



# CBILR: Camera Bi-Directional LiDAR-Radar Fusion for Robust Perception in Autonomous Driving

---

Arthur Nigmatzyanov, Gonzalo Ferrer and Dzmitry Tsetserukou

EasyChair preprints are intended for rapid dissemination of research results and are integrated with the rest of EasyChair.

June 14, 2024

# CBILR: Camera Bi-directional LiDAR-Radar Fusion for Robust Perception in Autonomous Driving

Arthur Nigmatzyanov  
Skoltech University  
Moscow, Russia  
Artur.Nigmatzyanov@skoltech.ru

Gonzalo Ferrer  
Skoltech University  
Moscow, Russia  
G.Ferrer@skoltech.ru

Dzmitry Tsetserukou  
Skoltech University  
Moscow, Russia  
D.Tsetserukou@skoltech.ru

## Abstract

Safe and reliable autonomous driving hinges on robust perception under challenging environments. Multi-sensor fusion, particularly camera-LiDAR-Radar integration, plays a pivotal role in achieving this goal. Different sensors have specific advantages and disadvantages. Existing pipelines are often constrained by adverse weather conditions, where cameras and LiDAR suffer significant degradation. This paper introduces the Camera Bi-directional LiDAR-Radar (CBILR) fusion pipeline, which leverages the strengths of sensors to enhance LiDAR and Radar point clouds. CBILR innovates with a bi-directional pre-fusion step between LiDAR and Radar, leading to richer feature representations. First, pre-fusion combines LiDAR and Radar points to compensate for individual sensor weaknesses. Next, the pipeline fuses the pre-fused features with camera features in the bird's eye view (BEV) space, resulting in a comprehensive multi-modal representation. Experiments have demonstrated that CBILR outperforms state-of-the-art pipelines, achieving superior robustness in challenging weather scenarios.

## Keywords

Fusion, Self-driving, Autonomous vehicle, Camera, LiDAR, Radar, Weather Conditions

### ACM Reference Format:

Arthur Nigmatzyanov, Gonzalo Ferrer, and Dzmitry Tsetserukou. 2024. CBILR: Camera Bi-directional LiDAR-Radar Fusion for Robust Perception in Autonomous Driving. In *Proceedings of 9th International Workshop on Sensor-Based Activity Recognition and Artificial Intelligence (iWOAR2024)*. ACM, New York, NY, USA, 7 pages.

## 1 Introduction

For self-driving systems, it is crucial to develop a fast and accurate 3D object detector that predicts the bounding boxes and categories of road objects. Nowadays, cameras, LiDARs, and Radars are often used in advanced systems such as drones, robots and autonomous vehicles. Many authors only use particular sensors to solve perception problems. This can lead to a generalization problem, because there is a high probability that one type of sensor will be more relevant than others for certain real-world scenarios. Each sensor has advantages and disadvantages. From cameras, we can only obtain

color and texture information about objects after projective transformation of captured 3D scene into a 2D plane and long stages of post-processing raw images [23], which is the field of color science. For this reason, cameras cannot provide accurate depth information (especially in low light conditions) compared to Radars and LiDARs that operate directly in 3D space [8, 12]. However, researchers continue to develop perception algorithms that rely only on cameras because it is a more cost-effective approach [10, 29].

### 1.1 Fusion approaches

Sensor fusion is an essential topic in many perception systems. A lot of papers [28, 31] are devoted to LiDAR-camera fusion because LiDARs have higher resolution, are less sparse than Radars and can provide accurate measurements at close range. Since Radar antennas are often installed horizontally, they cannot capture sufficient vertical height information [26]. For voxel representation, a highly sparse point cloud means that some voxels contain too few points for processing.

Although LiDARs can provide accurate geometric information about a scene, they do not perform as well as Radars at long distances and can introduce noise when the object is moving [7]. In [18, 19] the authors use Radar-camera fusion for 3D object detection and tracking. Since such sensors in many cases have opposite advantages and disadvantages, it is ideal to use multiple sensors [3, 13] for robust performance in a variety of scenarios and conditions. We have developed a fusion pipeline focused on improving sensors that can withstand adverse weather conditions.

There are several strategies for sensor fusion. Early fusion directly combines sensor inputs before feeding them into shared feature extractors. Late fusion processes sensor inputs independently and then combines the output results. Mid-level fusion [11] provides an intermediate representation for each sensor before the final fusion step.

