

# EPMF: Efficient Perception-aware Multi-sensor Fusion for 3D Semantic Segmentation

Zhuangwei Zhuang, Sitao Chen, Rong Li, Kui Jia, Qicheng Wang, Yuanqing Li, Mingkui Tan<sup>†</sup>

**Abstract**—We study multi-sensor fusion for 3D semantic segmentation that is important to scene understanding for many applications, such as autonomous driving and robotics. For example, for autonomous cars equipped with RGB cameras and LiDAR, it is crucial to fuse complementary information from different sensors for robust and accurate segmentation. Existing fusion-based methods, however, may not achieve promising performance due to the vast difference between the two modalities. In this work, we investigate a collaborative fusion scheme called perception-aware multi-sensor fusion (PMF) to effectively exploit perceptual information from two modalities, namely, appearance information from RGB images and spatio-depth information from point clouds. To this end, we first project point clouds to the camera coordinate using perspective projection. In this way, we can process both inputs from LiDAR and cameras in 2D space while preventing the information loss of RGB images. Then, we propose a two-stream network that consists of a LiDAR stream and a camera stream to extract features from the two modalities, separately. The extracted features are fused by effective residual-based fusion modules. Moreover, we introduce additional perception-aware losses to measure the perceptual difference between the two modalities. Last, we propose an improved version of PMF, *i.e.*, EPMF, which is more efficient and effective by optimizing data pre-processing and network architecture under perspective projection. Specifically, we propose cross-modal alignment and cropping to obtain tight inputs and reduce unnecessary computational costs. We then explore more efficient contextual modules under perspective projection and fuse the LiDAR features into the camera stream to boost the performance of the two-stream network. Extensive experiments on benchmark data sets show the superiority of our method. For example, on nuScenes test set, our EPMF outperforms the state-of-the-art method, *i.e.*, RangeFormer, by **0.9%** in mIoU. Compared to PMF, EPMF also achieves **2.06×** acceleration with **2.0%** improvement in mIoU. Our source code is available at <https://github.com/ICEORY/PMF>.

**Index Terms**—Multi-Sensor Fusion, 3D Semantic Segmentation, Scene Understanding, Deep Neural Networks, Autonomous Driving.

## 1 INTRODUCTION

SEMANTIC scene understanding is a fundamental task for many applications, such as autonomous driving and robotics [18], [44], [58], [59]. Specifically, in the scenes of autonomous driving, it provides fine-grained environmental information for high-level motion planning and improves the safety of autonomous cars [3], [20]. One of the important tasks in semantic scene understanding is semantic segmentation, which assigns a class label to each data point in the input data, and helps autonomous cars to better understand the environment.

According to the sensors used by semantic segmentation methods, recent studies can be divided into three categories: camera-only methods [2], [9], [10], [45], [73], LiDAR-only methods [1], [15], [31], [66], [77] and multi-sensor fusion methods [37], [47], [49], [64], [74]. Camera-only methods have achieved great progress with the help of a massive amount of open-access data sets [6], [14], [16]. Since images obtained by a camera are rich in appearance information (*e.g.*, texture and color), camera-only methods can

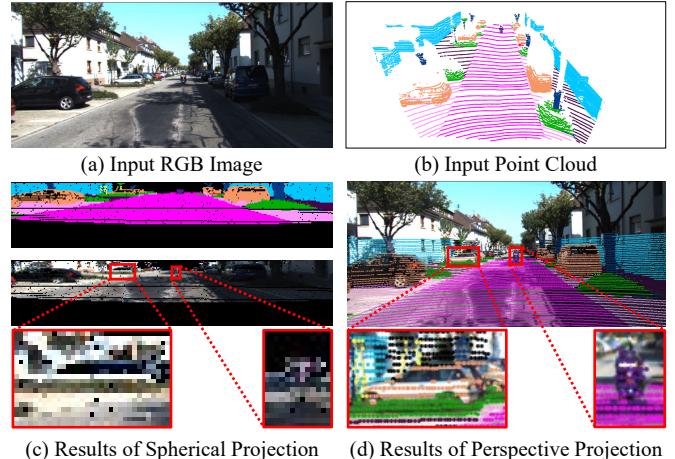


Fig. 1. Comparisons of spherical projection [50], [66] and perspective projection. With spherical projection, most of the appearance information from RGB images is lost. Instead, we preserve the information of images with perspective projection. To distinguish different classes, we colorize the point clouds using semantic labels from SemanticKITTI.

