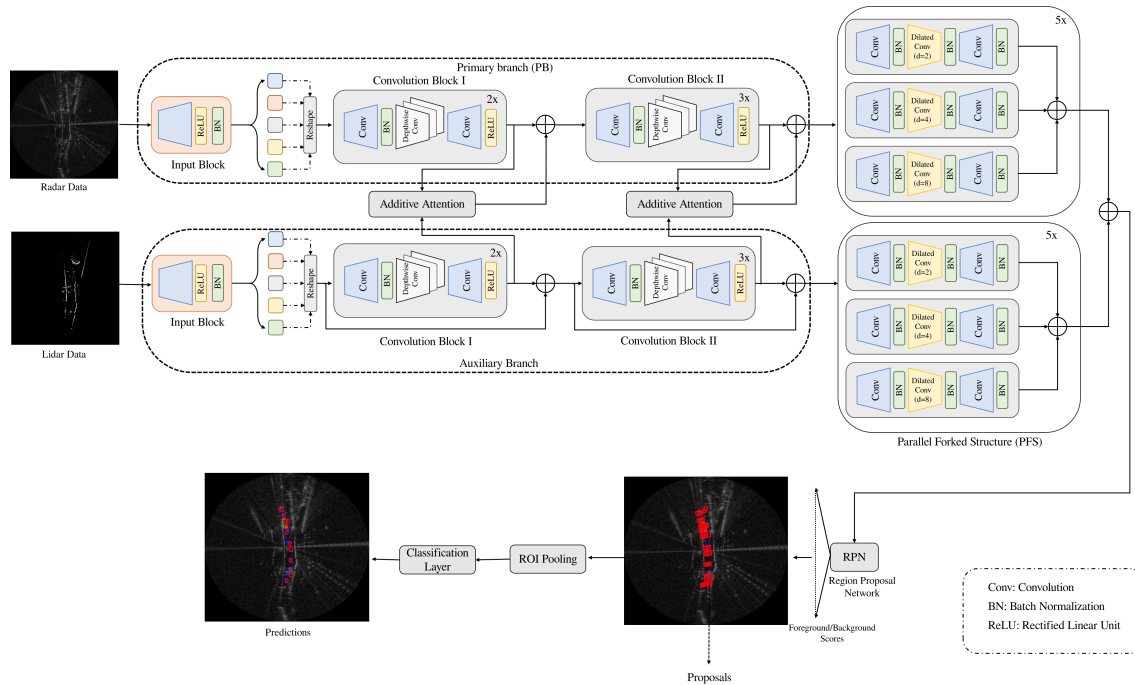


Fig. 1: The proposed method for radar-based object detection by fusing the radar and Lidar data. The architecture is composed of primary (PB) and auxiliary branch having the same configurations. Both the branches include input block, convolution block I and II for the features extraction by doing the depth-wise convolution, followed by parallel forked structure (PFS). The data is fused between the two branches using additive attention. Finally, the embedded features from the PFS are combined and used by the region-proposal detection for the object detection.

In the related work, a comprehensive analysis has been conducted on various data fusion strategies. Based on this analysis, the mid-fusion approach has been adopted in this particular work. In this work, a novel framework has been specifically designed for fusing radar and Lidar data for the task of object detection. The radar data allows to perceive the environment in adverse weather conditions, while Lidar data provides a dense point cloud of the environment, thus providing a unified framework for object detection. The proposed methodology, depicted in Fig. 1, consists of a pair of branched networks. The primary branch (PB) is responsible for extracting features from the radar data, while the secondary branch leverages Lidar data to provide supplementary information in the form of extracted features. The architectural configuration of both branches remains consistent: an initial input representation block is employed to partition the input data into patches, succeeded by multiple convolutional blocks and a parallel forked structure (PFS). Within the convolutional blocks, we employ depth-wise convolution to extract features from each channel. This strategy is chosen due to the inherent data sparsity [15]. The features extracted from these two branches are fused at an intermediate stage by using additive attention. Subsequently, the embedded features from both branches are concatenated and passed to a region-based proposal detection head for generating detection results. Similar to work, [9] have also used radar-Lidar fusion for object detection, yet the main differences to our work are that they have opted to use

late fusion and only consider one specific adverse weather condition: foggy weather. The main contributions of our work are:

- 1) We designed a novel framework for radar-based object detection by fusing the Lidar information at an intermediate level.
- 2) We employed additive attention to fuse the data between radar and Lidar, where the Lidar data provide auxiliary information that helps in radar-based object detection.
- 3) The proposed network introduces the parallel fork structure to effectively handle the scale variations by changing the receptive field size using the dilated convolution.

The rest of the paper is organized as follows: Section II covers the literature review. Section III presents the proposed framework, and Section IV presents experimentation, analysis and results. Finally, Section V concludes the paper with possible directions for future work.

