

# Pay “Attention” to Adverse Weather: Weather-aware Attention-based Object Detection

Saket S. Chaturvedi

College of Computing  
Michigan Technological University  
Houghton, MI, USA  
schatur2@mtu.edu

Lan Zhang

Dept of Electrical and Computer Engineering  
Michigan Technological University  
Houghton, MI, USA  
lanzhang@mtu.edu

Xiaoyong Yuan

College of Computing  
Michigan Technological University  
Houghton, MI, USA  
xyyuan@mtu.edu

**Abstract**—Despite the recent advances of deep neural networks, object detection for adverse weather remains challenging due to the poor perception of some sensors in adverse weather. Instead of relying on one single sensor, multimodal fusion has been one promising approach to provide redundant detection information based on multiple sensors. However, most existing multimodal fusion approaches are ineffective in adjusting the focus of different sensors under varying detection environments in dynamic adverse weather conditions. Moreover, it is critical to simultaneously observe local and global information under complex weather conditions, which has been neglected in most early or late-stage multimodal fusion works. In view of these, this paper proposes a Global-Local Attention (GLA) framework to adaptively fuse the multi-modality sensing streams, i.e., camera, gated camera, and lidar data, at two fusion stages. Specifically, GLA integrates an early-stage fusion via a local attention network and a late-stage fusion via a global attention network to deal with both local and global information, which automatically allocates higher weights to the modality with better detection features at the late-stage fusion to cope with the specific weather condition adaptively. Experimental results demonstrate the superior performance of the proposed GLA compared with state-of-the-art fusion approaches under various adverse weather conditions, such as light fog, dense fog, and snow.

**Index Terms**—Object Detection, Adverse Weather, Multimodal Fusion, Attention Neural Network

## I. INTRODUCTION

Object detection is a fundamental task in autonomous driving. Reliable and robust detection is crucial to various downstream autonomous driving tasks, such as object tracking [1], path planning [2], and motion control [3]. Although recent advances in deep neural networks have significantly improved object detection performance in normal weather, their performance in adverse weather such as light fog, dense fog, and snow has been challenging [4].

Object detection in adverse weather is significantly hindered by the poor performance of some sensors. For example, cameras alone are unreliable in low-light conditions as they do not provide depth information and are highly affected by adverse weather [4]. Similarly, the effective range of lidar is highly impaired in fog and snow due to backscatter [5]. To provide robust perception in adverse weather, multimodal sensor fusion is one of the most promising approaches since it leverages the benefits from multiple sensors to provide redundant detection information.



Fig. 1: The single modality model detections (Fig. 1 (a), 1 (b), and 1 (c)) failed to detect the Passenger Cars in the Adverse weather condition, which can be due to severe fog backscatter. The proposed Global Local Attention framework (Fig. 1 (d)) adapts to adverse weather by collaborative learning from Global and Local Attention features using the camera, gated, and lidar modalities.

Most existing multimodal fusion methods [3], [6]–[8] use simple concatenation or element-wise addition methods to fuse multiple modalities for object detection. These methods fail to adjust the focus of different sensors to deal with the dynamic conditions in adverse weather. In addition, the existing works [3], [6]–[10] use early or late-stage multimodal fusion and only consider low-level or high-level features from different sensors, which can be insufficient to deal with complex adverse weather conditions. For example, the scales of the objects vary during object detection. Using only local information (low-level features) or global information (high-level features) does not allow the model to fuse multimodality data in different object scales, which, however, is critical in adverse weather conditions.

To address the above-mentioned challenges, we leverage the attention method to adjust the focus of different sensors to deal with the dynamic adverse weather conditions. Additionally, we sequentially extract the local and global information from multiple sensors to adjust to the complex adverse weather conditions. We propose a Global-Local Attention (GLA) framework to adaptively fuse data from multiple sensors in two stages. In particular, the GLA framework learns

the local and global information using the Local Attention Network and Global Attention network and thus adaptively leverages the best camera, gated, and lidar features to deal with specific weather conditions. Each Attention Network learns to allocate higher weights to the modality with better detection features. To evaluate the performance of the GLA framework, we conduct experiments on the DENSE Dataset [11]. We compare our GLA framework with the state-of-art research [7], [11]. To further analyze the key components in GLA, we investigate the impact of the number of partitions in the Global Attention Network and the impact of using only Global or Local Attention.

Our main contributions are as follows.

- We propose a Global-Local Attention (GLA) framework to tackle object detection in adverse weather using the camera, gated, and lidar sensors to perform the dynamic fusion.
- We perform fusion at the two stages in our proposed GLA framework. First, the Local Attention features are used as an early-stage fusion, and second, the Global Attention features as a later stage fusion method. The collaborative learning from local and global features makes our model adaptive to deal with the highly-variable adverse weather conditions.
- We rigorously evaluate the proposed GLA framework. The proposed GLA framework achieves an average improvement of 19.83% mAP and 12.645% mAP compared with the state-of-the-art SSD-VGG model [7] and SSD Deep Entropy fusion model [11].
- We provide an ablation study to investigate the effectiveness of the attention map and the partition number's impact. We further demonstrate the effectiveness of Global-Local attention by comparing it with the global or local only attention methods.