- Zhuangwei Zhuang, Sitao Chen, Rong Li, and Mingkui Tan are with the School of Software Engineering, South China University of Technology. Zhuangwei Zhuang and Mingkui Tan are also with the Pazhou Laboratory, Guangzhou, China. E-mail:{z.zhuangwei, mechenst, selirong}@mail.scut.edu.cn, mingkuitan@scut.edu.cn.
- Kui Jia is with the School of Electronic and Information Engineering, South China University of Technology. E-mail: kuijia@scut.edu.cn.
- Yuanqing Li is with the Pazhou Laboratory, Guangzhou, China. E-mail: auyqli@scut.edu.cn.
- Qicheng Wang is with Department of Mathematics at the Hong Kong University of Science and Technology and is also with Minieye, Shenzhen, Guangdong, China. E-mail: wangqicheng@minieye.cc.
- <sup>†</sup> Corresponding author.

provide fine-grained and accurate semantic segmentation results. However, as passive sensors, cameras are susceptible to changes in lighting conditions and are thus unreliable [61].<sup>1</sup> To address this problem, researchers conduct semantic segmentation on point

1. See Section 4.7 for more details.

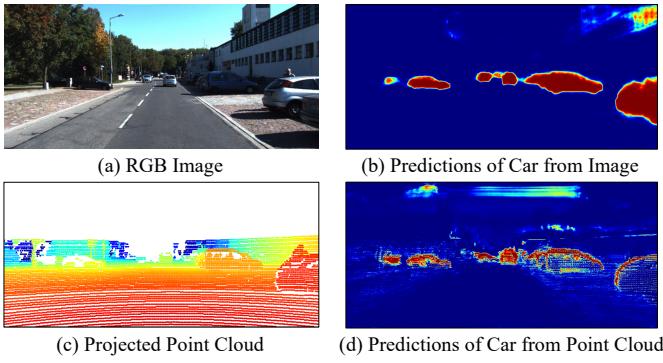


Fig. 2. Comparisons of the predictions from images and point clouds. Deep neural networks capture different perceptual information from RGB images and point clouds. Red indicates predictions with higher scores.

clouds from LiDAR. Compared with camera-only approaches, LiDAR-only methods are more robust in different light conditions, as LiDAR provides reliable and accurate spatio-depth information on the physical world. Unfortunately, LiDAR-only semantic segmentation is challenging due to the sparse and irregular distribution of point clouds. In addition, point clouds lack texture and color information, resulting in high classification error in the fine-grained segmentation task of LiDAR-only methods. A straightforward solution for addressing both drawbacks of camera-only and LiDAR-only methods is to fuse the multimodal data from both sensors, *i.e.*, multi-sensor fusion methods. Nevertheless, due to the large domain gap between RGB cameras and LiDAR, multi-sensor fusion is still a nontrivial task.

In multi-sensor fusion methods, fusing multimodal data from different sensors is an important problem. Existing fusion-based methods [47], [64] mainly lift the dense 2D image features to the 3D LiDAR coordinates using spherical projection [50] and conduct feature fusion in the sparse LiDAR domain. However, these methods suffer from a critical limitation: as the point clouds are very sparse, most of the appearance information from the RGB images is missing after un-projecting it to the LiDAR coordinates. For example, as shown in Figure 1 (c), the car and motorcycle in the image become distorted with spherical projection. As a result, existing fusion-based methods have difficulty capturing the appearance information from the projected RGB images.

In this paper, we aim to exploit an effective multi-sensor fusion method. Unlike existing methods [47], [64], we assume and highlight that the perceptual information from both RGB images and point clouds, *i.e.*, appearance information from images and spatio-depth information from point clouds, is important in fusion-based semantic segmentation. Based on this intuition, we propose a perception-aware multi-sensor fusion (PMF) scheme that conducts a collaborative fusion of perceptual information from two modalities of data in three aspects. **First**, we propose perspective projection to project the point clouds to the camera coordinate system to obtain additional spatio-depth information for RGB images. **Second**, we propose a two-stream network (TSNet) that contains a camera stream and a LiDAR stream to extract perceptual features from multi-modal sensors separately. Considering that the information from images is unreliable in an outdoor environment, we fuse the image features to the LiDAR stream by effective residual-based fusion (RF) modules, which are designed to learn the complementary features of the original LiDAR modules. **Third**, we propose perception-aware losses to measure the vast perceptual