### 1.2 BEV Perception

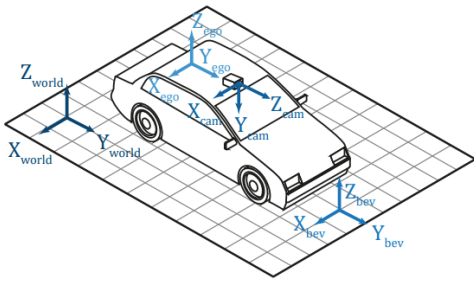
A unified representation is necessary to make it easier to transfer knowledge and combine features from different modalities [9]. The vast majority of modern perception methods use a bird's eye view (BEV) representation to describe a 3D scene [22, 33]. BEV is an informal perception standard for autonomous driving scenarios [12]. The BEV coordinate system is a rotation of the camera coordinate system, such that the  $Z$ -axis is aligned with the camera's negative  $Y$  direction, and is placed a fixed distance below the camera as shown in the Figure 1. Data from different modalities are used to provide complementary knowledge such as precise locations from point clouds and rich context from images. For example, fusion

---

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

*iWOAR2024, Sep 26–27, 2024, Potsdam, Germany*

© 2024 Copyright held by the owner/author(s). Publication rights licensed to ACM.



**Figure 1: Four main coordinate systems: world, ego - vehicle, camera and bird's eye view [22].**

algorithms translate features from different sensors into the BEV representation and then combine them [13].

Cameras are typically mounted on vehicles parallel to the ground and facing outward. For this reason, images are captured in a Perspective View (PV), which is orthogonal to BEV. Objects of the same shape and size in 3D space can have very different representation in the image plane because of their distance from the camera. The BEV representation does not have scale and occlusion problems compared to PV representation [8]. The transformation from PV to BEV is the inverse perspective map problem, and it can have more than one solution. Before the deep learning era, many works tackled this problem by using a homography transformation matrix because of its computational efficiency. Inverse Perspective Mapping (IPM) has been proposed to address this challenging mapping problem [16, 17]. IPM-based methods assume that all points are on the ground plane, sacrificing height variation. In complex real-world scenarios, 3D objects like vehicles possess *height* and such transformations can cause noticeable artifacts.

In recent years, data-driven methods have been widely used in complex systems such as self-driving vehicles. Data-driven PV-BEV transformation methods can be divided into three main groups: depth-based, MLP-based, and transformer-based approaches [16]. Depth-based methods estimate the depth distribution of the each image pixel along the ray (coming from the camera) that intersects objects in the environment. This allows to elevate the 2D features to 3D, and then obtain the BEV representations from 3D through dimensionality reduction. Depth-based PV-to-BEV methods can be divided into two classes depending on the using representation: point-based and voxel-based methods. Point-based methods are straightforward, they directly utilize depth estimation to convert pixels into point clouds. Examples: Pseudo-LiDAR [24], Pseudo-LiDAR++ [27], AM3D [15], PatchNet [14]. Voxel-based method discretize the 3D space to build a regular structure for feature transformation. The disadvantage of this approach is the loss of detailed local spatial information within each voxel. The advantage is that voxels are more effective at covering large-scale scene structure, they are more efficient for 3D scene understanding.

Another approach is to utilize a variational encoder-decoder or MLP to learn implicit representations of camera calibrations to project PV features to BEV. MLP plays the role of a universal approximator of the mapping function from PV to BEV [16]. MLP-based methods focus primarily on working with a single image. The drawback of MLP-based methods is that the learned weights

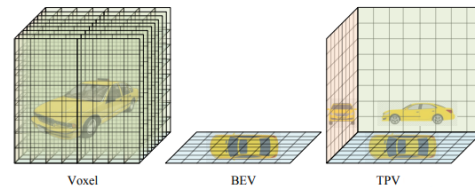
are fixed and not data dependent:

$$Y = WX, W \neq X$$

Transformer-based methods employ a top-down strategy constructing BEV queries and searching corresponding features in perspective images through cross-attention mechanism. These methods are more expressive, but hard to train.

### 1.3 BEV representation vs voxel-based

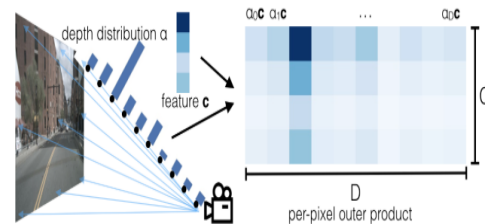
A voxel-based scene representation cannot provide computational efficiency because such representation describes a 3D scene with dense cubic features  $V \in \mathbb{R}^{H \times W \times D \times C}$  where  $H, W, D$  are the spatial resolution of the voxel space and  $C$  is the feature dimension. BEV provides the 3D scene with a 2D feature map  $B \in \mathbb{R}^{H \times W \times C}$  which encodes the top view of the scene. This represents the positional information of the ground plane by accumulating voxel features along the vertical  $z$ -axis Figure 2. The height dimension contains less information than the other two dimensions [5]. It is important to note that some researches do not directly use the BEV representation. In [5] for semantic prediction task due to the lack of  $z$ -axis information authors propose Tri-Perspective View (TPV) representation Figure 2.



