# SaViD: Spectravista Aesthetic Vision Integration for Robust and Discerning 3D Object Detection in Challenging Environments

Tanmoy Dam[1], Sanjay Bhargav Dharavath[2], Sameer Alam[1], Nimrod Lilith[1], Aniruddha Maiti [3]
Supriyo Chakraborty[2] and Mir Feroskhan[1]

*Abstract*—The fusion of LiDAR and camera sensors has demonstrated significant effectiveness in achieving accurate detection for short-range tasks in autonomous driving. However, this fusion approach could face challenges when dealing with long-range detection scenarios due to disparity between sparsity of LiDAR and high-resolution camera data. Moreover, sensor corruption introduces complexities that affect the ability to maintain robustness, despite the growing adoption of sensor fusion in this domain. We present SaViD, a novel framework comprised of a three-stage fusion alignment mechanism designed to address long-range detection challenges in the presence of natural corruption. The SaViD framework consists of three key elements: the Global Memory Attention Network (GMAN), which enhances the extraction of image features through offering a deeper understanding of global patterns; the Attentional Sparse Memory Network (ASMN), which enhances the integration of LiDAR and image features; and the KNNnectivity Graph Fusion (KGF), which enables the entire fusion of spatial information. SaViD achieves superior performance on the long-range detection Argoverse-2 (AV2) dataset with a performance improvement of 9.87% in AP value and an improvement of 2.39% in mAPH for L2 difficulties on the Waymo Open dataset (WOD). Comprehensive experiments are carried out to showcase its robustness against 14 natural sensor corruptions. SaViD exhibits a robust performance improvement of 31.43% for AV2 and 16.13% for WOD in RCE value compared to other existing fusion-based methods while considering all the corruptions for both datasets. Our code is available at SaViD.

*Index Terms*—GMAN, ASMN, KGF, Multi-modal fusion, 3D object detection

## I. INTRODUCTION

Autonomous driving uses LiDAR and cameras for 3D object detection [1]–[6]. Cameras offer high-resolution details; LiDAR adds depth and shape data. The fusion of both is crucial for accuracy but challenging [1], [7]. This study addresses these fusion challenges with a robust solution.

*Challenge 1: Weak fusion or alignment between LiDAR and camera features for long-range detection.* While camera features are often fused with LiDAR data [8], most methods [9]–[11] use mid-level fusion, such as additive fusion in

a. Long-range objects     b. Corruption-induced

Fig. 1: Argoverse-2 dataset two challenges

AVOD [12] or cross-attention in DeepFusion [9] and TransFusion [10]. These approaches may miss long-range details due to low-resolution camera inputs and sparse LiDAR data, highlighting the need for strong fusion alignment for accurate 3D detection in long-range scenarios [13].

*Challenge 2: Robustness in long-range detection.*

Data-driven deep learning models struggle to generalize on corrupted data from adverse weather, sensor noise, and other factors [14]–[16]. This limits the reliability of autonomous driving. Recent robustness assessments have developed datasets focused on adverse conditions [16], but evaluations are mainly on small-range datasets like KITTI and nuScenes. Achieving robustness for long-range detection remains a significant challenge, requiring dedicated benchmark analysis.

**Our contribution.** We introduce SaViD, a novel method for robust long-range detection, distinct from traditional LiDAR-camera fusion techniques. SaViD achieves strong modality alignment by integrating sparse LiDAR point cloud data with camera features through local-global view representations [17], and incorporates natural robustness to effectively handle long-range detection and adverse conditions (Figure 1). To summarize, our main contributions in this paper are described as follows:

- GMAN: A memory-based vision transformer that extracts image features using depth as a global query.
- ASMN: A single-stage method for aligning sparse point cloud features with global image features.
- KGF: A parameter-free fusion alignment technique for accurate integration of pseudo-point clouds and images.
- SaViD achieves state-of-the-art performance on long-range detection in both clean and corrupted Argoverse-2 and Waymo Open Dataset.

## II. RELATED WORKS

**LiDAR Point Clouds for 3D Object Detection.** LiDAR-only 3D detection aims to predict 3D bounding boxes within
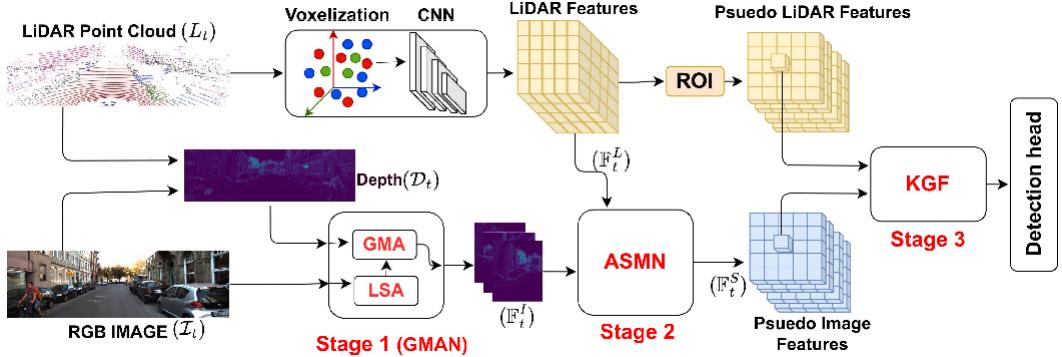
Fig. 2: SaViD Pipeline: The pipeline integrates multiple modalities through three essential components: GMAN, ASMN, and KGF. GMAN considers the image($\mathcal{I}_t$) as a local feature and the Depth($\mathcal{D}_t$) information as a global query, enabling an understanding of the scene through its vision-based Transformer in Stage 1. In Stage 2, a single-stage fusion integration of the ASMN module aligns the extracted image feature ($\mathbb{F}_t^{\mathcal{I}}$) by GMAN with the voxelized LiDAR feature, resulting in a cohesive global feature alignment. Finally, a parameter-free KGF alignment LiDAR RoI feature ($\mathbb{F}_t^L$) and the extracted ASMN Image feature ($\mathbb{F}_t^S$) with local importance are in Stage 3.

raw, unprocessed point clouds. Detectors often project these points onto grids like pillars [18], range images [13], or 3D voxels [19] to handle their irregular structure. Neural networks such as PointNet [20] and PointNet++ [21] extract features in the BEV plane, preserving object dimensions. Other methods use high-density range images for depth data [13], [22], [23]. These approaches can be classified as fully-dense, semi-dense, or point-dense [24], [25]. Due to the sparse nature of point clouds, achieving strong performance in long-range detection with single-modality detectors remains challenging.