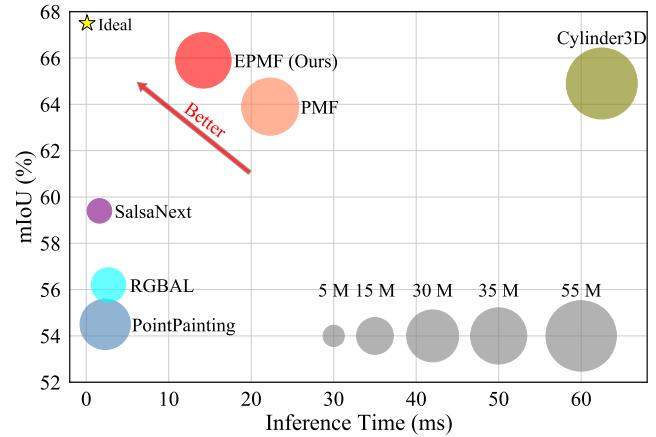


Fig. 3. Comparisons of efficiency and performance of different methods on SemanticKITTI-FV.

difference between the two data modalities and boost the fusion of different perceptual information. Specifically, as shown in Figure 2, the perceptual features captured by the camera stream and LiDAR stream are different. Therefore, we use the predictions with higher confidence to supervise those with lower confidence. Since model efficiency is also an essential factor for real-world applications, we further exploit the efficiency of PMF and propose an improved version, *i.e.*, EPMF.

Our contributions are summarized as follows. First, we propose a perception-aware multi-sensor fusion (PMF) scheme to effectively fuse the perceptual information from RGB images and point clouds. Second, by fusing the spatio-depth information from point clouds and appearance information from RGB images, PMF is able to address segmentation with undesired light conditions and sparse point clouds. More critically, PMF is robust in adversarial samples of RGB images by integrating the information from point clouds. Third, we introduce perception-aware losses into the network and force the network to capture the perceptual information from two different-modality sensors. As demonstrated in Figure 3, on top of PMF, we further propose EPMF, which reduces the model complexity of PMF while improving the model performance by a large margin. The extensive experiments on three benchmark data sets including SemanticKITTI-FV [3], nuScenes [7], and A2D2 [23], demonstrate the superior performance of our method. For example, on nuScenes test set, our EPMF outperforms the state-of-the-art methods, *i.e.*, SphereFormer [38] and RangeFormer [36] by 1.1% and 0.9% in mIoU, respectively, without additional fine-tuning or test-time augmentation.

This paper extends our prior version [82] from following aspects. 1) We propose cross-modal alignment and cropping (CAC) to address the miss-alignment issue of point clouds and RGB images. 2) We explore the impact of the different resolutions of point clouds and improve the efficiency of our method without performance degradation. 3) We adopt the proposed EPMF on more benchmark data sets and show the superior performance of our method on extremely sparse point clouds. 4) We provide more ablation studies to investigate the effectiveness of our method.

## 2 RELATED WORK

In this section, we revisit the existing literature on 2D and 3D semantic segmentation, *i.e.*, camera-only methods, LiDAR-only methods and multi-sensor fusion methods.

**Camera-only methods** Camera-only semantic segmentation aims to predict the pixel-wise labels of 2D images. FCN [45] is a fundamental work in semantic segmentation, which proposes an end-to-end fully convolutional architecture based on image classification networks. In addition to FCN, recent works have achieved significant improvements via exploring multi-scale information [9], [40], [78], dilated convolution [10], [48], [65], and attention mechanisms [32], [73]. More recently, transformer-based methods have been proposed for robust and accurate segmentation [11], [68]. Specifically, SegFormer [68] designs a positional-encoding-free and hierarchical transformer encoder that generates both high-resolution fine features and low-resolution coarse features. Moreover, it proposes a lightweight All-MLP decoder to aggregate the multiscale features for robust semantic segmentation. Mask2Former [11] proposes a universal architecture to address any image segmentation tasks, including panoptic, instance, or semantic segmentation. Although these methods show strong robustness to the corruptions and perturbations in autonomous driving scenes, their performance under poor lighting conditions (*i.e.*, at night-time) is unsatisfactory compared to LiDAR-based methods.

**LiDAR-only methods** To address the drawbacks of cameras, LiDAR is an important sensor on an autonomous car, as it is robust in more complex scenes. According to the preprocessing pipeline, existing methods for point clouds mainly contain two categories, including direct methods [31], [55], [56], [81] and projection-based methods [15], [66], [67], [69].