## II. RELATED WORK

### A. Multimodal Fusion

Multimodal fusion aims to provide robust prediction in object detection using multi-modality inputs. Simon *et al.* [3] proposed concatenation or element-wise addition for fusing low-level radar and camera features. Similarly, Xiaozhi *et al.* [8] proposed the concatenation for early/late-stage fusion and element-wise addition for the deep fusion of lidar and camera features in the architecture. Recently, attention-based fusion methods have attracted great interest in multimodal fusion. Zhang *et al.* [9] proposed a hybrid attention-aware fusion network (HAFNet) based on a cross-modal attention mechanism for HRI and Lidar features fusion using ATT-AFBlock. They used a global pooling layer followed by several FC layers and a sigmoid layer to perform fusion by interpreting channel-wise sigmoid weights of each modality. Dai *et al.* [12] introduced the idea of attentional feature fusion for fusing single modality features of different scales using sigmoid layers in the Global and Local Attention aggregation for a classification task. In comparison to their method, the focus

of our paper is to perform object detection in adverse weather by performing multimodal adaptive fusion by considering both local and global features. In the proposed Global Attention Network, we extract global partition-level features, i.e., global features for different partitions of an image, instead of focusing on image-level global features. Li *et al.* [10] integrated Global Attention with the YOLOv3 model for object detection. The attention block included the Global Max Pooling and Global Average Pooling layers to generate two attention tensor types and concatenate them before passing through FC layers and the softmax layer to map the output to a probability distribution. However, the global max pooling or global average pooling extracts only image-level features instead of local and partition-level features, which can be insufficient for detecting objects in complex adverse weather conditions.

### B. Object Detection in Adverse Weather

The major challenge in the multimodal fusion for object detection in adverse weather is developing a fusion method that can adapt to highly-variable adverse weather conditions. Mees *et al.* [13] showed that a novel mixture of deep network experts could be used to adaptively weight the prediction of different modality classifiers in harsh light conditions. Pfeuffer and Dietmayer [6] focused on finding the optimal fusion method between early, mid, and late fusion methods for object detection in adverse weather and used concatenation to fuse the camera and lidar features. Moreover, Debrigode and Guillem [7] developed the SSD-VGG halfway fusion model on a Dense dataset for Adversarial weather conditions. They used a simple concatenation method to fuse the Camera and Gated features. However, these studies use the conventional fusion method, such as concatenation that cannot adapt to dynamic adverse weather conditions. Additionally, early or late fusion allows the model to consider only low-level features or high-level features during fusion. Bijelic *et al.* [11] focused on weighted continuous multimodal fusion by calculating a multiplication matrix using entropy on inputs to scale each sensor based on available information. Their method performed an adaptive fusion of camera, gated, and lidar sensor features but lacked consideration of local (low-level) and global (partition-level) attention features during fusion. Therefore, our proposed GLA framework aims to learn from local and global features that can help detect partially visible objects in adverse weather conditions.

## III. GLOBAL-LOCAL ATTENTION FRAMEWORK

### A. Overview

The proposed Global Local Attention framework depicted in Figure 2 consists of three branches corresponding to the camera, gated, and lidar streams. The input is processed using the Inception block for each stream, then the Local Attention Network calculates local attention feature weights and performs the first-stage multimodal fusion. Further, the Global Attention Network takes Inception block processed input (same as Local Attention Network) and the local attention feature map from first-stage fusion after processing

