

# RoboFusion: Towards Robust Multi-Modal 3D Object Detection via SAM

Ziying Song<sup>1</sup>, Guoxing Zhang<sup>2</sup>, Lin Liu<sup>1</sup>, Lei Yang<sup>3</sup>, Shaoqing Xu<sup>4</sup>, Caiyan Jia<sup>1\*</sup>,  
Feiyang Jia<sup>1</sup>, Li Wang<sup>5</sup>

<sup>1</sup>School of Computer Science and Technology & Beijing Key Lab of Traffic Data Analysis and Mining, Beijing Jiaotong University

<sup>2</sup>Hebei University of Science and Technology <sup>3</sup>Tsinghua University <sup>4</sup>University of Macau <sup>5</sup>Beijing Institute of Technology

{songziying, cyjia, feiyangjia}@bjtu.edu.cn

## Abstract

Multi-modal 3D object detectors are dedicated to exploring secure and reliable perception systems for autonomous driving (AD). Although achieving state-of-the-art (SOTA) performance on clean benchmark datasets, they tend to overlook the complexity and harsh conditions of real-world environments. With the emergence of visual foundation models (VFM), opportunities and challenges are presented for improving the robustness and generalization of multi-modal 3D object detection in AD. Therefore, we propose **RoboFusion**, a robust framework that leverages VFMs like SAM to tackle out-of-distribution (OOD) noise scenarios. We first adapt the original SAM for AD scenarios named **SAM-AD**. To align SAM or **SAM-AD** with multi-modal methods, we then introduce **AD-FPN** for upsampling the image features extracted by SAM. We employ wavelet decomposition to denoise the depth-guided images for further noise reduction and weather interference. At last, we employ self-attention mechanisms to adaptively reweight the fused features, enhancing informative features while suppressing excess noise. In summary, RoboFusion significantly reduces noise by leveraging the generalization and robustness of VFMs, thereby enhancing the resilience of multi-modal 3D object detection. Consequently, RoboFusion achieves SOTA performance in noisy scenarios, as demonstrated by the KITTI-C and nuScenes-C benchmarks. Code is available at <https://github.com/adept-thu/RoboFusion>.

## 1 Introduction

Multi-modal 3D object detection plays a pivotal role in autonomous driving (AD) [Wang *et al.*, 2023a; Song *et al.*, 2024a]. Different modalities often provide complementary information. For instance, images contain richer semantic representations, yet lack depth information. In contrast, point clouds offer geometric and depth details, but they are sparse

and lack semantic information. Therefore, effectively leveraging the advantages of multi-model while mitigating their limitations contributes to enhancing the robustness and accuracy of perception systems [Song *et al.*, 2023].

With the emergence of AD datasets [Geiger *et al.*, 2012; Caesar *et al.*, 2020; Zhang *et al.*, 2023c], state-of-the-art (SOTA) methods [Liu *et al.*, 2023; Bai *et al.*, 2022; Chen *et al.*, 2022; Huang *et al.*, 2020; Li *et al.*, 2023; Song *et al.*, 2024b] on ‘clean’ datasets [Geiger *et al.*, 2012; Caesar *et al.*, 2020] have achieved record-breaking performance. However, they overlook the exploration of robustness and generalization in out-of-distribution (OOD) scenarios [Dong *et al.*, 2023]. For example, the KITTI dataset [Geiger *et al.*, 2012] lacks severe weather conditions. When SOTA methods [Chen *et al.*, 2022; Li *et al.*, 2023; Liu *et al.*, 2023] learn from these sunny weather datasets, can they truly generalize and maintain robustness in severe weather conditions like snow and fog?

The answer is ‘No’, as shown in Fig. 1 and verified in Table 3. People often utilize domain adaptation (DA) techniques to address these challenges [Wang *et al.*, 2023b; Tsai *et al.*, 2023; Peng *et al.*, 2023; Hu *et al.*, 2023]. Although DA techniques improve the robustness of 3D object detection and reduce the need for annotated data, they have some profound drawbacks, including domain shift limitations, label shift issues, and overfitting risks [Oza *et al.*, 2023]. For instance, DA techniques may be constrained if the differences between two domains are significant, leading to performance degradation on the target domain.

Recently, both Natural Language Processing (NLP) and Computer Vision (CV) have witnessed the appearance and the power of a series of foundation models [Kirillov *et al.*, 2023; OpenAI, 2023; Zhao *et al.*, 2023; Zhang *et al.*, 2023a], resulting in the emergence of new paradigms in deep learning. For example, a series of novel visual foundation models (VFMs) [Kirillov *et al.*, 2023; Zhao *et al.*, 2023; Zhang *et al.*, 2023a] have been developed. Thanks to their extensive training on huge datasets, these models exhibit powerful generalization capabilities. These developments have inspired new ideas, leveraging the robustness and generalization abilities of VFMs to achieve generalization in OOD noisy scenarios, much like how adults generalize knowledge when encountering new situations, without relying on DA techniques [Wang *et al.*, 2023b; Tsai *et al.*, 2023].

\*Corresponding author

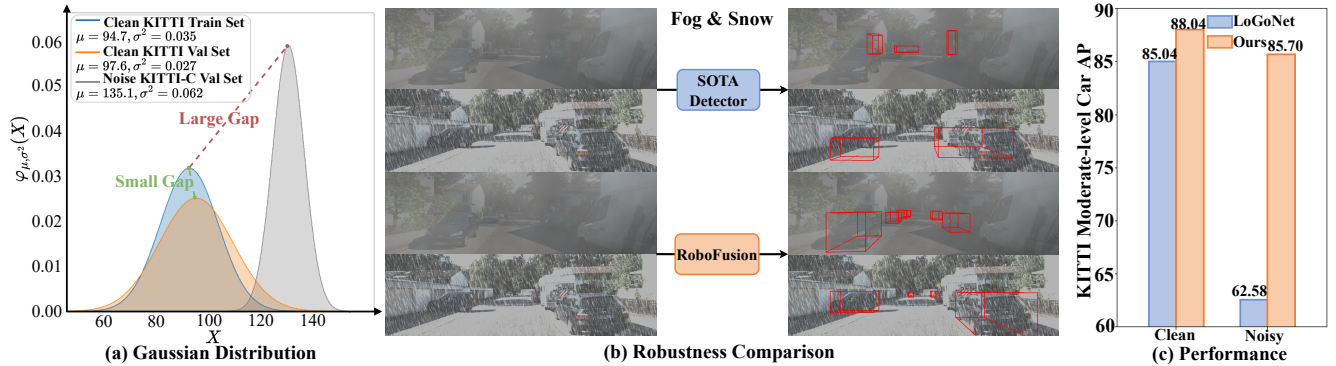


Figure 1: **(a)** We employ Gaussian distributions to represent the distributional disparities among the datasets. Indeed, there exists a large gap in data distribution between an OOD noise validation set and a clean validation set. Where the X-axis represents the set of mean pixel values in a dataset,  $X = \{x_i | i = 1, 2, \dots, N\}$ , with  $x_i = \frac{1}{H \times W \times 3} \sum_{i=1}^H \sum_{j=1}^W \sum_{k=1}^3 (I_{ijk})$ , where  $N$  is the number of the dataset,  $H$  is the height,  $W$  is the width, and  $I_{ijk}$  denotes the pixel values for each image. **(b)** Visual foundation models (VFMs) like SAM [Kirillov *et al.*, 2023], show robust performance in many noisy scenarios. Yet, the current methods are not robust enough to predict 3D tasks for autonomous driving perception. **(c)** To this end, we propose a robust framework, RoboFusion, which employs VFMs at the SOTA multi-modal 3D object detection. Empirical results reveal that our method surpasses the Top-performing LoGoNet [Li *et al.*, 2023] on the KITTI Leaderboard by a margin of 23.12% mAP (Weather) on KITTI-C [Dong *et al.*, 2023] noisy scenarios. Notably, our RoboFusion shows better performance with LoGoNet [Li *et al.*, 2023] in clean KITTI [Geiger *et al.*, 2012] dataset.

Inspired by the success of VFMs in CV tasks, in this work, we intend to use these models to tackle the challenges of multi-modal 3D object detectors in OOD noise scenarios. Therefore, we propose a robust framework, RoboFusion, which leverages VFMs like SAM to adapt a 3D multi-modal object detector from clean scenarios to OOD noise scenarios. In particular, the adaptation strategies for SAM are as follows. 1) We utilize features extracted from SAM rather than inference segmentation results. 2) We propose **SAM-AD**, which is a pre-trained SAM for AD scenarios. 3) We introduce a novel **AD-FPN** to address the issue of feature upsampling for aligning VFMs with multi-modal 3D object detector. 4) To further reduce noise interference and retain essential signal features, we design a **Depth-Guided Wavelet Attention (DGWA)** module that effectively attenuates both high-frequency and low-frequency noises. 5) After fusing point cloud features and image features, we propose **Adaptive Fusion** to further enhance feature robustness and noise resistance through self-attention to re-weight the fused features adaptively. We validate RoboFusion’s robustness against OOD noise scenarios in KITTI-C and nuScenes-C datasets [Dong *et al.*, 2023], achieving SOTA performance amid noise, as shown in Fig. 1.

## 2 Related Work

### 2.1 Multi-Modal 3D Object Detection

Currently, multi-modal 3D object detection has received considerable attention on popular datasets [Geiger *et al.*, 2012; Caesar *et al.*, 2020]. BEVFusion [Liu *et al.*, 2023] fuse multi-modal representations in a unified 3D or BEV space. TransFusion [Bai *et al.*, 2022] builds a two-stage pipeline where proposals are generated based on LiDAR features and further refined using query image features. DeepInteraction [Yang *et al.*, 2022] and SparseFusion [Xie *et al.*, 2023] further optimize the camera branch on top of TransFusion. Previous

methods are highly optimized to achieve the best performance on clean datasets. However, they ignore common factors in the real world (*e.g.*, bad weather and sensor noise). In this work, we consider a real-world robustness perspective and design a robust multi-modal 3D perception framework, RoboFusion.

### 2.2 Visual Foundation Models for 3D Object Detection

Motivated by the success of Large Language Models (LLMs) [OpenAI, 2023], VFMs start to be explored in CV community. SAM [Kirillov *et al.*, 2023] leverages ViT [Dosovitskiy *et al.*, 2020] to train on the huge SA-1B dataset, containing 11 million samples, which enables SAM to be generalized to many scenes. Currently, there have been a few research endeavors aiming at integrating 3D object detectors with SAM. For instance, SAM3D [Zhang *et al.*, 2023b], as a LiDAR-only method, solely transforms LiDAR’s 3D perspective into a BEV (Bird’s Eye View) 2D space to harness the generalization capabilities of SAM, yielding sub-optimal performance on ‘clean’ datasets. Another in progress work, 3D-Box-Segment-Anything<sup>1</sup>, tries to utilize SAM for 3D object detection. This indicates the highly attention of SAM like foundation models in 3D scenes in the literature. Our RoboFusion, as a multi-modal method, gives clear strategies to leverage the generalization capabilities of VFMs to address the OOD noise challenges inherent in existing 3D multi-modal object detection methods.