**LiDAR-camera Integration for 3D Object Detection.** LiDAR-camera fusion is challenging due to their differing data: LiDAR provides 3D depth, while cameras capture 2D visuals. Effective integration requires robust algorithms. Prior works like DeepFusion [9], TransFusion [10], and BEVFusion [11] address these complexities. DeepFusion uses self-attention but struggles with long-range detection due to resolution mismatches between LiDAR and camera data. TransFusion's Multi-head Attention faces issues with data density variations. BEVFusion's simple concatenation of both modalities may hinder detection of distant objects due to blurred camera images and sparse LiDAR points. These limitations suggest that current attention-based techniques may not fully ensure effective fusion.

## III. SaViD Pipeline

**Problem Definition.** The objective of this paper is to develop a robust 3D object detection approach using multi-modal sensors, achieving effective performance in challenging conditions. We consider multi-modal input-output sequences defined as $(\mathcal{X}_t, \mathcal{Y}_t) = \{(\mathcal{I}_t, L_t), (\mathcal{I}_{(t-1)}, L_{(t-1)}), \ldots\}$, where $\mathcal{I}_t \in \mathbb{R}^{H \times W \times 3}$ is the $t$-th camera image and $L_t \in \mathbb{R}^{N \times 3}$ is the LiDAR point cloud. The output $\mathcal{Y}_t$ consists of 3D bounding boxes associated with classes $\mathcal{M} = \{1, 2, \ldots, M\}$. Our framework, denoted as $\Phi$, predicts $\hat{\mathcal{Y}}_t = \Phi(\mathcal{I}_t, L_t)$ with high-confidence 3D bounding

boxes, closely resembling the ground truth $\mathcal{Y}_t$. We assume the use of depth information $\mathcal{D}_t$, which can be estimated from $\mathcal{I}_t$ and $L_t$, as an additional input for $\Phi$.

**Depth Estimation** ($\mathcal{D}_t$). To generate a high-resolution depth map from sparse LiDAR data $L_t$, combined with correlated RGB imagery $\mathcal{I}_t$, we extract global features beneficial for image processing. Therefore, a frame of point clouds $L_t$ can be transformed into a sparse depth map $\mathcal{D}_t \in \mathbb{R}^{H \times W \times 3}$ using a projection function $\mathcal{T}_{L_t, \mathcal{I}_t} \to \mathcal{D}_t$. In this context, the mapping function $\mathcal{T}$, implemented as a depth neural network, uses both the RGB image $\mathcal{I}_t$ and point clouds $L_t$ to produce the high-resolution depth map $\mathcal{D}_t$.

### A. LiDAR Feature Extraction through Voxelization

The point clouds $L_t$ are sparse and unevenly distributed. We preprocess by voxelizing the $t$-th point cloud with dimensions $H_v \times W_v \times C_v$ and compute voxel features by averaging point-wise features in non-empty voxels [26]. Key-points are identified using Furthest Point Sampling (FPS) [26], selecting $\mathcal{K} = 4096$ key-points ($L_t^{\mathcal{K}}$) for experiments. The characteristics of non-empty voxels are obtained by averaging 3D coordinates and reflectance values of contained points. Feature volumes are transformed through $3 \times 3 \times 3$ 3D sparse convolutions, resulting in downsampled resolutions of $1\times, 2\times, 4\times$, and $8\times$. These volumes are represented as feature vectors assigned to individual voxels, with the final voxel feature vectors denoted as $\mathbb{F}_t^L \in \mathbb{R}^{H \times W \times C}$.

### B. Multi-Modal Feature Fusion Alignment

**Stage 1: GMAN through $\mathcal{D}_t$.** We introduce GMAN, which combines local-global attention with frequency domain information via FFT-iFFT layers. As shown in Figure 3, the Local-Spectral Attention is a vision-based transfer architecture [27]. The local image tensor $\mathcal{I}_t$ serves as the query with dimensions $(B, H/P \times W/P, P \times P, C)$, where $B$ is the batch size, $P \times P$ represents the local patch window size, and $C$ is the channel count. The aggregated batch size is $B^* = H/P \times W/P$. This paper proposes a novel

transformer block that extracts features from $\mathcal{I}_t$ and $\mathcal{D}_t$ using two key modules: Local Spectral Attention (LSA) and Global Memory Attention (GMA). LSA and GMA capture local and global depth-relevant features to handle varying object scales within the same window. The local feature extraction is akin to the Swin Transformer [28] and Global Vision Transformer [27]. $\mathcal{I}_t$ passes through a local query generator utilizing the LSA module for detail-focused feature extraction.
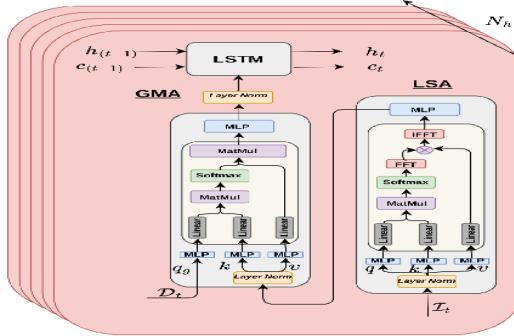


Fig. 3: The GMAN Architecture is composed of both LSA and GMA blocks, employing a total of $N_h$ attention heads. Moreover, the LSTM block operates across consecutive time frames.