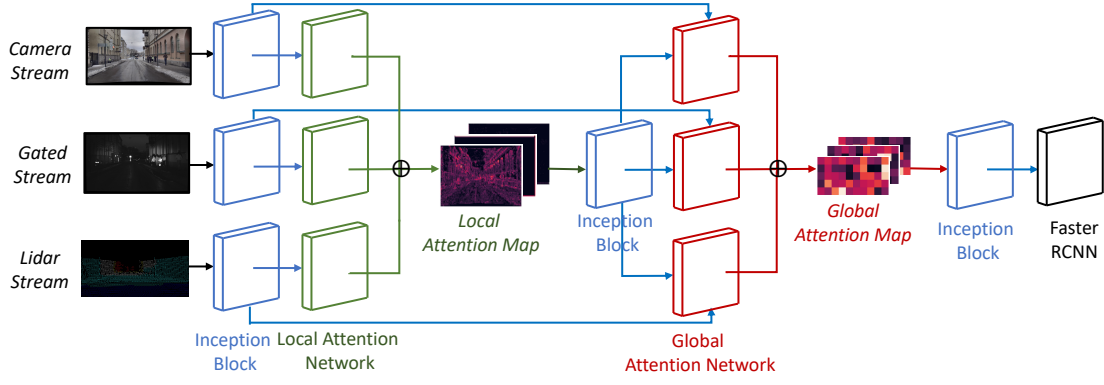


Fig. 2: The proposed Global Local Attention framework consists of three branches: camera, gated, and lidar streams. For each stream, the Local Attention Network (green) and Global Attention Network (red) take input from inception blocks (blue) which are either Camera, Gated, or Lidar features. The Global Attention Network has an additional input as a local attention feature. The Local/Global Attention Map is derived from the Local/Global Attention Network to perform multimodal weighted fusion represented with  $\oplus$ . Finally, the FasterRCNN analyzes the fused feature maps for object detection and outputs the final prediction.

with Inception blocks to perform the second-stage multimodal fusion. The final feature map is sent to the FasterRCNN's Region Proposal Network [14], a widely used object detection framework. The Region Proposal Network generates anchors of different sizes and scales. Finally, the object classification and bounding boxes are generated after performing non-max suppression on these anchors.

### B. Local Attention Network

We develop a Local Attention Network to select the best local features from different modality inputs representing each pixel. Given an intermediate Camera, Gated, Lidar features (extracted using the first Inception block)  $\mathbf{F}_{CA}, \mathbf{F}_G, \mathbf{F}_L \in \mathbf{R}^{H \times W \times C}$  with  $H \times W$  feature map and  $C$  channels. The local attention weight (LA) for camera, gated, and lidar streams can be computed as:

$$\mathbf{LA}_I = \text{Conv}(\text{BN}(\text{Conv}(\sigma(\text{BN}(\text{Conv}(\mathbf{F}_I))))), \quad (1)$$

where  $\mathbf{F}_I = \text{IB}(\mathbf{X})$ , IB denotes Inception Block, X denotes input feature for camera, gated, and lidar streams,  $\mathbf{LA}_I \in \mathbf{R}^{H \times W \times C}$  denotes the local attention weights for camera ( $LA_{CA}$ ), gated ( $LA_G$ ), and lidar ( $LA_L$ ) streams,  $\mathbf{F}_I$  represents corresponding camera ( $F_{CA}$ ), gated ( $F_G$ ), and lidar ( $F_L$ ) intermediate features. Conv denotes conv2d layer, BN denotes Batch Normalization layer, and  $\sigma$  denotes ReLU activation function. We then apply a channel-wise softmax operation to derive the local attention weights for the camera, gated, and lidar streams. Leveraging the local attention weights, the first-stage multimodal feature fusion calculates a local attention feature map  $F1$  to select the best feature for each pixel from different modality inputs. The calculation is summarized as follows:

$$\mathbf{LA}'_{CA}, \mathbf{LA}'_G, \mathbf{LA}'_L = \text{Softmax}(\mathbf{LA}_{CA}, \mathbf{LA}_L, \mathbf{LA}_G) \quad (2)$$

$$\mathbf{F1} = \mathbf{F}_{CA} * \mathbf{LA}'_{CA} + \mathbf{F}_G * \mathbf{LA}'_G + \mathbf{F}_L * \mathbf{LA}'_L. \quad (3)$$

### C. Global Attention Network

We develop a Global Attention Network, which performs partitions of an input feature to select the best global features from different modality inputs for each partition.

Moreover, given the camera stream, gated stream, lidar stream features (extracted using the first Inception block)  $\mathbf{F}_{CA}, \mathbf{F}_G, \mathbf{F}_L \in \mathbf{R}^{H \times W \times C}$  with  $H \times W$  feature map and  $C$  channels. In the proposed Global Network, we partition an input feature among 5 rows and 10 columns, resulting in a total of 50 partitions before using the average pooling layer over each partition to focus on the features of partition-level features instead of image-level features. Similar to Local Attention Network, we use a conv2d, batch normalization layers, followed by ReLU activation to extract features from inputs.

$$\mathbf{GA}_I = \text{Conv}(\text{BN}(\text{Conv}(\sigma(\text{BN}(\text{Conv}(\text{AP}(\mathbf{F}_I)))))), \quad (4)$$

where  $\mathbf{GA}_I \in \mathbf{R}^{H \times W \times C}$  denote the global attention feature weights for camera ( $GA_{CA}$ ), gated ( $GA_G$ ), and lidar ( $GA_L$ ) streams,  $\mathbf{F}_I$  represents corresponding camera ( $F_{CA}$ ), gated ( $F_G$ ), and lidar ( $F_L$ ) intermediate Features. Further, we apply channel-wise Softmax to compute the global attention weight (GA) for camera, gated, and lidar streams as follows:

$$\mathbf{GA}'_{CA}, \mathbf{GA}'_G, \mathbf{GA}'_L = \text{Softmax}(\mathbf{GA}_{CA}, \mathbf{GA}_G, \mathbf{GA}_L) \quad (5)$$

