

Article

# Dual-Camera LiDAR Fusion for Occlusion-Robust 3D Detection in Urban Driving Simulation

Xingnan Zhou <sup>1</sup> and Ciprian Alecsandru <sup>1,\*</sup>

<sup>1</sup> Department of Building, Civil and Environmental Engineering, Concordia University, Montreal, QC H3G 1M8, Canada

\* Correspondence: ciprian.alecsandru@concordia.ca (C.A.)

Version February 10, 2026 submitted to Sustainability

**Abstract:** Three-dimensional object detection from LiDAR point clouds is a cornerstone of autonomous driving perception, yet single-sensor systems remain vulnerable to occlusion in complex urban environments such as roundabouts and dense intersections. This paper proposes an asymmetric dual-camera LiDAR fusion framework that combines a PointPillar-based 3D LiDAR detector with YOLOv8-based 2D detections from two complementary camera viewpoints: a drone (top-down, 40 m altitude) and a subject-vehicle forward camera. The fusion operates at the decision level (late fusion), where camera-confirmed LiDAR detections receive confidence boosts— $\times 1.15$  for single-camera confirmation and  $\times 1.30$  for dual-camera agreement—while unconfirmed low-confidence detections within the drone’s field of view are suppressed ( $\times 0.75$ ). Crucially, the design is asymmetric: the drone camera applies both boost and suppress operations, whereas the limited-field-of-view forward camera applies boost-only to avoid penalizing legitimate detections outside its narrow coverage. Evaluated on a CARLA Town10HD dataset comprising 650 frames and 12,308 annotations across Car and Pedestrian classes, with five-seed repeated random sub-sampling validation (all at epoch 80), the drone-fused system improves mAP@0.5 by +1.48 percentage points (+5.8% relative; sign test  $p = 0.031$ ) and the full dual-camera fusion achieves +2.88 pp (+11.2% relative; sign test  $p = 0.031$ ), with all five seeds showing positive improvement in every configuration. Critically, domain-adapted YOLOv8 models fine-tuned on CARLA-rendered images are essential—the forward camera contributes a meaningful +0.62 pp mAP improvement (sign test  $p = 0.031$ ) only after domain adaptation. While baseline performance is constrained by the small training set (520 frames per seed), the consistent directional improvements across all seeds and configurations demonstrate the potential of asymmetric multi-viewpoint fusion for occlusion-robust perception.

**Keywords:** 3D object detection; LiDAR-camera fusion; intelligent transport system; late fusion; drone-assisted perception; sustainable traffic management; PointPillars; YOLOv8; CARLA simulation; autonomous driving safety

## 1. Introduction

Reliable three-dimensional (3D) object detection is a fundamental requirement for safe autonomous driving. LiDAR-based detectors have become the dominant paradigm for 3D perception in self-driving systems, offering precise depth measurements and geometric representations of the driving environment [1,2]. However, single-viewpoint LiDAR systems suffer from well-documented limitations: occlusion by foreground objects, sparse point density at long range, and blind spots caused by the sensor’s mounting position [3]. These limitations are particularly acute in complex urban geometries such as roundabouts, dense intersections, and narrow streets, where vehicles, pedestrians, and infrastructure elements frequently occlude one another from the ego vehicle’s perspective.

34 Camera-LiDAR fusion has emerged as a promising strategy to mitigate single-sensor limitations  
35 by combining the geometric precision of LiDAR with the rich semantic and texture information  
36 provided by cameras [4,5]. Early fusion methods such as PointPainting [6] and Frustum PointNets [7]  
37 augment point clouds with camera features at the data level, while deep fusion approaches like  
38 BEVFusion [8,9] and TransFusion [10] learn joint representations in a unified bird's-eye view (BEV)  
39 space. Although these methods have achieved impressive results on benchmarks such as nuScenes [11]  
40 and KITTI [12], they predominantly rely on cameras co-located with the ego vehicle, inheriting the  
41 same viewpoint limitations that constrain the LiDAR sensor.

42 A fundamentally different approach to overcoming single-viewpoint occlusion is to  
43 introduce cameras at *complementary* vantage points. Drone-assisted perception [13,14] and  
44 vehicle-to-infrastructure (V2I) cooperative systems [15,16] have demonstrated that elevated or offset  
45 viewpoints can observe objects hidden from the street-level perspective. However, most cooperative  
46 perception research has focused on sharing intermediate features between multiple LiDAR-equipped  
47 agents [17–19], which requires expensive sensor suites on all cooperating platforms. In contrast,  
48 cameras are lightweight, inexpensive, and easily deployable on drones or infrastructure poles, making  
49 camera-based viewpoint augmentation a practical and cost-effective alternative to full multi-LiDAR  
50 cooperative perception.

51 This paper proposes an asymmetric dual-camera LiDAR fusion framework that addresses  
52 occlusion in urban driving by combining a vehicle-mounted LiDAR 3D detector with 2D object  
53 detections from two complementary cameras: (1) a drone camera providing a top-down perspective  
54 at 40 m altitude, and (2) the subject vehicle's forward-facing camera. The key insight motivating the  
55 asymmetric design is that these two cameras have fundamentally different coverage characteristics.  
56 The drone camera provides a wide, overhead field of view (FOV) that can observe nearly the entire  
57 scene, including objects fully occluded from street level. The forward camera, by contrast, covers only a  
58 narrow frontal sector but provides high-resolution appearance information for objects directly ahead of  
59 the vehicle. Treating these cameras identically in the fusion pipeline would be suboptimal: the drone's  
60 comprehensive coverage justifies both boosting confirmed detections *and* suppressing unconfirmed  
61 ones, while the forward camera's limited FOV means that the absence of a camera detection does not  
62 reliably indicate the absence of an object.

63 The fusion operates at the decision level (late fusion), modifying the confidence scores of the 3D  
64 LiDAR detector's output based on spatial agreement with 2D camera detections projected into the  
65 BEV plane. Camera-confirmed LiDAR detections receive a confidence boost of  $\times 1.15$  for single-camera  
66 confirmation or  $\times 1.30$  for dual-camera agreement. Unconfirmed low-confidence Car detections within  
67 the drone's FOV are suppressed by a factor of  $\times 0.75$ , while no suppression is applied for detections  
68 outside the drone's coverage or within the forward camera's limited FOV alone. This late-fusion  
69 approach has several practical advantages: it is modular (any 3D detector and any 2D detector  
70 can be used), computationally lightweight (no retraining required), and transparent (the confidence  
71 adjustments are interpretable).

72 From a sustainable transportation perspective, improved 3D perception directly supports safer  
73 autonomous driving—a critical enabler of sustainable urban mobility. Traffic accidents remain a leading  
74 cause of preventable death [20], and occlusion-related collisions at intersections and roundabouts  
75 disproportionately affect vulnerable road users such as pedestrians and cyclists. By leveraging  
76 cost-effective drone cameras rather than expensive multi-LiDAR infrastructure, the proposed approach  
77 offers a scalable pathway to enhanced perception that aligns with sustainable smart-city transportation  
78 strategies [21].

79 We evaluate the proposed framework using the CARLA driving simulator [22], which provides  
80 precise ground-truth annotations and full control over sensor placement. The dataset comprises 650  
81 frames collected in Town10HD with 12,308 object annotations spanning Car and Pedestrian classes.  
82 To ensure statistical rigor, we conduct five-seed repeated random sub-sampling validation with all

83 models trained to epoch 80, and report both parametric (paired *t*-test) and non-parametric (sign test)  
84 significance measures. Our contributions are as follows:

- 85 1. We propose an *asymmetric* late-fusion framework that differentiates between cameras based on  
86 their coverage characteristics, applying boost-and-suppress operations for the wide-FOV drone  
87 camera and boost-only operations for the narrow-FOV forward camera.
- 88 2. We demonstrate that drone-camera fusion improves mAP@0.5 by +1.48 percentage points (+5.8%  
89 relative) over the LiDAR-only baseline (sign test  $p = 0.031$ ), with the full dual-camera system  
90 achieving +2.88 pp (+11.2% relative). All three fusion configurations achieve 5/5 positive seeds  
91 ( $p = 0.031$ ), with domain-adapted YOLOv8 detectors proving essential for effective fusion.
- 92 3. We conduct five-seed repeated random sub-sampling validation with both parametric (*t*-test) and  
93 non-parametric (sign test) significance measures, providing transparent reporting of statistical  
94 power and its limitations in small-sample detection experiments.
- 95 4. We provide the complete data-collection and fusion pipeline built on CARLA and OpenPCDet [23],  
96 with code to be released upon publication, enabling reproducible research on multi-viewpoint  
97 fusion for autonomous driving perception.

98 The remainder of this paper is organized as follows. Section 2 reviews related work on  
99 LiDAR-based 3D detection, camera-LiDAR fusion, and drone-assisted perception. Section 3 describes  
100 the proposed asymmetric fusion methodology. Section 4 details the experimental setup, including the  
101 CARLA data-collection pipeline and training configuration. Section 5 presents quantitative results  
102 with statistical analysis. Section 6 discusses findings, limitations, and practical implications. Section 7  
103 concludes the paper.

## 104 2. Related Work

### 105 2.1. LiDAR-Based 3D Object Detection

106 LiDAR-based 3D object detection has progressed through several architectural paradigms.  
107 Point-based methods such as PointNet [24] and PointNet++ [25] operate directly on raw point  
108 clouds, learning per-point features through shared multi-layer perceptrons and set abstraction layers.  
109 PointRCNN [26] extended this paradigm to two-stage 3D detection by generating proposals from  
110 point-level features. While point-based methods preserve fine geometric detail, their computational  
111 cost scales unfavorably with point cloud density.

112 Voxel-based methods discretize the point cloud into a regular 3D grid and apply sparse 3D  
113 convolutions. VoxelNet [27] pioneered end-to-end voxel-based detection, and SECOND [28] introduced  
114 spatially sparse convolutions that dramatically reduced computation by only processing occupied  
115 voxels. CenterPoint [29] further advanced voxel-based detection with center-heatmap prediction and  
116 a two-stage refinement module. These methods achieve strong accuracy but require careful voxel  
117 resolution tuning to balance precision and efficiency.