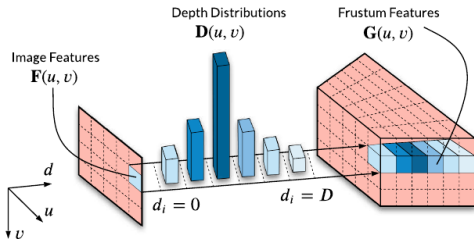
**Figure 2: Voxel, BEV and TPV representations. Voxel representation is more informative, but cannot provide computation efficiency [5].**

### 1.4 Camera-to-BEV View Transform

Transforming from a camera view to a bird's eye view is complex because the depth associated with each camera feature pixel can be ambiguous. The idea of camera-to-BEV transformation is based on projective geometry (see Figure 3). The process of monocular depth estimation involves generating a unique depth value for each pixel in an image. The state-of-the-art approach involves predicting a categorical distribution of depth for each pixel in the image [13, 20, 21]. This technique is known as *feature lifting* [20].



**Figure 3: The important stage of Camera-to-BEV View Transform is estimation of a categorical depth distribution [20].**



**Figure 4: Frustum generating: each feature pixel  $F(u, v)$  is weighted by its depth distribution probabilities  $D(u, v)$  of belonging to  $D$  discrete depth bins to generate frustum features  $G(u, v)$  [21].**

In [21], the model utilize the estimated categorical depth distributions to “lift” an input image in 3D, generating a frustum-shaped point cloud of contextual features. The frustum feature grid is then transformed into a voxel grid using specific camera calibration parameters, and then collapsed into a BEV feature grid. All steps are well-illustrated in the paper [21]. By associating image features with estimated depths, image information can be projected into 3D space using a frustum feature network.

The input to the frustum feature network is an image  $I \in \mathbb{R}^{W_I \times H_I \times 3}$ , where  $W_I, H_I$  are the image width and height. The network output is a frustum feature grid  $G \in \mathbb{R}^{W_F \times H_F \times D \times C}$ , where  $W_F, H_F$  are the width and height of the image feature representation,  $D$  is the number of discretized depth bins, and  $C$  is the number of feature channels. If we have  $N$  cameras, the full size of the frustum features is  $N \times W_F \times H_F \times D$ .

Let’s denote  $(u, v, c)$  as a coordinate in image features  $F$  and  $(u, v, d_i)$  as a coordinate in categorical depth distributions  $D$ , where  $(u, v)$  is the location of feature pixel,  $c$  is the channel index, and  $d_i$  is the depth bin index. In order to create a frustum feature grid  $G$ , each feature pixel  $F(u, v)$  is weighted by its associated depth bin probabilities in  $D(u, v)$ . It adds a new depth axis  $d_i$ , as shown in Figure 4. The outer product can be used to weight feature pixels:

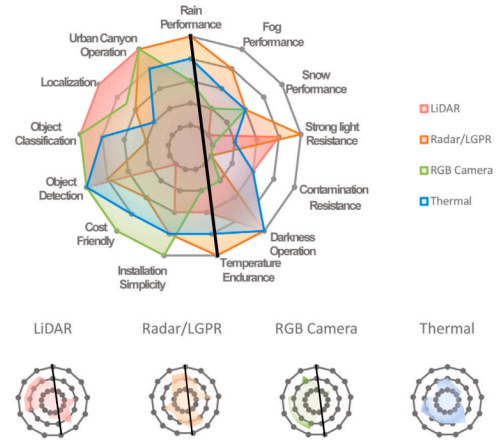
$$G(u, v) = D(u, v) \otimes F(u, v) \quad (1)$$

where  $D(u, v)$  is the predicted depth distribution and  $G(u, v)$  is a matrix  $D \times C$ . The outer product is calculated for each pixel to generate frustum features  $G \in \mathbb{R}^{W_F \times H_F \times D \times C}$ . The next steps are voxel transformation using the camera calibration matrix [21] and collapsing to BEV.

For example, Bevfusion [13] converts camera features into a point cloud, aggregates it with BEV pooling and flattens it along the  $z$ -axis. Such algorithms can be related to the Lift-Splat category [20, 21, 32].

## 1.5 Motivation

In [30] authors made a detailed review of how autonomous vehicles perceive the environment under adverse weather conditions. They summarized the strengths and weaknesses of each sensor in the chart 5. As we can see, camera sensors are the most sensitive to environmental conditions. But, not all parts of an image typically contain destructive information. For example, in the Figure 6 certain



**Figure 5: Sensor performance and characteristics [30].**



**Figure 6: Camera in rain condition.**

regions of the image provide crucial details about the objects in the scene.