Finally, the second-stage fusion was performed using  $\mathbf{F1}'$ , derived from processing first-stage feature fusion  $F1$  with inception block and Global Attention Weights ( $\mathbf{GA}'_{CA}, \mathbf{GA}'_G, \mathbf{GA}'_L$ ). The global feature map  $F2$  is calculated as follows.

$$\mathbf{F2} = \mathbf{F1}' * \mathbf{GA}'_{CA} + \mathbf{F1}' * \mathbf{GA}'_G + \mathbf{F1}' * \mathbf{GA}'_L, \quad (6)$$

where  $\mathbf{F1}' = \text{IB}(\mathbf{F1})$ , IB denotes Inception Block and  $F1$  denotes first-stage feature fusion. The proposed fusion model

TABLE I: The performance comparison of baseline models with the proposed Global-Local Attention method over PassengerCar class. We report mAP (%) on the DENSE dataset over easy (E), moderate (M), and hard (H) difficulty levels.

Models	Day/ Night	Weather Conditions											
		Clear			LightFog			DenseFog			Snow		
		E (%)	M (%)	H (%)	E (%)	M (%)	H (%)	E (%)	M (%)	H (%)	E (%)	M (%)	H (%)
Camera	Day	92.23	57.17	38.76	97.86	62.37	20.62	90.37	22.99	35.11	92.78	51.40	36.82
	Night	87.49	49.02	36.45	91.57	56.11	45.92	90.28	52.56	22.03	94.54	48.55	29.96
Gated	Day	86.81	55.52	14.78	95.12	55.52	14.78	91.55	31.35	22.67	91.00	49.84	32.03
	Night	88.24	44.63	41.26	88.24	44.63	41.26	86.42	24.39	18.37	88.52	34.15	29.06
Lidar	Day	85.62	46.13	27.55	97.65	59.16	16.97	94.42	23.31	39.25	94.43	55.37	31.59
	Night	88.63	46.81	33.47	92.62	55.75	39.83	90.12	48.35	14.23	94.29	47.52	25.87
Camera-Gated	Day	<b>95.40</b>	71.79	54.22	98.65	71.75	39.95	<b>97.90</b>	53.12	57.83	<b>96.97</b>	<b>70.44</b>	55.62
	Night	93.69	63.26	52.18	96.62	<b>69.53</b>	60.17	<b>94.63</b>	63.11	38.23	96.39	<b>67.74</b>	50.55
Camera-Lidar	Day	95.33	72.26	<b>55.21</b>	99.03	77.38	<b>41.78</b>	97.86	48.08	56.55	96.22	68.97	57.40
	Night	<b>94.19</b>	63.43	53.37	96.82	66.36	65.25	93.38	<b>63.53</b>	37.62	96.43	66.44	50.73
Camera-Gated-Lidar (Concat)	Day	93.87	63.96	43.63	97.75	60.77	22.32	96.79	36.35	39.70	95.66	62.09	46.73
	Night	91.35	56.42	47.47	95.12	60.43	55.81	93.32	47.87	24.40	95.47	57.32	36.54
<b>Proposed Method</b>	Day	<b>95.40</b>	<b>73.09</b>	<b>55.00*</b>	<b>99.05</b>	<b>79.69</b>	<b>40.85*</b>	97.82	<b>56.47</b>	<b>58.16</b>	<b>96.86*</b>	68.39	<b>53.65</b>
	Night	93.47	<b>63.44</b>	<b>54.15</b>	<b>97.12</b>	<b>68.11*</b>	<b>66.67</b>	<b>94.36*</b>	60.61	<b>43.96</b>	<b>97.48</b>	<b>67.28*</b>	<b>51.19</b>

\* represents the second-best model results for the specified Day/Night weather condition.

for each partition assigns higher weights to the modality with better features and the weights for each corresponding partition in the global attention map sum up to one.

#### IV. EVALUATION

##### A. Dataset Description

We evaluate the proposed GLA framework on the DENSE dataset [11]. The dataset contains multiple modality inputs, including camera, gated, and lidar sensors for clear and adverse weather conditions such as light fog, dense fog, and snow. Following [11], we project the three-dimensional lidar data to front-view considering depth, height, and pulse-intensity values. The original Gated images in the DENSE dataset of dimension  $1280 \times 768$  were projected to the RGB camera plane by adding black padding of 210, 46, 280, and 360 to the top, bottom left, and right borders, respectively.