II. RELATED WORK

A. Radar Object Detection

Even though cameras and Lidars are fundamental for perception in autonomous driving in recent years we can observe an influx of works focusing on the radar as a single modality for object detection. Earlier work is focused on classical radar processing techniques involved in detecting the peaks in the radar sensor data. The resulting peaks are clustered and tracked over time and then classified by a single-layer neural network [16]. Similarly, [17] have accumulated

the radar peaks in a 2D grid and then adopted a random forest and convolution neural network to classify them in an ensemble network. Furthermore, deep learning approaches have also been studied for radar-based object detection. [18] have proposed a deep learning-based approach for vehicle detection using range-azimuth-doppler measurements. [12] have adopted the channel boosting method in an ensemble learning framework for extracting the features from radar data and then employed the transformer network for object detection. Similarly, [10] have first transformed the radar data into ego-centric radar image frames and then exploited the temporal relations for object detection. All of the aforementioned object detection algorithms have used only radar data, but our work focuses on fusion of Lidar and radar data.

B. Fusion algorithms

Neural networks allow for diverse sensor data fusion strategies due to their hierarchical feature representation [19]. Commonly, these strategies include early, mid, and late fusion. Early fusion combines data before feature extraction but can face misalignment issues. Late fusion works at the decision level, but may not capture the full potential of joint sensor data representations. Mid-fusion strikes a balance between the two, effectively learning joint representations of sensor modalities [19]. Some works have fused the radar data with Lidar and camera to complement the sensor’s individual limitations. [14] have fused the radar and camera by first finding the centre points in image data using the centre point detection network. Then, it solves the data association problem between radar detections and the object’s centre using a frustum-based method. Later, the associated radar features are used for radar object detection. Another work, RODNet [20] is a cross-supervised camera-radar fusion framework for radar object detection. It uses the teacher-student network, where the teacher network fuses both the radar images and RGB images to obtain the object classes and location in the radar images. Finally, the teacher network supervises the student network that only uses radar images for object detection. MT-DETR [13] has proposed a multi-modal fusion network by introducing the residual and confidence fusion modules to fuse the radar, camera and Lidar. In addition, a residual enhancement module is designed to enhance each unimodal branch’s feature extraction. Similarly, [21] have also fused the Lidar, camera and radar for object detection. MVDNet [9] has fused the Lidar and radar by first generating the proposal from two sensor data and then adopting the region-wise features to improve the object detections. Similarly, ST-MVDNet [11] is built on top of MVDNet [9] by introducing the teacher-student network for radar-based object detection. In their framework, the teacher model generates the predictions to train the student network using a consistency constraint while the student network passes the parameters it learns to the teacher network via an exponential moving average.

MVDNet [9] and ST-MVDNet [11] are more similar to our work employing radar and Lidar fusion for object detection. However, both of these methods only consider

foggy weather conditions. In addition, MVDNet [9] adopts a late fusion approach that limits learning the joint distribution of sensor modalities in contrast to our mid-fusion approach. However, ST-MVDNet [11] uses the knowledge distillation approach, which has a problem of large discrepancy of predictive distribution between teacher and student networks. To address these issues, we designed a novel fusion network that uses additive attention for sensor fusion and depth-wise convolution to extract useful features from the sparse data of radar and Lidar, respectively. In addition, we introduce a parallel fork structure to handle the scale variations that solve the predictive discrepancy between the learned features.

III. PROPOSED METHOD

In this section, we explain in detail the proposed multi-modal fusion framework for object detection using Lidar and radar data, as illustrated in Fig. 1.

A. Input Block

Effective feature representation is crucial for deep learning-based object detection. The key idea is to explore the occurrence of features from multi-modal data by designing the neural network backbone for learning robust feature representations. In our proposed method settings, the inputs are represented in an image format. The radar image is acquired from the sensor, capturing information about the surrounding objects. Additionally, the Lidar data is preprocessed to obtain an image representation. The Lidar point-cloud data is preprocessed to a bird-eye-view (BEV) representation to make it more applicable to be used by 2D convolution neural networks in our proposed method. The majority of the works in object detection use a standard convolutional stem with a large kernel size as the input layer, which down-samples the input to acquire a proper feature map [22] [23]. However, we have opted for a patch-based representation of the input data, where the input data is split into fixed-size patches and obtained linear embeddings, which preserve locality and reduces computations complexity [24] [25]. The input image pair is given by $P = (I_i, J_i)$, where $I \in \mathbb{R}^{H \times W \times C}$ represent the radar range-azimuth images in cartesian coordinate system and $J \in \mathbb{R}^{H \times W \times C}$ shows the Lidar BEV image. The image I is partitioned into $p \times p$ non-overlapping patches (x_p) given by $x_p \in \mathbb{R}^{N \times (P^2 \cdot C)}$, where (H, W) is the height and width of the original image, C represent the input channels, p is the patch-size and $N = HW/p^2$ are the number of patches. The Eq.1 gives the linear embedding, which is further given to the convolution blocks, where ζ denotes nonlinear activation function.