## 3 RoboFusion

In this section, we present RoboFusion, a framework that harnesses the robustness and generalization capabilities of VFMs

<sup>1</sup><https://github.com/dvlab-research/3D-Box-Segment-Anything>

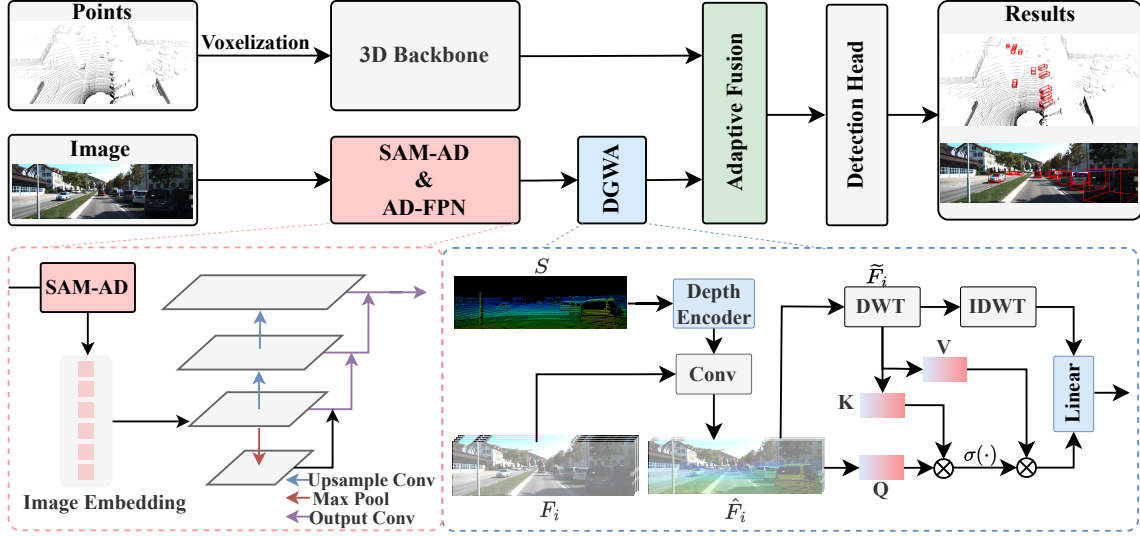


Figure 2: The framework of RoboFusion. The LiDAR branch follows the baselines [Chen *et al.*, 2022; Bai *et al.*, 2022] to generate LiDAR features. In the camera branch, first, we extract robust image features using a highly optimized SAM-AD and acquire multi-scale features using AD-FPN. Second, the sparse depth map  $S$  is generated by the raw points and fed into a depth encoder to obtain depth features and fused with multi-scale image features  $F_i$  to obtain depth-guided image features  $\hat{F}_i$ . Then wave attention is used to remove the mutation noise. Finally, adaptive Fusion integrates point cloud features with robust image features with depth information via self-attention mechanism.

such as SAM [Kirillov *et al.*, 2023] for multi-modal 3D object detection. The overall architecture is depicted in Fig. 2 and comprises the following components: 1) **SAM-AD & AD-FPN** module which obtains robust multi-scale image features, 2) **Depth-Guided Wavelet Attention (DGWA)** module which employs wavelet decomposition to denoise depth-guided image features, 3) **Adaptive Fusion** module which adaptively fuses point cloud features with image features.

### 3.1 SAM-AD & AD-FPN

**Preliminaries.** SAM [Kirillov *et al.*, 2023], a VFM, achieves generalization across diverse scenes due to its extensive training on the large-scale SA-1B dataset—with over 11 million samples and 1 billion high-quality masks. Currently, SAM family [Kirillov *et al.*, 2023; Zhao *et al.*, 2023; Zhang *et al.*, 2023a] primarily support 2D tasks. However, directly extending VFMs like SAM to 3D tasks presents a gap. To address this, we combine SAM with multi-modal 3D models, merging 2D robust feature representations with 3D point cloud features to achieve robust fused features.

**SAM-AD.** To further adapt SAM with AD (autonomous driving) scenarios, we perform pre-training on SAM to obtain SAM-AD. Specifically, we curate an extensive collection of image samples from well-established datasets (*i.e.*, KITTI [Geiger *et al.*, 2012] and nuScenes [Caesar *et al.*, 2020]), forming the foundational AD dataset. Following DMAE [Wu *et al.*, 2023], we perform pre-training on SAM to obtain SAM-AD in AD scenarios, as shown in Fig. 3. We denote  $x$  as a clean image from the AD dataset (*i.e.* KITTI [Geiger *et al.*, 2012] and nuScenes [Caesar *et al.*, 2020]) and  $\eta$  as a set of noise images generated by [Dong *et al.*, 2023] based on  $x$ . And the noise type and the severity are randomly chosen from the four weather (*i.e.*, rain, snow, fog, and strong sunlight)

and the five severities from 1 to 5, respectively. We employ the image encoder of SAM [Kirillov *et al.*, 2023], MobileSAM [Zhang *et al.*, 2023a] as our encoder while the decoder and the reconstruction loss are the same as DMAE [Wu *et al.*, 2023]. For FastSAM [Zhao *et al.*, 2023], we adopt YOLOv8<sup>2</sup> to pre-train FastSAM on the AD dataset. To avoid overfitting, we use random resizing and cropping as data augmentation. We also set the mask ratio as 0.75 and have trained 400 epochs on 8 NVIDIA A100 GPUs.

**AD-FPN.** As a promptable segmentation model, SAM has three components: image encoder, prompt encoder and mask decoder. Generally, the image encoder can provide high-quality and highly robust image embedding for downstream models, while the mask decoder is only designed to provide decoding services for semantic segmentation. Furthermore, what we require are robust image features rather than the processing of prompting information by the prompt encoder. Therefore, we employ SAM’s image encoder to extract robust image features. However, SAM utilizes the ViT series [Dosovitskiy *et al.*, 2020] as its image encoder, which excludes multi-scale features and provides only high-dimensional low-resolution features. To generate the multi-scale features required for object detection, inspired by [Li *et al.*, 2022a], we design an AD-FPN that offers ViT-based multi-scale features. Specifically, leveraging height-dimensional image embedding with stride 16 (scale=1/16) provided by SAM, we produce a series of multi-scale features  $F_{ms}$  with stride of  $\{32, 16, 8, 4\}$ . Sequentially, we acquire multi-scale feature  $F_i \in \mathbb{R}^{\frac{H}{4} \times \frac{W}{4} \times C_i}$  by integrate  $F_{ms}$  in a bottom-up manner similar to FPN [Lin *et al.*, 2017].

<sup>2</sup><https://github.com/ultralytics/ultralytics>

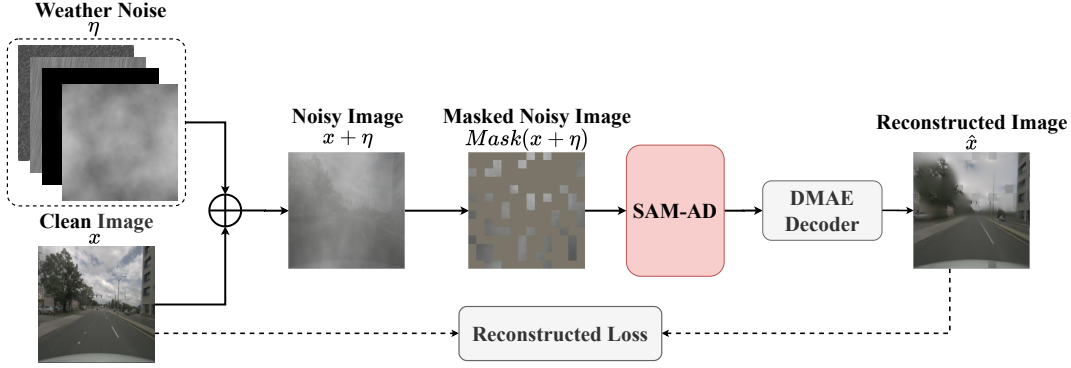


Figure 3: An illustration of the pre-training framework. We corrupt a clean image  $x$  by  $\eta$  which contains multiple weather noises and then randomly masking several patches on a noisy image  $x + \eta$  to obtain a masked noisy image  $Mask(x + \eta)$ . The SAM-AD and DMAE decoder are trained to reconstruct the clean image  $\hat{x}$  from  $Mask(x + \eta)$ .

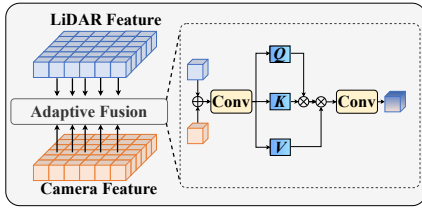


Figure 4: The architecture of **Adaptive Fusion**, which involves adaptively re-weighting the fused features using self-attention.

### 3.2 Depth-Guided Wavelet Attention

Although SAM-AD or SAM has the capability to extract robust image features, the gap between 2D and 3D domains still persists and cameras lacking geometric information in a corrupted environment often amplify noise and give rise to negative migration issues. To mitigate this problem, we propose the Depth-Guided Wavelet Attention (DGWA) module, which can be split into two steps. 1) A depth-guided network is designed, that adds geometry prior to image features by combining image features and depth features from a point cloud. 2) The features of an image are decomposed into four wavelet subbands using the Haar wavelet transform [Liu *et al.*, 2020a], then attention mechanism allows to denoise informative features in the subbands.

Formally, given image features  $F_i \in \mathbb{R}^{\frac{H}{4} \times \frac{W}{4} \times C_i}$  and raw points  $P \in \mathbb{R}^{N, C_p}$  as input. We project  $P$  onto the image plane to acquire a sparse depth map  $S \in \mathbb{R}^{H \times W \times 2}$ . Next, we feed  $S$  into the depth encoder  $DE(\cdot)$ , which consists of several convolution and max pooling blocks, to acquire depth features  $F_d \in \mathbb{R}^{\frac{H}{4} \times \frac{W}{4} \times C_i}$ . Afterward, we leverage convolution encode  $(F_i, F_d)$  to acquire depth-guided image features  $\hat{F}_i \in \mathbb{R}^{\frac{H}{4} \times \frac{W}{4} \times 16}$ , given by

$$\hat{F}_i = \text{Conv}(\text{Concat}(F_i, DE(S))). \quad (1)$$

Subsequently, we employ discrete wavelet transform (DWT), a reversible operator, to partition the input  $\hat{F}_i$  into four subbands. Specifically, we encode the rows and columns of the input separately into one low-frequency band  $\tilde{f}_i^{LL} \in \mathbb{R}^{\frac{H}{8} \times \frac{W}{8} \times 4}$  and three high-frequency bands

$(\tilde{f}_i^{LH}, \tilde{f}_i^{HL}, \tilde{f}_i^{HH}) \in \mathbb{R}^{\frac{H}{8} \times \frac{W}{8} \times 4}$ , with the low-filter  $f_L = (\frac{1}{\sqrt{2}}, \frac{1}{\sqrt{2}})$  and the high-filter  $f_H = (\frac{1}{\sqrt{2}}, -\frac{1}{\sqrt{2}})$ . In this state, the low-frequency band retains coarse-grained information while the high-frequency band retains fine-grained information. In other words, it is easier to capture the mutation signal, so as to filter the noise information. We concatenate the four-subband features along channel dimension to acquire wavelet features  $\tilde{F}_i = [\hat{f}_i^{LL}, \hat{f}_i^{LH}, \hat{f}_i^{HL}, \hat{f}_i^{HH}] \in \mathbb{R}^{\frac{H}{8} \times \frac{W}{8} \times 16}$ . Next, we perform wave-attention  $Att_\omega$  to query informative features in the wavelet features. Concretely, we employ  $\hat{F}_i$  as a Query and  $\tilde{F}_i$  as a Key/Value given by

$$F_{att} = Att_\omega(\hat{F}_i, \tilde{F}_i) = \sigma\left(\frac{\hat{F}_i W^q (\tilde{F}_i W^k)^T}{\sqrt{C_i}}\right) \tilde{F}_i W^v. \quad (2)$$

Finally, we leverage the IDWT (inverse DWT) to convert  $\tilde{F}_i$  back to  $\hat{F}_i$  and integrate this converted  $\hat{F}_i$  and  $F_{att}$  to obtain denoise features  $F_{out} \in \mathbb{R}^{\frac{H}{16} \times \frac{W}{16} \times 16}$  by

$$F_{out} = MLP(\text{Concat}(F_{att}, \hat{F}_i)), \quad (3)$$

where  $F_{out}$  preserves informative features and restrains redundant mutation noise in the frequency domain.