In our work, we considered PassengerCar, Pedestrian, LargeVehicle, and Ridable Vehicle classes among the total 14 classes available in the dataset, similar to [7]. We constructed the training dataset consisting of samples from all weather conditions. For clear weather samples, the information about the training and test dataset samples was taken from the GitHub repository [11]. For other weather conditions, we randomly split the available samples in the training and test datasets such that around 20 percent of Day and Night samples are in the test dataset.

##### B. Experimental Settings

The proposed GLA method is a general framework for object detection in adverse weather, which can be used with both single-stage and two-stage object detection models. In this work, we have used the GLA framework in conjunction with the Faster RCNN model (one of the most commonly used object detection models) and implemented using Tensorflow Object Detection API. We used the FasterRCNN-InceptionV2 pre-trained model weights on the COCO dataset to perform

training with a constant learning rate of 0.00002 using the Adam optimizer by setting a batch size to 2. We use the Weighted Smooth L1 Localization Loss (Huber Loss) function for both the first and second stages of the FasterRCNN model. However, the Classification Loss Function was changed from Softmax Cross-Entropy in the first stage to Sigmoid Focal loss in the second stage to address the class imbalance problem in the dataset. We also used Non-Maximum Suppression (NMS) [15] to filter the predicted bounding box above the 0.7 IoU threshold. Finally, the performance evaluation of the proposed and baseline models was conducted using mean average precision (mAP) with IoU 0.5 using the Pascal VOC Detection Metrics across easy (E), moderate (M), and hard (H) difficulty cases following [16].

##### C. Evaluation Results

This section evaluates our proposed fusion model on the individual test dataset for each weather. Similar to the state-of-art method [11], we evaluated the proposed GLA framework over PassengerCar class because these are the most prevalent class in the DENSE dataset.

We compare our proposed GLA framework with previous works [7], [11] on the DENSE dataset, consisting of Day/Night samples for Clear, light fog, dense fog, and snow weather conditions. The study [11] showed the superiority of their method over [17]–[20] object detection methods for adverse weather on the DENSE dataset. Hence, we also show the effectiveness of our proposed GLA framework for adverse weather over [17]–[20] methods by comparing it with [11].

We perform a comparison of the proposed GLA framework with three single modality baseline models: Camera-Only, Gated-Only, Lidar-Only, and two fusion baseline models: Camera-Gated, Camera-Lidar on the test dataset. Also, we compare the proposed GLA framework with the simple concatenation model, which has the same architecture except using concatenation to perform fusion instead of Global-Local

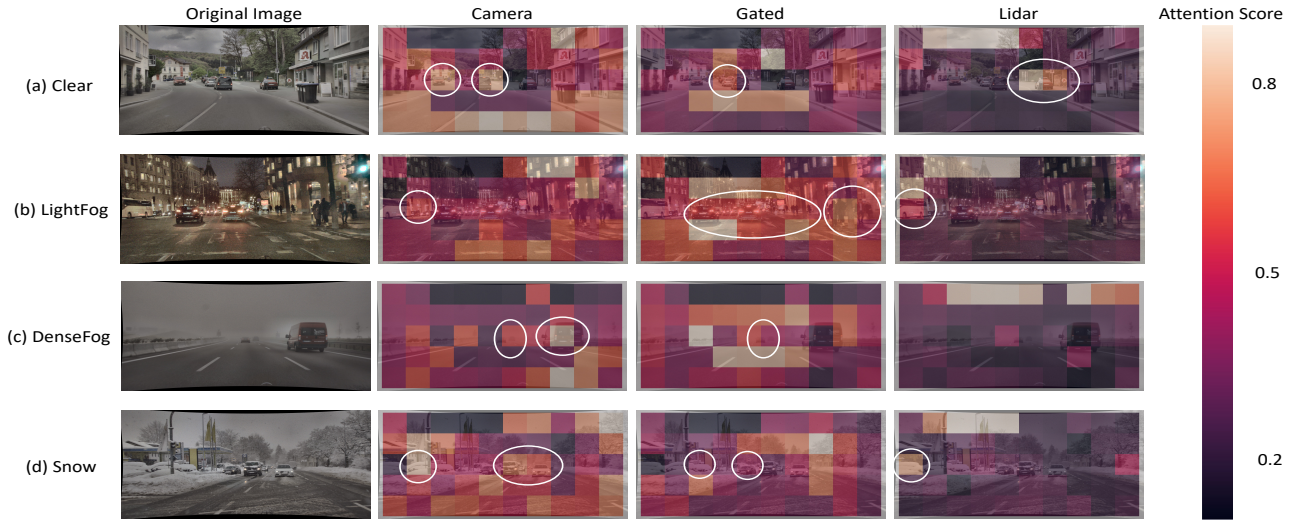


Fig. 3: The camera, gated, lidar attention maps demonstrate the effectiveness of our proposed GLA framework in the clear weather (Fig. 3(a)) and adverse weather conditions, light fog (Fig. 3(b)), dense fog (Fig. 3(c)), snow (Fig. 3(d)). The attention maps show the attention score corresponding to different colors in the color map. The *white* Ellipse in the Attention map highlights each modality’s best detection features (i.e., higher attention score features representing PassengerCar, Pedestrian, LargeVehicle, or Ridable Vehicle).