Direct methods perform semantic segmentation by processing the raw 3D point clouds directly. PointNet [55] is a pioneering work in this category that extracts point cloud features by multi-layer perception. A subsequent extension, *i.e.*, PointNet++ [56], further aggregates a multi-scale sampling mechanism to aggregate global and local features. However, these methods do not consider the varying sparsity of point clouds in outdoor scenes. Cylinder3D [81] addresses this issue by using 3D cylindrical partitions and asymmetrical 3D convolutional networks. SphereFormer [38] directly aggregates information from dense close points to sparse distant ones through radial window partitions and proposes dynamic feature selection to select local neighbor features or radial contextual features. However, direct methods have a high computational complexity, which limits their applicability in autonomous driving. PVKD [30] transfers both point-level and voxel-level hidden knowledge from a large LiDAR semantic segmentation model to a slim network to achieve model compression. WaffleIron [54] uses standard MLPs and dense 2D convolutions to build a 3D backbone for point cloud semantic segmentation, which does not rely on sparse 3D convolution. In addition, one can easily improve the efficiency of networks by existing neural architecture search [8], [26], [52] and model compression techniques [27], [43], [71].

Projection-based methods are more efficient because they convert 3D point clouds to a 2D grid. In projection-based methods, researchers focus on exploiting effective projection methods, such as spherical projection [50], [66] and bird's-eye projection [77]. Such 2D representations allow researchers to investigate efficient network architectures based on existing 2D convolutional networks [1], [15], [25]. AMVNet [42] utilizes the late fusion to combine the advantages of both the range view and bird's-eye view networks. RPVNet [70] proposes a range-point-voxel fusion network to synergize all the three view's representations to alleviate the above problem. RangeFormer [36] formulates the segmentation of range view grids as a seq2seq problem and adopts several standard Transformer blocks to capture the rich

contextual information, then uses the MLP heads to decode the multi-scale features. Unlike uni-modal methods, we focus on fusing information from both the camera and LiDAR to achieve accurate and robust 3D semantic segmentation for autonomous driving.

**Multi-sensor fusion methods** To leverage the benefits of both camera and LiDAR, recent work has attempted to fuse information from two complementary sensors to improve the accuracy and robustness of the 3D semantic segmentation algorithm [37], [47], [49], [64]. RGBAL [47] converts RGB images to a polar-grid mapping representation and designs early and mid-level fusion strategies. PointPainting [64] obtains the segmentation results of images and projects them to the LiDAR space by using bird's-eye projection [77] or spherical projection [50]. The projected segmentation scores are concatenated with the original point cloud features to improve the performance of LiDAR networks. 2DPASS [72] enhances the representation learning of 3D semantic segmentation network by distilling multi-modal knowledge to single point cloud modality. In this way, 2DPASS can use LiDAR-only input in test-time. However, the model performance is unsatisfactory under the scene with sparser point clouds (*e.g.*, A2D2 with 16-beam LiDARs). In contrast, our EPMF can achieve promising performance by fusing 2D images and 3D point clouds during inference. Besides, unlike existing methods that perform feature fusion in the LiDAR domain, PMF [82] exploits a collaborative fusion of multimodal data in camera coordinates. In this work, we further extend PMF to improve its efficiency and performance.

### 3 PROPOSED METHOD

In this work, we propose an efficient perception-aware multi-sensor fusion (EPMF) scheme to perform an effective fusion of the perceptual information from both RGB images and point clouds. Specifically, as shown in Figure 4, EPMF contains three components: (1) perspective projection with cross-modal alignment and cropping; (2) a two-stream network (TSNet) with residual-based fusion modules; (3) perception-aware losses. The general scheme of EPMF is shown in Algorithm 1. We first project the point clouds to the camera coordinate system by using perspective projection. Then, we use a two-stream network that contains a camera stream and a LiDAR stream to extract perceptual features from the two modalities, separately. The features from the camera stream are fused into the LiDAR stream by residual-based fusion modules. Finally, we introduce perception-aware losses into the optimization of the network.

#### 3.1 Pipeline of data pre-processing

Existing methods [47], [64] mainly project images to the LiDAR coordinate system using spherical projection. However, due to the sparse nature of point clouds, most of the appearance information from the images is lost with spherical projection (see Figure 1). To address this issue, we propose perspective projection to project the sparse point clouds to the camera coordinate system.

**Formulation of perspective projection.** Let  $\{\mathbf{P}, \mathbf{X}, \mathbf{y}\}$  be one of the training samples from a given data set, where  $\mathbf{P} \in \mathbb{R}^{4 \times N}$  indicates a point cloud from LiDAR and  $N$  denotes the number of points. Each point  $\mathbf{P}_i$  in point cloud  $\mathbf{P}$  consists of 3D coordinates  $(x, y, z)$  and a reflectance value ( $r$ ). Let  $\mathbf{X} \in \mathbb{R}^{3 \times H \times W}$  be an image from an RGB camera, where  $H$  and  $W$  represent the height and width of the image, respectively.  $\mathbf{y} \in \mathbb{R}^N$  is the set of semantic labels for point cloud  $\mathbf{P}$ .