### 3.3 Adaptive Fusion

Following the incorporation of image depth features within the DGWA module, we propose the **Adaptive Fusion** technique to combine point cloud attributes with robust image features enriched with depth information. Specifically, different types of noise affect LiDAR and images to different degrees, which raises a corruption imbalance problem. Therefore, considering the distinct influences of various noises on LiDAR and camera, we employ self-attention to re-weight the fused features adaptively as shown in Fig. 4. The corruption degree of modality-specificity is dynamic, and self-attention mechanism allows adaptive re-weighting features to enhance informative features and suppress redundant noise.

## 4 Experiments

### 4.1 Datasets

We perform experiments on both the clean public benchmarks (KITTI [Geiger *et al.*, 2012] and nuScenes [Caesar *et al.*,



Table 1: Comparison with SOTA methods on **KITTI validation and test** sets for car class with AP of  $R_{40}$ .

Method	AP <sub>3D</sub> (%) (validation set)				AP <sub>3D</sub> (%) (test set)			
	mAP	Easy	Mod.	Hard	mAP	Easy	Mod.	Hard
Voxel R-CNN	86.84	92.38	85.29	82.86	83.19	90.90	81.62	77.06
VFF	86.91	92.31	85.51	82.92	83.62	89.50	82.09	79.29
CAT-Det	83.58	90.12	81.46	79.15	82.62	89.87	81.32	76.68
LoGoNet	87.13	92.04	85.04	84.31	85.87	91.80	<b>85.06</b>	<b>80.74</b>
Focals Conv-F	-	-	-	-	83.47	90.55	82.28	77.59
Baseline*	86.75	92.05	85.51	82.70	-	-	-	-
RoboFusion-L	<b>88.87</b>	<b>93.30</b>	<b>88.04</b>	<b>85.27</b>	<b>85.58</b>	91.75	84.08	80.71
RoboFusion-B	88.45	93.22	87.87	84.27	85.32	<b>91.98</b>	83.76	80.23
RoboFusion-T	88.08	93.28	87.60	83.36	85.09	91.68	83.70	79.89

\* denotes our reproduced results based on the officially released codes.

2020]) and the noisy public benchmarks (KITTI-C[Dong *et al.*, 2023] and nuScenes-C [Dong *et al.*, 2023]).

### KITTI

The KITTI dataset provides synchronized LiDAR point clouds and front-view camera images, consists of 3,712 training samples, 3,769 validation samples and 7,518 test samples. The standard evaluation metric for object detection is the mean Average Precision (mAP), computed using recall at 40 positions (R40).

### nuScenes

The nuScenes dataset is a large-scale 3D detection benchmark consisting of 700 training scenes, 150 validation scenes, and 150 testing scenes. The data are collected using six multi-view cameras and a 32-channel LiDAR sensor. It includes 360-degree object annotations for 10 object classes. To evaluate the detection performance, the primary metrics used are the mean Average Precision (mAP) and the nuScenes detection score (NDS).

### KITTI-C and nuScenes-C

In terms of data robustness, [Dong *et al.*, 2023] has designed 27 types of common corruptions for both LiDAR and camera, with the aim of benchmarking the corruption robustness of existing 3D object detectors. [Dong *et al.*, 2023] has established corruption robustness benchmarks<sup>3</sup>, including **KITTI-C** and **nuScenes-C**, by synthesizing corruptions on public datasets. Specifically, we utilize **KITTI-C** and **nuScenes-C** in our work. It is worth noting that Ref. [Dong *et al.*, 2023] has only added noise to the validation dataset and kept the train and test datasets clear.

## 4.2 Experimental Settings

### Network Architecture.

Our RoboFusion consists of three variants: RoboFusion-L, RoboFusion-B, and RoboFusion-T, which utilize the models SAM-B [Kirillov *et al.*, 2023], FastSAM [Zhao *et al.*, 2023], and MobileSAM [Zhang *et al.*, 2023a], respectively. It is noteworthy that due to the convolutional operations of FastSAM in RoboFusion-B which is capable of generating multi-scale features, the AD-FPN module is not employed.

<sup>3</sup>[https://github.com/thu-ml/3D\\_Corruptions\\_AD](https://github.com/thu-ml/3D_Corruptions_AD)

Table 2: Comparison with SOTA methods on **nuScenes validation and test** sets.

Method	LiDAR	Camera	validation set		test set	
			NDS	mAP	NDS	mAP
FUTR3D	VoxelNet	ResNet-101	68.3	64.5	-	-
BEVFusion-mit	VoxelNet	Swin-T	71.4	68.5	72.9	70.2
DeepInteraction	VoxelNet	ResNet-50	72.6	69.9	73.4	70.8
CMT	VoxelNet	ResNet-50	72.9	70.3	74.1	72.0
SparseFusion	VoxelNet	ResNet-50	72.8	70.4	73.8	72.0
TransFusion	VoxelNet	ResNet-50	71.3	67.5	71.6	68.9
Baseline*	VoxelNet	ResNet-50	70.8	67.3	-	-
RoboFusion-L	VoxelNet	SAM	72.1	69.9	72.0	69.9
RoboFusion-B	VoxelNet	FastSAM	71.9	69.4	71.8	69.4
RoboFusion-T	VoxelNet	MobileSAM	71.3	69.1	71.5	69.1

\* denotes our reproduced results based on the officially released codes.

Table 3: Comparison with SOTA methods on **KITTI-C validation** set. The results are evaluated based on the car class with AP of  $R_{40}$  at moderate difficulty. ‘S.L.’, ‘D.’, ‘C.O.’, and ‘C.T.’ denotes Strong Sunlight, Density, Cutout, and Crosstalk, respectively.

Rai	Clean	Weather					Sensor		
		mAP	Snow	Rain	Fog	S.L.	D.	C.O.	C.T.
SECOND <sup>†</sup>	81.59	64.33	52.34	52.55	74.10	78.32	80.18	73.59	80.24
PointPillars <sup>†</sup>	78.41	49.80	36.47	36.18	64.28	62.28	76.49	70.28	70.85
PointRCNN <sup>†</sup>	80.57	59.14	50.36	51.27	72.14	62.78	80.35	73.94	71.53
PV-RCNN <sup>†</sup>	84.39	65.83	52.35	51.58	79.47	79.91	82.79	76.09	82.34
SMOKE <sup>†</sup>	7.09	4.51	2.47	3.94	5.63	6.00	-	-	-
ImVoxelNet <sup>†</sup>	11.49	3.22	0.22	1.24	1.34	10.08	-	-	-
EPNet <sup>†</sup>	82.72	46.21	34.58	36.27	44.35	69.65	82.09	76.10	82.10
Focals Conv-F <sup>†</sup>	85.88	50.40	34.77	41.30	44.55	80.97	84.95	78.06	<b>85.82</b>
LoGoNet*	85.04	62.58	51.45	55.80	67.53	75.54	83.68	77.17	82.00
RoboFusion-L	<b>88.04</b>	<b>85.70</b>	<b>85.29</b>	<b>86.48</b>	<b>85.53</b>	<b>85.50</b>	<b>85.71</b>	<b>83.17</b>	84.12
RoboFusion-B	87.87	84.70	84.11	85.54	84.00	85.15	84.34	81.30	82.45
RoboFusion-T	87.60	84.60	84.67	84.79	84.17	84.75	84.11	81.21	83.07

<sup>†</sup>: Results from Ref. [Dong *et al.*, 2023].

\* denotes re-implement result.

Since KITTI and nuScenes are distinct datasets with varying evaluation metrics and characteristics, we provide a detailed description of our RoboFusion settings for each dataset.

**RoboFusion in KITTI and KITTI-C.** We validate our RoboFusion on the KITTI dataset using Focals Conv [Chen *et al.*, 2022] as the baseline. The input voxel size is set to (0.05m, 0.05m, 0.1m), with anchor sizes for cars at [3.9, 1.6, 1.56] and anchor rotations at [0, 1.57]. We adopt the same data augmentation solution as Focals Conv-F.

**RoboFusion with nuScenes and nuScenes-C.** We validate our RoboFusion on the nuScenes dataset using TransFusion [Bai *et al.*, 2022] as the baseline. The detection range for the X and Y axis is set at [-54m, 54m] and [-5m, 3m] for the Z axis. The input voxel size is set at (0.075m, 0.075m, 0.2m), and the maximum number of point clouds contained in each voxel is set to 10. It is noteworthy that the Adaptive Fusion module is applied exclusively to Focals Conv rather than TransFusion, while TransFusion uses its own fusion module.

### Training and Testing Details.

Our RoboFusion is meticulously trained from scratch using the Adam optimizer and incorporates several foundation models as image encoders including SAM, FastSAM and

Table 4: Comparison with SOTA methods on **nuScenes-C validation** set with mAP. ‘S.L.’, ‘D.’, ‘C.O.’, and ‘C.T.’ denotes Strong Sunlight, Density, Cutout, and Crosstalk, respectively.

Method	Clean	Weather					Sensor			
		mAP	Snow	Rain	Fog	S.L.	D.	C.O.	C.T.	
PointPillars <sup>†</sup>	27.69	25.87	27.57	27.71	24.49	23.71	27.27	24.14	25.92	
SSN <sup>†</sup>	46.65	43.70	46.38	46.50	41.64	40.28	46.14	40.95	44.08	
CenterPoint <sup>†</sup>	59.28	52.49	55.90	56.08	43.78	54.20	58.60	56.28	56.64	
FCOS3D <sup>†</sup>	23.86	11.44	2.01	13.00	13.53	17.20	-	-	-	
PGD <sup>†</sup>	23.19	12.85	2.30	13.51	12.83	22.77	-	-	-	
DETR3D <sup>†</sup>	34.71	22.00	5.08	20.39	27.89	34.66	-	-	-	
BEVFormer <sup>†</sup>	41.65	26.29	5.73	24.97	32.76	41.68	-	-	-	
FUTR3D <sup>†</sup>	64.17	55.50	52.73	58.40	53.19	57.70	63.72	62.25	62.66	
TransFusion <sup>†</sup>	66.38	58.87	63.30	63.35	53.67	55.14	65.77	63.66	64.67	
BEVFusion <sup>†</sup>	68.45	61.87	62.84	66.13	54.10	64.42	67.79	66.18	67.32	
DeepInteraction <sup>*</sup>	69.90	62.14	62.36	66.48	54.79	64.93	68.15	66.23	68.12	
CMT <sup>*</sup>	<b>70.28</b>	63.46	62.56	61.44	66.26	63.59	<b>69.65</b>	68.70	68.26	
RoboFusion-L	69.91	<b>67.24</b>	<b>67.12</b>	<b>67.58</b>	<b>67.01</b>	<b>67.24</b>	69.48	<b>69.18</b>	<b>68.68</b>	
RoboFusion-B	69.40	66.33	66.07	67.01	65.54	66.71	69.02	69.01	68.04	
RoboFusion-T	69.09	65.82	65.96	66.45	64.34	66.54	68.58	68.20	68.17	

<sup>†</sup>: Results from Ref. [Dong *et al.*, 2023].