Attention. Table I presents the results of our proposed method and other associated baseline model results over PassengerCar class. The single modality models were trained with training hyperparameters following Kitti Config in TensorFlow Object Detection API, and two baseline fusion models were trained with similar training hyperparameters as the proposed fusion model.

The camera-only model performed best across all weather conditions among the single modality models. The proposed GLA framework has an average performance improvement of 13.72% over the camera-only model and an average performance improvement of 8.79% over a simple concatenation-based method keeping the modalities same; this shows the effectiveness of our proposed method. The adaptive feature fusion of camera, gated, and lidar modalities using Local Attention and Global Attention at two stages helps perform collaborative learning from each sensor by feature exchange. The proposed GLA framework, other than single modality models, also mostly outperformed or has the second-best results compared to the camera-gated and camera-lidar baseline models.

Table II compares the proposed GLA framework with the state-of-art multimodal fusion models SSD Deep Entropy Fusion [11] and SSD-VGG [7] over PassengerCar class. In the study [7], the proposed SSD-VGG model’s overall difficulty evaluation in different weather conditions is conducted instead of individual difficulty evaluation with easy, moderate, and hard weather conditions. So, we also compare the overall difficulty evaluation results of our proposed GLA framework with the state-of-art models used in this study. Overall, our proposed GLA framework outperformed the SSD-VGG [7] model by a margin of 23.05%-19.72%, 20.74%-

TABLE II: Performance comparison of proposed GLA framework with the state-of-art methods in terms of mAP (%) over PassengerCar class.

Model	Day/Night	Weather Conditions			
		Clear (%)	LightFog (%)	DenseFog (%)	Snow (%)
SSD-VGG [7]	Day	56.09	67.18	69.71	58.61
	Night	58.60	61.01	70.18	62.62
SSD Deep [11]	Day	67.13	75.47	74.75	69.24
	Night	62.20	71.96	68.72	72.04
Entropy Fusion	Day	<b>79.14</b>	<b>87.92</b>	<b>87.99</b>	<b>78.79</b>
	Night	<b>78.32</b>	<b>84.50</b>	<b>82.90</b>	<b>83.08</b>

23.49%, 18.28%-12.72%, and 20.18%-20.46% mAP during Day-Night time for clear, light fog, dense fog, and snow weather, respectively. For reproducing the SSD Deep Entropy Fusion [11] model, as the authors have not released their model codes, experimental settings, and training hyperparameters, we follow the experimental settings of this study. To perform a fair comparison, we reproduce their fusion model logic using SSD-Inception and perform training using a pretrained model since we have also used an inception feature extractor and a pretrained model for training. *Compared with their approach, the proposed GLA framework improved the object detection mAP by 12.04%-16.12%, 12.45%-12.54%, 13.24%-14.18%, and 9.55%-11.04% mAP during Day-Night time for clear, light fog, dense fog, and snow weather, respectively.*

## V. ABLATION STUDY

This section provides an ablation study to investigate the effectiveness of the attention map and the impact of partition numbers in the Global Attention Network. We further demon-



strate the effectiveness of Global-Local attention by comparing it with the global or local only attention methods.

#### A. Attention Maps.

To demonstrate the effectiveness of our proposed model in Adverse weather conditions, we visualize the camera, gated, and lidar attention map using the heatmap in Fig. 3 over clear, light fog, dense fog, and snow weathers. In the attention maps, light color represents high attention weights (0.8), violet color represents moderate attention weights (0.5), and dark color represents low attention weights (0.2). Our method builds attention weights for 50 partitions of an image, where different modalities can get attended for a different partition of an image. This approach makes our method adaptive to adverse weather conditions, where different sensors can be preferred depending on weather conditions.

This paragraph demonstrates the effectiveness of the proposed GLA framework in clear and light fog weather during both day and night time. The attention maps for a sample clear weather image in Fig. 3(a) show that the gated modality weights focused on a few distant passenger cars and the lidar attention map focused on the passenger car on the right side. At the same time, camera modality weights focused on the multiple passenger cars in two partitions of the attention map. This shows that each camera, gated, and lidar modality contributed to the passenger cars detection in the clear day weather condition. Fig. 3(b) shows the attention maps for a sample light fog weather image. Since the performance of the camera and lidar sensor can be affected during nighttime or due to fog backscatter, only a few passenger cars and large vehicles are focused on their attention map. On the other hand, the gated attention map focused on almost every center passenger car and the pedestrians walking on the right side of the road.