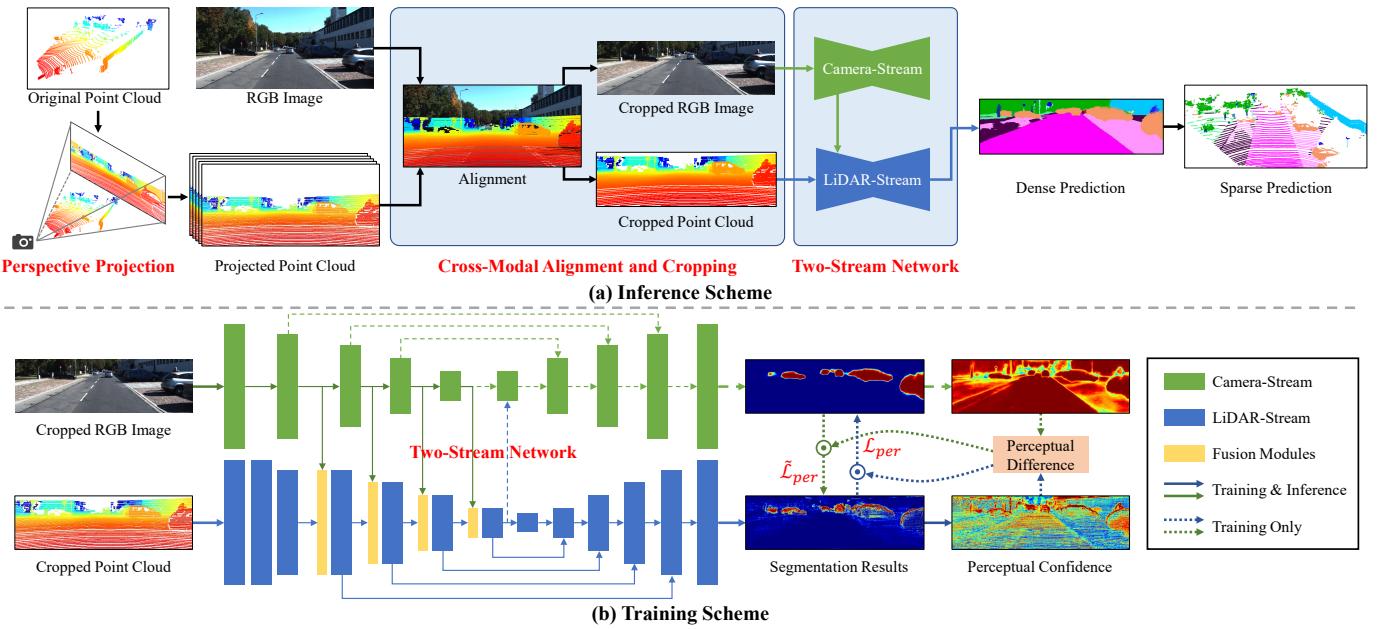


Fig. 4. Illustration of the training and inference schemes of EPMF. EPMF consists of three components: (1) perspective projection with cross-modal alignment and crop; (2) a two-stream network (TSNet) with feature fusion modules; and (3) perception-aware losses  $\mathcal{L}_{per}$ ,  $\tilde{\mathcal{L}}_{per}$  w.r.t. the camera stream and the LiDAR stream. We first project the point clouds to the camera coordinate with perspective projection and learn the features from both the RGB images and point clouds using TSNet. The image features are fused into the LiDAR stream network by fusion modules. In the training procedure, we use perception-aware losses to help the network focus on the perceptual features of both images and point clouds. In the inference procedure, we apply dense-to-sparse mapping to obtain 3D segmentation results of point clouds.

#### Algorithm 1 General Scheme of EPMF

**Input:** Training data  $\{\mathbf{P}, \mathbf{X}, \mathbf{y}\}$ , TSNet with submodels  $M, \tilde{M}$ , hyperparameters  $\tau, \lambda, \gamma$ .