<sup>\*</sup> denotes re-implement result.

Table 5: Performance of different VFMs on RoboFusion. ‘RCE’ denotes Relative Corruption Error [Dong *et al.*, 2023]. ‘mAP (Weather)’ denotes the average value across four types of weather corruptions, Snow, Rain, Fog, and Strong Sunlight.

Method	Model Size	FPS (A100)	mAP (Weather)	mAP (Clean)	RCE (%)
RoboFusion-L	97.54M	3.1	67.24	69.91	0.04
RoboFusion-B	81.01M	3.5	66.33	69.40	0.04
RoboFusion-T	13.94M	6.0	65.82	69.09	0.05
DeepInteraction	57.82M	4.9	62.14	69.90	0.10
TransFusion	36.96M	6.2	58.37	66.38	0.12

MobileSAM. To enable effective training on the KITTI and nuScenes datasets, we utilize 8 NVIDIA A100 GPUs for network training. Additionally, the runtime is evaluated on an NVIDIA A100 GPU. Specifically, for KITTI, our RoboFusion based on Focals Conv[Chen *et al.*, 2022] involves training for 80 epochs. For nuScenes, our RoboFusion based on TransFusion [Bai *et al.*, 2022] has 20 epochs of training. During the model inference stage, we employ a non-maximal suppression (NMS) operation in the Region Proposal Network (RPN) with an IoU threshold of 0.7. We select the top 100 region proposals to serve as inputs for the detection head. After refinement, we apply NMS again with an IoU threshold of 0.1 to eliminate redundant predictions. For additional details regarding our method, please refer to OpenPCDet<sup>4</sup>.

### 4.3 Comparing with state-of-the-art

We conduct evaluations on the clean datasets KITTI and nuScenes, as well as the noisy datasets KITTI-C and nuScenes-C. While SOTA methods are primarily focused on achieving high accuracy, we place greater emphasis on the robustness and generalization of the methods. These factors are crucial for the practical deployment of 3D object detection in AD scenarios, making the evaluation on the noisy datasets more important in our perspective.

<sup>4</sup><https://github.com/open-mmlab/OpenPCDet>

### Results on the clean benchmark.

As shown in Table 1, we compare our RoboFusion with SOTA methods, including Voxel R-CNN [Deng *et al.*, 2021], VFF [Li *et al.*, 2022b], CAT-Det [Zhang *et al.*, 2022], Focals Conv-F [Chen *et al.*, 2022], and LoGoNet [Li *et al.*, 2023] on the KITTI validation and test sets. As shown in Table 2, we also compare our RoboFusion with SOTA methods, including FUTR3D [Chen *et al.*, 2023], TransFusion [Bai *et al.*, 2022], BEVFusion [Liu *et al.*, 2023], DeepInteraction [Yang *et al.*, 2022], CMT [Yan *et al.*, 2023] and SparseFusion [Xie *et al.*, 2023], on the nuScenes test and validation sets. Our RoboFusion has achieved SOTA performance on the clean benchmarks (KITTI and nuScenes).

### Results on the noisy benchmark.

In the real-world AD scenarios, the distribution of data often differs from that of training or testing data, as shown in Fig. 1 (a). Specifically, Ref. [Dong *et al.*, 2023] provides a novel noisy benchmark that includes KITTI-C and nuScenes-C, which we primarily use to evaluate the weather and sensor noise corruptions, including rain, snow, fog, and strong sunlight, density, cutout, and so on. In addition, comparisons of our RoboFusion with SOTA methods in other settings are presented in the Appendix<sup>5</sup>.

As shown in Table 3, SOTA methods, including SECOND [Yan *et al.*, 2018], PointPillars [Lang *et al.*, 2019], PointR-CNN [Shi *et al.*, 2019], PV-RCNN [Shi *et al.*, 2020], SMOKE [Liu *et al.*, 2020b], ImVoxelNet [Rukhovich *et al.*, 2022], EpNet [Huang *et al.*, 2020], Focals Conv-F [Chen *et al.*, 2022], and LoGoNet [Li *et al.*, 2023], experience a significant decrease in performance on the noisy scenarios, particularly for weather conditions such as snow and rain. It can be attributed to the fact that the ‘clean’ KITTI dataset does not include examples in snowy or rainy weather. On the other hand, VFMs like SAM-AD have been trained on a diverse range of data and exhibit robustness and generalization to OOD scenarios, leading to higher performance on our RoboFusion metric. Furthermore, multi-modal methods like LoGoNet, and Focals Conv-F demonstrate better robustness and generalization in sensor noise scenarios, while LiDAR-only methods like PV-RCNN [Shi *et al.*, 2020] are more robust in weather noise scenarios. This observation motivates our research on adaptive fusion schemes for point cloud and image features. Overall, in the KITTI-C [Dong *et al.*, 2023] dataset, our RoboFusion’s performance is nearly on par with the clean scene, indicating high level of robustness and generalization.

As shown in Table 4, SOTA methods including PointPillars [Lang *et al.*, 2019], SSN [Zhu *et al.*, 2020], CenterPoint [Yin *et al.*, 2021], FCOS3D [Wang *et al.*, 2021], PGD [Wang *et al.*, 2022a], DETR3D [Wang *et al.*, 2022b], BEVFormer [Li *et al.*, 2022c], FUTR3D [Chen *et al.*, 2023], TransFusion [Bai *et al.*, 2022], BEVFusion [Liu *et al.*, 2023], DeepInteraction [Yang *et al.*, 2022] and CMT [Yan *et al.*, 2023] in nuScenes-C show relatively higher robustness than in KITTI-C when faced with weather noise. However, BEVFusion performs well in the presence of snow, rain, and strong sunlight noise but experiences a significant performance drop in foggy scenarios. In contrast, our method exhibits strong robustness and

<sup>5</sup><https://arxiv.org/abs/2401.03907>

Table 6: Impacts of different SAM usages on **KITTI** and **KITTI-C validation** sets for car class with AP of  $R_{40}$ . ‘S.L.’ denotes Strong Sunlight.

Solution	AP <sub>3D</sub> (%)				AP <sub>Weather</sub> (%)			
	mAP	Easy	Mod.	Hard	Snow	Rain	Fog	S.L.
Offline	80.41	88.76	77.38	75.11	-	-	-	-
No optim	86.45	91.86	84.80	82.71	45.11	47.77	63.10	79.21
Optim	88.00	92.41	86.77	84.81	57.43	54.27	68.81	82.07

Table 7: Influence of pre-training on SAM at **KITTI-C validation** set for car class with AP of  $R_{40}$  at moderate difficulty. ‘S.L.’, ‘D.’, ‘C.O.’, and ‘C.T.’ denotes Strong Sunlight, Density, Cutout, and Crosstalk, respectively.

VFM	Weather				Sensor		
	Snow	Rain	Fog	S.L.	D.	C.O.	C.T.
SAM	57.43	54.27	68.81	82.07	84.21	83.04	84.06
SAM-AD	80.68	81.68	81.67	83.48	84.71	84.17	84.12

generalization in both weather and sensor noise scenarios in nuScenes-C.

#### 4.4 Ablation Study

##### Performance of Different VFMs on RoboFusion.

In order to analyze the noise robustness and FPS performance of different-sized VFMs, SAM, FastSAM and MobileSAM, we conduct comparative experiments of RoboFusion-L, RoboFusion-B and RoboFusion-T with SOTA methods, DeepInteraction [Yang *et al.*, 2022] and TransFusion [Bai *et al.*, 2022], on the nuScenes-C [Dong *et al.*, 2023] validation set, as shown in Table 5. Specifically, our RoboFusion exhibits remarkable robustness to weather noise scenarios. Furthermore, our RoboFusion-T has a similar FPS to TransFusion [Bai *et al.*, 2022]. Overall, we have presented a viable application of SAM in 3D object detection tasks.

##### Impacts of Different SAM usages.

As shown in Table 6, our RoboFusion-L is experimented upon. Specifically, the first row is the offline usage, which involves loading pre-saved image features during training. It implies that certain online data augmentation cannot be utilized. The second (No optim) and the third (Optim) rows are online usages, where the former omits fine-tuning and keeps the model parameters fixed, the latter follows fine-tuning and updating. Therefore, offline usage perform worse than online usages. Additionally, fine-tuning the weights of SAM has demonstrated superior performance, resulting in a performance improvement in the presence of snow, rain, and fog noise scenarios.

##### Influence of Pre-training on SAM.

As shown in Table 7, to investigate the scientific value of pre-trained VFMs like SAM, FastSAM, and MobileSAM in AD scenarios, we conduct our RoboFusion-L with SAM evaluation on SAM and SAM-AD. Through pre-training, SAM-AD has gained a better understanding of AD scenarios than the original SAM. The pre-training strategy effectively improves the performance of our RoboFusion, demonstrating a significant improvement in the snow, rain, and fog noise scenarios.

Table 8: Roles of SAM3DFusion-L modules on **KITTI-C validation** set for car class with AP of  $R_{40}$  at moderate difficulty. ‘A.F.’ denotes **Adaptive Fusion** module. ‘S.L.’ denotes strong sunlight.

Method	SAM-AD	AD-FPN	DGWA	A.F.	Snow	Rain	Fog	S.L.	FPS(A100)
a)					34.77	41.30	44.55	80.97	10.8
b)	✓				80.68	81.68	81.67	83.48	4.0
c)	✓	✓			82.32	83.60	82.39	83.98	3.6
d)	✓	✓	✓		83.99	85.63	84.01	84.81	3.4
e)	✓	✓	✓	✓	85.29	86.48	85.53	85.50	3.1

##### Roles of Different Modules in RoboFusion.

As shown in Table 8, we present ablation experiments for different modules of our RoboFusion-L, built upon SAM-AD, including AD-FPN, DGWA, and Adaptive Fusion. Leveraging the strong capabilities of SAM-AD in AD scenarios, SAM-AD has a significant improvement from baseline Focals Conv [Chen *et al.*, 2022] (34.77%, 41.30%, 44.55%, 80.97%) to (80.68%, 81.68%, 81.67%, 83.48%). Subsequently, AD-FPN, DGWA, and Adaptive Fusion achieve even higher performance on the foundation of SAM-AD. This further highlights the substantial contributions of diverse modules within our RoboFusion framework in addressing OOD noise scenarios in AD.

## 5 Conclusions

In this work, we propose a robust framework RoboFusion to enhance the robustness and generalization of multi-modal 3D object detectors using VFMs like SAM, FastSAM, and MobileSAM. Specifically, we pre-train SAM for AD scenarios, yielding SAM-AD. To align SAM or SAM-AD with multi-modal 3D object detectors, we introduce AD-FPN for feature upsampling. To further mitigate noise and weather interference, we apply wavelet decomposition for depth-guided image denoising. Subsequently, we utilize self-attention mechanisms to adaptively reweight fused features, enhancing informative attributes and suppressing excess noises. Extensive experiments demonstrate that our RoboFusion effectively integrates VFMs to boost feature robustness and address OOD noise challenges. We anticipate this work to lay a strong foundation for future research on building robust and dependable foundation AD models.

**Limitation and Future Work.** First, RoboFusion has a heavy reliance on the representation capability of VFMs. This raises the baseline models’ generalization ability, but increases their complexities. Second, the inference speed of RoboFusion-L and RoboFusion-B is relatively slow due to the limitations of SAM and FastSAM. However, the inference speed of RoboFusion-T is competitive with some SOTA methods (e.g. TransFusion) without VFMs. In the future, for improving the real-time application ability of VFMs, we will attempt to incorporate SAM only in the training phase to guide a fast-speed student model, meanwhile explore more noise scenarios.