**GMA:** While LSA focuses on local patches within the current frame ($\mathcal{I}_t$), GMA operates on a unified framework, utilizing depth modality from depth maps ($\mathcal{D}_t$). Unlike single modality methods [29], GMA's global query computation is predefined, using globally extracted query tokens from $\mathcal{D}_t$ and interacting with local key and value representations from LSA. This allows GMA to integrate both local and global information. The GMA module leverages global context by applying attention across $\mathcal{I}_t$ and $\mathcal{D}_t$, correlating $\mathcal{I}_t$-derived key-value pairs with $\mathcal{D}_t$ as the query. Since $\mathcal{D}_t$ integrates LiDAR ($L_t$) and RGB ($\mathcal{I}_t$) data, GMA effectively attends to various locations within $L_t$, enhancing contextual understanding. GMA is formulated as:

$$\text{GMA}(\mathcal{I}_t, \mathcal{D}_t) = \mathbb{LN}(\alpha v), \quad \alpha = \text{Softmax}(g(q_g, k)) \quad (1)$$

where $q_g \in \mathcal{D}_t$, $k, v \in \mathcal{I}_t$ are the query, key, and value, respectively, and $g(\cdot)$ is the global attention function. Temporal feature accumulation is managed using ReLU-activated LSTM cells [30] after GMA to capture sequential patterns. Algorithm 1 provides pseudo-code for the GMA module.

**Stage 2: ASMN.** ASMN introduces a novel temporal fusion mechanism between feature extractors $\mathbb{F}_t^{\mathcal{I}} \in \mathbb{R}^{H \times W \times C}$ and $\mathbb{F}_t^L \in \mathbb{R}^{H \times W \times C}$. Unlike GMAN, ASMN uses a single-stage integration that handles sequential information from both modalities. To address the sparsity of LiDAR voxel features ($\mathbb{F}_t^L$), we incorporate sparse attention [31] combined with LSTM cells to enhance fusion with $\mathbb{F}_t^{\mathcal{I}}$. Following GMA principles, voxel features act as global queries, interacting with key-value pairs from image features through sparse attention. This interaction generates a correspondence map

---

**Algorithm 1** Global Memory Attention (GMA)

**Input/Output:** (B*, N, C) where B* = B × N*, $N_h$: Attention Heads

**Initialization:**
$\quad k_m, q_m, v_m = nn.\text{Linear}(C, C)$
$\quad \text{softmax} = nn.\text{Softmax}(dim = -1)$
$\quad \text{LSTM} = nn.\text{LSTM}(.)$

**Forward**($\mathcal{I}_t, \mathcal{D}_t$):
$\quad k = k_m(\mathcal{I}_t), \ q\_g = q_m(\mathcal{D}_t), \ v = v_m(\mathcal{I}_t)$
$\quad k, q\_g, v = k.\text{view}(B*, N, N_h, -1).\text{permute}(0, 2, 1, 3)$
$\quad q\_k = \text{matmul}(q\_g, k.\text{transpose}(-2, -1))$
$\quad \text{attn} = \text{softmax}(q\_k)$
$\quad \text{attn}\_v = \text{matmul}(\text{attn}, v.\text{transpose}(-2, -1))$
$\quad$**return** $\text{LSTM}(\text{attn}\_v).\text{reshape}(B*, N, C)$

---

linking $\mathbb{F}_t^L$ to regions in $\mathbb{F}_t^{\mathcal{I}}$, with LSTM states adapting to both modalities for effective sequence integration.
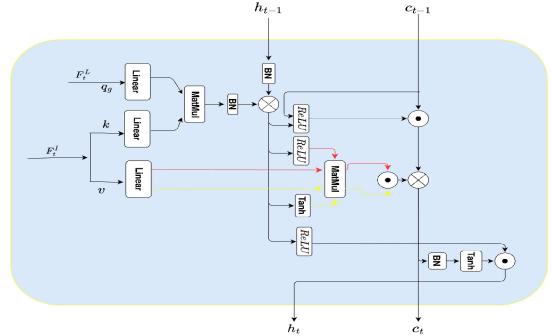


Fig. 4: ASMN Architecture: The LiDAR feature extractor ($\mathbb{F}_t^L$) serves as a global query in a unified stage with $\mathbb{F}_t^{\mathcal{I}}$. Sequential data is processed using an LSTM cell with sparse attention and ReLU activation.

$$\text{ASMN}(\mathbb{F}_t^{\mathcal{I}}, \mathbb{F}_t^L) = (\beta_c v) \cdot \text{Tanh}(\beta_h v) \quad (2)$$

where $\beta_h = f_h(q_g, k, h_{t-1})$ and $\beta_c = f_c(\beta_h, c_{t-1})$, with $q_g \in \mathbb{F}_t^L$, $k, v \in \mathbb{F}_t^{\mathcal{I}}$. $f_h$ and $f_c$ update key-query pairs using previous hidden and cell states with batch normalization and ReLU activation. Therefore, LSTM updates are:

$$c_t = \beta_c \cdot ((\beta_c v) \cdot \text{Tanh}(\beta_h v)), \quad h_t = \text{ReLU}(\beta_h) \cdot \text{Tanh}(c_t) \quad (3)$$

For ASMN, $\beta_h$ and $\beta_c$ are defined as:

$$\beta_h = BN(\updownarrow(q_g) \cdot \updownarrow(k)) \cdot BN(h_{t-1}), \quad \beta_c = \text{ReLU}(\beta_h \cdot c_{t-1}) \quad (4)$$

$$\mathbb{F}_t^S = \text{ReLU}(\beta_c \cdot \updownarrow v) \cdot \text{Tanh}(\beta_h \cdot \updownarrow v) \quad (5)$$

Here, $\updownarrow$ is a linear function, and $h_{t-1}$, $c_{t-1}$ are hidden and cell states from the previous time frame.