Recent works [3, 13] have used a mid-level fusion approach to aggregate features from all modalities. Combining the representations of different modalities allows to solve perception problems in adverse weather conditions (see the Table 1).

**Table 1: Sensor fusion and target weather conditions. ”L”, ”C” and ”R” represent LiDAR, Camera, and Radar modalities respectively.**

Sensor fusion	Configuration	Target weather
Bi-LRFusion (2023)	R + L	Fog
RadarNet (2020)	R + L	Rain
MVDNet (2021)	R + L	Fog
Liu (2021)	R + C	Rain, fog, nighttime
Rawashdeh (2021)	C + L + R	Snow
SLS-Fusion (2021)	L + C	Fog
Radecki (2016)	L + R + C	Wet conditions

In last time LiDARs and Radars sensors were significantly improved in terms of spatial resolution, accuracy, velocity measurement and resistance to adverse weather conditions [1].

Because images received from cameras may have artifacts and be overlapped for any reason, using visual transformers may not be as efficient as convolutional neural networks. Transformers take and process every patch of an image, even areas that may not be

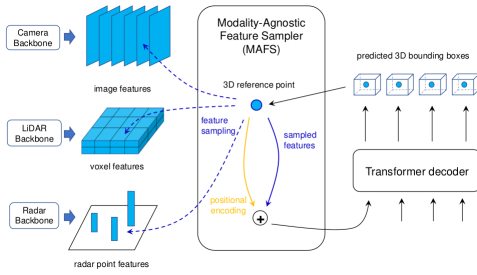


Figure 7: The illustration of the FUTR3D pipeline [3].

relevant for the specific task. For this reason, authors prefer to use convolutional layers first in neural networks as preprocessing part [3, 12, 13]. Also ViTs require *enormous amounts of data and computation to train*, and in some cases have longer inference time. For this reason, researchers avoid the unwise use of ViTs [2, 3, 12, 13].

## 2 Related Work

In the mid-fusion approach, features are combined after *feature extraction* [4, 11].

**FUTR3D.** Particularly, in [3] every modality is encoded in its own coordinate. This framework (see Figure 7) does not assume any particular modalities and their model architectures. For this reason FUTR3D can work with any selected feature encoders. Researches used three types of data: LiDAR point cloud, Radar point cloud, and multi-view camera images.

VoxelNet was used to encode LiDAR point clouds as multi-scale Bird’s-eye view (BEV) feature maps  $\{\mathcal{F}_{\text{lid}}^j \in \mathbb{R}^{C \times H_j \times W_j}\}_{j=1}^m$ , where  $H_i \times W_i$  is the size of the  $i$ -th BEV feature map,  $m$  is the count of feature maps. Radar points  $\{r_j\}_{j=1}^N \in \mathbb{R}^{C_{ri}}$  are pillarized into 0.8 m pillars.

Then MLP  $\Phi_{\text{rad}}$  is used to achieve per-pillar features  $\mathcal{F}_{\text{rad}}^j = \Phi_{\text{rad}}(r_j) \in \mathbb{R}^{C_{ro}}$ , where  $C_{ro}$  is the number of encoded Radar features. In this way the Radar BEV feature map  $\mathcal{F}_{\text{rad}} \in \mathbb{R}^{C_{ro} \times H \times W}$  is obtained. It is also assumed that there are  $N$  surrounding cameras installed in the car. It is supposed that each camera has taken  $m$  images. For image feature extraction ResNet is used. It outputs multi-scale features for each image, denoted as  $\mathcal{F}_{\text{cam}}^k = \{\mathcal{F}_{\text{cam}}^{kj} \in \mathbb{R}^{C \times H_j \times W_j}\}_{j=1}^m$  for the  $k$ -th camera. So, after camera backbone there are  $m$  image feature maps for each camera.

A transformer decoder uses queries to predict 3D bounding boxes. The predicted boxes can be repeatedly sent back into the transformer decoder and MAFS to refine the predictions. Modality-Agnostic Feature Sampler (MAFS) creates and aggregate features from each modality based on the 3D reference point (initial position) of each query. The 3D reference point is ground to collect features from multiple sources. The input of detection head is a set of *object queries*  $Q = \{q_i\}_{i=1}^{N_q} \subset \mathbb{R}^C$ , and features from all sensors, where  $C$  is the output channel of BEV feature map after processing LiDAR point clouds with VoxelNet. MAFS updates each query by sampling features from each sensor feature and fusing them. Object queries are updated using self-attention modules and FFN.