118 Pillar-based methods, led by PointPillars [30], offer a compelling trade-off between speed and  
119 accuracy by collapsing the vertical dimension into a single “pillar” per horizontal grid cell. The  
120 resulting 2D pseudo-image can be processed by standard 2D convolutional backbones and detection  
121 heads, enabling real-time inference. PointPillars remains a widely used baseline in both academic  
122 benchmarks and practical deployments due to its simplicity, speed, and competitive accuracy. Wang et  
123 al. [31] further explored pillar-based architectures with improved feature encoding. Hybrid approaches  
124 such as PV-RCNN [32] combine voxel and point-based processing for state-of-the-art accuracy at the  
125 cost of increased complexity.

126 In this work, we adopt PointPillars as the LiDAR 3D detector due to its favorable speed–accuracy  
127 trade-off and its suitability as a baseline for evaluating fusion strategies. The modular nature of  
128 our late-fusion approach means that the 3D detector can be replaced with any alternative without  
129 modifying the fusion logic.

130 2.2. *Camera-LiDAR Fusion for 3D Detection*

131 Camera-LiDAR fusion methods can be broadly categorized by the stage at which sensor  
132 information is combined: early (data-level), deep (feature-level), and late (decision-level) fusion [4,5].

133 2.2.1. Early Fusion

134 Early fusion methods augment the LiDAR point cloud with camera-derived features before  
135 detection. PointPainting [6] projects LiDAR points onto the camera image and appends per-point  
136 semantic segmentation scores to the point features, enabling the 3D detector to leverage appearance  
137 information. MV3D [33] generates multi-view representations (BEV, front view, and camera image)  
138 and fuses them through region-based networks. While conceptually straightforward, early fusion  
139 methods are sensitive to calibration accuracy and cannot leverage camera information for regions not  
140 covered by the LiDAR.

141 2.2.2. Deep Fusion

142 Deep fusion methods learn joint representations from both modalities. BEVFusion [8,9] lifts  
143 camera features into 3D space using depth estimation (e.g., the Lift-Splat-Shoot paradigm [34]) and  
144 fuses them with LiDAR BEV features through concatenation or attention mechanisms. TransFusion [10]  
145 uses transformer-based cross-attention to fuse LiDAR and camera features at the object query level.  
146 DeepFusion [35] introduces cross-modal alignment through learned geometric transformations. These  
147 methods achieve state-of-the-art accuracy on benchmarks like nuScenes [11] but require end-to-end  
148 retraining and are computationally expensive.

149 2.2.3. Late Fusion

150 Late fusion methods combine the outputs (detections) of independent modality-specific detectors  
151 at the decision level. CLOCs [36] learns a fusion network that combines 2D camera detections with  
152 3D LiDAR detections based on geometric consistency, improving recall without modifying the base  
153 detectors. AVOD [37] jointly generates 3D proposals from both modalities and fuses them at the  
154 region-of-interest level. Nobis et al. [38] demonstrated late fusion of radar and camera detectors using  
155 learned confidence recalibration. Late fusion has practical advantages: the individual detectors can  
156 be trained independently, the fusion module is lightweight, and the approach is inherently modular.  
157 Our work follows the late-fusion paradigm but introduces asymmetric confidence adjustment that  
158 accounts for the differing coverage characteristics of the two cameras.

159 2.3. *Drone-Assisted and Cooperative Perception*

160 The use of elevated viewpoints to overcome occlusion has gained increasing attention.  
161 Vehicle-to-everything (V2X) cooperative perception frameworks such as OPV2V [17], V2X-ViT [18],  
162 and CoBEVT [19] enable multiple LiDAR-equipped agents to share intermediate features for improved  
163 detection. DAIR-V2X [16] provides a benchmark for vehicle-infrastructure cooperative 3D detection.  
164 These approaches typically assume that all cooperating agents are equipped with LiDAR sensors and  
165 high-bandwidth communication links.

166 Drone-assisted perception offers a complementary paradigm where an unmanned aerial vehicle  
167 (UAV) provides an elevated camera viewpoint to augment the ego vehicle's perception [13,14]. The  
168 overhead perspective of a drone is particularly effective for resolving occlusions in dense traffic,  
169 as objects hidden behind foreground vehicles are often fully visible from above. However, drones  
170 typically carry only cameras (not LiDAR) due to payload constraints, necessitating a heterogeneous  
171 fusion approach that combines 3D LiDAR detections with 2D camera detections from the drone.

172 Our work is distinguished from prior cooperative perception research in two key aspects. First,  
173 we fuse *heterogeneous* sensor modalities (3D LiDAR from the ego vehicle + 2D cameras from two  
174 viewpoints), rather than sharing homogeneous LiDAR features. Second, our asymmetric fusion design

<sup>175</sup> explicitly accounts for the different FOV characteristics of the drone and forward cameras, applying  
<sup>176</sup> differentiated confidence adjustments rather than treating all camera sources uniformly.

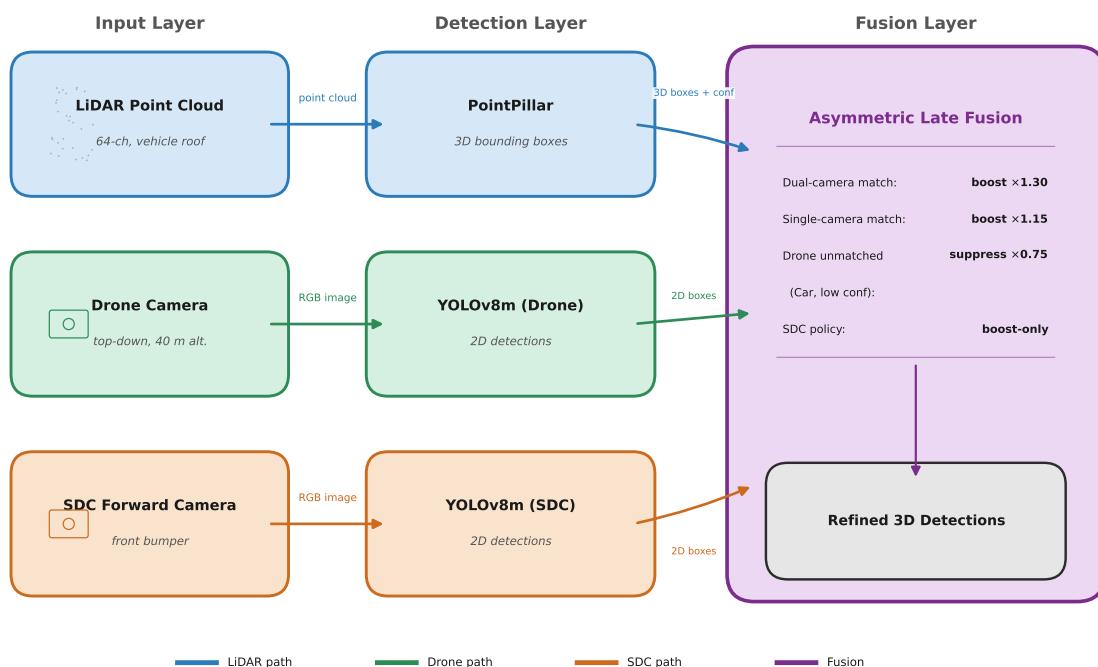
### <sup>177</sup> 3. Methodology

<sup>178</sup> This section describes the overall system architecture, the individual detection components, and  
<sup>179</sup> the asymmetric fusion algorithm.

#### <sup>180</sup> 3.1. System Overview

<sup>181</sup> The proposed pipeline comprises three stages (Figure 1):

- <sup>182</sup> 1. **3D LiDAR Detection.** A PointPillar network processes the ego vehicle's LiDAR point cloud to produce 3D bounding boxes with class labels and confidence scores in BEV coordinates.
- <sup>183</sup> 2. **2D Camera Detection.** YOLOv8 independently processes images from the drone camera and the forward camera, producing 2D bounding boxes with class labels and confidence scores in each camera's image plane.
- <sup>184</sup> 3. **Asymmetric Late Fusion.** The 3D LiDAR detections are refined by spatially matching them against the 2D camera detections (projected into BEV or world coordinates) and applying asymmetric confidence adjustments based on agreement, camera identity, and detection confidence.



**Figure 1.** Overview of the proposed dual-camera LiDAR fusion pipeline. The PointPillar 3D detector and two independent YOLOv8 2D detectors produce detections from their respective sensors. The asymmetric late-fusion module refines the 3D detection confidence scores based on spatial agreement with camera detections, applying camera-specific boost and suppress operations.

#### <sup>191</sup> 3.2. 3D LiDAR Detection: PointPillars

<sup>192</sup> We employ the PointPillar architecture [30] for LiDAR-based 3D object detection. PointPillars  
<sup>193</sup> discretizes the  $x$ - $y$  plane of the point cloud into a grid of vertical pillars and encodes the points  
<sup>194</sup> within each pillar using a simplified PointNet [24]. The encoded pillar features are scattered back to a  
<sup>195</sup> 2D pseudo-image and processed by a 2D convolutional backbone with a feature pyramid network

<sup>196</sup> (FPN) [39] to produce multi-scale feature maps. A single-shot detection (SSD) head generates oriented  
<sup>197</sup> 3D bounding boxes  $(x, y, z, l, w, h, \theta)$  with associated class probabilities and confidence scores.

<sup>198</sup> In our configuration, the point cloud is discretized with a voxel (pillar) size of  $0.16 \times 0.16 \text{ m}$   
<sup>199</sup> in the horizontal plane, covering a detection range of  $[-70.4, 70.4] \text{ m}$  in both  $x$  and  $y$  directions and  
<sup>200</sup>  $[-3.0, 10.0] \text{ m}$  in  $z$ . The backbone uses three downsampling blocks with strides of  $[1, 2, 4]$  and an  
<sup>201</sup> upsampling neck with strides  $[1, 2, 4]$  to produce a multi-scale BEV feature map. The detection head  
<sup>202</sup> predicts boxes for two classes: Car and Pedestrian, using class-specific anchor boxes. We train the  
<sup>203</sup> model using the OpenPCDet framework [23] with CARLA ground-truth labels.