**Stage 3: KGF.** KGF is a parameter-free fusion alignment technique that captures local correlations between $\mathbb{F}_t^S$ and $\mathbb{F}_t^L$ by identifying similar attributes across modalities and emphasizing channel importance. Given ASMN-extracted features $\mathbb{F}_t^S \in \mathbb{R}^{H \times W \times C}$ and voxelized LiDAR features

$\mathbb{F}_t^L \in \mathbb{R}^{H \times W \times C}$, KGF correlates pixel features $(\tau, \epsilon)$ with LiDAR points $(\xi, \gamma)$ using cosine distance, factoring in neighboring points. The minimum cosine distance $\mathbb{V}$ between corresponding features is calculated as:

$$\text{Cosine}(\mathbb{F}_t^S, \mathbb{F}_t^L)(\tau, \epsilon) = \min \left[ \frac{\mathbb{F}_t^S \cdot \mathbb{F}_t^L}{\sqrt{(\mathbb{F}_t^S)^2 + (\mathbb{F}_t^L)^2}} \right]. \quad (6)$$

The final KGF value accumulates weighted sums across channels and pseudo code is given

$$KGF(\tau, \epsilon, \kappa) = \sum_{\kappa=1}^{C} 2^{-\kappa} \cdot \mathbb{V}(\tau, \epsilon, \kappa) + \mathbb{F}_t^S(\tau, \epsilon, \kappa). \quad (7)$$

Algorithm 2 for the pseudo-code of the KGF module.

---

**Algorithm 2** KGF

---

**Input:** $\mathbb{F}_t^L, \mathbb{F}_t^S$ : (H, W, C)
**Output:** Projected features
**def cosine_dist**$(a, b)$:
    **return** $\frac{a \cdot b}{\sqrt{a^2 + b^2}}$
**def project**$(\mathbb{F}_t^L, \mathbb{F}_t^S, \tau, \epsilon)$:
  $H\_range, W\_range, C = \text{shape}(\mathbb{F}_t^L)$
  $count = 0$
  **for** $\kappa$ **in** $C$:
    $\text{val\_nei} = [\text{val} \mid \text{val} \in \text{neighbors}]$
    $\text{min\_dist} = \min(\textbf{cosine\_dist}(\mathbb{F}_t^L, \mathbb{F}_t^S[\tau][\epsilon]) \mid (\tau, \epsilon) \in$
  $\text{val\_nei})$
    $count{+}= (2^{-\kappa}) \times \text{min\_dist}$
  **return** $count$
**def KGF**$(\mathbb{F}_t^L, \mathbb{F}_t^S)$:
  $output = \text{zeros\_like}(\mathbb{F}_t^S)$
  **for** $\kappa$ **in** $C$:
    **for** $\tau$ **in** $H$, $\epsilon$ **in** $W$:
      $project\_value = \textbf{project}(\mathbb{F}_t^L, \mathbb{F}_t^S[\kappa], \tau, \epsilon)$
      $output[\kappa][\tau][\epsilon] = \mathbb{F}_t^S + project\_value$
  **return** $output$

---

### C. Loss function

SaViD uses Voxel R-CNN [32] for RPN and RoI loss, in addition to using Fusion Loss [31] and LSTM loss [33].

## IV. EXPERIMENTS

### A. Dataset details

Our goal is to conduct robust, long-range experiments with multi-modal fusion using the Argoverse2 (AV2) [24] and Waymo Open Dataset (WOD).

**AV2:** includes 1000 sequences: 700 for training, 150 for validation, and 150 for testing. It has a perception range of 200 meters and covers a 400m × 400m area, making it more extensive than other standard benchmarks like Waymo [9] and nuScenes [10]. AV2 features 30 object categories with a long-tail distribution; we focus on the top 20 classes, excluding the 10 tail classes. SaViD is tested in two scenarios:

AV2-C for clean data and AV2-R (AV2-Robust) for corrupted sensor data $(\mathcal{L}_t, \mathcal{I}_t)$.

**WOD** is the leading benchmark for LiDAR-based 3D object detection, known for its large and complex dataset of 1,150 sequences with over 200,000 frames, including LiDAR points, camera images, and 3D bounding boxes. The dataset is split into 798 training, 202 validation, and 150 testing sequences. The clean WOD-C detection range is 75 meters, covering a 150m x 150m area. Our evaluation focuses on long-range performance using LEVEL_2 (L2) difficulty, excluding LEVEL_1 (L1) for small-range detection. We also introduce WOD-R for corrupted sensors $(\mathcal{L}_t, \mathcal{I}_t)$, similar to AV2-R.

### B. Natural Robustness

Natural robustness addresses real-world corruptions in autonomous driving, categorized into weather-induced and sensor-induced corruptions. We identify 14 common corruptions relevant to AV2 and WOD for long-range detection [16].

**Weather-Induced Corruptions.** These include Snow, Rain, Fog, and Sunlight, significantly affecting LiDAR and camera data. Weather effects are simulated on LiDAR using physics-based methods [14], [42], [43] and visually augmented for cameras [15].

**Sensor-Induced Corruptions.** We introduce 10 sensor-level corruptions: seven for LiDAR (e.g., Density Decrease, Cutout, LiDAR Crosstalk, FOV Lost, various noise types) and three for images (Gaussian, Uniform, and Impulse Noise) to simulate visual disturbances from lighting or camera faults [15].

### C. Implementation Details

**Network Architecture.** SaViD employs a three-stage strategy to integrate features from LiDAR and image modalities, using pseudo feature extraction. For LiDAR, SaViD builds on the Voxel-RCNN framework [44] with dynamic voxelization and feature dimensions of 16, 32, 64, and 64 to manage sparse point clouds. The image stream feature extractor relies on depth information $(\mathcal{D}_t)$ from a pretrained Twise network [45]. In the GMA module, a dropout rate of 30% is applied to the attention affinity matrix during training, with parameters: dimension = 64, $N_h = 8$, and $P = 7$. The MLP layer after GMA is a fully connected layer with 64 filters.

**Training and Inference Details.** SaViD is trained from scratch using the ADAM optimizer on 32 GTX 1080 Tesla T4 GPUs with a cosine annealing learning rate. The proposal refinement stage samples 128 proposals, maintaining a 1:1 ratio of positive (IoU $\geq 0.55$) to negative proposals for enhanced long-range detection. Data augmentation techniques are employed during training [7], [26]. For inference, non-maximum suppression (NMS) is used in the RPN with IoU thresholds of 0.7 and 0.1 to filter redundant predictions [44].