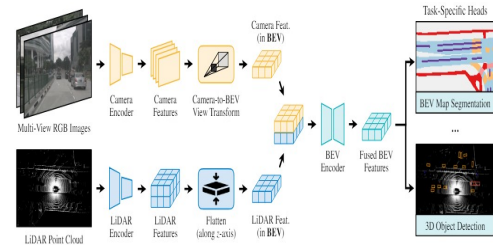


Figure 8: The illustration of the BEVFusion pipeline [13].

**BEVFusion.** BEVFusion [13] is the state-of-the-art fusion pipeline on the nuScenes dataset. It fuses camera and LiDAR sensors in BEV space to perform 3D detection and tracking simultaneously. BEVFusion uses an effective method of transforming camera images into a BEV representation and combining them with LiDAR BEV features using convolutional layers.

Like FUTR3D, BEVFusion provides independent camera and LiDAR streams (see Figure 8).

**Bi-LRFusion.** Radar provides long range detection and velocity hints, while LiDAR is better at capturing the object’s 3D shape [26]. To fully utilize the advantages of combining LiDAR and Radar, the authors enhance the Radar features to make them more powerful before the final fusion.

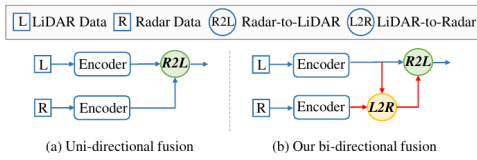
*Encoding of LiDAR Features.* The LiDAR encoding layer is used to extract voxelwise features. This process consists of the following steps:

- dividing LiDAR points into voxels.
- taking all points in the same voxel as input and using a multi-layer perception (MLP) to extract pointwise features.
- using elementwise max pooling to obtain the locally aggregated features for each voxel.
- 3D Voxel Backbone composed of 3D sparse convolutional layers and 3D sub-manifold convolutional layers [25].
- producing a LiDAR BEV feature map by stacking volume features along the Z-axis.

*Radar Feature Encoding.* By utilizing the Pillar Feature Backbone [6], the Radar point cloud is converted into a series of pillars. It is important to note that the Radar point’s value on the z-axis is set to the height of the Radar sensor by default.

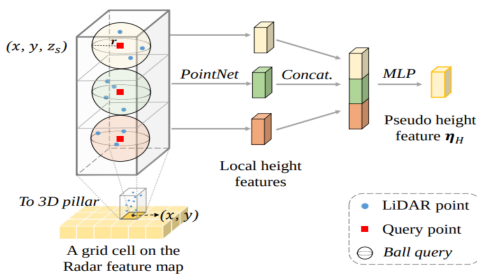
The Bi-LRFusion algorithm consists of the following steps:

- encoding BEV features: the LiDAR feature stream and the Radar feature stream receive the input LiDAR and Radar points to create BEV features.
- enhancing Radar features by combining LiDAR raw points and Radar features through the LiDAR-to-Radar (L2R) fusion module: due to the lack of height information and sparsity, the Radar’s local features are enriched by learning important details from the LiDAR points. To acquire more comprehensive Radar features, for each valid (non-empty) grid cell on the Radar feature map, the nearby LiDAR data (height and BEV perspectives) are queried and grouped with the Radar features.

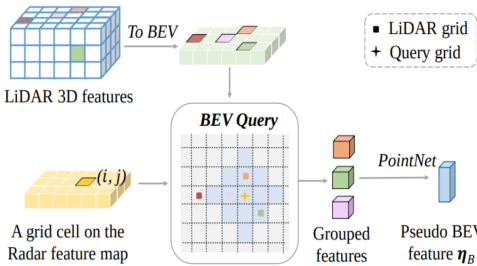


**Figure 9: The difference between Bi-LRFusion and standard uni-directional fusion [26].**

- the Radar-to-LiDAR (R2L) fusion step: combining LiDAR features with the enhanced Radar features in a *unified BEV representation*.
- predicting 3D bounding boxes for dynamic objects using the obtained BEV features.