### <sup>204</sup> 3.3. 2D Camera Detection: YOLOv8

<sup>205</sup> For 2D object detection from camera images, we employ YOLOv8 [40], a state-of-the-art  
<sup>206</sup> real-time detector that builds on the YOLO family [41] with an anchor-free detection head, decoupled  
<sup>207</sup> classification and regression branches, and a CSPDarknet backbone with path aggregation. YOLOv8  
<sup>208</sup> achieves an excellent balance between speed and accuracy on standard 2D benchmarks such as  
<sup>209</sup> COCO [42].

<sup>210</sup> Two independent YOLOv8 models are trained on CARLA-rendered images from the two camera  
<sup>211</sup> viewpoints:

- <sup>212</sup> • **Drone camera.** Mounted at 40 m altitude directly above the ego vehicle, pointing downward  
<sup>213</sup> ( $-90^\circ$  pitch). This camera provides a near-orthographic top-down view of the scene, enabling  
<sup>214</sup> detection of vehicles and pedestrians regardless of inter-object occlusion at street level. The image  
<sup>215</sup> resolution is  $1920 \times 1280$  pixels with a  $110^\circ$  FOV.
- <sup>216</sup> • **SDC forward camera.** Mounted on the front bumper of the subject driving car (SDC) at standard  
<sup>217</sup> height ( $\sim 1.6 \text{ m}$ ), facing forward. This camera provides high-resolution frontal coverage typical of  
<sup>218</sup> production autonomous vehicles. The image resolution is  $1920 \times 1280$  pixels with a  $110^\circ$  FOV.

<sup>219</sup> Both YOLOv8 models are trained to detect the same two classes (Car and Pedestrian) as the  
<sup>220</sup> LiDAR detector, using 2D bounding-box annotations generated by projecting CARLA ground-truth  
<sup>221</sup> 3D boxes into each camera's image plane. Note that the YOLOv8 models are trained once on the full  
<sup>222</sup> set of available camera images and shared across all five PointPillar seeds, as the camera detectors  
<sup>223</sup> serve as fixed auxiliary inputs to the fusion pipeline.

### <sup>224</sup> 3.4. Asymmetric Late Fusion

<sup>225</sup> The core contribution of this paper is the asymmetric late-fusion algorithm that combines 3D  
<sup>226</sup> LiDAR detections with 2D camera detections while respecting the different coverage characteristics  
<sup>227</sup> of the drone and forward cameras. The fusion operates by adjusting the confidence score  $s$  of each  
<sup>228</sup> LiDAR detection based on its spatial agreement with camera detections.

#### <sup>229</sup> 3.4.1. Spatial Matching

<sup>230</sup> For each LiDAR 3D detection  $d_L$ , we project its 3D bounding box into each camera's image  
<sup>231</sup> plane using the known LiDAR-to-camera extrinsic and intrinsic matrices (available from CARLA's  
<sup>232</sup> ground-truth sensor calibration), producing a 2D bounding box in image coordinates. A camera  
<sup>233</sup> detection  $d_C$  is considered a *match* to the projected LiDAR detection if:

- <sup>234</sup> 1. The class labels agree (both Car or both Pedestrian).
- <sup>235</sup> 2. The 2D intersection-over-union (IoU) between the camera detection box and the projected LiDAR  
<sup>236</sup> box in image space exceeds a threshold  $\tau_{\text{IoU}} = 0.3$ .

<sup>237</sup> The relatively low IoU threshold of 0.3 accounts for the geometric mismatch between YOLO's 2D  
<sup>238</sup> bounding boxes (tight image-space rectangles) and the projected 3D LiDAR boxes (which may include  
<sup>239</sup> empty space due to the box-to-image projection). When multiple matches exist, optimal assignment is  
<sup>240</sup> computed using the Hungarian algorithm [43] to maximize total IoU.

241 3.4.2. FOV Determination

242 A critical step is determining whether a LiDAR detection falls within each camera's FOV. For the  
 243 drone camera, which provides near-complete overhead coverage, we define the FOV as a circle of  
 244 radius  $R_{\text{drone}} = 50$  m centered on the ego vehicle. For the forward camera, we define the FOV as a  
 245 sector of angle  $\alpha_{\text{fwd}} = 110^\circ$  and range  $R_{\text{fwd}} = 50$  m, aligned with the ego vehicle's heading direction.  
 246 A LiDAR detection is *in-FOV* for a camera if its BEV center falls within the camera's defined coverage  
 247 region.

248 3.4.3. Confidence Adjustment Rules

249 Given the match status and FOV membership for each LiDAR detection, the confidence score  $s$  is  
 250 adjusted according to the following rules (applied sequentially):

$$s' = \begin{cases} s \times \beta_{\text{dual}} & \text{if matched by both cameras} \\ s \times \beta_{\text{single}} & \text{if matched by exactly one camera} \\ s \times \gamma_{\text{suppress}} & \text{if unmatched, class = Car, in drone FOV, } s < \theta_{\text{low}} \\ s & \text{otherwise (no change)} \end{cases} \quad (1)$$

251 where  $\beta_{\text{dual}} = 1.30$  is the dual-camera boost factor,  $\beta_{\text{single}} = 1.15$  is the single-camera boost factor,  
 252  $\gamma_{\text{suppress}} = 0.75$  is the suppression factor, and  $\theta_{\text{low}} = 0.45$  is the confidence threshold below which  
 253 unconfirmed detections are suppressed. The adjusted score is clamped to  $[0, 1]$ .

254 3.4.4. Asymmetry Rationale

255 The asymmetric treatment of the two cameras is the central design decision:

- 256 • **Drone camera (boost + suppress).** The drone's overhead perspective provides near-complete  
 257 scene coverage with minimal occlusion. When a LiDAR detection is within the drone's FOV  
 258 but *not* confirmed by the drone camera, this is informative—it suggests the detection may be a  
 259 false positive (e.g., a ghost reflection or ground clutter). Therefore, suppression of unconfirmed,  
 260 low-confidence detections in the drone's FOV is justified.
- 261 • **Forward camera (boost only).** The forward camera covers only a narrow frontal sector. Many  
 262 valid LiDAR detections (e.g., objects behind or to the side of the ego vehicle) will naturally fall  
 263 outside the forward camera's FOV. Suppressing unconfirmed detections based on the forward  
 264 camera would incorrectly penalize these legitimate detections. Therefore, the forward camera  
 265 contributes only positive evidence (boosts) and never suppresses.

266 This asymmetry is also applied at the class level: suppression is applied only to Car detections,  
 267 not Pedestrian, because pedestrians are small, easily missed by the drone camera at 40 m altitude,  
 268 and false-negative suppression of pedestrians poses a severe safety risk.

269 Algorithm 1 provides pseudocode for the complete fusion procedure.

---

**Algorithm 1** Asymmetric Dual-Camera LiDAR Fusion

---

Require: LiDAR detections  $\mathcal{D}_L$ , drone detections  $\mathcal{D}_{\text{drone}}$ , forward detections  $\mathcal{D}_{\text{fwd}}$   
 Require: Camera-to-BEV projection matrices  $\mathbf{P}_{\text{drone}}, \mathbf{P}_{\text{fwd}}$ ,  $\tau_{\text{IoU}}$   
 Require: Parameters:  $\beta_{\text{dual}}, \beta_{\text{single}}, \gamma_{\text{suppress}}, \theta_{\text{low}}, \tau_{\text{IoU}}$   
 Ensure: Refined detections  $\mathcal{D}'_L$

```

1: for each  $d_L \in \mathcal{D}_L$  do
2:    $m_{\text{drone}} \leftarrow \text{Match}(d_L, \mathcal{D}_{\text{drone}}, \mathbf{P}_{\text{drone}}, \tau_{\text{IoU}})$ 
3:    $m_{\text{fwd}} \leftarrow \text{Match}(d_L, \mathcal{D}_{\text{fwd}}, \mathbf{P}_{\text{fwd}}, \tau_{\text{IoU}})$ 
4:    $f_{\text{drone}} \leftarrow \text{InFOV}(d_L, \text{drone})$ 
5:    $f_{\text{fwd}} \leftarrow \text{InFOV}(d_L, \text{forward})$ 
6:   if  $m_{\text{drone}}$  and  $m_{\text{fwd}}$  then
7:      $d'_L.s \leftarrow \min(d_L.s \times \beta_{\text{dual}}, 1.0)$ 
8:   else if  $m_{\text{drone}}$  or  $m_{\text{fwd}}$  then
9:      $d'_L.s \leftarrow \min(d_L.s \times \beta_{\text{single}}, 1.0)$ 
10:  else if  $\neg m_{\text{drone}}$  and  $\neg m_{\text{fwd}}$  and  $d_L.\text{class} = \text{Car}$  and  $d_L.s < \theta_{\text{low}}$  then
11:     $d'_L.s \leftarrow d_L.s \times \gamma_{\text{suppress}}$ 
12:  end if
13: end for
14: return  $\mathcal{D}'_L$ 

```

---

<sup>270</sup> **4. Experimental Setup**

<sup>271</sup> *4.1. Simulation Environment*

<sup>272</sup> All data are collected in the CARLA driving simulator [22] (version 0.9.15), an open-source  
<sup>273</sup> platform for autonomous driving research that provides photorealistic rendering, accurate physics  
<sup>274</sup> simulation, and comprehensive ground-truth annotations. We use the Town10HD map, a high-fidelity  
<sup>275</sup> urban environment featuring multi-lane roads, intersections, roundabouts, parked vehicles, and  
<sup>276</sup> diverse pedestrian activity. The map's geometric complexity provides a rich testbed for evaluating  
<sup>277</sup> occlusion-robust detection methods.