### D. Performance on AV2-C and WOD-C

We evaluated single and multi-modal fusion methods on AV2-C (Table I). Initial results using CenterPoint improved

| Methods | Vehicle | Bus | Pedestrian | Stop Sign | Box Truck | Bollard | C-Barrel | Motorcyclist | MPC-Sign | Motorcycle | Bicycle | A-Bus | School Bus | Truck Cab | C-Cone | V-Trailer | Sign | Large Vehicle | Stroller | Bicyclist | AP |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| *Precision* | | | | | | | | | | | | | | | | | | | | | |
| CenterPoint | 61.0 | 36.0 | 33.0 | 28.0 | 26.0 | 25.0 | 22.5 | 16.0 | 16.0 | 12.5 | 9.5 | 8.5 | 7.5 | 8.0 | 8.0 | 7.0 | 6.5 | 3.0 | 2.0 | 14 | 17.5 |
| CenterPoint+ | 67.6 | 38.9 | 46.5 | 16.9 | 37.4 | 40.1 | 32.2 | 28.6 | 27.4 | 33.4 | 24.5 | 8.7 | 25.8 | 22.6 | 29.5 | 22.4 | 6.3 | 3.9 | 0.5 | 20.1 | 26.67 |
| FSD | 67.1 | 39.8 | 57.4 | 21.3 | 38.3 | 38.3 | 38.1 | 30.0 | 23.6 | 38.1 | 25.5 | 15.6 | 30.0 | 20.1 | 38.9 | 23.9 | 7.9 | 5.1 | 5.7 | 27.0 | 29.58 |
| BEVFusion # | 67.2 | 39.8 | 58.1 | 31.9 | 36.3 | 35.2 | 36.7 | 34.1 | 26.1 | 46.8 | 33.6 | 21.2 | 22.2 | 16.9 | 31.2 | 22.8 | 13.2 | 5.4 | 9.6 | 32.6 | 31.05 |
| TransFusion # | 67.6 | 40.5 | 58.4 | 32.6 | 38.5 | 36.1 | 38.6 | 34.3 | 26.8 | 48.3 | 37.3 | 21.7 | 22.9 | 18.5 | 33.8 | 23.2 | 13.5 | 6.2 | 9.8 | 33.1 | 32.09 |
| DeepFusion # | 70.7 | 42.3 | 62.1 | 32.8 | 40.8 | 40.0 | 42.2 | 42.6 | 28.3 | 50.1 | 40.1 | 21.7 | 29.7 | 17.6 | 40.2 | 25.3 | 14.7 | 7.9 | 10.7 | 35.1 | 34.74 |
| SaViD (t=1) | 78.2 | 48.6 | 67.8 | 38.6 | 40.7 | 42.8 | 45.3 | 42.4 | 30.8 | 53.2 | 41.5 | 25.9 | 30.9 | 22.6 | 41.3 | 30.9 | 19.6 | 10.8 | 12.8 | 38.8 | 38.17 (+3.43) |
| SaViD (t=7) | 79.7 | 49.5 | 68.7 | 40.1 | 41.9 | 43.8 | 47.2 | 44.1 | 33.3 | 55.4 | 42.1 | 25.9 | 32.3 | 25.1 | 44.9 | 31.6 | 20.3 | 12.9 | 13.7 | 40.5 | 39.65(+1.48) |

TABLE I: The table shows AP results on AV2-C validation for categories like C-Barrel, MPC-Sign, A-Bus, C-Cone, and V-Trailer. Bolded AP values indicate significant gains over single-frame ($t = 1$) performance. #: Simulated in the same environments.
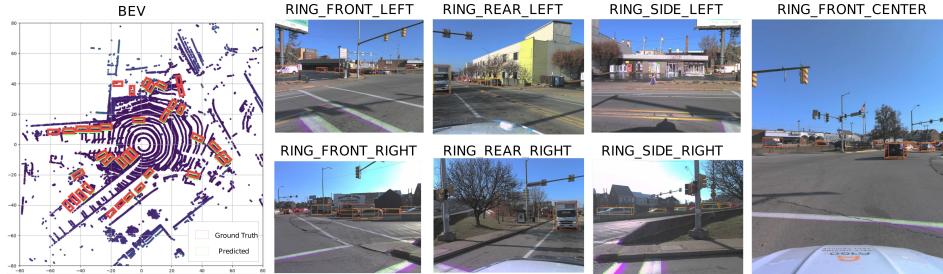


Fig. 5: A qualitative comparison of long-range 3D object detection using multi-modal fusion methods on AV2-C validation set. BEV maps on left, 2D image space on right. Red: Ground truth, Green: Predicted boxes.
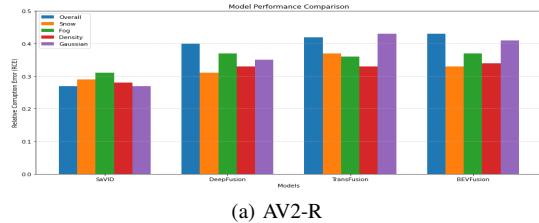
TABLE II: Comparison of Model Performance for 3D Detection on the WOD Test Set. In the table, 'L' and 'I' denote LiDAR and camera sensors, respectively. 'TTA' and 'Ens' represent test-time augmentation and ensemble model outputs, indicated by #