## Acknowledgments

This work was supported by the Fundamental Research Funds for the Central Universities (2023YJS019), the National Key R&D Program of China (2018AAA0100302).



## A Appendix

### A.1 Broader Impacts

Our work aims to develop a robust framework to address out-of-distribution (OOD) noise scenarios in autonomous driving (AD). To the best of our knowledge, RoboFusion is the first method that leverages the generalization capabilities of visual foundation models (VFM) like SAM [Kirillov *et al.*, 2023], FastSAM [Zhao *et al.*, 2023], and MobileSAM [Zhang *et al.*, 2023a] for multi-modal 3D object detection. Although existing multi-model 3D object detection methods achieve the state-of-the-art (SOTA) performance of ‘clean’ datasets, they overlook the robustness of real-world scenarios [Song *et al.*, 2024a]. Therefore, we believe it is valuable to combine VFMs and multi-modal 3D object detection to mitigate the impact of OOD noise scenarios.

### A.2 More Results

#### Specific classes AP on the nuScenes-C validation set.

As shown in Table 9, we present a comparison of Specific classes AP between TransFusion and our RoboFusion-L on the nuScenes-C validation set, encompassing scenarios with snow, rain, fog, and strong sunlight noise. It is evident from the results that RoboFusion-L exhibits superior performance.

Table 9: Comparison with TransFusion on nuScenes validation ‘Snow, Rain, Fog, and Strong Sunlight’ noisy scenarios. ‘T.F.’, ‘R.F.’, ‘S.L.’, ‘C.V.’, ‘Motor.’, ‘Ped.’, and ‘T.C.’ are short for **TransFusion**, **RoboFusion-L**, Strong Sunlight, construction vehicle, motorcycle, pedestrian, and traffic cone, respectively.

		mAP	Car	Truck	C.V.	Bus	Trailer	Barrier	Motor.	Bike	Ped.	T.C.
Snow	T.F.	63.30	84.55	58.41	25.50	62.31	56.00	70.19	69.98	43.69	84.24	78.16
	R.F.	67.12	87.21	60.88	29.47	67.45	58.99	75.12	71.45	48.28	86.23	86.12
		+3.82				+5.14		+4.93		+4.59		+7.96
Rain	T.F.	63.35	85.37	56.87	25.12	64.65	55.10	71.99	68.21	44.13	83.87	78.14
	R.F.	67.58	86.79	60.44	30.21	65.41	58.12	75.47	71.39	50.87	88.91	88.23
		+4.23			+5.09					+6.74	+5.04	+10.09
Fog	T.F.	53.67	80.23	48.51	18.04	50.69	53.03	62.24	54.53	25.27	80.63	66.48
	R.F.	67.01	87.56	59.03	29.36	66.10	57.23	74.33	72.01	49.50	87.08	87.91
		+13.34		+10.52	+11.32	+15.41		+12.09	+17.48	+24.23		+21.43
S.L.	T.F.	55.14	81.99	48.07	19.78	51.09	52.57	63.68	55.09	26.98	82.68	69.49
	R.F.	67.24	87.67	57.74	31.00	64.29	58.94	75.23	70.23	50.82	88.70	87.82
		+12.10			+11.22	+13.20		+11.55	+15.14	+23.30		+18.33

#### Roles of Different Modules in RoboFusion.

To assess the roles of different modules in RoboFusion, we conduct an ablation study on the original SAM rather than SAM-AD, as shown in Table 10, where a) is the results of the baseline [Chen *et al.*, 2022], b)-e) shows the performance of our RoboFusion-L under different modules. According to Table 10, SAM and AD-FPN modules significantly improve the performance in OOD noisy scenarios. It is worth noticing that DGWA module significantly improves the performance, especially in snow noisy scenarios. By Table 10, the impact of fog noise on point clouds is relatively minor. But, using A.F. (Adaptive Fusion) module to dynamically aggregate point cloud features and image features exhibits significant enhancements in fog-noise scenarios.

Table 10: Roles of RoboFusion modules on **KITTI-C validation** set for car class with AP of  $R_{40}$  at moderate difficulty. ‘A.F.’ denotes **Adaptive Fusion** module. ‘S.L.’ denotes Strong Sunlight.

Method	SAM	AD-FPN	DGWA	A.F.	Snow	Rain	Fog	S.L.
a)					34.77	41.30	44.55	80.97
b)	✓				57.43	54.27	68.81	82.07
c)	✓	✓			59.81	56.59	69.68	83.20
d)	✓	✓	✓		66.45	58.11	70.53	84.01
e)	✓	✓	✓	✓	68.47	59.07	74.38	84.07

#### More Results on the KITTI-C validation set.

Besides the experimental results mentioned in the main text, we test our RoboFusion on KITTI-C and nuScenes-C [Dong *et al.*, 2023] to extend our work to a wider range of noise scenarios, including Gaussian, Uniform, Impulse, Moving Object, Motion Blur, Local Density, Local Cutout, Local Gaussian, Local Uniform, and Local Impulse, as shown in Tables 11, 12, and 13. From these Tables, compared with LiDAR-only methods including SECOND [Yan *et al.*, 2018], PointPillars [Lang *et al.*, 2019], PointRCNN [Shi *et al.*, 2019] and PV-RCNN [Shi *et al.*, 2020], Camera-Only methods including Smoke [Liu *et al.*, 2020b], ImVoxelNet [Rukhovich *et al.*, 2022], and multi-modal methods including EPNet [Huang *et al.*, 2020], Focals Conv [Chen *et al.*, 2022], and LoGoNet [Li *et al.*, 2023], our RoboFusion-L, RoboFusion-B, and RoboFusion-T consistently outperform across various noise scenarios and achieve the best overall performance. Overall, our RoboFusion demonstrates superior performance in weather-noisy (*i.e.* Snow, Rain, Fog, and Strong Sunlight) scenarios and exhibits better results across a broader range of scenarios, which shows remarkable robustness and generalizability.

#### Performance Comparison Analysis with the LoGoNet.

In addition, to provide a clearer analysis of performance across different noise scenarios, we present a more detailed comparative study of our RoboFusion-L and LoGoNet [Li *et al.*, 2023] on the KITTI-C validation dataset, as shown in Table 14. It is worth noting that LoGoNet is a SOTA multi-modal 3D detector known for its exceptional robustness and high accuracy. [Dong *et al.*, 2023] provides noise at varying levels, with the KITTI-C dataset including 5 severities. It is evident that our method demonstrates a high degree of robustness, exhibiting the most stable results with the variance of noise severities. For instance, when considering snow conditions, the performance of our RoboFusion-L shows a marginal variation from 86.69% to 83.67% across severities from 1 to 5. In contrast, LoGoNet’s performance drops from 55.07% to 45.02% over the same severity range. Furthermore, in the presence of moving object noise, our method outperforms LoGoNet. In summary, our RoboFusion exhibits remarkable robustness and generalization capabilities, making it well-suited to diverse noise scenarios.

#### More Results on the nuScenes-C validation set.

As depicted in Table 15, compared with LiDAR-only methods including PointPillars [Lang *et al.*, 2019], and CenterPoint [Yin *et al.*, 2021], Camera-Only methods FCOS3D



Table 11: Comparison with SOTA methods on **KITTI-C validation** set. The results are evaluated based on the **car** class with AP of  $R_{40}$  at **moderate** difficulty. The best one is highlighted in **bold**. ‘S.L.’ denote Strong Sunlight. ‘RCE’ denotes Relative Corruption Error from Ref.[Dong *et al.*, 2023].

Corruptions		LiDAR-Only				Camera-Only		LC Fusion			RoboFusion (Ours)		
		SECOND <sup>†</sup>	PointPillars <sup>†</sup>	PointRCNN <sup>†</sup>	PV-RCNN <sup>†</sup>	SMOKE <sup>†</sup>	ImVoxelNet <sup>†</sup>	EPNet <sup>†</sup>	Focals Conv <sup>†</sup>	LoGoNet <sup>*</sup>	L	B	T
None(AP <sub>clean</sub> )		81.59	78.41	80.57	84.39	7.09	11.49	82.72	85.88	85.04	<b>88.04</b>	87.87	87.60
Weather	Snow	52.34	36.47	50.36	52.35	2.47	0.22	34.58	34.77	51.45	<b>85.29</b>	84.70	84.60
	Rain	52.55	36.18	51.27	51.58	3.94	1.24	36.27	41.30	55.80	<b>86.48</b>	85.54	84.79
	Fog	74.10	64.28	72.14	79.47	5.63	1.34	44.35	44.55	67.53	<b>85.53</b>	84.00	84.17
	S.L.	78.32	62.28	62.78	79.91	6.00	10.08	69.65	80.97	75.54	<b>85.50</b>	85.15	84.75
Sensor	Density	80.18	76.49	80.35	82.79	-	-	82.09	84.95	83.68	<b>85.71</b>	84.34	84.11
	Cutout	73.59	70.28	73.94	76.09	-	-	76.10	78.06	77.17	<b>83.17</b>	81.30	81.21
	Crosstalk	80.24	70.85	71.53	82.34	-	-	82.10	<b>85.82</b>	82.00	84.12	82.45	83.07
	Gaussian (L)	64.90	74.68	61.20	65.11	-	-	60.88	<b>82.14</b>	61.85	76.56	78.32	76.52
	Uniform (L)	79.18	77.31	76.39	81.16	-	-	79.24	<b>85.81</b>	82.94	85.05	83.04	84.11
	Impulse (L)	81.43	78.17	79.78	82.81	-	-	81.63	85.01	84.66	85.26	85.06	<b>85.46</b>
	Gaussian (C)	-	-	-	-	1.56	2.43	80.64	80.97	84.29	82.16	<b>84.63</b>	82.17
	Uniform (C)	-	-	-	-	2.67	4.85	81.61	83.38	84.45	83.30	<b>85.20</b>	83.30
	Impulse (C)	-	-	-	-	1.83	2.13	81.18	80.83	84.20	83.51	<b>84.55</b>	82.91
Motion	Moving Obj.	52.69	50.15	50.54	54.60	1.67	5.93	<b>55.78</b>	49.14	14.44	49.30	49.12	49.90
	Motion Blur	-	-	-	-	3.51	4.19	74.71	81.08	84.52	84.17	<b>84.56</b>	84.18
Object	Local Density	75.10	69.56	74.24	77.63	-	-	76.73	80.84	78.63	83.21	82.53	<b>83.22</b>
	Local Cutout	68.29	61.80	67.94	72.29	-	-	69.92	76.64	64.88	<b>77.22</b>	75.27	76.23
	Local Gaussian	72.31	76.58	69.82	70.44	-	-	75.76	<b>82.02</b>	55.66	79.02	78.32	78.33
	Local Uniform	80.17	78.04	77.67	82.09	-	-	81.71	84.69	79.94	<b>84.69</b>	83.70	84.37
	Local Impulse	81.56	78.43	80.26	84.03	-	-	82.21	<b>85.78</b>	84.29	85.26	85.08	85.06
Average(AP <sub>cor</sub> )		71.68	66.34	68.76	73.41	3.25	3.60	70.35	74.43	71.89	<b>81.72</b>	81.31	81.12
RCE (%) ↓		12.14	15.38	14.65	13.00	54.11	68.65	14.94	13.32	15.46	<b>7.17</b>	7.46	7.38

<sup>†</sup>: Results from Ref. [Dong *et al.*, 2023].