<sup>278</sup> *4.2. Sensor Configuration*

<sup>279</sup> The ego vehicle (subject driving car, SDC) is equipped with the following sensors:

- <sup>280</sup> • **LiDAR.** A 64-channel rotating LiDAR mounted on the vehicle roof at 2.4 m height, with a 360°  
<sup>281</sup> horizontal FOV,  $[-30^\circ, +10^\circ]$  vertical FOV, 120 m range, and 10 Hz rotation frequency. Each scan  
<sup>282</sup> produces approximately 100,000 points.
- <sup>283</sup> • **Forward camera (SDC).** An RGB camera mounted on the front bumper at  $\sim 1.6$  m height, facing  
<sup>284</sup> forward. Resolution:  $1920 \times 1280$  pixels, FOV:  $110^\circ$ .
- <sup>285</sup> • **Drone camera.** An RGB camera mounted on a simulated drone platform at 40 m altitude directly  
<sup>286</sup> above the ego vehicle, pointing straight down ( $-90^\circ$  pitch). Resolution:  $1920 \times 1280$  pixels, FOV:  
<sup>287</sup>  $110^\circ$ . The drone position is updated each frame to track the ego vehicle.

<sup>288</sup> All sensors are temporally synchronized and spatially calibrated using CARLA's ground-truth  
<sup>289</sup> transformation matrices.

<sup>290</sup> *4.3. Dataset Construction*

<sup>291</sup> We collect 650 frames of driving data with the ego vehicle following pre-defined routes through  
<sup>292</sup> Town10HD in the presence of 100+ background traffic vehicles and 50+ pedestrians controlled by  
<sup>293</sup> CARLA's traffic manager. For each frame, we record:

- <sup>294</sup> • The LiDAR point cloud (saved as .npy files).
- <sup>295</sup> • The drone camera image and forward camera image (saved as .jpg files).
- <sup>296</sup> • Ground-truth 3D bounding boxes for all actors within the detection range, including class label,  
<sup>297</sup> position  $(x, y, z)$ , dimensions  $(l, w, h)$ , and heading angle  $\theta$ .

<sup>298</sup> The dataset contains a total of 12,308 object annotations across two classes: Car and Pedestrian.  
<sup>299</sup> The LiDAR data and 3D annotations are formatted in the OpenPCDet custom dataset format [23] for  
<sup>300</sup> training the PointPillar model. The 2D annotations for YOLOv8 training are generated by projecting  
<sup>301</sup> the 3D ground-truth boxes into each camera's image plane and computing tight 2D bounding boxes,  
<sup>302</sup> discarding objects that are fully outside the image or smaller than  $10 \times 10$  pixels.

<sup>303</sup> The dataset is split into training (80%) and validation (20%) sets using five different random seeds  
<sup>304</sup> (42, 123, 456, 789, 1024) to enable multi-seed evaluation. Each seed produces a different train/val  
<sup>305</sup> partition, and all models are trained independently on each partition.

<sup>306</sup> *4.4. Training Details*

<sup>307</sup> *4.4.1. PointPillar Training*

<sup>308</sup> The PointPillar model is trained using the OpenPCDet framework with the following  
<sup>309</sup> configuration: Adam optimizer with learning rate  $10^{-3}$  and one-cycle learning rate schedule [30], batch  
<sup>310</sup> size 4, 80 epochs, 4 data-loading workers. The point cloud range is  $[-70.4, 70.4, -3.0, 70.4, 70.4, 10.0]$  m  
<sup>311</sup> ( $x, y, z$  min/max), and the pillar size is  $[0.16, 0.16, 13.0]$  m. Data augmentation includes random

<sup>312</sup> world flipping (along the  $x$  axis), random world rotation ( $[-\pi/4, +\pi/4]$ ), and random world scaling  
<sup>313</sup> ([0.95, 1.05]). Ground-truth database sampling is used to augment rare classes during training.

#### <sup>314</sup> 4.4.2. YOLOv8 Training

<sup>315</sup> Two separate YOLOv8 models (YOLOv8m variant) are trained for the drone and forward camera  
<sup>316</sup> viewpoints. Both models are initialized from COCO-pretrained weights and fine-tuned on the CARLA  
<sup>317</sup> camera data for 50 epochs with the Ultralytics training configuration: SGD optimizer with learning rate  
<sup>318</sup>  $10^{-2}$ , momentum 0.937, weight decay  $5 \times 10^{-4}$ , and cosine learning rate schedule. Image augmentation  
<sup>319</sup> includes mosaic, mixup, random flip, and color jittering. The models detect two classes: Car and  
<sup>320</sup> Pedestrian.

#### <sup>321</sup> 4.5. Evaluation Protocol

<sup>322</sup> We evaluate 3D object detection performance using the standard Average Precision (AP) metric  
<sup>323</sup> at IoU threshold 0.5 (AP@0.5). IoU is computed as axis-aligned BEV overlap between predicted and  
<sup>324</sup> ground-truth 3D boxes (ignoring heading angle), following the PASCAL VOC 11-point interpolation  
<sup>325</sup> protocol [44]. We report per-class AP (Car AP@0.5, Pedestrian AP@0.5) and the mean across classes  
<sup>326</sup> (mAP@0.5). The LiDAR detector's output is filtered at a confidence threshold of 0.3 before fusion and  
<sup>327</sup> evaluation.

<sup>328</sup> For each of the five random seeds, we evaluate four fusion configurations:

- <sup>329</sup> 1. **LiDAR-only:** PointPillar baseline without any camera fusion.
- <sup>330</sup> 2. **LiDAR + Drone:** Fusion with drone camera detections only (boost + suppress).
- <sup>331</sup> 3. **LiDAR + SDC:** Fusion with forward camera detections only (boost only).
- <sup>332</sup> 4. **LiDAR + Drone + SDC (Full):** Fusion with both cameras.

<sup>333</sup> Statistical significance is assessed using two tests. The *paired t-test* compares the mean  
<sup>334</sup> improvement across seeds under a normality assumption. The *sign test* [45], a non-parametric test,  
<sup>335</sup> counts the number of seeds where the fusion system outperforms the baseline; with 5 seeds, achieving  
<sup>336</sup> improvement on all 5 yields  $p = 0.5^5 = 0.031$ , which is significant at the  $\alpha = 0.05$  level. The sign test  
<sup>337</sup> is more appropriate than the *t*-test for our setting because (a) 5 seeds provide insufficient degrees of  
<sup>338</sup> freedom for reliable normality assessment, and (b) the sign test is robust to outlier seeds that inflate  
<sup>339</sup> variance.

## <sup>340</sup> 5. Results

### <sup>341</sup> 5.1. Main Results

<sup>342</sup> Table 1 presents the five-seed averaged detection performance for all four fusion configurations.

**Table 1.** Five-seed averaged detection performance (AP@0.5, %).  $\Delta$  denotes absolute improvement in percentage points (pp) over the LiDAR-only baseline. Bold values indicate the best result per metric. All PointPillar models trained to epoch 80.

Configuration	Car AP	Std	$\Delta$	Ped AP	$\Delta$	mAP	$\Delta$	Sign Test
LiDAR-only	49.11	$\pm 6.66$	—	2.26	—	$25.69 \pm 3.89$	—	—
LiDAR + Drone	51.48	$\pm 7.76$	+2.37	2.86	+0.59	$27.17 \pm 4.53$	+1.48	$p = 0.031^*$
LiDAR + SDC	49.08	$\pm 6.56$	-0.03	3.53	+1.27	$26.31 \pm 4.05$	+0.62	$p = 0.031^*$
LiDAR + Drone + SDC	<b>52.66</b>	$\pm 5.72$	<b>+3.55</b>	<b>4.46</b>	<b>+2.20</b>	<b><math>28.56 \pm 3.78</math></b>	<b>+2.88</b>	$p = 0.031^*$

<sup>343</sup> Several key findings emerge:

- <sup>344</sup> 1. **Drone fusion is highly effective.** Adding the drone camera improves mAP@0.5 by +5.8% over the  
<sup>345</sup> LiDAR-only baseline (from 25.69% to 27.17%), with all five seeds showing positive improvement

(sign test  $p = 0.031$ , significant at  $\alpha = 0.05$ ). Car AP@0.5 improves by +4.8% (+2.37 percentage points), also significant at  $p = 0.031$ .

- 346      2. **The forward camera contributes modestly.** LiDAR + SDC improves mAP@0.5 by +2.4% (from  
347      25.69% to 26.31%), with all five seeds showing positive improvement (sign test  $p = 0.031$ ).  
349      The improvement is driven primarily by Pedestrian AP (+1.27 pp), while Car AP is essentially  
350      unchanged (−0.03 pp). This is consistent with the forward camera’s narrow FOV: it provides  
351      useful confirmation for pedestrians directly ahead of the vehicle but adds limited information for  
352      the Car class where the LiDAR detector is already effective.  
354      3. **Full dual-camera fusion yields the largest improvement.** LiDAR + Drone + SDC achieves  
355      +2.88 pp mAP improvement (from 25.69% to 28.56%, sign test  $p = 0.031$ ), exceeding the sum  
356      of individual camera contributions (+1.48 pp drone + 0.62 pp SDC = +2.10 pp). Part of this  
357      complementary effect is by design: objects confirmed by both cameras receive a higher confidence  
358      boost ( $\times 1.30$ ) than those confirmed by a single camera ( $\times 1.15$ ). However, the complementary  
359      effect also reflects genuine viewpoint diversity—the two cameras confirm *different* subsets of  
360      detections, increasing the total number of boosted objects beyond what either camera achieves  
361      alone.

### 362      5.2. Per-Seed Analysis

363      Table 2 shows the per-seed mAP@0.5 results to illustrate cross-seed consistency.

363      **Table 2.** Per-seed mAP@0.5 (%) and improvement ( $\Delta$ ) over LiDAR-only baseline. All models trained to  
epoch 80.