- 1: **while** not convergent **do**
- 2: Project the point clouds  $\mathbf{P}$  by using perspective projection to obtain  $\tilde{\mathbf{X}}$ .
- 3: Use  $\{\tilde{\mathbf{X}}, \mathbf{X}\}$  as the inputs of TSNet and compute the output probabilities  $\{\tilde{\mathbf{O}}, \mathbf{O}\}$  with Eq. (2).
- 4: Compute the perceptual confidence  $\tilde{\mathbf{C}}$  and  $\mathbf{C}$ .
- 5: Construct perception-aware losses to measure the perceptual difference with Eqs. (7) and (9).
- 6: Update  $\tilde{M}$  and  $M$  by minimizing the objective in Eq. (17).
- 7: **end while**

In perspective projection, we aim to project the point cloud  $\mathbf{P}$  from LiDAR coordinate to the camera coordinate to obtain the 2D LiDAR features  $\tilde{\mathbf{X}} \in \mathbb{R}^{C \times H \times W}$ . Here,  $C$  indicates the number of channels w.r.t. the projected point cloud. Following [19], we obtain  $\mathbf{P}_i = (x, y, z, 1)^\top$  by appending a fourth column to  $\mathbf{P}_i$  and compute the projected point  $\tilde{\mathbf{P}}_i = (\tilde{x}, \tilde{y}, \tilde{z})^\top$  in the camera coordinates by

$$\tilde{\mathbf{P}}_i = \mathbf{T} \mathbf{R} \mathbf{P}_i, \quad (1)$$

where  $\mathbf{T} \in \mathbb{R}^{3 \times 4}$  is the projection matrix from LiDAR coordinates to camera coordinates.  $\mathbf{R} \in \mathbb{R}^{4 \times 4}$  is expanded from the rectifying rotation matrix  $\mathbf{R}^{(0)} \in \mathbb{R}^{3 \times 3}$  by appending a fourth zero row and column and setting  $\mathbf{R}(4, 4) = 1$ . The calibration parameters  $\mathbf{T}$  and  $\mathbf{R}^{(0)}$  can be obtained by the approach in [21]. Subsequently, the corresponding pixel  $(h, w)$  in the projected image  $\tilde{\mathbf{X}}$  w.r.t. the point  $\mathbf{P}_i$  is computed by  $h = \tilde{x}/\tilde{z}$  and  $w = \tilde{y}/\tilde{z}$ .

**Cross-modal alignment and cropping.** As shown in Figure 4(a), since we only focus on the segmentation of point clouds, directly projecting point clouds to the view of cameras leads to unnecessary computational costs. To address this issue, we introduce cross-modal alignment and cropping (CAC). First, we align the RGB image and the projected point clouds to find the overlap of the multi-modal inputs. Then, we crop both RGB images and projected point clouds to obtain compact inputs: For RGB images, we only keep the area that contains point clouds. For the projected point clouds, as the area outside the horizontal field of view (FOV) of the camera is covered by other cameras, we only keep the points within the horizontal FOV of the camera. In the case that the LiDAR sensor has a larger vertical FOV, we can keep the point clouds outside the image.<sup>2</sup>

After applying CAC, we compute the features of the projected point clouds. Because the point cloud is very sparse, each pixel in the projected  $\tilde{\mathbf{X}}$  may not have a corresponding point  $\mathbf{p}$ . Therefore, we first initialize all pixels in  $\tilde{\mathbf{X}}$  to 0. Following [15], [50], we compute 5-channel LiDAR features, i.e.,  $(d, x, y, z, r)$ , for each pixel  $(h, w)$  in the projected 2D image  $\tilde{\mathbf{X}}$ , where  $d = \sqrt{x^2 + y^2 + z^2}$  represents the range value of each point. Thus, we set the number of channels  $C$  to 5 in this work.

## 3.2 Architecture design of EPMF

As images and point clouds are different-modality data, it is difficult to handle both types of information from the two modalities by using a single network [37]. Motivated by [17], [60], we propose a two-stream network (TSNet) that contains a camera stream and a LiDAR stream to process the features from camera and LiDAR,

2. More discussions of CAC can be found in Section 5.3

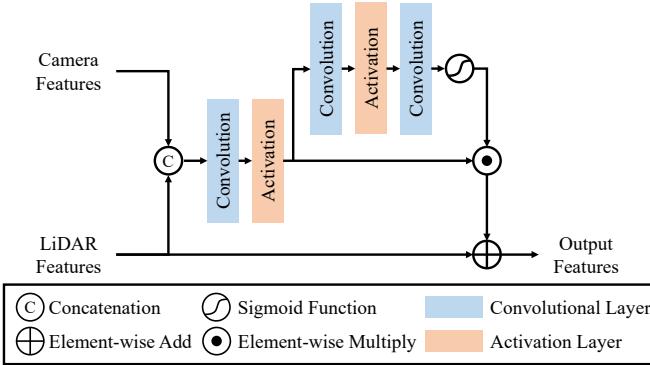


Fig. 5. Illustration of the residual-based fusion (RF) module. RF fuses features from both the camera and LiDAR to generate the complementary information of the original LiDAR features.

separately, as illustrated in Figure 4. In this way, we can use the network architectures designed for images and point clouds as the backbones of each stream in TSNet.