\* denotes re-implement result.

Table 12: Comparison with SOTA methods on **KITTI-C validation** set. The results are evaluated based on the **car** class with AP of  $R_{40}$  at **easy** difficulty. The best one is highlighted in **bold**. ‘S.L.’ denotes Strong Sunlight. ‘RCE’ denotes Relative Corruption Error from Ref.[Dong *et al.*, 2023].

Corruptions		Lidar-Only				Camera-Only		LC Fusion			RoboFusion (Ours)		
		SECOND <sup>†</sup>	PointPillars <sup>†</sup>	PointRCNN <sup>†</sup>	PV-RCNN <sup>†</sup>	SMOKE <sup>†</sup>	ImVoxelNet <sup>†</sup>	EPNet <sup>†</sup>	Focals Conv <sup>†</sup>	LoGoNet <sup>*</sup>	L	B	T
None(AP <sub>clean</sub> )		90.53	87.75	91.65	92.10	10.42	17.85	92.29	92.00	92.04	<b>93.30</b>	93.22	93.28
Weather	Snow	73.05	55.99	71.93	73.06	3.68	0.30	48.03	53.80	74.24	<b>88.77</b>	88.18	88.31
	Rain	73.31	55.17	70.79	72.37	5.66	1.77	50.93	61.44	75.96	88.12	<b>88.57</b>	87.75
	Fog	85.58	74.27	85.01	<b>89.21</b>	8.06	2.37	64.83	68.03	86.60	88.96	88.16	88.09
	S.L.	88.05	67.42	64.90	87.27	8.75	15.72	81.77	90.03	80.30	89.79	89.23	<b>90.36</b>
Sensor	Density	90.45	86.86	91.33	91.98	-	-	91.89	91.14	91.85	<b>92.90</b>	92.08	92.12
	Cutout	81.75	78.90	83.33	83.40	-	-	84.17	83.84	84.20	<b>85.94</b>	85.75	84.75
	Crosstalk	89.63	78.51	77.38	90.52	-	-	91.30	92.01	88.15	91.71	91.54	<b>92.07</b>
	Gaussian (L)	73.21	86.24	74.28	74.61	-	-	66.99	<b>88.56</b>	64.62	80.96	84.30	83.23
	Uniform (L)	89.50	87.49	89.48	90.65	-	-	89.70	91.77	90.75	<b>92.89</b>	91.28	91.63
	Impulse (L)	90.70	87.75	90.80	91.91	-	-	91.44	92.10	91.66	91.90	91.95	<b>92.30</b>
	Gaussian (C)	-	-	-	-	2.09	3.74	91.62	89.51	91.64	91.94	<b>92.08</b>	91.57
	Uniform (C)	-	-	-	-	3.81	7.66	91.95	91.20	91.84	92.01	92.14	<b>92.93</b>
	Impulse (C)	-	-	-	-	2.57	3.35	91.68	89.90	91.65	91.96	<b>92.04</b>	91.33
Motion	Moving Obj.	62.64	58.49	59.29	63.36	2.69	9.63	<b>66.32</b>	54.57	16.83	53.09	51.94	51.70
	Motion Blur	-	-	-	-	5.39	6.75	89.65	91.56	91.96	91.99	<b>92.09</b>	92.06
Object	Local Density	87.74	82.90	88.37	89.60	-	-	89.40	89.60	89.00	92.02	<b>92.42</b>	92.42
	Local Cutout	81.29	75.22	83.30	84.38	-	-	82.40	85.55	77.57	87.30	87.49	<b>87.79</b>
	Local Gaussian	82.05	87.69	82.44	77.89	-	-	85.72	<b>89.78</b>	60.03	89.56	89.41	89.62
	Local Uniform	90.11	87.83	89.30	90.63	-	-	91.32	91.88	88.51	91.59	91.53	<b>91.75</b>
	Local Impulse	90.58	87.84	90.60	91.91	-	-	91.67	92.02	91.34	<b>92.09</b>	91.97	90.69
Average(AP <sub>cor</sub> )		83.10	77.41	80.78	83.92	4.74	5.69	81.63	83.91	80.93	<b>88.27</b>	88.20	88.12
RCE(%) ↓		8.20	11.78	11.85	8.87	54.46	68.07	11.54	8.78	12.07	<b>5.39</b>	<b>5.39</b>	5.53

<sup>†</sup>: Results from Ref. [Dong *et al.*, 2023].

\* denotes re-implement result.

Table 13: Comparison with SOTA methods on **KITTI-C validation** set. The results are evaluated based on the **car** class with AP of  $R_{40}$  at **hard** difficulty. The best one is highlighted in **bold**. ‘S.L.’ denotes Strong Sunlight. ‘RCE’ denotes Relative Corruption Error from Ref.[Dong *et al.*, 2023].

Corruptions		Lidar-Only				Camera-Only		LC Fusion			RoboFusion (Ours)		
		SECOND <sup>†</sup>	PointPillars <sup>†</sup>	PointRCNN <sup>†</sup>	PV-RCNN <sup>†</sup>	SMOKE <sup>†</sup>	ImVoxelNet <sup>†</sup>	EPNet <sup>†</sup>	Focals Conv <sup>†</sup>	LoGoNet *	L	B	T
None(AP <sub>clean</sub> )		78.57	75.19	78.06	82.49	5.57	9.20	80.16	83.36	84.31	<b>85.27</b>	84.27	83.36
Weather	Snow	48.62	32.96	45.41	48.62	1.92	0.20	32.39	30.41	45.57	<b>64.26</b>	62.49	62.74
	Rain	48.79	32.65	45.78	48.20	3.16	0.99	34.69	35.71	50.12	<b>66.07</b>	64.89	63.18
	Fog	68.93	58.19	68.05	75.05	4.56	1.03	38.12	39.50	60.47	<b>80.03</b>	78.37	77.29
	S.L.	74.62	58.69	61.11	78.02	4.91	8.24	66.43	78.06	73.62	80.02	77.52	<b>81.61</b>
Sensor	Density	77.04	72.85	77.58	81.15	-	-	79.77	82.38	81.98	<b>83.06</b>	83.03	83.05
	Cutout	70.79	67.32	71.57	74.60	-	-	73.95	76.69	76.18	76.96	77.00	<b>77.38</b>
	Crosstalk	76.92	67.51	69.41	80.98	-	-	79.54	83.22	80.36	82.94	<b>83.22</b>	83.08
	Gaussian (L)	61.09	71.12	56.73	62.70	-	-	56.88	<b>77.15</b>	59.98	74.45	75.03	73.81
	Uniform (L)	75.61	74.09	72.25	78.93	-	-	75.92	81.62	80.68	81.74	81.79	<b>82.44</b>
	Impulse (L)	78.33	74.65	76.88	81.79	-	-	79.14	<b>83.28</b>	82.51	83.13	83.16	83.24
	Gaussian (C)	-	-	-	-	1.18	1.96	78.20	79.01	82.22	82.86	<b>83.05</b>	81.32
	Uniform (C)	-	-	-	-	2.19	3.90	79.14	81.39	82.37	<b>83.22</b>	83.03	82.06
	Impulse (C)	-	-	-	-	1.52	1.71	78.51	78.87	82.16	82.75	<b>83.00</b>	81.59
Motion	Moving Obj.	48.02	45.47	46.23	50.75	1.40	4.63	<b>50.97</b>	45.34	13.66	43.56	42.62	42.89
	Motion Blur	-	-	-	-	2.95	3.32	72.49	77.75	82.50	<b>83.12</b>	83.06	82.92
Object	Local Density	71.45	65.70	71.09	75.39	-	-	74.36	77.30	76.83	<b>81.71</b>	81.24	81.15
	Local Cutout	63.25	56.69	63.50	68.58	-	-	66.53	72.40	60.62	71.95	72.07	<b>73.78</b>
	Local Gaussian	68.16	73.11	65.65	68.03	-	-	72.71	<b>78.52</b>	54.02	76.38	76.41	76.26
	Local Uniform	76.67	74.68	74.37	80.17	-	-	78.85	81.99	77.44	82.04	82.06	<b>82.33</b>
	Local Impulse	78.47	75.18	77.38	82.33	-	-	79.79	<b>83.20</b>	82.21	82.99	83.16	82.99
Average(AP <sub>cor</sub> )		67.92	62.55	65.18	70.95	2.64	2.88	67.41	71.18	69.27	<b>77.16</b>	76.81	76.75
RCE(%)↓		13.55	16.80	16.49	13.98	52.54	68.62	15.89	14.59	17.83	9.51	9.71	<b>7.93</b>

<sup>†</sup>: Results from Ref. [Dong *et al.*, 2023].

\* denotes re-implement result.

Table 14: Performance comparison of our **RoboFusion-L** with **LoGoNet** on KITTI-C with 5 noise severities. The results are reported based on the **car** with AP of  $R_{40}$  at **moderate** difficulty. ‘S.L.’ denotes Strong Sunlight. The better one is marked in **bold**.

Corruptions		Severity					AP <sub>s</sub>
		1	2	3	4	5	
Weather	Snow	55.07 / <b>86.69</b>	52.98 / <b>86.55</b>	53.08 / <b>85.94</b>	51.14 / <b>83.61</b>	45.02 / <b>83.67</b>	51.45 / <b>85.29</b>
	Rain	57.29 / <b>87.84</b>	56.90 / <b>87.75</b>	56.76 / <b>86.49</b>	55.05 / <b>85.24</b>	53.01 / <b>85.07</b>	55.80 / <b>86.48</b>
	Fog	75.93 / <b>87.31</b>	69.69 / <b>86.58</b>	64.77 / <b>84.71</b>	64.69 / <b>84.56</b>	62.58 / <b>84.51</b>	67.53 / <b>85.53</b>
	S.L.	82.03 / <b>87.26</b>	80.53 / <b>86.53</b>	76.75 / <b>84.66</b>	71.12 / <b>84.61</b>	67.31 / <b>84.46</b>	75.54 / <b>85.50</b>
Sensor	Density	86.60 / <b>86.81</b>	84.59 / <b>86.59</b>	84.05 / <b>85.60</b>	82.74 / <b>85.27</b>	82.42 / <b>84.30</b>	83.68 / <b>85.71</b>
	Cutout	82.18 / <b>87.64</b>	80.02 / <b>86.21</b>	77.41 / <b>83.25</b>	74.66 / <b>80.81</b>	71.59 / <b>77.94</b>	77.17 / <b>83.17</b>
	Crosstalk	84.22 / <b>84.41</b>	83.38 / <b>84.38</b>	81.41 / <b>84.13</b>	80.78 / <b>83.79</b>	80.22 / <b>83.90</b>	82.00 / <b>84.12</b>
	Gaussian (L)	84.69 / <b>85.41</b>	82.52 / <b>84.66</b>	77.43 / <b>81.39</b>	47.28 / <b>73.58</b>	17.31 / <b>57.79</b>	61.85 / <b>76.56</b>
	Uniform (L)	84.77 / <b>85.77</b>	84.64 / <b>85.42</b>	84.39 / <b>85.47</b>	82.32 / <b>85.00</b>	78.59 / <b>83.59</b>	82.94 / <b>85.05</b>
	Impulse (L)	84.45 / <b>84.95</b>	84.73 / 82.88	84.92 / 82.20	84.63 / 80.51	84.56 / 80.29	84.66 / 82.16
	Gaussian (C)	84.53 / <b>85.77</b>	84.47 / <b>85.42</b>	84.31 / <b>85.47</b>	84.18 / <b>85.32</b>	83.96 / <b>84.32</b>	84.29 / <b>85.26</b>
	Uniform (C)	84.74 / <b>85.57</b>	84.57 / <b>85.08</b>	84.54 / 82.96	84.36 / 82.53	84.05 / 80.36	84.45 / 83.30
	Impulse (C)	84.53 / <b>85.70</b>	84.26 / 83.63	84.38 / 83.54	83.95 / 82.42	83.86 / 82.28	84.20 / 83.51
Motion	Moving Obj.	58.89 / <b>78.46</b>	12.78 / <b>67.86</b>	0.43 / <b>41.07</b>	0.06 / <b>36.28</b>	0.07 / <b>22.85</b>	14.44 / <b>49.30</b>
	Motion Blur	84.64 / <b>85.23</b>	84.53 / <b>84.98</b>	84.56 / <b>84.72</b>	84.45 / 83.00	84.43 / 82.96	84.52 / 84.17
Object	Local Density	82.31 / <b>85.23</b>	81.66 / <b>84.87</b>	80.15 / <b>82.70</b>	76.53 / <b>82.08</b>	72.52 / <b>81.21</b>	78.63 / <b>83.21</b>
	Local Cutout	76.77 / <b>82.94</b>	72.46 / <b>81.31</b>	65.87 / <b>78.14</b>	59.14 / <b>74.12</b>	50.17 / <b>69.61</b>	64.88 / <b>77.22</b>
	Local Gaussian	84.45 / <b>86.81</b>	81.12 / <b>86.25</b>	67.13 / <b>82.72</b>	33.33 / <b>76.01</b>	12.27 / <b>63.31</b>	55.66 / <b>79.02</b>
	Local Uniform	84.51 / <b>85.91</b>	84.35 / <b>85.65</b>	81.95 / <b>85.23</b>	79.62 / <b>84.66</b>	69.25 / <b>81.99</b>	79.94 / <b>84.68</b>
	Local Impulse	84.53 / <b>85.65</b>	84.47 / <b>85.13</b>	84.32 / <b>85.18</b>	84.40 / <b>85.16</b>	83.72 / <b>85.16</b>	84.29 / <b>85.25</b>
AP <sub>c</sub>		79.35 / <b>85.56</b>	75.73 / <b>84.38</b>	72.93 / <b>81.77</b>	68.22 / <b>79.92</b>	63.34 / <b>76.97</b>	71.81 / <b>81.72</b>
Clean							85.04 / <b>88.04</b>