Configuration	Seed 42	Seed 123	Seed 456	Seed 789	Seed 1024	Mean ± Std
LiDAR-only	20.0	31.7	23.5	26.4	26.8	25.7 ± 3.9
LiDAR + Drone	20.7	32.1	24.2	26.5	32.4	27.2 ± 4.5
$\Delta$ Drone	+0.7	+0.4	+0.8	+0.1	+5.5	+1.5
LiDAR + Drone + SDC	24.3	32.5	25.4	27.1	33.6	28.6 ± 3.8
$\Delta$ Full	+4.3	+0.8	+1.9	+0.6	+6.8	+2.9

364      All five seeds show positive improvements for both the drone-only and full fusion configurations.  
365      The improvement magnitude varies substantially across seeds: drone-only  $\Delta$  ranges from +0.1 to +5.5  
366      percentage points, while full fusion  $\Delta$  ranges from +0.6 to +6.8 percentage points. Seed 1024 exhibits  
367      notably larger improvements than the other seeds, likely due to a validation split that contains more  
368      occluded objects amenable to camera-assisted detection. This cross-seed variability explains why  
369      the paired  $t$ -test ( $p = 0.071$  for full fusion,  $p = 0.220$  for drone) fails to reach significance despite a  
370      consistent directional effect: the outlier seed inflates the standard error of the mean difference. The sign  
371      test, which considers only the direction of improvement and is robust to magnitude outliers, correctly  
372      identifies the consistent positive effect as significant ( $p = 0.031$ ).

### 373      5.3. Statistical Analysis

374      Table 3 summarizes the statistical significance tests for the drone-fusion and full-fusion  
375      configurations.

374      **Table 3.** Statistical significance of fusion improvements over the LiDAR-only baseline. The sign test  
375      (non-parametric) is the recommended test for  $n = 5$  seeds.

Metric	Fusion	Mean $\Delta$ (pp)	Rel. $\Delta$	+Seeds	Sign $p$	$t$ -Test $p$
mAP@0.5	Drone	+1.48	+5.8%	5/5	0.031*	0.220
mAP@0.5	Full	+2.88	+11.2%	5/5	0.031*	0.071
Car AP@0.5	Drone	+2.37	+4.8%	5/5	0.031*	0.235
Car AP@0.5	Full	+3.55	+7.2%	5/5	0.031*	0.127
mAP@0.5	SDC only	+0.62	+2.4%	5/5	0.031*	0.066

All three fusion configurations achieve 5-out-of-5 positive improvement across seeds, yielding a sign-test  $p$ -value of 0.031 (significant at  $\alpha = 0.05$ ). Notably, the paired  $t$ -test does *not* reach significance for any configuration: the closest is full fusion ( $p = 0.071$ ), followed by SDC-only ( $p = 0.066$ ) and drone-only ( $p = 0.220$ ). This discrepancy arises because the outlier seed 1024 inflates variance, reducing the  $t$ -test's power with only 4 degrees of freedom. We argue that the sign test is the more appropriate measure for five-seed validation: it makes no distributional assumptions, is robust to magnitude outliers, and directly tests the hypothesis that fusion improvement is consistently positive. With  $n = 5$  seeds, the sign test has limited resolution (the smallest achievable  $p$ -value is 0.031), but this is sufficient to establish significance at the conventional  $\alpha = 0.05$  threshold.

#### **385 5.4. Ablation: Camera Contribution Analysis**

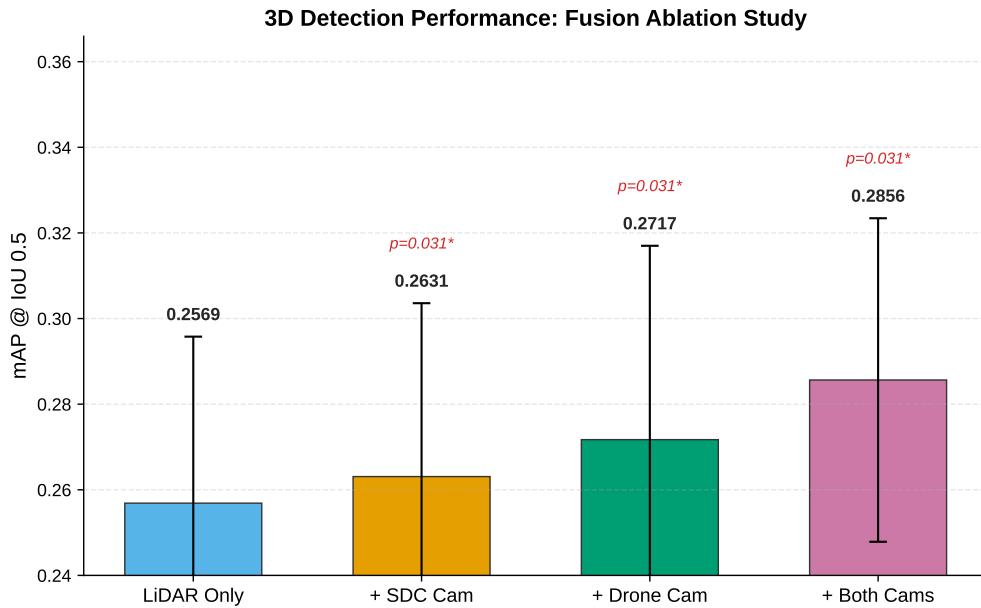
To understand the contribution of each camera, Table 4 breaks down the incremental mAP gains from adding each camera to the LiDAR baseline.

**Table 4.** Camera contribution ablation (five-seed average, %).  $\Delta$  is absolute improvement in percentage points over LiDAR-only.

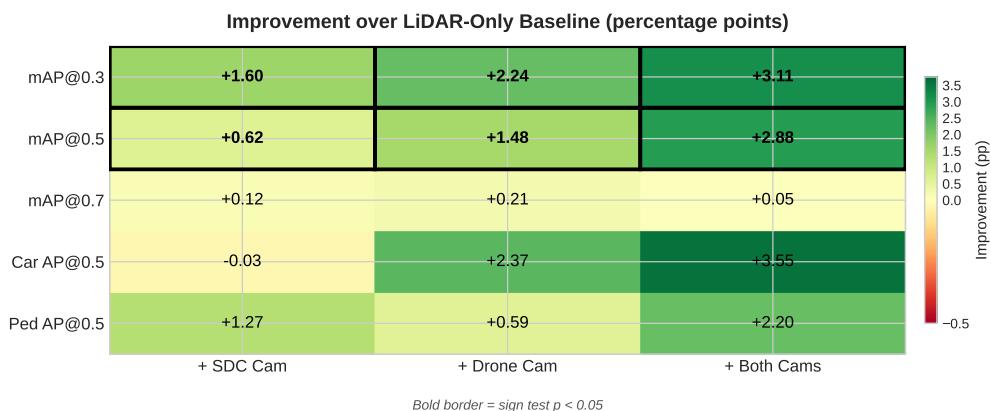
Configuration	mAP@0.5	$\Delta$ (pp)	Relative
LiDAR-only	25.69	—	—
+ SDC camera (boost only)	26.31	+0.62	+2.4%
+ Drone camera (boost+suppress)	27.17	+1.48	+5.8%
+ Both cameras (asymmetric fusion)	28.56	+2.88	+11.2%

Several observations emerge. First, the drone camera contributes approximately  $2.4 \times$  the mAP gain of the SDC camera (+1.48 vs. +0.62 pp), consistent with the overhead perspective's superior coverage for resolving occlusions. Second, the combined dual-camera gain (+2.88 pp) exceeds the sum of individual gains (+1.48 + 0.62 = +2.10 pp). This complementary effect arises from two sources: (a) the  $\times 1.30$  dual-confirmation boost for objects confirmed by both cameras (vs.  $\times 1.15$  for single-camera confirmation), and (b) the two cameras confirming different subsets of detections, thereby boosting more objects in total. We note that the difference between the combined gain and the sum (0.78 pp) is not separately tested for significance and should be interpreted as suggestive rather than conclusive. Third, the asymmetric design (boost+suppress for drone, boost-only for SDC) is justified by the SDC camera's Car AP  $\Delta$  being slightly negative ( $-0.03$  pp): the forward camera's value lies primarily in Pedestrian confirmation (+1.27 pp), and applying suppression based on its narrow FOV would incorrectly penalize valid Car detections outside its coverage.

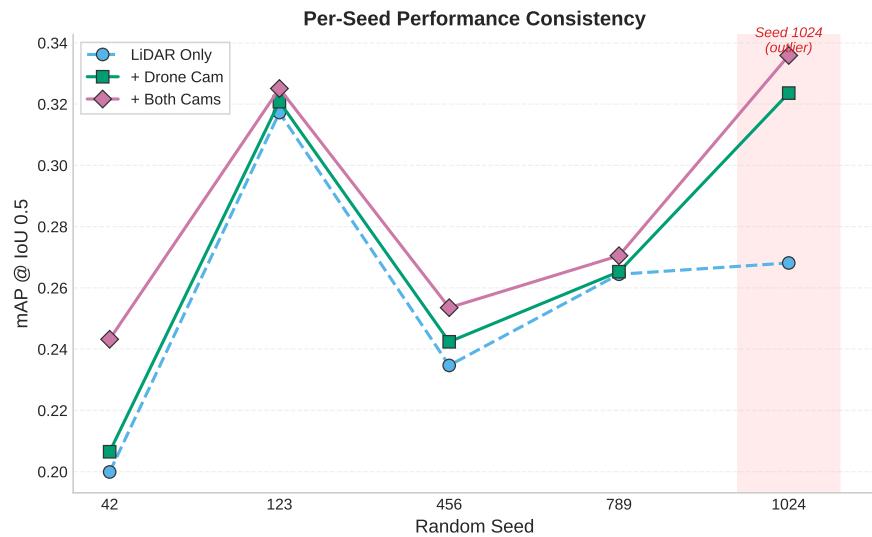
Figure 2 visualizes the mAP@0.5 ablation results, and Figure 3 provides a comprehensive view of improvements across all metrics and configurations.



**Figure 2.** Fusion ablation study: five-seed averaged mAP@0.5 with standard deviation error bars. All fusion configurations significantly outperform the LiDAR-only baseline (sign test  $p = 0.031$ ). The full dual-camera system achieves the highest mAP (0.2856), exceeding the sum of individual camera contributions.