| Method | Modality | ALL (mAPH) | VEH (APH) | PED (APH) | CYC (APH) |
|---|---|---|---|---|---|
| SaViD (t=7) (Ours) | L+I | 82.96 (+1.94) | 82.94 | 84.15 | 81.78 |
| SaViD (t=1) (Ours) | L+I | 82.16 | 82.37 | 83.51 | 80.59 |
| LoGoNet Ens# [34] | L+I | 81.02 | 81.72 | 81.28 | 80.06 |
| BEVFusion TTA# [29] | L+I | 79.97 | 80.92 | 79.65 | 79.33 |
| LidarMultiNet TTA# [35] | L | 79.94 | 80.36 | 79.86 | 79.59 |
| MPPNet Ens# [36] | L | 79.60 | 80.93 | 80.14 | 77.73 |
| MT-Net Ens# [37] | L | 78.45 | 80.11 | 78.08 | 77.17 |
| DeepFusion Ens# [9] | L+I | 78.41 | 79.09 | 78.57 | 77.58 |
| AFDetV2 Ens# [38] | L | 77.64 | 78.34 | 76.75 | 77.83 |
| INT Ens# [39] | L | 77.21 | 78.73 | 76.36 | 76.54 |
| HorizonLiDAR3D Ens# [40] | L+I | 77.11 | 77.83 | 76.50 | 76.98 |
| LoGoNet [34] | L+I | 77.10 | 79.30 | 78.91 | 73.10 |
| BEVFusion [29] | L+I | 76.33 | 77.48 | 76.41 | 75.09 |
| CenterFormer [41] | L | 76.29 | 78.28 | 77.42 | 73.17 |
| MPPNet [36] | L | 75.67 | 76.91 | 75.93 | 74.18 |
| DeepFusion [9] | L+I | 75.54 | 75.69 | 76.40 | 74.51 |

by 52.4% in AP with the modified CenterPoint+. The FSD method, with its sparse attention mechanism, enhanced performance by 10.91% over CenterPoint+. Two-modality 3D detection outperformed single-modality approaches, consistent with AV2-C results. BEVFusion improved AP by 4.73% over FSD, and TransFusion surpassed BEVFusion by 3.34%. DeepFusion, with cross-former feature alignment, further improved AP by 7.62% over TransFusion. SaViD, using three-stage feature fusion, achieved the highest AP of 38.17, a 9.87% gain over DeepFusion. Adding temporal alignment in SaViD increased AP to 39.65, up 3.73% from the single-frame model. Figure 5 highlights SaViD's qualitative performance using BEV maps and front camera views.

In Table II, we compare model performance for 3D detection on the WOD test set for clean data. SaViD (t=7) achieved the highest mAPH of 82.96, improving L2 difficulties by 1.94 points. It excelled across all classes, with APH scores of 82.94 for vehicles, 84.15 for pedestrians, and 81.78 for cyclists. The single-frame SaViD (t=1) also outperformed LoGoNet Ens, with improvements of 1.14 mAPH, 1.38 APH for vehicles, 1.34 APH for pedestrians, and 1.56 APH for cyclists. Compared to BEVFusion [29], SaViD (t=7) showed a 2.99 mAPH increase, with specific gains of 2.02 for vehicles, 4.50 for pedestrians, and 2.45 for cyclists. The single-frame SaViD (t=1) also surpassed BEVFusion with a 2.19 mAPH boost. SaViD (t=7) outperformed LidarMultiNet TTA [35] by 3.02 mAPH, and SaViD (t=1) achieved a 2.22 mAPH increase compared to LidarMultiNet TTA. These results underscore SaViD (t=7)'s superior performance in 3D object detection across various categories, establishing its effectiveness compared to other methods, especially when leveraging sequential frame information.
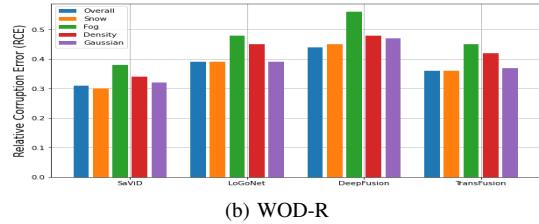
(a) AV2-R



(b) WOD-R

Fig. 6: The RCE depicts the overall results under all levels of corruption as well as the outcomes under each level of corruption for AV2-R and WOD-R dataset

### E. Performance on AV2-R and WOD-R

Robustness performance is assessed by measuring individual and relative corruption effects at various severity levels. Fusion-based methods' clean performance on AV2-C is denoted as $AP_{cln}$, while corrupted performances are denoted $AP_{r,s}$ for each corruption type (r) and severity level (s) [16]. The average corruption performance is given by:

$$AP_{corr} = \frac{1}{|\nu|} \sum_{r \in \nu} \frac{1}{5} \sum_{s=1}^{5} AP_{r,s}, \qquad (8)$$

where $\nu$ represents the set of corruptions. The Relative Corruption Error (RCE) quantifies robustness by evaluating performance degradation under clean conditions:

$$RCE = \left[\frac{AP_{cln} - AP_{corr}}{AP_{cln}}\right], \quad RCE_{r,s} = \left[\frac{AP_{cln} - AP_{r,s}}{AP_{cln}}\right]. \qquad (9)$$

Figure 6 shows individual and overall RCE for AV2-R and WOD-R datasets. SaViD demonstrates the most robust performance with a 27.24% overall RCE on AV2-R, a 31.43% improvement over DeepFusion. BEVFusion shows the lowest robustness with a 48.82% RCE decline. On WOD-R, SaViD's RCE drop is 31.24%, compared to 36.36% for LoGoNet, and 44.15% for Transfusion. SaViD's superior robustness is due to its 3-stage fusion alignment mechanism.

### V. ABLATION STUDY ON AV2-C

**Necessities of three-stage Fusions.** We conducted extensive experiments, as shown in Table III, to assess the impact of the three components on SaViD's performance: (1) Missing ASMN: while considering the only local-global image feature extractor GMAN and KGF, there is a 7.2% decrease in AP value compared to considering all the components. The observed performance drop can be attributed to the limitations of the sparse LiDAR feature extractor($\mathbb{F}_t^L$).

| GMAN | ASMN | KGF | AP |
|:---:|:---:|:---:|:---:|
| ✓ | | ✓ | 35.4 |
| | ✓ | ✓ | 24.8 |
| ✓ | ✓ | ✓ | 38.17 |

TABLE III: Performance of Each Component in SaViD (t=1)

Despite estimating depth ($\mathcal{D}_t$) using LiDAR on Image, it fails to ensure optimal fusion alignment, leading to sub-optimal results in the fusion process. (2) Missing GMAN: while excluding the image feature extractor results in a 35% decrease in the AP value compared to using all the components. Despite providing projected image data($\mathcal{I}_t$) to align with the channels of the LiDAR feature extractor($\mathbb{F}_t^L$), it practically lacks information due to the sparsity of the LiDAR data. Consequently, it behaves like a conventional voxelized LiDAR-based detector.