Table 15: Comparison with SOTA methods on **nuScenes-C validation** set with **mAP**. ‘D.I.’ refers to DeepInteraction [Yang *et al.*, 2022]. The best one is highlighted in **bold**. ‘S.L.’ denotes Strong Sunlight. ‘RCE’ denotes Relative Corruption Error from Ref.[Dong *et al.*, 2023].

Corruptions		Lidar-Only		Camera-Only			LC Fusion					RoboFusion (Ours)		
		PointPillars <sup>†</sup>	CenterPoint <sup>†</sup>	FCOS3D <sup>†</sup>	DETR3D <sup>†</sup>	BEVFormer <sup>†</sup>	FUTR3D <sup>†</sup>	TransFusion <sup>†</sup>	BEVFusion <sup>†</sup>	D.I.*		L	B	T
None(AP <sub>clean</sub> )		27.69	59.28	23.86	34.71	41.65	64.17	66.38	68.45	69.90		<b>69.91</b>	69.40	69.09
Weather	Snow	27.57	55.90	2.01	5.08	5.73	52.73	63.30	62.84	62.36		<b>67.12</b>	66.07	65.96
	Rain	27.71	56.08	13.00	20.39	24.97	58.40	65.35	66.13	66.48		<b>67.58</b>	67.01	66.45
	Fog	24.49	43.78	13.53	27.89	32.76	53.19	53.67	54.10	54.79		<b>67.01</b>	65.54	64.34
	S.L.	23.71	54.20	17.20	34.66	41.68	57.70	55.14	64.42	64.93		<b>67.24</b>	66.71	66.54
Sensor	Density	27.27	58.60	-	-	-	63.72	65.77	67.79	68.15		<b>69.48</b>	69.02	68.58
	Cutout	24.14	56.28	-	-	-	62.25	63.66	66.18	66.23		<b>69.18</b>	69.01	68.20
	Crosstalk	25.92	56.64	-	-	-	62.66	64.67	67.32	68.12		<b>68.68</b>	68.04	68.17
	FOV lost	8.87	20.84	-	-	-	26.32	24.63	27.17	<b>42.66</b>		39.48	39.30	39.43
	Gaussian (L)	19.41	45.79	-	-	-	58.94	55.10	<b>60.64</b>	57.46		57.77	57.07	56.00
	Uniform (L)	25.60	56.12	-	-	-	63.21	64.72	66.81	<b>67.42</b>		64.57	64.25	64.99
	Impulse (L)	26.44	57.67	-	-	-	63.43	65.51	<b>67.54</b>	67.41		65.64	65.45	65.44
	Gaussian (C)	-	-	3.96	14.86	15.04	54.96	64.52	64.44	66.52		66.73	<b>66.75</b>	66.53
	Uniform (C)	-	-	8.12	21.49	23.00	57.61	65.26	65.81	<b>65.90</b>		65.77	65.76	65.56
	Impulse (C)	-	-	3.55	14.32	13.99	55.16	64.37	64.30	<b>65.65</b>		64.82	64.75	64.56
Motion	Compensation	3.85	11.02	-	-	-	31.87	9.01	27.57	39.95		<b>41.88</b>	39.54	41.28
	Motion Blur	-	-	10.19	11.06	19.79	55.99	64.39	64.74	65.45		<b>67.21</b>	66.52	66.42
Obeject	Local Density	26.70	57.55	-	-	-	63.60	65.65	67.42	<b>67.71</b>		66.74	66.59	65.88
	Local Cutout	17.97	48.36	-	-	-	61.85	63.33	63.41	65.19		<b>66.82</b>	66.53	66.76
	Local Gaussian	25.93	51.13	-	-	-	62.94	63.76	64.34	64.75		65.08	<b>65.17</b>	64.77
	Local Uniform	27.69	57.87	-	-	-	64.09	66.20	<b>67.58</b>	66.44		66.71	66.19	65.40
	Local Impulse	27.67	58.49	-	-	-	64.02	66.29	<b>67.91</b>	67.86		66.53	66.87	66.67
Average(AP <sub>cor</sub> )		22.99	49.78	8.94	18.71	22.12	56.88	58.77	61.35	62.92		<b>63.90</b>	63.43	63.23
RCE (%) ↓		16.95	16.01	62.51	46.07	46.89	11.34	11.45	10.36	9.97		8.58	8.59	<b>8.47</b>

<sup>†</sup>: Results from Ref. [Dong *et al.*, 2023].

\* denotes re-implement result.

[Wang *et al.*, 2021], DETR3D [Wang *et al.*, 2022b], and BEVFormer [Li *et al.*, 2022c] and multi-modal methods including FUTR3D [Chen *et al.*, 2023], TransFusion [Bai *et al.*, 2022], BEVFusion [Liu *et al.*, 2023] and DeepInteraction [Yang *et al.*, 2022], our RoboFusion demonstrates superior performance across more noise scenarios in AD on average. For instance, our RoboFusion-L excels in 10 noise scenarios, including Weather (Snow, Rain, Fog, Strong Sunlight), Sensor (Density, Cutout, Crosstalk), Motion (Compensation, Motion Blur), and Object (Local Cutout), outperforming DeepInteraction [Yang *et al.*, 2022] which achieves the best performance only in 5 of these noise scenarios. Overall, our method exhibits not only exceptional robustness in weather-induced noise scenarios, but also shows remarkable resilience across a broader noise include sensor, motion and object noise.

### A.3 Visualization

As shown in Fig. 5, we provide visualization results between our RoboFusion-L and LoGoNet on the KITTI-C dataset. Overall, compared to SOTA methods like LoGoNet [Li *et al.*, 2023], our method enhances the robustness of multi-modal 3D object detection by leveraging the generalization capability and robustness of VFMs to mitigate OOD noisy scenarios in AD.

### A.4 More Limitations

Although we have mentioned the two main limitations in the ‘Conclusions’ section of the main text, our RoboFusion still

has other limitations. Our method does not achieve the best performance in all noisy scenarios. For instance, as shown in Table 11, our method does not show the best in ‘Moving Object’ noisy scenarios. Furthermore, we conduct experiments only on the corruption datasets [Dong *et al.*, 2023] rather than real-world datasets. It is valuable to construct a real-world corruption dataset, but it must be an expensive work.

## References

- [Bai *et al.*, 2022] Xuyang Bai, Zeyu Hu, Xinge Zhu, Qingqiu Huang, Yilun Chen, Hongbo Fu, and Chiew-Lan Tai. TransFusion: Robust LiDAR-camera fusion for 3D object detection with transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1090–1099, 2022.
- [Caesar *et al.*, 2020] 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. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11621–11631, 2020.
- [Chen *et al.*, 2022] Yukang Chen, Yanwei Li, Xiangyu Zhang, Jian Sun, and Jiaya Jia. Focal sparse convolutional networks for 3D object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5428–5437, 2022.
- [Chen *et al.*, 2023] Xuanyao Chen, Tianyuan Zhang, Yue Wang, Yilun Wang, and Hang Zhao. FUTR3D: A unified