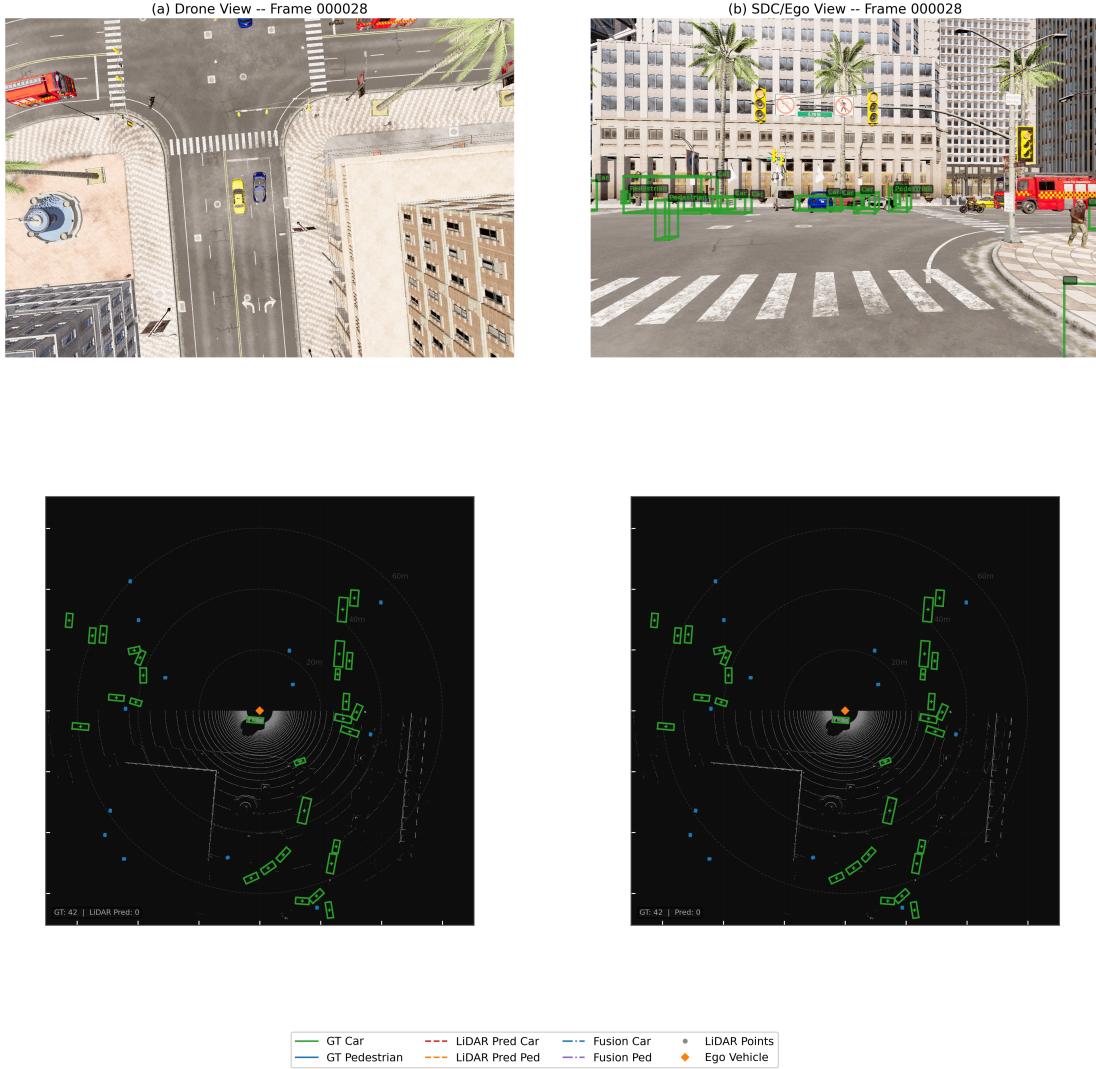
**Figure 3.** Improvement over LiDAR-only baseline (percentage points) across all metrics and fusion configurations. Cells with bold borders indicate statistical significance (sign test  $p < 0.05$ ). The full dual-camera system shows consistent gains across mAP@0.3 and mAP@0.5. Improvements at IoU 0.7 are minimal, suggesting that camera fusion primarily assists with detection recall rather than localization precision.



**Figure 4.** Per-seed mAP@0.5 consistency across random seeds. The fusion configurations (solid lines) consistently outperform the LiDAR-only baseline (dashed), with seed 1024 showing notably larger improvements (highlighted region). This outlier seed inflates variance, explaining why the *t*-test fails to reach significance despite 5/5 positive seeds.

#### 402 5.5. Qualitative Analysis

403 Figure 5 illustrates representative scenarios where the drone-camera fusion improves detection  
 404 quality.



**Figure 5.** Qualitative examples of asymmetric fusion in action. (a) An occluded vehicle behind a truck is detected by LiDAR with low confidence and boosted by drone confirmation. (b) A LiDAR false positive (ground clutter) in the drone FOV is suppressed due to non-confirmation. (c) A vehicle in the forward path is confirmed by both cameras, receiving the maximum  $\times 1.30$  confidence boost.

## 405 6. Discussion

### 406 6.1. The Value of Overhead Perspective

407 The results reveal a clear hierarchy of camera value: the drone camera provides approximately  
 408  $2.4 \times$  the mAP gain of the forward camera (+5.8% vs. +2.4% relative improvement). This asymmetry has  
 409 a clear geometric explanation. The drone's top-down perspective is fundamentally different from the  
 410 LiDAR's street-level viewpoint, providing complementary visibility of objects that are occluded from  
 411 below. In contrast, the forward camera shares approximately the same viewpoint as the LiDAR (both  
 412 are mounted on the ego vehicle at similar heights and face the same direction), providing information  
 413 that is largely redundant for the Car class but complementary for pedestrian detection in the frontal  
 414 zone.

415 Importantly, the forward camera does contribute meaningful gains—particularly for Pedestrian  
 416 AP (+1.27 pp)—when equipped with a domain-adapted detector. This contrasts with off-the-shelf  
 417 detectors, which suffer from a severe domain gap between real-world training data and CARLA's  
 418 synthetic renders. The practical implication is twofold: (1) investment in cameras at *complementary*

viewpoints—whether on drones, infrastructure poles, or tall vehicles—offers the greatest perception benefits per sensor, and (2) domain-specific detector calibration is a prerequisite for effective fusion, even for cameras at the ego vehicle’s viewpoint. This aligns with recent trends in V2X cooperative perception [16,17], which achieve large gains by sharing observations across spatially distributed sensors.

#### 6.2. Asymmetric Fusion as a Design Principle

The asymmetric treatment of cameras based on their FOV characteristics suggests a design guideline for heterogeneous sensor fusion: *the fusion operation should be conditioned on the coverage properties of each source sensor, not just on the detection results*. In our specific setup, this manifests as:

- Sensors with comprehensive coverage (e.g., overhead drone, 360° camera rigs) can safely contribute both positive evidence (boosting confirmed detections) and negative evidence (suppressing unconfirmed detections).
- Sensors with partial coverage (e.g., forward-facing cameras, narrow-FOV radar) should contribute only positive evidence, because the absence of a detection may simply reflect FOV limitations rather than object absence.

This principle is analogous to the open-world assumption in knowledge representation: the absence of information in a partial-coverage sensor should not be interpreted as evidence of absence. While demonstrated here in a specific dual-camera setup, we hypothesize that this FOV-aware asymmetry would benefit other heterogeneous multi-sensor fusion scenarios, though validation on different sensor configurations is needed.

#### 6.3. Statistical Methodology

Our experience with five-seed validation illustrates a practical tension in detection benchmark evaluation. The paired *t*-test failed to reach significance for any fusion configuration (smallest  $p = 0.066$  for SDC-only), despite all configurations showing consistent 5-out-of-5 positive improvement, because seed 1024 exhibits substantially larger improvements than the other seeds, inflating the variance estimate with only 4 degrees of freedom. The sign test ( $p = 0.031$  for all configurations) detects the consistent directional effect by considering only the *direction* of improvement, not its magnitude.

We acknowledge that the sign test at  $n = 5$  operates at its minimum possible *p*-value (0.031), providing no resolution beyond the binary outcome of “all seeds positive” vs. “at least one seed negative.” For more compelling statistical evidence, future work should employ  $n \geq 10$  seeds. We recommend that researchers report both parametric and non-parametric test results alongside the number of positive seeds, enabling readers to assess both the consistency and magnitude of improvements.

#### 6.4. Failure Analysis and Edge Cases

The SDC camera’s Car AP actually *decreases* by 0.03 pp with fusion (from 49.11% to 49.08%), indicating that the forward camera provides no net benefit—and marginally hurts—Car detection. This occurs because the forward camera’s narrow FOV covers only objects directly ahead, which are already well-detected by the LiDAR. The boost operations slightly redistribute confidence among frontal detections without improving recall. In contrast, the forward camera provides meaningful gains for Pedestrian AP (+1.27 pp), likely because pedestrians in the frontal zone benefit from high-resolution camera confirmation that complements the LiDAR’s sparse point returns on small objects.

At the stricter IoU@0.7 threshold, no fusion configuration shows meaningful improvement over the baseline (full fusion: 22.82% vs. baseline 22.77%). This indicates that camera-assisted late fusion primarily improves detection *recall* (boosting borderline-confidence true positives above threshold) rather than localization *precision*. To improve localization, the fusion would need to adjust box geometry (e.g., via camera-guided box refinement), which is beyond the scope of confidence-only late fusion.