**Figure 10: The pseudo height feature formation [26].**

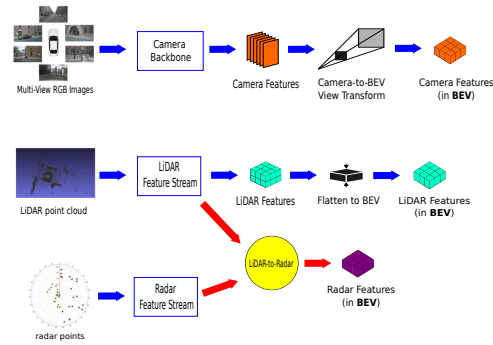


**Figure 11: The pseudo BEV feature formation [26].**

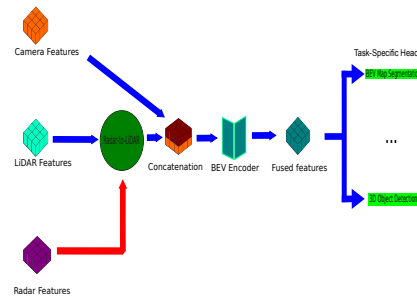
As shown in the Figure 9, the authors propose a bidirectional fusion scheme. L2R Fusion Module consists of two submodules: Query-based Height feature Fusion block and Query-based BEV feature Fusion block, as shown in Figures 10 and 11 respectively. To form *pseudo-Radar height features* the LiDAR raw points are aggregated, and the LiDAR BEV features are aggregated into *pseudo-Radar BEV*. Then the pseudo-Radar height and pseudo-Radar BEV are concatenated to the Radar BEV features. The next step is the Radar-to-LiDAR (R2L) fusion in a unified BEV.

### 3 Method

Since both LiDARs and Radars operate in 3D space and they are more reliable than cameras under adverse environmental conditions, we first do their prefusion [26]. The figures 12 and 13 illustrate



**Figure 12: This is the first part of the pipeline. Transformation of raw LiDAR, Radar points, and images into a BEV representation.**



**Figure 13: This is the second part of the pipeline. We fuse all the BEV representations together, encode the result and then send it to specific heads.**

the concept of our pipeline. We have split the illustration of it into two parts. First of all, we want to make the feature extraction like in [13]. Then, we use a specific transformation for a particular sensor to represent the extracted feature in the BEV. The next steps are bilateral LiDAR-Radar prefusion and image feature concatenation for final fusion.

The LiDAR-to-Radar step enriches the Radar points. Similar to [26], we mix Radar points with LiDAR points before encoding. This eliminates the lack of Radar points per voxel, especially in the height direction. LiDAR points are encoded with SECOND [25], Radar points are encoded with PillarFeatureNet [6]. As a BEVFusion, this pipeline can be used for different tasks such as segmentation and 3D object detection. This article includes a link to GitHub for more information. It is recommended to match the configuration file with the pipeline Figures 12, 13.

### 4 Experiments

The Nuscenes Dataset is widely used dataset for vision-centric perception with six calibrated cameras covering a 360-degree horizontal FOV, 1 LiDAR and 5 Radars. The camera image resolution is 1600×900. Nuscenes consists of 1000 scenes, each one of them is 20 seconds long. 850 scenes are for training/validation and 150 for testing.

The most commonly used criterion for BEV Detection is average precision (AP) and the mean average precision (mAP) over different classes. For BEV Segmentation, IoU for each class and mIoU over all classes. The Average Precision (AP) metric is extended from 2D to the 3D space:

$$AP = \int_0^1 \max \{p(r' \mid r' \geq r)\} dr \quad (3)$$

where  $p(r)$  is the precision-recall curve. The difference between 2D AP and 3D AP is the matching criteria between ground truth and predictions when calculating precision and recall.

Instead of IoU to select TP, NuScenes proposes  $AP_{\text{center}}$  where a predicted object is matched to a ground truth object if the distance of their center locations on the ground (BEV) plane is below a certain threshold  $d$ . The  $AP_{\text{center}}$  is calculated under different distance thresholds:  $\mathbb{D} = \{0.5, 1, 2, 4\}$  meters. The mAP is computed by averaging the  $AP_{\text{center}}$  over all matching thresholds and all classes  $\mathbb{C}$ :  $\text{mAP} = \frac{1}{|\mathbb{C}||\mathbb{D}|} \sum_{c \in \mathbb{C}} \sum_{d \in \mathbb{D}} AP_{c,d}$ . NuScenes Detection Score (NDS) is further proposed to take both  $AP_{\text{center}}$  and the error of other parameters, i.e. size, heading, velocity, into consideration.

In our experiments we compared BEVFusion [13] and BiFusion [26] with our method (see the Table 2).

**Table 2: Results of experiments. "L", "C" and "R" represent LiDAR, Camera, and Radar modalities respectively.**

Model	mAP	NDS
Bi-LRFusion (R + L)	62.3	65.54
BEVFusion (C + L)	68.57	71.40
CBILR (C + R + L)	<b>71.09</b>	<b>73.36</b>

Experiments show that it is important to use all modalities in a clever way. Combining different modalities helps to overcome the limitations of individual sensors.

## 5 Conclusion

This work has demonstrated CBILR, a promising multi-sensor fusion framework that aims to improve perception robustness for autonomous vehicles. It has addressed the critical challenge of limited sensor performance in adverse weather conditions, a significant hurdle on the path to achieving truly autonomous navigation. CBILR aims to overcome the limitations of existing fusion methods by using the Bi-LRFusion module. This module promotes a mutually beneficial LiDAR/radar relationship, allowing each to benefit from the other's strengths.