$$x_0 = BN(\zeta\{Conv2D(x, stride = P, Kernel = P)\}) \quad (1)$$

B. Convolution Blocks I & II

Multi-modal object detection encounters challenges in answering the central question, which is the best way to fuse information from different modalities. Although, in literature, early, mid and late fusion techniques have been explored; it has been shown, in studies [26] [27], mid-fusion techniques

are most efficient in determining characteristics that co-occur and assist in acquiring rich features for object detection. Therefore, we also follow mid-fusion architecture, where we use two feature extractor branches, one for each modality, consisting of convolutional blocks I & II, and then leverage additive attention to fuse the features as illustrated in Fig. 1. Table I gives the details of each convolution block, which consists of point-wise convolutions, BatchNorm, depth-wise convolution and nonlinear activation. The depth-wise convolution is adopted to learn the spatial locality of features from the sparse data, which assists in learning better local and global feature representation from each data channel and using large kernel size provides robustness on out-of-distribution samples [25] [28]. The output from the convolutional block of the radar (x_r) and Lidar (x_l) branches is fused using additive attention. [29]. The output from additive attention is added to the PB since we are using ground-truth labels in the radar reference frame. Using additive attention enhances the feature representation and localization, reducing the feature response in unrelated background parts without the need to extract the area of interest.

Fig. 2 shows the operation of additive attention. The fusion of features through attention provides flexibility regarding what to focus on from a regional basis when fusing the information with the gating signal. The additive attention coefficient $\alpha \in [0, 1]$ distinguishes prominent feature map regions and prunes features for task-specific activations [30]. The output of the additive attention module is the element-wise multiplication of the fused feature map and input feature map described in Eq. 2.

$$\hat{x} = x_r \cdot \alpha \quad (2)$$

The α is given by Eq.3.

$$\begin{aligned} x_{att} &= \psi^T(\sigma(W_r^T x_r + W_l^T x_l)) + b_\psi, \\ \alpha &= \sigma_2(x_{att}(x_r, x_l, i; \Theta_{att})). \end{aligned} \quad (3)$$

where $\sigma_2(x) = \frac{1}{1+\exp(-x)}$ is the sigmoid activation function. Θ_{att} characterizes the set of learning parameters including W_r , W_l , and ψ , and the bias term, b_g . These parameters are trained with standard backpropagation.

TABLE I: The layer detail of the architecture for the proposed framework.

Blocks	Input/Output channel size	Layers
input block	3/64	$\begin{bmatrix} 16 \times 16, 64 \\ \text{RELU} \\ \text{BN} \end{bmatrix}$
convolution block I	64/256	$\begin{bmatrix} 1 \times 1, 256 \\ \text{BN} \\ 11 \times 11, 256 \\ \text{RELU} \end{bmatrix} \times 3$
Convolution block II	256/512	$\begin{bmatrix} 1 \times 1, 512 \\ \text{BN} \\ 11 \times 11, 512 \\ 1 \times 1, 512 \\ \text{RELU} \end{bmatrix} \times 4$
Parallel Forked Structure	512/1024	$\begin{bmatrix} 1 \times 1, 512 \\ \text{BN} \\ \{11 \times 11, d=2, 512\} \{11 \times 11, d=4, 512\} \\ \{11 \times 11, d=8, 512\} \\ \text{BN} \\ 1 \times 1, 1024 \\ \text{RELU} \end{bmatrix} \times 5$

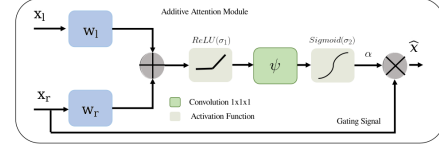


Fig. 2: The additive attention working operation to fuse radar and Lidar feature maps.

C. Parallel Forked Structure (PFS) Block

This section explains the parallel fork structure, which handles scale variation for object detection by varying the receptive field while having the same transformation parameters. As a result, the PFS renders scale-specific feature maps with a consistent representation of features. Contrary to our work, the most popular techniques to handle scale variation in object detection are image pyramids [31] and feature pyramids. The image pyramid introduces computational inefficiency but does provide representational power to deal with objects of all scales equally. However, the feature pyramid sacrifices the feature consistency across different scales leading to a higher risk of overfitting [32].