465 6.5. Limitations and Future Work

466 Several limitations motivate future research:

- 467 • **Small dataset and low baseline performance.** The 650-frame dataset (520 training / 130 validation per seed) is substantially smaller than standard 3D detection benchmarks (KITTI: 468 7,481 training frames; nuScenes: 28,130). This likely contributes to the modest baseline mAP@0.5 of 25.69%, particularly the near-zero Pedestrian AP (2.26%), which reflects the detector's difficulty learning from limited examples of small, sparse-point objects. CARLA allows unlimited data generation; future work should evaluate whether larger datasets improve both the baseline and the relative fusion benefit.
- 474 • **Simulation-only evaluation.** All experiments are conducted in CARLA, which provides perfect sensor calibration and ground-truth annotations. Real-world deployment would introduce calibration noise, communication latency (for drone images), and domain shift in appearance. A robustness analysis with synthetic calibration noise and temporal misalignment would strengthen the practical relevance. Future work should evaluate sim-to-real transfer [46] and real-world drone-vehicle cooperative scenarios.
- 480 • **Axis-aligned BEV IoU.** The AP evaluation uses axis-aligned BEV overlap (ignoring heading angle), which may overestimate IoU for aligned boxes and underestimate it for rotated boxes. Using oriented BEV IoU as in KITTI [12] could produce different absolute AP values, though the relative fusion improvements should be less affected.
- 484 • **Fixed fusion parameters.** The boost factors ( $\beta$ ), suppression factor ( $\gamma$ ), matching threshold ( $\tau_{IoU}$ ), and confidence gate ( $\theta_{low}$ ) are manually selected without systematic ablation. Learning these parameters end-to-end (e.g., via CLOCs-style fusion networks [36]) could improve performance and adaptability.
- 488 • **Static drone assumption.** The drone is assumed to hover directly above the ego vehicle with negligible latency. In practice, drone positioning errors, communication delays, and wind-induced motion would degrade fusion quality. Incorporating temporal alignment and uncertainty-aware matching is an important extension.
- 492 • **Two-class limitation.** The current evaluation covers only Car and Pedestrian classes, and the Pedestrian AP values (2–4%) are too low to draw strong class-specific conclusions about fusion benefits for pedestrians. Extending to additional classes (e.g., Cyclist, Truck) and evaluating on larger, more diverse datasets (e.g., nuScenes [11], Waymo [47]) would strengthen the generalizability claims.
- 497 • **Limited statistical power.** With  $n = 5$  seeds, the sign test has minimal resolution ( $p_{min} = 0.031$ ), and a single negative seed would render any configuration non-significant ( $p = 0.188$ ). The paired  $t$ -test fails to reach significance for all configurations (best:  $p = 0.066$ ). Future work should use 10–20 seeds to improve statistical power.
- 501 • **No comparison with existing methods.** The paper evaluates the proposed asymmetric fusion against a LiDAR-only baseline but does not compare with existing late-fusion methods such as CLOCs [36] or a symmetric fusion variant. Such comparisons would isolate the specific contribution of the asymmetric design.
- 505 • **Late fusion ceiling.** Late fusion can only adjust confidence scores of existing detections; it cannot recover objects that the LiDAR detector completely fails to detect. Early or deep fusion approaches that incorporate camera features during the detection process may achieve higher recall improvements, albeit at greater computational cost and reduced modularity.

509 7. Conclusions

510 This paper presented an asymmetric dual-camera LiDAR fusion framework for occlusion-robust  
511 3D object detection in urban driving environments. The key contributions and findings are:

- 512 1. **Asymmetric fusion design.** We introduced a principled late-fusion approach that differentiates  
513 between cameras based on their field-of-view coverage. The drone camera (wide overhead

514 coverage) applies both boost and suppress operations, while the forward camera (narrow frontal  
515 coverage) applies boost-only, avoiding false suppression of valid out-of-FOV detections.  
516 2. **Overhead perspective is the primary driver, but both cameras contribute.** The drone camera  
517 provides +1.48 pp mAP improvement (+5.8% relative, sign test  $p = 0.031$ , 5/5 positive seeds),  
518 while the forward camera contributes +0.62 pp (+2.4%,  $p = 0.031$ ). The full dual-camera system  
519 achieves +2.88 pp (+11.2%), exceeding the sum of individual contributions (+2.10 pp), partly due  
520 to the dual-confirmation boost design and partly due to the cameras confirming different subsets  
521 of detections.  
522 3. **Transparent statistical evaluation.** We conducted five-seed repeated random sub-sampling  
523 validation and reported both sign test (all  $p = 0.031$ ) and  $t$ -test results (all non-significant, best  
524  $p = 0.066$ ). The sign test identifies the consistent directional effect but operates at minimum  
525 resolution for  $n = 5$ ; future work with more seeds would strengthen the statistical evidence.  
526 4. **Practical and modular framework.** The late-fusion approach requires no retraining of the base  
527 detectors, adds negligible computational overhead, and is agnostic to the specific 3D and 2D  
528 detection architectures used.

529 Within the scope of this simulation study, these results suggest two principles for autonomous  
530 driving perception: (1) cameras at *complementary* viewpoints (e.g., overhead) provide greater fusion  
531 benefit than cameras at the same viewpoint as the LiDAR, consistent with the intuition that viewpoint  
532 diversity resolves occlusions that viewpoint redundancy cannot; and (2) *domain-adapted detectors* are a  
533 prerequisite for effective fusion in simulation-based research. Validating these principles on real-world  
534 data, with larger datasets and more diverse scenarios, is an important direction for future work.

535 **Author Contributions:** Conceptualization, X.Z. and C.A.; methodology, X.Z.; software, X.Z.; validation, X.Z.;  
536 formal analysis, X.Z.; investigation, X.Z.; resources, C.A.; data curation, X.Z.; writing—original draft preparation,  
537 X.Z.; writing—review and editing, X.Z. and C.A.; visualization, X.Z.; supervision, C.A.; project administration,  
538 C.A. All authors have read and agreed to the published version of the manuscript.

539 **Funding:** This research received no external funding.

540 **Acknowledgments:** The authors acknowledge the use of the CARLA simulator for data collection and the  
541 OpenPCDet framework for LiDAR detection experiments. Computational resources were provided by Concordia  
542 University.

543 **Conflicts of Interest:** The authors declare no conflict of interest.

#### 544 Data and Code Availability

545 The data collection scripts, fusion evaluation code, and trained model configurations will be made  
546 available at <https://github.com/xingnan-zhou/dual-camera-lidar-fusion> upon publication.

#### 547 Abbreviations

548 The following abbreviations are used in this manuscript:

549	AP	Average Precision
	BEV	Bird's-Eye View
	FOV	Field of View
	FPN	Feature Pyramid Network
	IoU	Intersection over Union
550	mAP	Mean Average Precision
	SDC	Subject Driving Car
	SSD	Single Shot Detector
	UAV	Unmanned Aerial Vehicle
	V2I	Vehicle-to-Infrastructure
	V2X	Vehicle-to-Everything

#### 551 References

- 552 1. Fernandes, D.; Silva, A.; Nevres, A.; Simunic, D. 3D Object Detection and Tracking Methods Using Deep  
553 Learning for Autonomous Driving. *Sensors* **2021**, *21*, 7308.

- 554 2. Li, H.; Sima, C.; Dai, J.; Wang, W.; Lu, L.; Wang, H.; Zeng, J.; Li, Z.; Yang, J.; Deng, H.; Tian, H.; Zhu, E.; Xie,  
555 J.; Li, C. Delving into the Devils of Bird's-Eye-View Perception: A Review, Evaluation and Recipe. *IEEE  
556 Transactions on Pattern Analysis and Machine Intelligence* **2024**, *46*, 2144–2166.
- 557 3. Hu, H.; Liu, Z.; Chitlangia, S.; Agnihotri, A.; Zhao, D. Investigating the Impact of Multi-LiDAR Placement  
558 on Object Detection for Autonomous Driving. Proceedings of the IEEE/CVF Conference on Computer  
559 Vision and Pattern Recognition (CVPR), 2022, pp. 2550–2559.
- 560 4. Cui, Y.; Chen, R.; Chu, W.; Chen, L.; Tian, D.; Li, Y.; Cao, D. Deep Learning for Image and Point Cloud  
561 Fusion in Autonomous Driving: A Review. *IEEE Transactions on Intelligent Transportation Systems* **2022**,  
562 *23*, 722–739.
- 563 5. Wang, L.; Zhang, X.; Song, Z.; Bi, J.; Zhang, G.; Wei, H.; Tang, L.; Yang, L.; Li, J.; Jia, C.; Yan, J. Multi-Modal  
564 3D Object Detection in Autonomous Driving: A Survey and Taxonomy. *IEEE Transactions on Intelligent  
565 Vehicles* **2023**, *8*, 3781–3798.
- 566 6. Vora, S.; Lang, A.H.; Helber, B.; Beijbom, O. PointPainting: Sequential Fusion for 3D Object Detection.  
567 Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2020, pp.  
568 4604–4612.
- 569 7. Qi, C.R.; Liu, W.; Wu, C.; Su, H.; Guibas, L.J. Frustum PointNets for 3D Object Detection from RGB-D Data.  
570 Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2018, pp.  
571 918–927.
- 572 8. Liu, Z.; Tang, H.; Amini, A.; Yang, X.; Mao, H.; Rus, D.L.; Han, S. BEVFusion: Multi-Task Multi-Sensor  
573 Fusion with Unified Bird's-Eye View Representation. Proceedings of the IEEE International Conference on  
574 Robotics and Automation (ICRA), 2023, pp. 2774–2781.
- 575 9. Liang, T.; Xie, H.; Yu, K.; Xia, Z.; Lin, Z.; Wang, Y.; Tang, T.; Wang, B.; Tang, Z. BEVFusion: A Simple  
576 and Robust LiDAR-Camera Fusion Framework. *Advances in Neural Information Processing Systems* **2022**,  
577 *35*, 10421–10434.
- 578 10. Bai, X.; Hu, Z.; Zhu, X.; Huang, Q.; Chen, Y.; Fu, H.; Tai, C.L. TransFusion: Robust LiDAR-Camera Fusion  
579 for 3D Object Detection with Transformers. Proceedings of the IEEE/CVF Conference on Computer Vision  
580 and Pattern Recognition (CVPR), 2022, pp. 1090–1099.
- 581 11. Caesar, H.; Bankiti, V.; Lang, A.H.; Vora, S.; Lioing, V.E.; Xu, Q.; Krishnan, A.; Pan, Y.; Baldan, G.; Beijbom,  
582 O. nuScenes: A Multimodal Dataset for Autonomous Driving. Proceedings of the IEEE/CVF Conference  
583 on Computer Vision and Pattern Recognition (CVPR), 2020, pp. 11621–11631.
- 584 12. Geiger, A.; Lenz, P.; Urtasun, R. Are We Ready for Autonomous Driving? The KITTI Vision Benchmark  
585 Suite. Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR),  
586 2012, pp. 3354–3361.
- 587 13. Liu, Y.; Jia, T.; Fan, J.; Liu, S.; Zhang, X.; Sun, J. When Autonomous Vehicles Meet Drones: Challenges and  
588 Opportunities for 3D Perception. *arXiv preprint arXiv:2202.07588* **2022**.
- 589 14. Shi, S.; Jiang, C.; Guo, D.; Chen, Z. Drone-Vehicle Cooperative Perception for Autonomous Driving: A  
590 Survey. *IEEE Transactions on Intelligent Transportation Systems* **2024**, *25*, 4784–4800.
- 591 15. Arnold, E.; Dianati, M.; de Temple, R.; Fallah, S. Cooperative Perception for 3D Object Detection in Driving  
592 Scenarios Using Infrastructure Sensors. *IEEE Transactions on Intelligent Transportation Systems*, 2022,  
593 Vol. 23, pp. 1852–1864.
- 594 16. Yu, H.; Luo, Y.; Shu, M.; Huo, Y.; Yang, Z.; Shi, Y.; Guo, Z.; Li, H.; Hu, X.; Yuan, J.; Nie, Z. DAIR-V2X:  
595 A Large-Scale Dataset for Vehicle-Infrastructure Cooperative 3D Object Detection. Proceedings of the  
596 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2022, pp. 21361–21370.
- 597 17. Xu, R.; Xiang, H.; Xia, X.; Han, X.; Li, J.; Ma, J. OPV2V: An Open Benchmark Dataset and Fusion Pipeline  
598 for Perception with Vehicle-to-Vehicle Communication. Proceedings of the IEEE International Conference  
599 on Robotics and Automation (ICRA), 2022, pp. 2583–2589.
- 600 18. Xu, R.; Xiang, H.; Tu, Z.; Xia, X.; Yang, M.H.; Ma, J. V2X-ViT: Vehicle-to-Everything Cooperative Perception  
601 with Vision Transformer. Proceedings of the European Conference on Computer Vision (ECCV), 2022, pp.  
602 107–124.
- 603 19. Xu, R.; Tu, Z.; Xiang, H.; Shao, W.; Zhou, B.; Ma, J. CoBEVT: Cooperative Bird's Eye View Semantic  
604 Segmentation with Sparse Transformers. *Proceedings of the Conference on Robot Learning (CoRL)* **2023**.
- 605 20. Choi, E.H. Crash Factors in Intersection-Related Crashes: An On-Scene Perspective. Technical Report DOT  
606 HS 811 366, National Highway Traffic Safety Administration, 2010.