Bi-LRFusion module enrich the sparse Radar point cloud with Lidar original points in two directions (R2L submodule) and then add enhanced Radar feature representation to Lidar feature representation (L2R submodule). This enriched representation significantly enhances the overall perception accuracy, especially under challenging weather scenarios.

The experiments show that using multiple sensors for fusion increases reliability in challenging weather conditions. Previous works uniformly combine all sensors together. They do not consider the weaknesses of different sensors. By utilizing Bi-LRFusion and promoting a thorough understanding of the environment, CBILR

strives to lead the way into a new era of strong and adaptable perception. This effort aims to bring autonomous vehicles closer to the ultimate goal of safe and reliable operation in all conditions.

## References

- [1] [n. d.]. 4D LiDARs vs 4D RADARS: Why the LiDAR vs RADAR comparison is more relevant today than ever. <https://www.thinkautonomous.ai/blog/fmcw-lidars-vs-imaging-radars/>. Accessed: 2024-02-21.
- [2] [n. d.]. Transforming Computer Vision: The Rise of Vision Transformers And Its Impact. <https://clck.ru/38qgMw>. Accessed: 2023-9-13.
- [3] Xuanyao Chen, Tianyuan Zhang, Yue Wang, Yilun Wang, and Hang Zhao. 2023. Futr3d: A unified sensor fusion framework for 3d detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 172–181.
- [4] Tengfeng Huang, Zhe Liu, Xiwu Chen, and Xiang Bai. 2020. EPNet: Enhancing Point Features with Image Semantics for 3D Object Detection. (July 2020).
- [5] Yuanhui Huang, Wenzhao Zheng, Yunpeng Zhang, Jie Zhou, and Jiwen Lu. 2023. Tri-perspective view for vision-based 3d semantic occupancy prediction. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 9223–9232.
- [6] Alex H Lang, Sourabh Vora, Holger Caesar, Lubing Zhou, Jiong Yang, and Oscar Beijbom. 2019. Pointpillars: Fast encoders for object detection from point clouds. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 12697–12705.
- [7] Huadong Li, Minhao Jing, Jiajun Liang, Haoqiang Fan, and Rehne Ji. 2023. Sparse Beats Dense: Rethinking Supervision in Radar-Camera Depth Completion. *arXiv:2312.00844* (2023).
- [8] Hongyang Li, Chonghao Sima, Jifeng Dai, Wenhao Wang, Lewei Lu, Huijie Wang, Jia Zeng, Zhiqi Li, Jiazhi Yang, Hanming Deng, et al. 2023. Delving into the devils of bird's-eye-view perception: A review, evaluation and recipe. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2023).
- [9] Yanwei Li, Yilun Chen, Xiaojuan Qi, Zeming Li, Jian Sun, and Jiaya Jia. 2022. Unifying voxel-based representation with transformer for 3d object detection. *Advances in Neural Information Processing Systems* 35 (2022), 18442–18455.
- [10] Y. Li, Z. Ge, G. Yu, J. Yang, Z. Wang, Y. Shi, J. Sun, and Z. Li. 2023. BEVDepth: Acquisition of Reliable Depth for Multi-View 3D Object Detection. In *AAAI Technical Track on Computer Vision II*, Vol. 37, 1477–1485.
- [11] Ming Liang, Bin Yang, Shenlong Wang, and Raquel Urtasun. 2018. Deep continuous fusion for multi-sensor 3d object detection. In *Proceedings of the European conference on computer vision (ECCV)*, 641–656.
- [12] Tingting Liang, Hongwei Xie, Kaicheng Yu, Zhongyu Xia, Zhiwei Lin, Yongtao Wang, Tao Tang, Bing Wang, and Zhi Tang. 2022. BEVFusion: A simple and robust lidar-camera fusion framework. *Advances in Neural Information Processing Systems* 35 (2022), 10421–10434.
- [13] Zhijian Liu, Haotian Tang, Alexander Amini, Xinyu Yang, Huihui Mao, Daniela L Rus, and Song Han. 2023. BEVFusion: Multi-task multi-sensor fusion with unified bird's-eye view representation. In *2023 IEEE international conference on robotics and automation (ICRA)*. IEEE, 2774–2781.
- [14] Xinzhu Ma, Shinan Liu, Zhiyi Xia, Hongwen Zhang, Xingyu Zeng, and Wanli Ouyang. 2020. Rethinking pseudo-lidar representation. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XIII 16*. Springer, 311–327.
- [15] Xinzhu Ma, Zhihui Wang, Haojie Li, Pengbo Zhang, Wanli Ouyang, and Xin Fan. 2019. Accurate monocular 3d object detection via color-embedded 3d reconstruction for autonomous driving. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 6851–6860.
- [16] Yuexin Ma, Tai Wang, Xuyang Bai, Huitong Yang, Yuenan Hou, Yaming Wang, Yu Qiao, Ruigang Yang, Dinesh Manocha, and Xinge Zhu. 2022. Vision-centric bev perception: A survey. *arXiv preprint arXiv:2208.02797* (2022).
- [17] Hanspeter A Mallot, Heinrich H Bülthoff, JJ Little, and Stefan Bohrer. 1991. Inverse perspective mapping simplifies optical flow computation and obstacle detection. *Biological cybernetics* 64, 3 (1991), 177–185.
- [18] Ramin Nabati, Landon Harris, and Hairong Qi. 2021. CFTrack: Center-based Radar and Camera Fusion for 3D Multi-Object. *arXiv preprint arXiv:2107.05150* (2021).
- [19] Ramin Nabati and Hairong Qi. 2020. CenterFusion: Center-based Radar and Camera Fusion for 3D Object Detection. *arXiv preprint arXiv:2011.04841* (2020).
- [20] Jonah Philion and Sanja Fidler. 2020. Lift, splat, shoot: Encoding images from arbitrary camera rigs by implicitly unprojecting to 3d. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XIV 16*. Springer, 194–210.
- [21] Cody Reading, Ali Harakeh, Julia Chae, and Steven L Waslander. 2021. Categorical depth distribution network for monocular 3d object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 8555–8564.
- [22] Thomas Roddick. 2021. *Learning Birds-Eye View Representations for Autonomous Driving*. Ph. D. Dissertation.