The PFS block consists of three parallel forked branches, each with a different receptive field, implemented using dilated convolution [33]. Dilation convolutions enhance feature expression and reduce computation by selectively spacing convolutional kernel elements, enabling a wider receptive field with fewer parameters. This helps to effectively captures intricate patterns and contextual information, showcasing the superior efficiency and effectiveness of this convolutional technique. The dilated convolution for one-dimensional input data is given by Eq. 4, where input x_i , output y_i , kernel filter w_k of the length of k and d is the dilated parameter that corresponds to stride through the input feature.

$$y_i = \sum_{k=1}^K x_{i+d \cdot k} \cdot w_k, \quad (4)$$

The standard convolution is dilated convolution with $d = 1$. i variable shows the location on output feature y_i , when kernel filter w_k is applied on input feature x_i . k denotes the indices of the dilated convolution kernel. The dilated convolution with rate d introduces $d - 1$ zeroes in the consecutive filter values, increasing the kernel size of the $k \times k$ filter to $k_d = k + (k - 1)(d - 1)$ with no increase in computation power and the same number of parameters. Therefore, it provides an effective method for handling scale variation and finding the best concession between context assimilation and finding the region of interest.

D. Detection Head

The outputs from the multi-modal branches are fused and passed through 1×1 convolution to learn the combined features for the downstream object detection task. The output of this convolution is passed to the region proposal network

[22], which predicts the bounding box and class for each object. Eq. 5 gives the region proposal network loss.

$$L = \frac{1}{N_{cls}} \sum_j L_{cls}(p_j, p_j^*) + \lambda \frac{1}{N_{reg}} \sum_j p_j^*(b_j, b_j^*) \quad (5)$$

where, j denotes the anchor index, p_j is predicted probability of objectness in an anchor, p_j^* is the ground truth. b_i denotes the bounding box coordinates and b_i^* represent the ground truth values. The classification loss L_{class} is log loss, and regression loss L_{reg} is defined by smooth $L1$ loss. N_{cls} and N_{reg} are the normalization term and λ is weighted term which is set to $\lambda = 10$.

IV. EXPERIMENTATION AND RESULTS

A. Dataset

In this study, we have used the RADIATE (RADar Dataset In Adverse weaThEr) dataset [34]. The RADIATE dataset is collected in adverse weather conditions and covers numerous traffic scenarios, facilitating object detection research in dynamic environments under adverse weather. The perception sensor data is collected using the Navtech CTS350-X radar and 32-channel Velodyne Lidar. The radar generates 360° range-azimuth images with a 1152×1152 resolution. However, Lidar renders a point cloud that has 360° coverage. The sensors are calibrated by finding extrinsic parameters for Lidar and radar. The radar is used as the reference frame to calculate the parameters. The Lidar point-cloud data is pre-processed to obtain bird-eye-view (BEV) two-dimensional images by using binary occupancy grid mapping. For this, the BEV grid of $100m \times 100m$ grid is considered, where the block size of the grid is $0.174m \times 0.174m$, which results in a BEV image of 576×576 pixels. To further enhance the BEV images without data loss, the resulting BEV image is up-sampled to give the resolution of 1152×1152 pixels having 3 channels, representing occupancy, height, and intensity information obtained from the Lidar point cloud. We opted for BEV representation for input data compared to 3D voxelization of point cloud data, as it is more computationally effective and preserves the metric space, allowing the framework to explore features about the size and shape of objects.

The dataset contains data from different weather conditions, for instance, fog, rain, snow, night, sunny and cloudy. Moreover, it has various driving scenarios like urban, suburban and motorway. The annotated data consists of 40K images with an estimated 200K labelled vehicles. The vehicles are annotated using a bounding box defined by $(x, y, w, h, angle)$. Table II shows the training and test split of the dataset.

B. Experimentation Details

The proposed framework is trained on two RTX3090 GPUs using the PyTorch library. The input training data is augmented, scaled, randomly flipped and randomly cropped to provide more generalized datasets. The network is trained for 100K iterations using a batch size of 4. A stochastic gradient descent (SGD) optimizer is used with a learning rate

TABLE II: The RADIATE dataset partition topology for training and testing the proposed approach in adverse weather conditions

	No. of Images	No. of vehicles
Training data <i>Good Weather</i>	23091	106931
Training data <i>Good & Bad Weather</i>	9760	39647
Testing data	11289	147005

TABLE III: Quantitative analysis of proposed method using only radar and fusion of radar and Lidar is performed. The analysis illustrates the mAP scores of proposed method and state-of-the-art methods in both good, and adverse weather conditions.