**Necessities of $\mathcal{D}_t$ in SaViD.** In our experiment, we examine two variations of vision transformers, namely SwinV2 [28] and GCVIT [27]. We integrate them with the two proposed fusion stages, ASMN and KGF, while excluding the depth information ($\mathcal{D}_t$). Table IV presents the results. When using SwinV2 to extract $\mathbb{F}_t^I$, the obtained AP value is 33.4. However, when considering fused conv2D in GCVIT, the performance improved to 34.9 in AP. Nevertheless, both cases exhibit sub-optimal performance due to the lack of alignment between spatial feature extraction and sparse LiDAR information. Despite the inclusion of other two proposed alignment methods, the absence of global context LiDAR information (disparity with image) leads to a notable decline in long-range detection performance. Therefore, this ablation study highlights the importance of using $\mathcal{D}_t$ as a global query to minimize the disparity with $\mathcal{I}_t$ and improve performance.

| Model | AP |
|:---|:---:|
| SwinV2\$\mathcal{D}_t$ + ASMN + KGF | 33.4 |
| GCVIT\$\mathcal{D}_t$ + ASMN + KGF | 34.9 |
| GMAN + ASMN + KGF (t=1) | 38.17 |

TABLE IV: Performance Without $\mathcal{D}_t$ on Varying Vision Transformers

### VI. CONCLUSION

This paper introduces SaViD, a novel three-stage robust fusion alignment method incorporating a local-global perspective for 3D object detection. The first stage uses a vision transformer-based GMAN to extract image features, considering local and global depth information. It then introduces a ASMN to align sparse LiDAR features with extracted image features. Lastly, a parameter-free KGF achieves final fusion. SaViD achieves notable performance gains on AV2-C, shows resilience to corruptions on AV2-R, and excels on WOD, especially in L2 difficulties. With its long-range detection capability, SaViD's potential extends beyond AVs to Digital Airport Tower Control, enhancing operational efficacy and safety in complex airport environments.

## References

[1] E. Arnold, O. Y. Al-Jarrah, M. Dianati, S. Fallah, D. Oxtoby, and A. Mouzakitis, "A survey on 3d object detection methods for autonomous driving applications," *IEEE Transactions on Intelligent Transportation Systems*, vol. 20, no. 10, pp. 3782–3795, 2019.

[2] J. Xiao, P. Pisutsin, and M. Feroskhan, "Collaborative target search with a visual drone swarm: An adaptive curriculum embedded multi-stage reinforcement learning approach," *IEEE Transactions on Neural Networks and Learning Systems*, 2023.

[3] A. K. Kamath, S. G. Anavatti, and M. Feroskhan, "A physics-informed neural network approach to augmented dynamics visual servoing of multirotors," *IEEE Transactions on Cybernetics*, 2024.

[4] J. Xiao, J. H. Chee, and M. Feroskhan, "Real-time multi-drone detection and tracking for pursuit-evasion with parameter search," *IEEE Transactions on Intelligent Vehicles*, 2024.

[5] L. Chen, J. Xiao, Y. Zheng, N. A. Alagappan, and M. Feroskhan, "Design, modeling, and control of a coaxial drone," *IEEE Transactions on Robotics*, vol. 40, pp. 1650–1663, 2024.

[6] T. Dam, "Developing generative adversarial networks for classification and clustering: Overcoming class imbalance and catastrophic forgetting," Ph.D. dissertation, University of New South Wales (Australia), 2022.

[7] T. Dam, S. B. Dharavath, S. Alam, N. Lilith, S. Chakraborty, and M. Feroskhan, "Aydiv: Adaptable yielding 3d object detection via integrated contextual vision transformer," in *2024 IEEE International Conference on Robotics and Automation (ICRA)*, 2024, pp. 10 657–10 664.

[8] S. Vora, A. H. Lang, B. Helou, and O. Beijbom, "Pointpainting: Sequential fusion for 3d object detection," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 4604–4612.

[9] Y. Li, A. W. Yu, T. Meng, B. Caine, J. Ngiam, D. Peng, J. Shen, Y. Lu, D. Zhou, Q. V. Le *et al.*, "Deepfusion: Lidar-camera deep fusion for multi-modal 3d object detection," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 17 182–17 191.

[10] X. Bai, Z. Hu, X. Zhu, Q. Huang, Y. Chen, H. Fu, and C.-L. 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*, 2022, pp. 1090–1099.

[11] T. Liang, H. Xie, K. Yu, Z. Xia, Z. Lin, Y. Wang, T. Tang, B. Wang, and Z. Tang, "Bevfusion: A simple and robust lidar-camera fusion framework," 2022.

[12] J. Ku, M. Mozifian, J. Lee, A. Harakeh, and S. L. Waslander, "Joint 3d proposal generation and object detection from view aggregation," in *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2018, pp. 1–8.

[13] T. Yin, X. Zhou, and P. Krahenbuhl, "Center-based 3d object detection and tracking," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 11 784–11 793.

[14] V. Kilic, D. Hegde, V. Sindagi, A. B. Cooper, M. A. Foster, and V. M. Patel, "Lidar light scattering augmentation (lisa): Physics-based simulation of adverse weather conditions for 3d object detection," *arXiv preprint arXiv:2107.07004*, 2021.

[15] D. Hendrycks and T. Dietterich, "Benchmarking neural network robustness to common corruptions and perturbations," *arXiv preprint arXiv:1903.12261*, 2019.

[16] Y. Dong, C. Kang, J. Zhang, Z. Zhu, Y. Wang, X. Yang, H. Su, X. Wei, and J. Zhu, "Benchmarking robustness of 3d object detection to common corruptions," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 1022–1032.