We also demonstrate the effectiveness of our proposed GLA framework in dense fog and snow weather conditions. In the dense fog weather, each camera and gated attention map focused on the center distant passenger cars and a right side passenger car in Fig. 3(c). In the snow weather (Fig. 3(d)), the camera attention map has a higher attention weight for center passenger cars and left-side passenger cars, almost hidden in the snow. The gated and lidar sensor also contributed to detecting a few centers and left-side passenger cars.

#### B. Impact of the Number of Partitions in Global Attention Network.

To study the impact of the number of partitions in our method, we construct different variants of the Global-Local Attention Networks. The Local Attention Network and all training hyperparameters were kept the same for all model variants in this experiment. The only difference is the number of partitions in the Global Network, which were set to 18 (3 rows, 6 columns, 32 (4 rows, 8 columns), 50 (5 rows, 10 columns), and 72 (6 rows, 12 columns) in the different variants.

TABLE III: Impact of Number of Partitions in Global Attention Network for camera-gated-lidar fusion in terms of mAP (%) compared to proposed method over PassengerCar class.

Model	Day/Night	Weather Conditions			
		Clear (%)	LightFog (%)	DenseFog (%)	Snow (%)
18 partitions	Day	78.92	87.10	86.08	78.50
	Night	78.01	84.58	81.57	81.42
32 partitions	Day	79.15	86.68	86.33	<b>79.40</b>
	Night	78.27	83.97	82.09	82.62
<b>50 partitions</b>	Day	<b>79.15</b>	<b>87.92</b>	<b>87.99*</b>	78.79
	Night	<b>78.32*</b>	<b>84.50</b>	<b>82.90</b>	<b>83.08</b>
72 partitions	Day	78.93	85.56	<b>88.42</b>	79.21
	Night	<b>78.75</b>	84.39	80.96	82.90

\* represents the second-best model results for the specified Day/Night weather condition.

TABLE IV: Impact of using Global Attention only or Local Attention only for camera-gated-lidar fusion in terms of mAP (%) compared to the proposed Global Local Attention method over PassengerCar class.

Model	Day/Night	Weather Conditions			
		Clear (%)	LightFog (%)	DenseFog (%)	Snow (%)
Global Attention	Day	78.01	87.46	85.89	<b>79.11</b>
	Night	<b>78.80</b>	82.47	81.24	82.56
Local Attention	Day	78.79	86.23	86.93	78.64
	Night	78.04	84.02	82.73	81.69
<b>Proposed Method</b>	Day	<b>79.14</b>	<b>87.92</b>	<b>87.99</b>	<b>78.79*</b>
	Night	<b>78.32*</b>	<b>84.50</b>	<b>82.90</b>	<b>83.08</b>

\* represents the second-best model results for the specified Day/Night weather condition.

Table III represents the results of four partition variants of the Global-Local Attention Network on the test dataset for overall weather conditions. In most cases, the Global Network with 50 partitions outperformed the 18 and 32 partition variants of Global Network, with exceptions for daytime snow weather cases. In this experiment, We noted the general trend of increase in mAP with an increase in partitions from 18 to 32 to 50. However, when the number of partitions in the Global Network was set to 72, the mAP mainly dropped, except for a few cases. Hence, the results suggest that the number of partitions in the Global Attention Network has an important role in working the proposed Global Local Attention model, which gives the best results with 50 partitions.

#### C. Impact of using only Global or Local Attention method.

Further, we investigate the impact of using Global Attention only or Local Attention only instead of the proposed Global-Local Attention method for the Adverse weather conditions. We keep the model's architecture the same except using Local Attention only for early-stage multimodal fusion or Global Attention only for late-stage multimodal fusion.

Table IV summarizes the results for the Global Attention model, Local Attention model, and the proposed Global-Local

Attention method. The Local Attention model showed better performance than the Global Attention model in terms of mAP for dense fog weather conditions, and the Global Attention model yielded better results for snow weather conditions. However, the superiority between the Global Attention and Local Attention model for clear and light fog weather is not clear. Notably, the proposed Global-Local Attention model outperformed the Local Attention model and Global Attention model in almost every weather condition, except daytime snow and clear nighttime weather, where the proposed method yielded second-best results. This helps to understand that although the roles of Local Attention and Global Attention are different, they are complementary and work better in combination. Hence, we propose using the Global Local Attention model for object detection in adverse weather.

## VI. CONCLUSION

Recent advances in deep learning have greatly improved object detection for autonomous driving in normal weather; however, progress remains limited in adverse weather conditions. The current state-of-the-art methods are ineffective in dealing with the dynamic nature of the adverse weather by using conventional fusion methods and only considering local or global information for multimodal fusion. This work proposes Global-Local Attention (GLA), a general framework for object detection in adverse weather conditions to address these issues. The GLA framework extracts the global and local attention feature map at two stages to adaptively leverage the best of the camera, gated, and lidar features in the multimodal fusion, resulting in superior performance compared with the state-of-the-art methods. The proposed GLA framework outperformed the state-of-the-art research models, leading to an average improvement of 19.83% mAP and 12.645% mAP over the SSD-VGG and SSD Deep Entropy Fusion models.