**Formulation of Two Stream Network.** Let  $\tilde{M}$  and  $M$  be the LiDAR stream and the camera stream in TSNet, respectively. Let  $\tilde{\mathbf{O}} \in \mathbb{R}^{S \times H \times W}$  and  $\mathbf{O} \in \mathbb{R}^{S \times H \times W}$  be the output probabilities w.r.t. each network, where  $S$  indicates the number of semantic classes. The outputs of TSNet are computed by

$$\begin{cases} \mathbf{O} = M(\mathbf{X}), \\ \tilde{\mathbf{O}} = \tilde{M}(\tilde{\mathbf{X}}). \end{cases} \quad (2)$$

**Design of Residual-based Fusion Module.** Since the features of images contain many details of objects, we then introduce a residual-based fusion module, as illustrated in Figure 5, to fuse the image features to the LiDAR stream<sup>3</sup>. Let  $\{\mathbf{F}_l \in \mathbb{R}^{C_l \times H_l \times W_l}\}_{l=1}^L$  be a set of image features from the camera stream, where  $l$  indicates the layer in which we obtain the features.  $C_l$  indicates the number of channels of the  $l$ -th layer in the camera stream.  $H_l$  and  $W_l$  indicate the height and width of the feature maps from the  $l$ -th layer, respectively. Let  $\{\tilde{\mathbf{F}}_l \in \mathbb{R}^{\tilde{C}_l \times H_l \times W_l}\}_{l=1}^L$  be the features from the LiDAR stream, where  $\tilde{C}_l$  indicates the number of channels of the  $l$ -th layer in the LiDAR stream. To obtain the fused features, we first concatenate the features from each network and use a convolutional layer to reduce the number of channels of the fused features. The fused features  $\mathbf{F}_l^{fuse} \in \mathbb{R}^{\tilde{C}_l \times H_l \times W_l}$  are computed by

$$\mathbf{F}_l^{fuse} = f_l([\tilde{\mathbf{F}}_l; \mathbf{F}_l]), \quad (3)$$

where  $[\cdot; \cdot]$  indicates the concatenation operation.  $f_l(\cdot)$  is the convolution operation w.r.t. the  $l$ -th fusion module.

Considering that the camera is easily affected by different lighting and weather conditions, the information from RGB images is not reliable in an outdoor environment. We use the fused features as the complement of the original LiDAR features and design the fusion module based on the residual structure [28]. Incorporating with the attention module [5], the output features  $\mathbf{F}_l^{out} \in \mathbb{R}^{\tilde{C}_l \times H_l \times W_l}$  of the fusion module are computed by

$$\mathbf{F}_l^{out} = \tilde{\mathbf{F}}_l + \sigma(g_l(\mathbf{F}_l^{fuse})) \odot \mathbf{F}_l^{fuse}, \quad (4)$$

3. We discuss different designs of fusion modules in Section 2 of the supplementary material of [82]

where  $\sigma(x) = 1/(1 + e^{-x})$  indicates sigmoid function.  $g_l(\cdot)$  indicates convolution operation in the attention module w.r.t. the  $l$ -th fusion module.  $\odot$  indicates element-wise multiplication operation.

### 3.3 Construction of perception-aware loss

The construction of perception-aware loss is very important in our method. As demonstrated in Figure 2, because the point clouds are very sparse, the LiDAR stream network learns only the local features of points while ignoring the shape of objects. In contrast, the camera stream can easily capture the shape and texture of objects from dense images. In other words, the perceptual features captured by the camera stream and LiDAR stream are different. With this intuition, we introduce a perception-aware loss to make the fusion network focus on the perceptual features from the camera and LiDAR.