Methods	Split: train good weather	Split: train good and bad weather
	mAP@0.5 scores	mAP@0.5 scores
RetinaNet-OB-ResNet18 [10]	37.83	31.57
RetinaNet-OB-ResNet34 [10]	35.61	31.10
RetinaNet-OB-ResNet34-T [10]	37.30	24.50
CenterPoint-OB-ResNet18 [10]	51.43	42.37
CenterPoint-OB-ResNet34 [10]	49.41	44.48
CenterPoint-OB-ResNet34 [10]	50.17	42.81
BBAVectors-ResNet18 [10]	50.53	45.43
BBAVectors-ResNet34 [10]	51.26	44.61
TRL-EfficientNetB4 [10]	50.98	43.05
TRL-ResNet18 [10]	53.11	46.42
TRL-ResNet34 [10]	54.00	43.98
Radiate-ResNet50 [34]	45.31	45.77
Radiate-ResNet101 [34]	45.84	46.55
Channel-Boosted-ResNet50 [12]	54.90	55.54
Channel-Boosted-ResNet101 [12]	58.39	59.03
Ours (Radar only)	63.45	60.95
MVDNet [9]	60.57	59.35
ST-MVDNet [11]	65.25	62.47
MT-DETR (w/o camera stream) [13]	64.95	61.89
Ours (Radar+Lidar)	67.14	65.08

of 0.02, and a linear scheduler is incorporated to decrease the learning rate by a factor of 0.1 after 1K iterations. The input patch-embedding size is 16, and the depthwise convolution kernel size is 11. The parallel fork structure consists of 3 branches with 2, 4 and 8 dilation rates. The region proposal network obtained about 12000 proposals from the parallel forked structure, and 128 ROI were sampled for training. The standard COCO evaluation metric of Average precision (AP) is used for evaluation [35]. The higher value of AP score corresponds to better object detection model.

C. Results

The proposed method is quantitatively and qualitatively evaluated on the Radiate dataset [34], having two different configurations of weather settings: good weather and adverse weather conditions (good and bad weather). In order to quantitatively evaluate the proposed method's efficacy, we first analyzed object detection using radar as a single modality.

TABLE IV: The table shows the result for the early and late fusion of multi-sensor data in comparison to mid-fusion.

Model	split: train good weather	Split: train good and bad weather
Early Fusion	61.59	59.10
Late Fusion	63.36	61.47
Mid-fusion (Proposed)	67.14	65.08

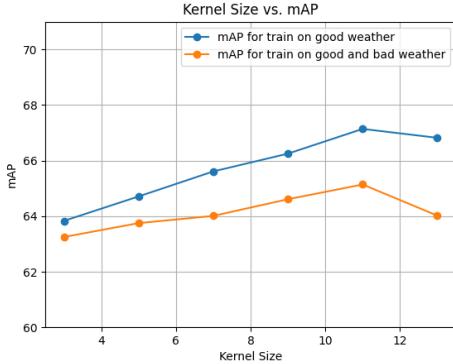


Fig. 3: The results on test data for using different kernel sizes in depthwise convolution.

This quantitative analysis forms a baseline for the proposed method to address how much auxiliary information from the Lidar branch provides robustness in object detection. To this end, the proposed method only with radar input is quantitatively evaluated with the state-of-the-art method that utilizes only the radar data. Table III illustrates the quantitative results for the proposed method (only radar) and state-of-the-art methods in both weather conditions settings. The proposed method with only radar has achieved the mAP score of 63.45 in good weather conditions and 60.95 in adverse weather conditions, respectively, thus outperforming the best among state-of-the-art method [12] by 5.06% in good weather and by a margin of 1.92% in adverse weather conditions respectively. It is to be noted here that we only used the methods that utilize the Radiate dataset for radar-based object detection.