## ACKNOWLEDGEMENT

This work is supported in part by MTU Research Excellence Fund (REF) Award.

## REFERENCES

- [1] Ankith Manjunath, Ying Liu, Bernardo Henriques, and Armin Engstle. Radar based object detection and tracking for autonomous driving. In *2018 IEEE MTT-S International Conference on Microwaves for Intelligent Mobility (ICMIM)*, pages 1–4. IEEE, 2018.
- [2] Vu Trieu Minh and John Pumwa. Feasible path planning for autonomous vehicles. *Mathematical Problems in Engineering*, 2014, 2014.
- [3] Shobit Sharma, Girma Tewolde, and Jaerock Kwon. Lateral and longitudinal motion control of autonomous vehicles using deep learning. In *2019 IEEE International Conference on Electro Information Technology (EIT)*, pages 1–5. IEEE, 2019.
- [4] Shizhe Zang, Ming Ding, David Smith, Paul Tyler, Thierry Rakotoarivelo, and Mohamed Ali Kaafar. The impact of adverse weather conditions on autonomous vehicles: how rain, snow, fog, and hail affect the performance of a self-driving car. *IEEE vehicular technology magazine*, 14(2):103–111, 2019.
- [5] Robin Heinzler, Florian Piewak, Philipp Schindler, and Wilhelm Stork. Cnn-based lidar point cloud de-noising in adverse weather. *IEEE Robotics and Automation Letters*, 5(2):2514–2521, 2020.
- [6] Andreas Pfeuffer and Klaus Dietmayer. Optimal sensor data fusion architecture for object detection in adverse weather conditions. In *2018 21st International Conference on Information Fusion (FUSION)*, pages 1–8. IEEE, 2018.

- [7] Guillem Tudela Debrigode. Design of an image fusion object detection algorithm for the autonomous driving during adversarial weather conditions. 2021.
- [8] Xiaozhi Chen, Huimin Ma, Ji Wan, Bo Li, and Tian Xia. Multi-view 3d object detection network for autonomous driving. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 1907–1915, 2017.
- [9] Peng Zhang, Peijun Du, Cong Lin, Xin Wang, Erzhu Li, Zhaohui Xue, and Xuyu Bai. A hybrid attention-aware fusion network (hafnet) for building extraction from high-resolution imagery and lidar data. *Remote Sensing*, 12(22):3764, 2020.
- [10] Wei Li, Kai Liu, Lizhe Zhang, and Fei Cheng. Object detection based on an adaptive attention mechanism. *Scientific Reports*, 10(1):1–13, 2020.
- [11] Mario Bijelic, Tobias Gruber, Fahim Mannan, Florian Kraus, Werner Ritter, Klaus Dietmayer, and Felix Heide. Seeing through fog without seeing fog: Deep multimodal sensor fusion in unseen adverse weather. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11682–11692, 2020.
- [12] Yimian Dai, Fabian Gieseke, Stefan Oehmcke, Yiquan Wu, and Kobus Barnard. Attentional feature fusion. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 3560–3569, 2021.
- [13] Oier Mees, Andreas Eitel, and Wolfram Burgard. Choosing smartly: Adaptive multimodal fusion for object detection in changing environments. In *2016 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 151–156. IEEE, 2016.
- [14] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. *Advances in neural information processing systems*, 28:91–99, 2015.
- [15] Jan Hosang, Rodrigo Benenson, and Bernt Schiele. Learning non-maximum suppression. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4507–4515, 2017.
- [16] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? the kitti vision benchmark suite. In *2012 IEEE conference on computer vision and pattern recognition*, pages 3354–3361. IEEE, 2012.
- [17] Eric Tzeng, Judy Hoffman, Kate Saenko, and Trevor Darrell. Adversarial discriminative domain adaptation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7167–7176, 2017.
- [18] Judy Hoffman, Eric Tzeng, Taesung Park, Jun-Yan Zhu, Phillip Isola, Kate Saenko, Alexei Efros, and Trevor Darrell. Cycada: Cycle-consistent adversarial domain adaptation. In *International conference on machine learning*, pages 1989–1998. PMLR, 2018.
- [19] Charles R Qi, Wei Liu, Chenxia Wu, Hao Su, and Leonidas J Guibas. Frustum pointnets for 3d object detection from rgb-d data. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 918–927, 2018.
- [20] Jason Ku, Melissa Mozifian, Jungwook Lee, Ali Harakeh, and Steven L Waslander. Joint 3d proposal generation and object detection from view aggregation. In *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 1–8. IEEE, 2018.