- 607 21. Yurtsever, E.; Lambert, J.; Carballo, A.; Takeda, K. A Survey of Autonomous Driving: Common Practices  
608 and Emerging Technologies. *IEEE Access* **2020**, *8*, 58443–58469.
- 609 22. Dosovitskiy, A.; Ros, G.; Codevilla, F.; Lopez, A.; Koltun, V. CARLA: An Open Urban Driving Simulator.  
610 Proceedings of the 1st Annual Conference on Robot Learning (CoRL), 2017, pp. 1–16.
- 611 23. OpenPCDet Development Team. OpenPCDet: An Open-source Toolbox for 3D Object Detection from  
612 Point Clouds. <https://github.com/open-mmlab/OpenPCDet>, 2020. Accessed: 2026-01-15.
- 613 24. Qi, C.R.; Su, H.; Mo, K.; Guibas, L.J. PointNet: Deep Learning on Point Sets for 3D Classification and  
614 Segmentation. Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition  
615 (CVPR), 2017, pp. 652–660.
- 616 25. Qi, C.R.; Yi, L.; Su, H.; Guibas, L.J. PointNet++: Deep Hierarchical Feature Learning on Point Sets in a  
617 Metric Space. Advances in Neural Information Processing Systems, 2017, pp. 5099–5108.
- 618 26. Shi, S.; Wang, X.; Li, H. PointRCNN: 3D Object Proposal Generation and Detection from Point Cloud.  
619 Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2019, pp.  
620 770–779.
- 621 27. Zhou, Y.; Tuzel, O. VoxelNet: End-to-End Learning for Point Cloud Based 3D Object Detection. Proceedings  
622 of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2018, pp. 4490–4499.
- 623 28. Yan, Y.; Mao, Y.; Li, B. SECOND: Sparsely Embedded Convolutional Detection. *Sensors* **2018**, *18*, 3337.
- 624 29. Yin, T.; Zhou, X.; Krähenbühl, P. Center-Based 3D Object Detection and Tracking. Proceedings of the  
625 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2021, pp. 11784–11793.
- 626 30. Lang, A.H.; Vora, S.; Caesar, H.; Zhou, L.; Yang, J.; Beijbom, O. PointPillars: Fast Encoders for Object  
627 Detection from Point Clouds. Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern  
628 Recognition (CVPR), 2019, pp. 12697–12705.
- 629 31. Wang, Y.; Fathi, A.; Kunze, A.; Ross, D.A.; Pantofaru, C.; Funkhouser, T.; Solomon, J. Pillar-based Object  
630 Detection for Autonomous Driving. Proceedings of the European Conference on Computer Vision (ECCV),  
631 2020, pp. 18–34.
- 632 32. Shi, S.; Guo, C.; Jiang, L.; Wang, Z.; Shi, J.; Wang, X.; Li, H. PV-RCNN: Point-Voxel Feature Set Abstraction  
633 for 3D Object Detection. Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern  
634 Recognition (CVPR), 2020, pp. 10529–10538.
- 635 33. Chen, X.; Ma, H.; Wan, J.; Li, B.; Xia, T. Multi-View 3D Object Detection Network for Autonomous Driving.  
636 Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2017, pp.  
637 1907–1915.
- 638 34. Philion, J.; Fidler, S. Lift, Splat, Shoot: Encoding Images from Arbitrary Camera Rigs by Implicitly  
639 Unprojecting to 3D. Proceedings of the European Conference on Computer Vision (ECCV), 2020, pp.  
640 194–210.
- 641 35. Li, Y.; Yu, A.W.; Meng, T.; Caine, B.; Ngiam, J.; Peng, D.; Shen, J.; Lu, Y.; Zhou, D.; Le, Q.V.; Yuille, A.; Tan,  
642 M. DeepFusion: LiDAR-Camera Deep Fusion for Multi-Modal 3D Object Detection. Proceedings of the  
643 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2022, pp. 17182–17191.
- 644 36. Pang, S.; Morris, D.; Radha, H. CLOCs: Camera-LiDAR Object Candidates Fusion for 3D Object Detection.  
645 Proceedings of the IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), 2020, pp.  
646 10386–10393.
- 647 37. Ku, J.; Mozifian, M.; Lee, J.; Harakeh, A.; Waslander, S.L. Joint 3D Proposal Generation and Object Detection  
648 from View Aggregation. Proceedings of the IEEE/RSJ International Conference on Intelligent Robots and  
649 Systems (IROS), 2018, pp. 1–8.
- 650 38. Nobis, F.; Geisslinger, M.; Weber, M.; Betz, J.; Lienkamp, M. A Deep Learning-Based Radar and Camera  
651 Sensor Fusion Architecture for Object Detection. *Sensors* **2021**, *21*, 2321.
- 652 39. Lin, T.Y.; Dollár, P.; Girshick, R.; He, K.; Hariharan, B.; Belongie, S. Feature Pyramid Networks for Object  
653 Detection. Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR),  
654 2017, pp. 2117–2125.
- 655 40. Jocher, G.; Chaurasia, A.; Qiu, J. Ultralytics YOLOv8. <https://github.com/ultralytics/ultralytics>, 2023.  
656 Accessed: 2026-01-15.
- 657 41. Redmon, J.; Divvala, S.; Girshick, R.; Farhadi, A. You Only Look Once: Unified, Real-Time Object Detection.  
658 Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2016, pp.  
659 779–788.

- 660 42. Lin, T.Y.; Maire, M.; Belongie, S.; Hays, J.; Perona, P.; Ramanan, D.; Dollár, P.; Zitnick, C.L. Microsoft COCO:  
661 Common Objects in Context. Proceedings of the European Conference on Computer Vision (ECCV), 2014,  
662 pp. 740–755.
- 663 43. Kuhn, H.W. The Hungarian method for the assignment problem. *Naval Research Logistics Quarterly* **1955**,  
664 2, 83–97.
- 665 44. Everingham, M.; Van Gool, L.; Williams, C.K.I.; Winn, J.; Zisserman, A. The Pascal Visual Object Classes  
666 (VOC) Challenge. *International Journal of Computer Vision* **2010**, 88, 303–338.
- 667 45. Conover, W. *Practical Nonparametric Statistics*, 3rd ed.; John Wiley & Sons: New York, 1999.
- 668 46. Yue, X.; Zhang, Y.; Zhao, S.; Sangiovanni-Vincentelli, A.; Keutzer, K.; Gong, B. Domain Randomization for  
669 Transferring Deep Neural Networks from Simulation to the Real World. *arXiv preprint arXiv:1903.11499*  
670 **2019**.
- 671 47. Sun, P.; Kretzschmar, H.; Dotiwalla, X.; Chouard, A.; Patnaik, V.; Tsui, P.; Guo, J.; Zhou, Y.; Chai, Y.; Caine,  
672 B.; others. Scalability in Perception for Autonomous Driving: Waymo Open Dataset. Proceedings of the  
673 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2020, pp. 2446–2454.

674 © 2026 by the authors. Submitted to *Sustainability* for possible open access publication  
675 under the terms and conditions of the Creative Commons Attribution (CC BY) license  
676 (<http://creativecommons.org/licenses/by/4.0/>).