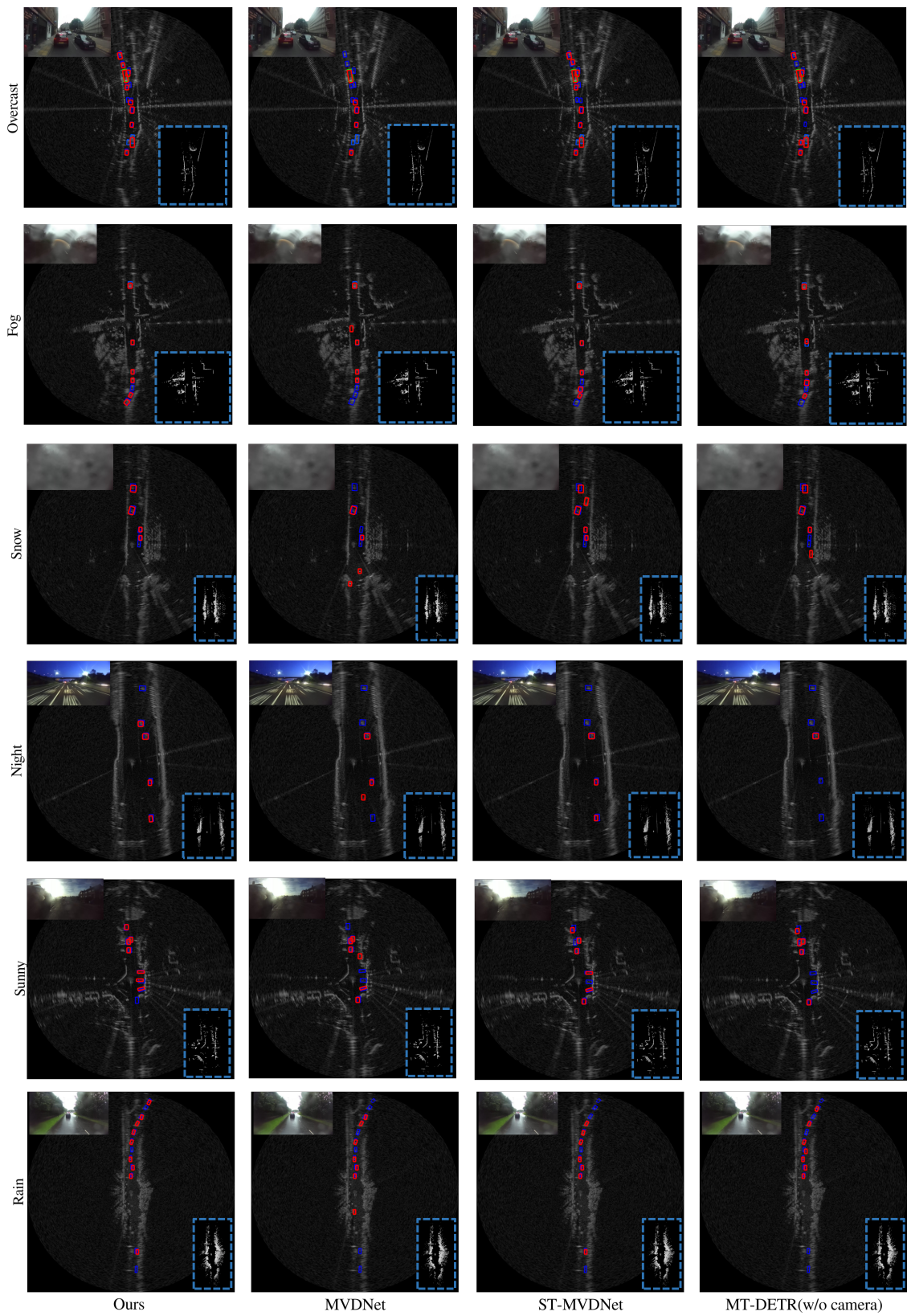
In order to quantitatively evaluate the proposed method in a sensor fusion setting, we employed the state-of-the-art methods developed using other datasets. To the best of our knowledge, no work has used the Radiate dataset [34] for fusing radar and Lidar data in sensor fusion settings. To this end, we re-implemented and refactored the codebase to be utilized by the proposed method using the Radiate dataset. The state-of-the-art methods that are used for the quantitative evaluation are MVDNet [9], ST-MVDNet [11] and MT-DETR (without camera stream) [13]. MVDNet [9] is a two-stage fusion network where region proposals are generated for both sensor modalities and fused together for object detection. MVDNet [9] uses Oxford Radar Robotcar (ORR) dataset [36] for the evaluation. Similarly, ST-MVDNet [11] is developed on top of MVDNet [9], which employs the teacher-student network for mutual learning of object detection features and is evaluated on the ORR dataset. MT-DETR [13] adopts a hierarchical fusion mechanism by designing a residual fusion module, a confidence fusion module a residual enhancement module to learn the features effectively. MT-DETR [13] uses the STF dataset in the experimental analysis. However, the original MT-DETR [13] uses the camera, Lidar and radar data in the sensor fusion settings, but in our proposed work, we refactored the original

MT-DETR only to accept two sensor streams: radar and Lidar. We opted for the same quantitative evaluation metrics of mAP scores and compared the proposed method with the state-of-the-art method in both good and adverse weather conditions. Table III shows the quantitative results where the proposed method has achieved the mAP score of 67.14 in good weather and 65.08 mAP score in adverse weather conditions, respectively. The proposed method in sensor fusion settings has outperformed ST-MVDNet by a margin of 1.89% in good weather conditions and by 2.61% in adverse weather conditions, respectively. The quantitative evaluation of the proposed method clearly shows the improvement when the Lidar auxiliary information is used for radar-based object detection. The addition of auxiliary information has increased the proposed method mAP score by 3.61% in good weather conditions and similarly by a margin of 4.13% mAP score improvement in adverse weather conditions. In order to visualize the radar object detection results, the proposed method and state-of-the-art methods qualitative results are illustrated in Fig.4 for both good and adverse weather conditions.