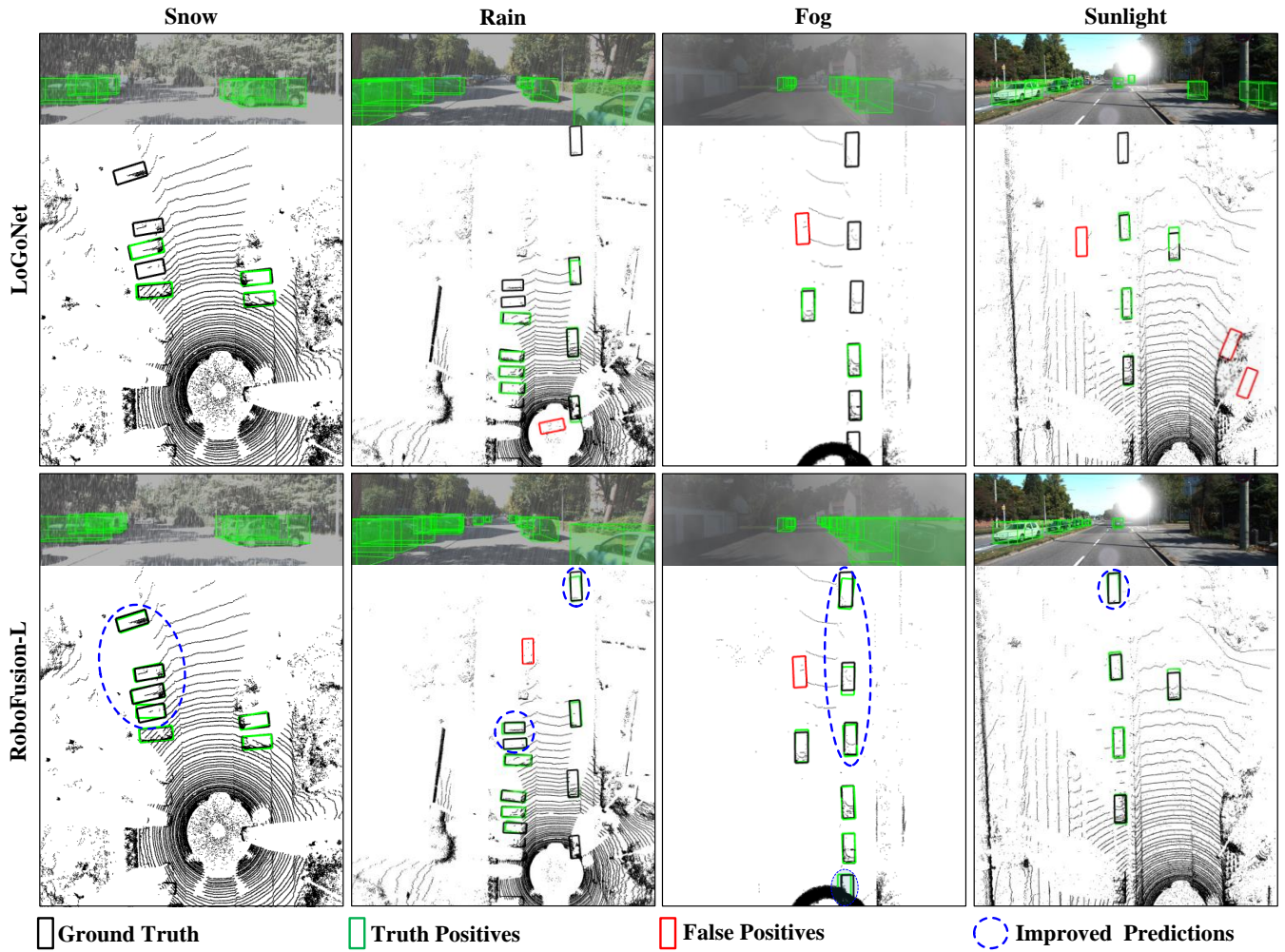


Figure 5: **Visualization Results of LoGoNet and our RoboFusion in KITTI-C dataset.** We use boxes in red to represent false positives, green boxes for truth positives, and black for the ground truth. We use blue dashed ovals to highlight the pronounced improvements in predictions.

sensor fusion framework for 3D detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 172–181, 2023.

[Deng et al., 2021] Jiajun Deng, Shaoshuai Shi, Peiwei Li, Wengang Zhou, Yanyong Zhang, and Houqiang Li. Voxel R-CNN: Towards high performance voxel-based 3D object detection. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 1201–1209, 2021.

[Dong et al., 2023] Yinpeng Dong, Caixin Kang, Jinlai Zhang, Zijian Zhu, et al. Benchmarking robustness of 3D object detection to common corruptions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1022–1032, 2023.

[Dosovitskiy et al., 2020] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *CoRR*, abs/2010.11929, 2020.

[Geiger et al., 2012] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? the kitti vision benchmark suite. In *IEEE conference on computer vision and pattern recognition*, pages 3354–3361, 2012.

[Hu et al., 2023] Qianjiang Hu, Daizong Liu, and Wei Hu. Density-insensitive unsupervised domain adaption on 3D object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 17556–17566, 2023.

[Huang et al., 2020] Tengpeng Huang, Zhe Liu, Xiwu Chen, and Xiang Bai. EPNet: Enhancing point features with image semantics for 3D object detection. In *European Conference on Computer Vision*, pages 35–52. Springer, 2020.

[Kirillov et al., 2023] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, et al. Segment anything. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4015–4026, 2023.



- [Lang et al., 2019] Alex H Lang, Sourabh Vora, et al. PointPillars: Fast encoders for object detection from point clouds. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 12697–12705, 2019.
- [Li et al., 2022a] Yanghao Li, Hanzi Mao, Ross Girshick, and Kaiming He. Exploring plain vision transformer backbones for object detection. In *European Conference on Computer Vision*, pages 280–296. Springer, 2022.
- [Li et al., 2022b] Yanwei Li, Xiaojuan Qi, Yukang Chen, Liwei Wang, Zeming Li, Jian Sun, and Jiaya Jia. Voxel field fusion for 3D object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1120–1129, 2022.
- [Li et al., 2022c] Zhiqi Li, Wenhai Wang, Hongyang Li, Enze Xie, et al. BEVFormer: Learning bird’s-eye-view representation from multi-camera images via spatiotemporal transformers. In *European conference on computer vision*, pages 1–18. Springer, 2022.
- [Li et al., 2023] Xin Li, Tao Ma, Yuenan Hou, Botian Shi, et al. LoGoNet: Towards accurate 3D object detection with local-to-global cross-modal fusion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 17524–17534, 2023.
- [Lin et al., 2017] 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.
- [Liu et al., 2020a] Lin Liu, Jianzhuang Liu, Shanxin Yuan, Gregory Slabaugh, Aleš Leonardis, Wengang Zhou, and Qi Tian. Wavelet-based dual-branch network for image demoiréing. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XIII 16*, pages 86–102. Springer, 2020.
- [Liu et al., 2020b] Zechen Liu, Zizhang Wu, and Roland Tóth. SMOKE: Single-stage monocular 3D object detection via keypoint estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 996–997, 2020.
- [Liu et al., 2023] Zhijian Liu, Haotian Tang, Alexander Amini, Xinyu Yang, Huizi Mao, Daniela L Rus, and Song Han. BEVFusion: Multi-task multi-sensor fusion with unified bird’s-eye view representation. In *2023 IEEE international conference on robotics and automation (ICRA)*, pages 2774–2781. IEEE, 2023.
- [OpenAI, 2023] OpenAI. Gpt-4 technical report. <https://cdn.openai.com/papers/gpt-4.pdf>, 2023.
- [Oza et al., 2023] Poojan Oza, Vishwanath A Sindagi, Vibashan Vishnukumar Sharmini, and Vishal M Patel. Unsupervised domain adaptation of object detectors: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023.
- [Peng et al., 2023] Xidong Peng, Xinge Zhu, and Yuexin Ma. CL3D: Unsupervised domain adaptation for cross-LiDAR 3D detection. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pages 2047–2055, 2023.
- [Rukhovich et al., 2022] Danila Rukhovich, Anna Vorontsova, and Anton Konushin. ImVoxelNet: Image to voxels projection for monocular and multi-view general-purpose 3D object detection. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 2397–2406, 2022.
- [Shi et al., 2019] Shaoshuai Shi, Xiaogang Wang, and Hongsheng Li. PointRCNN: 3D object proposal generation and detection from point cloud. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 770–779, 2019.
- [Shi et al., 2020] Shaoshuai Shi, Chaoxu Guo, Li Jiang, Zhe Wang, Jianping Shi, Xiaogang Wang, and Hongsheng Li. PV-RCNN: point-voxel feature set abstraction for 3D object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10529–10538, 2020.
- [Song et al., 2023] Ziyang Song, Haiyue Wei, Lin Bai, et al. GraphAlign: Enhancing accurate feature alignment by graph matching for multi-modal 3D object detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3358–3369, 2023.
- [Song et al., 2024a] Ziyang Song, Lin Liu, Feiyang Jia, et al. Robustness-aware 3D object detection in autonomous driving: A review and outlook. *arXiv preprint arXiv:2401.06542*, 2024.
- [Song et al., 2024b] Ziyang Song, Guoxin Zhang, Jun Xie, Lin Liu, et al. Voxelnexfusion: A simple, unified and effective voxel fusion framework for multi-modal 3D object detection. *arXiv preprint arXiv:2401.02702*, 2024.
- [Tsai et al., 2023] Darren Tsai, Julie Stephany Berrio, Mao Shan, Eduardo Nebot, and Stewart Worrall. Viewer-centred surface completion for unsupervised domain adaptation in 3D object detection. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 9346–9353, 2023.
- [Wang et al., 2021] Tai Wang, Xinge Zhu, Jiangmiao Pang, and Dahua Lin. FCOS3D: Fully convolutional one-stage monocular 3D object detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 913–922, 2021.
- [Wang et al., 2022a] Tai Wang, ZHU Xinge, Jiangmiao Pang, and Dahua Lin. Probabilistic and geometric depth: Detecting objects in perspective. In *Conference on Robot Learning*, pages 1475–1485. PMLR, 2022.
- [Wang et al., 2022b] Yue Wang, Vitor Campagnolo Guizilini, Tianyuan Zhang, et al. DERT3D: 3D object detection from multi-view images via 3D-to-2D queries. In *Conference on Robot Learning*, pages 180–191. PMLR, 2022.
- [Wang et al., 2023a] Li Wang, Xinyu Zhang, Ziyang Song, et al. Multi-modal 3d object detection in autonomous driv-

- ing: A survey and taxonomy. *IEEE Transactions on Intelligent Vehicles*, 2023.
- [Wang *et al.*, 2023b] Yan Wang, Junbo Yin, Wei Li, et al. SSDA3D: Semi-supervised domain adaptation for 3D object detection from point cloud. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pages 2707–2715, 2023.
- [Wu *et al.*, 2023] QuanLin Wu, Hang Ye, Yuntian Gu, et al. Denoising masked autoencoders help robust classification. In *The Eleventh International Conference on Learning Representations*, 2023.
- [Xie *et al.*, 2023] Yichen Xie, Chenfeng Xu, Marie-Julie Rakotosaona, Patrick Rim, et al. SparseFusion: Fusing multi-modal sparse representations for multi-sensor 3D object detection. *arXiv preprint arXiv:2304.14340*, 2023.
- [Yan *et al.*, 2018] Yan Yan, Yuxing Mao, and Bo Li. SECOND: Sparsely embedded convolutional detection. *Sensors*, 18(10):3337, 2018.
- [Yan *et al.*, 2023] Junjie Yan, Yingfei Liu, Jianjian Sun, Fan Jia, et al. Cross modal transformer via coordinates encoding for 3D object detection. *arXiv preprint arXiv:2301.01283*, 2023.
- [Yang *et al.*, 2022] Zeyu Yang, Jiaqi Chen, Zhenwei Miao, et al. DeepInteraction: 3D object detection via modality interaction. *Advances in Neural Information Processing Systems*, 35:1992–2005, 2022.
- [Yin *et al.*, 2021] Tianwei Yin, Xingyi Zhou, and Philipp Krahenbuhl. Center-based 3d object detection and tracking. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11784–11793, 2021.
- [Zhang *et al.*, 2022] Yanan Zhang, Jiaxin Chen, and Di Huang. CAT-Det: Contrastively augmented transformer for multi-modal 3D object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 908–917, 2022.
- [Zhang *et al.*, 2023a] Chaoning Zhang, Dongshen Han, Yu Qiao, Jung Uk Kim, et al. Faster segment anything: Towards lightweight SAM for mobile applications. *arXiv preprint arXiv:2306.14289*, 2023.
- [Zhang *et al.*, 2023b] Dingyuan Zhang, Dingkan Liang, Hongcheng Yang, Zhikang Zou, Xiaoqing Ye, Zhe Liu, and Xiang Bai. SAM3D: Zero-shot 3D object detection via segment anything model. *arXiv preprint arXiv:2306.02245*, 2023.
- [Zhang *et al.*, 2023c] Xinyu Zhang, Li Wang, Jian Chen, et al. Dual radar: A multi-modal dataset with dual 4d radar for autonomous driving. *arXiv preprint arXiv:2310.07602*, 2023.
- [Zhao *et al.*, 2023] Xu Zhao, Wenchao Ding, Yongqi An, Yinglong Du, et al. Fast segment anything, 2023.
- [Zhu *et al.*, 2020] Xinge Zhu, Yuexin Ma, Tai Wang, et al. SSN: Shape signature networks for multi-class object detection from point clouds. In *European Conference on Computer Vision*, pages 581–597. Springer, 2020.