[17] Z. Yang, Y. Sun, S. Liu, X. Shen, and J. Jia, "Ipod: Intensive point-based object detector for point cloud," 2018.

[18] A. H. Lang, S. Vora, H. Caesar, L. Zhou, J. Yang, and O. Beijbom, "Pointpillars: Fast encoders for object detection from point clouds," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 12 697–12 705.

[19] Y. Zhou and O. Tuzel, "Voxelnet: End-to-end learning for point cloud based 3d object detection," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 4490–4499.

[20] C. R. Qi, H. Su, K. Mo, and L. J. Guibas, "Pointnet: Deep learning on point sets for 3d classification and segmentation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 652–660.

[21] C. R. Qi, L. Yi, H. Su, and L. J. Guibas, "Pointnet++: Deep hierarchical feature learning on point sets in a metric space," *Advances in neural information processing systems*, vol. 30, 2017.

[22] A. Bewley, P. Sun, T. Mensink, D. Anguelov, and C. Sminchisescu, "Range conditioned dilated convolutions for scale invariant 3d object detection," *arXiv preprint arXiv:2005.09927*, 2020.

[23] S. B. Dharavath, T. Dam, S. Chakraborty, P. Roy, and A. Maiti, "Quantum inverse contextual vision transformers (q-icvt): A new frontier in 3d object detection for avs," in *Proceedings of the 33rd ACM International Conference on Information and Knowledge Management*, 2024, pp. 3724–3729.

[24] L. Fan, F. Wang, N. Wang, and Z. Zhang, "Fully sparse 3d object detection," 2022.

[25] X. Chen, H. Ma, J. Wan, B. Li, and T. Xia, "Multi-view 3d object detection network for autonomous driving," 2017.

[26] S. Shi, L. Jiang, J. Deng, Z. Wang, C. Guo, J. Shi, X. Wang, and H. Li, "Pv-rcnn++: Point-voxel feature set abstraction with local vector representation for 3d object detection," *International Journal of Computer Vision*, vol. 131, no. 2, pp. 531–551, 2023.

[27] A. Hatamizadeh, H. Yin, J. Kautz, and P. Molchanov, "Global context vision transformers," *arXiv preprint arXiv:2206.09959*, 2022.

[28] Z. Liu, H. Hu, Y. Lin, Z. Yao, Z. Xie, Y. Wei, J. Ning, Y. Cao, Z. Zhang, L. Dong *et al.*, "Swin transformer v2: Scaling up capacity and resolution," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp. 12 009–12 019.

[29] Z. Liu, H. Tang, A. Amini, X. Yang, H. Mao, D. L. Rus, and S. Han, "Bevfusion: Multi-task multi-sensor fusion with unified bird's-eye view representation," in *2023 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2023, pp. 2774–2781.

[30] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.

[31] B. Zhang, I. Titov, and R. Sennrich, "Sparse attention with linear units," *arXiv preprint arXiv:2104.07012*, 2021.

[32] J. Deng, S. Shi, P. Li, W. Zhou, Y. Zhang, and H. Li, "Voxel r-cnn: Towards high performance voxel-based 3d object detection," 2021.

[33] L. Fan, Y. Yang, F. Wang, N. Wang, and Z. Zhang, "Super sparse 3d object detection," 2023.

[34] X. Li, T. Ma, Y. Hou, B. Shi, Y. Yang, Y. Liu, X. Wu, Q. Chen, Y. Li, Y. Qiao *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*, 2023, pp. 17 524–17 534.

[35] D. Ye, Z. Zhou, W. Chen, Y. Xie, Y. Wang, P. Wang, and H. Foroosh, "Lidarmultinet: Towards a unified multi-task network for lidar perception," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 37, no. 3, 2023, pp. 3231–3240.

[36] X. Chen, S. Shi, B. Zhu, K. C. Cheung, H. Xu, and H. Li, "Mppnet: Multi-frame feature intertwining with proxy points for 3d temporal object detection," in *European Conference on Computer Vision*. Springer, 2022, pp. 680–697.

[37] S. Chen, Z. Jie, X. Wei, and L. Ma, "Mt-net submission to the waymo 3d detection leaderboard," *arXiv preprint arXiv:2207.04781*, 2022.

[38] Y. Hu, Z. Ding, R. Ge, W. Shao, L. Huang, K. Li, and Q. Liu, "Afdetv2: Rethinking the necessity of the second stage for object detection from point clouds," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 36, no. 1, 2022, pp. 969–979.

[39] J. Xu, Z. Miao, D. Zhang, H. Pan, K. Liu, P. Hao, J. Zhu, Z. Sun, H. Li, and X. Zhan, "Int: Towards infinite-frames 3d detection with an efficient framework," in *European Conference on Computer Vision*. Springer, 2022, pp. 193–209.

[40] Z. Ding, Y. Hu, R. Ge, L. Huang, S. Chen, Y. Wang, and J. Liao, "1st place solution for waymo open dataset challenge–3d detection and domain adaptation," *arXiv preprint arXiv:2006.15505*, 2020.

[41] Z. Zhou, X. Zhao, Y. Wang, P. Wang, and H. Foroosh, "Centerformer: Center-based transformer for 3d object detection," in *European Conference on Computer Vision*. Springer, 2022, pp. 496–513.

[42] M. Hahner, C. Sakaridis, D. Dai, and L. Van Gool, "Fog simulation on real lidar point clouds for 3d object detection in adverse weather," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 15 283–15 292.

[43] M. Hahner, C. Sakaridis, M. Bijelic, F. Heide, F. Yu, D. Dai, and L. Van Gool, "Lidar snowfall simulation for robust 3d object detection," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 16 364–16 374.

[44] S. Shi, C. Guo, L. Jiang, Z. Wang, J. Shi, X. Wang, and H. 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*, 2020, pp. 10 529–10 538.

[45] S. Imran, X. Liu, and D. Morris, "Depth completion with twin surface extrapolation at occlusion boundaries," 2021.