Furthermore, we investigate the effect of kernel size on the robustness of the proposed framework. Fig.3 shows the mAP scores at different settings of kernel size. The optimal scores are obtained at 11×11 , which provides a larger receptive field to capture the contextual understanding of the feature maps and robustness to the variation due to noise. We examine different fusion approaches for the radar and Lidar fusion; Table IV shows the mAP scores for the early, mid and late fusion schemes. The input images were stacked channel-wise and given to the network in early fusion. However, in late fusion, the feature maps were concatenated after the parallel forked structure, and there was no information exchange between the two branches. The mid-fusion approach outperforms the prior strategies, as it provides the network to learn the joint representation of features map between different sensor modalities and also provides robustness against spatial and temporal misalignment of data.

V. CONCLUSION

This study proposed a multimodal fusion network for object detection in adverse weather conditions using radar and Lidar data. A branched architecture is introduced, which has an input representation block, followed by convolution blocks I & II and the PFS block. The object detections are generated using a region proposal network. The efficacy of the algorithm is evaluated using the COCO metric. The mean average precision of 67.14 and 65.08 is achieved on test data for good and good-bad weather, respectively. The results show that the mid-fusion network performs better than the early or late fusion of features. Furthermore, using Lidar data improves the detection accuracy, and a large kernel size is adequate for capturing robust feature representation from images. In future work, we aim to extend this work to include images and extensively test the framework on other public datasets.



Groundtruth
 Predictions
 Lidar BEV

Fig. 4

ACKNOWLEDGMENT

This work is supported by the Academy of Finland Flagship program: Finnish Center for Artificial Intelligence (FCAI), and also by Culture, Sports and Tourism Research and Development Program through the Korea Creative Content Agency Grant funded by the Ministry of Culture, Sports and Tourism (R2022060001, “Development of Service Robot and Contents Supporting Children’s Reading Activities Based on Artificial Intelligence”).