- [23] Carlota Salinas, Roemi Fernandez, Hector Montes, and Armada Manuel. 2015. A New Approach for Combining Time-of-Flight and RGB Cameras Based on Depth-Dependent Planar Projective Transformations. *Sensors* 15 (2015), 24615–24643.
- [24] Yan Wang, Wei-Lun Chao, Divyansh Garg, Bharath Hariharan, Mark Campbell, and Kilian Q Weinberger. 2019. Pseudo-lidar from visual depth estimation: Bridging the gap in 3d object detection for autonomous driving. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 8445–8453.
- [25] Yan Yan, Yuxing Mao, and Bo Li. 2018. Second: Sparsely embedded convolutional detection. *Sensors* 18, 10 (2018), 3337.
- [26] Jiajun Yingjie Wang, Deng, Yao Li, Jinshui Hu, Cong Liu, Yu Zhang, Jianmin Ji, Wanli Ouyang, and Yanyong Zhang. 2023. Bi-LRFusion: Bi-Directional LiDAR-Radar Fusion for 3D Dynamic Object Detection. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [27] Yurong You, Yan Wang, Wei-Lun Chao, Divyansh Garg, Geoff Pleiss, Bharath Hariharan, Mark Campbell, and Kilian Q Weinberger. 2019. Pseudo-lidar++: Accurate depth for 3d object detection in autonomous driving. *arXiv preprint arXiv:1906.06310* (2019).
- [28] Ce Zhang, Chengjie Zhang, Yiluan Guo, Lingji Chen, and Michael Happold. 2023. Motiontrack: end-to-end transformer-based multi-object tracking with lidar-camera fusion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 151–160.
- [29] Tianyuan Zhang, Xuanyao Chen, Yue Wang, Yilun Wang, and Hang Zhao. 2022. Mutr3d: A multi-camera tracking framework via 3d-to-2d queries. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 4537–4546.
- [30] Yuxiao Zhang, Alexander Carballo, Hanting Yang, and Kazuya Takeda. 2023. Perception and sensing for autonomous vehicles under adverse weather conditions: A survey. *ISPRS Journal of Photogrammetry and Remote Sensing* 196 (2023), 146–177.
- [31] Huazan Zhong, Hao Wang, Zhengrong Wu, Chen Zhang, Yongwei Zheng, and Tao Tang. 2021. A survey of LiDAR and camera fusion enhancement. In *10th International Conference of Information and Communication Technology*.
- [32] Hongyu Zhou, Zheng Ge, Zeming Li, and Xiangyu Zhang. 2023. Matrixvt: Efficient multi-camera to bev transformation for 3d perception. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 8548–8557.
- [33] Zijian Zhu, Yichi Zhang, Hai Chen, Yinpeng Dong, Shu Zhao, Wenbo Ding, Jiachen Zhong, and Shibao Zheng. 2023. Understanding the Robustness of 3D Object Detection With Bird's-Eye-View Representations in Autonomous Driving. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 21600–21610.