REFERENCES

- [1] L.-H. Wen and K.-H. Jo, “Deep learning-based perception systems for autonomous driving: A comprehensive survey,” *Neurocomputing*, 2022.
- [2] S. Azam, F. Munir, A. M. Sheri, J. Kim, and M. Jeon, “System, design and experimental validation of autonomous vehicle in an unconstrained environment,” *Sensors*, vol. 20, no. 21, p. 5999, 2020.
- [3] G. Velasco-Hernandez, J. Barry, J. Walsh *et al.*, “Autonomous driving architectures, perception and data fusion: A review,” in *2020 IEEE 16th International Conference on Intelligent Computer Communication and Processing (ICCP)*. IEEE, 2020, pp. 315–321.
- [4] Y. Jiao, Z. Jie, S. Chen, J. Chen, L. Ma, and Y.-G. Jiang, “Msmdfusion: Fusing lidar and camera at multiple scales with multi-depth seeds for 3d object detection,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 21 643–21 652.
- [5] X. Chen, H. Ma, J. Wan, B. Li, and T. Xia, “Multi-view 3d object detection network for autonomous driving,” in *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, 2017, pp. 1907–1915.
- [6] H. Kuang, X. Liu, J. Zhang, and Z. Fang, “Multi-modality cascaded fusion technology for autonomous driving,” in *2020 4th International Conference on Robotics and Automation Sciences (ICRAS)*. IEEE, 2020, pp. 44–49.
- [7] Y. Zhang, A. Carballo, H. Yang, and K. Takeda, “Perception and sensing for autonomous vehicles under adverse weather conditions: A survey,” *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 196, pp. 146–177, 2023.
- [8] Y. Zhang, Q. Hu, G. Xu, Y. Ma, J. Wan, and Y. Guo, “Not all points are equal: Learning highly efficient point-based detectors for 3d lidar point clouds,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 18 953–18 962.
- [9] K. Qian, S. Zhu, X. Zhang, and L. E. Li, “Robust multimodal vehicle detection in foggy weather using complementary lidar and radar signals,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 444–453.
- [10] P. Li, P. Wang, K. Berntorp, and H. Liu, “Exploiting temporal relations on radar perception for autonomous driving,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 17 071–17 080.
- [11] Y.-J. Li, J. Park, M. O’Toole, and K. Kitani, “Modality-agnostic learning for radar-lidar fusion in vehicle detection,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 918–927.
- [12] S. Azam, F. Munir, and M. Jeon, “Channel boosting feature ensemble for radar-based object detection,” in *2021 IEEE Intelligent Vehicles Symposium (IV)*. IEEE, 2021, pp. 762–769.
- [13] S.-Y. Chu and M.-S. Lee, “Mt-detr: Robust end-to-end multimodal detection with confidence fusion,” in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2023, pp. 5252–5261.
- [14] R. Nabati and H. 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*, 2021, pp. 1527–1536.
- [15] Z. Wang, Y. Bai, Y. Zhou, and C. Xie, “Can cnns be more robust than transformers?” *arXiv preprint arXiv:2206.03452*, 2022.
- [16] J. Dickmann, N. Appenrodt, J. Klappstein, H.-L. Bloecher, M. Muntzinger, A. Sailer, M. Hahn, and C. Brenk, “Making bertha see even more: Radar contribution,” *IEEE Access*, vol. 3, pp. 1233–1247, 2015.
- [17] J. Lombacher, M. Hahn, J. Dickmann, and C. Wöhler, “Object classification in radar using ensemble methods,” in *2017 IEEE MTT-S International Conference on Microwaves for Intelligent Mobility (ICMIM)*. IEEE, 2017, pp. 87–90.
- [18] B. Major, D. Fontijne, A. Ansari, R. Teja Sukhavasi, R. Gowaikar, M. Hamilton, S. Lee, S. Grzechnik, and S. Subramanian, “Vehicle detection with automotive radar using deep learning on range-azimuth-doppler tensors,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops*, 2019, pp. 0–0.
- [19] D. Feng, C. Haase-Schütz, L. Rosenbaum, H. Hertlein, C. Glaeser, F. Timm, W. Wiesbeck, and K. Dietmayer, “Deep multi-modal object detection and semantic segmentation for autonomous driving: Datasets, methods, and challenges,” *IEEE Transactions on Intelligent Transportation Systems*, vol. 22, no. 3, pp. 1341–1360, 2020.
- [20] Y. Wang, Z. Jiang, Y. Li, J.-N. Hwang, G. Xing, and H. Liu, “Rodnet: A real-time radar object detection network cross-supervised by camera-radar fused object 3d localization,” *IEEE Journal of Selected Topics in Signal Processing*, vol. 15, no. 4, pp. 954–967, 2021.
- [21] F. Nobis, E. Shafiei, P. Karle, J. Betz, and M. Lienkamp, “Radar voxel fusion for 3d object detection,” *Applied Sciences*, vol. 11, no. 12, p. 5598, 2021.
- [22] S. Ren, K. He, R. Girshick, and J. Sun, “Faster r-cnn: Towards real-time object detection with region proposal networks,” *Advances in neural information processing systems*, vol. 28, 2015.
- [23] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, “Focal loss for dense object detection,” in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2980–2988.
- [24] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly *et al.*, “An image is worth 16x16 words: Transformers for image recognition at scale,” *arXiv preprint arXiv:2010.11929*, 2020.
- [25] Z. Liu, H. Mao, C.-Y. Wu, C. Feichtenhofer, T. Darrell, and S. Xie, “A convnet for the 2020s,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 11 976–11 986.
- [26] K. Chitta, A. Prakash, B. Jaeger, Z. Yu, K. Renz, and A. Geiger, “Transfuser: Imitation with transformer-based sensor fusion for autonomous driving,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022.
- [27] T. Meng, X. Jing, Z. Yan, and W. Pedrycz, “A survey on machine learning for data fusion,” *Information Fusion*, vol. 57, pp. 115–129, 2020.
- [28] S. Xie, R. Girshick, P. Dollár, Z. Tu, and K. He, “Aggregated residual transformations for deep neural networks,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 1492–1500.
- [29] D. Bahdanau, K. Cho, and Y. Bengio, “Neural machine translation by jointly learning to align and translate,” *arXiv preprint arXiv:1409.0473*, 2014.
- [30] T. Shen, T. Zhou, G. Long, J. Jiang, S. Pan, and C. Zhang, “Disan: Directional self-attention network for rnn/cnn-free language understanding,” in *Proceedings of the AAAI conference on artificial intelligence*, vol. 32, no. 1, 2018.
- [31] E. H. Adelson, C. H. Anderson, J. R. Bergen, P. J. Burt, and J. M. Ogden, “Pyramid methods in image processing,” *RCA engineer*, vol. 29, no. 6, pp. 33–41, 1984.
- [32] Z. Li, C. Peng, G. Yu, X. Zhang, Y. Deng, and J. Sun, “Detnet: Design backbone for object detection,” in *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 334–350.
- [33] F. Yu and V. Koltun, “Multi-scale context aggregation by dilated convolutions,” *arXiv preprint arXiv:1511.07122*, 2015.
- [34] M. Sheeny, E. De Pellegrin, S. Mukherjee, A. Ahrabian, S. Wang, and A. Wallace, “Radiate: A radar dataset for automotive perception in bad weather,” in *2021 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2021, pp. 1–7.
- [35] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, “Microsoft coco: Common objects in context,” in *Computer Vision—ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part V 13*. Springer, 2014, pp. 740–755.
- [36] D. Barnes, M. Gadd, P. Murcutt, P. Newman, and I. Posner, “The oxford radar robotcar dataset: A radar extension to the oxford robotcar dataset,” in *2020 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2020, pp. 6433–6438.