To measure the perceptual confidence of the predictions w.r.t. the LiDAR stream, we first compute the entropy map  $\tilde{\mathbf{E}} \in \mathbb{R}^{H \times W}$  by

$$\tilde{\mathbf{E}}_{h,w} = -\frac{1}{\log S} \sum_{s=1}^S \tilde{\mathbf{O}}_{s,h,w} \log(\tilde{\mathbf{O}}_{s,h,w}). \quad (5)$$

Following [57], we use  $\log S$  to normalize the entropy to  $(0, 1]$ . Then, the perceptual confidence map  $\tilde{\mathbf{C}}$  w.r.t. the LiDAR stream is computed by  $\tilde{\mathbf{C}} = \mathbf{1} - \tilde{\mathbf{E}}$ . For the camera stream, the confidence map is computed by  $\mathbf{C} = \mathbf{1} - \mathbf{E}$ .

Note that not all information from the camera stream is useful. For example, the camera stream is confident inside objects but may make mistakes at the boundary of the objects. In addition, the predictions with lower confidence scores are more likely to be wrong. Incorporating a confidence threshold, we measure the importance of perceptual information from the camera stream by

$$\tilde{\Omega}_{h,w} = \begin{cases} \max(\mathbf{C}_{h,w} - \tilde{\mathbf{C}}_{h,w}, 0), & \text{if } \mathbf{C}_{h,w} > \tau, \\ 0, & \text{otherwise.} \end{cases} \quad (6)$$

Here  $\tau$  indicates the confidence threshold.

Inspired by [29], [33], [76], to learn the perceptual information from the camera stream, we construct the perception-aware loss w.r.t. the LiDAR stream by

$$\tilde{\mathcal{L}}_{per} = \frac{1}{Q} \sum_{h=1}^H \sum_{w=1}^W \tilde{\Omega}_{h,w} D_{KL}(\tilde{\mathbf{O}}_{:,h,w} || \mathbf{O}_{:,h,w}), \quad (7)$$

where  $Q = H \cdot W$  and  $D_{KL}(\cdot || \cdot)$  indicates the Kullback-Leibler divergence [29].

For the camera stream, the importance of information from LiDAR stream is computed by

$$\Omega_{h,w} = \begin{cases} \max(\tilde{\mathbf{C}}_{h,w} - \mathbf{C}_{h,w}, 0), & \text{if } \tilde{\mathbf{C}}_{h,w} > \tau, \\ 0, & \text{otherwise.} \end{cases} \quad (8)$$

The perception-aware loss w.r.t. the camera stream is

$$\mathcal{L}_{per} = \frac{1}{Q} \sum_{h=1}^H \sum_{w=1}^W \Omega_{h,w} D_{KL}(\mathbf{O}_{:,h,w} || \tilde{\mathbf{O}}_{:,h,w}). \quad (9)$$

### 3.4 Objective functions

In addition to the perception-aware loss, we also use multi-class focal loss [41] and Lovász-softmax loss [4], which are commonly used in existing segmentation work [15], [81], to train the two stream network.

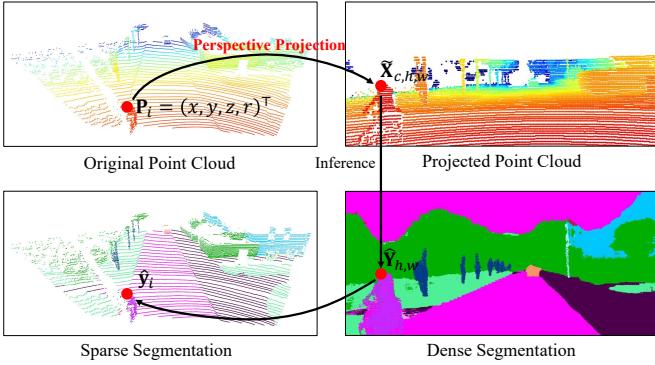


Fig. 6. Illustration of the pipeline to obtain the sparse segmentation from the dense prediction results.  $\bar{\mathbf{X}}$  indicates the projected point cloud.  $\hat{\mathbf{Y}}$  and  $\hat{\mathbf{y}}$  indicate the dense predictions and sparse predictions, respectively. For each point  $\mathbf{P}_i$ , we first compute the corresponding pixel  $(h, w)$  in the camera coordinate system by perspective projection. Second, we get the dense segmentation  $\hat{\mathbf{Y}}$  from the prediction results of PMF. Last, we obtain the corresponding sparse prediction  $\hat{\mathbf{y}}_i$  w.r.t. the point  $\mathbf{P}_i$  from the dense segmentation  $\hat{\mathbf{Y}}_{h,w}$ .

Let  $\mathbf{Y} \in \mathbb{R}^{H \times W}$  be the projected labels in the camera coordinates.  $H$  and  $W$  indicate the height and width, respectively. For each point  $\mathbf{P}_i$ , we project the 3D coordinates  $(x, y, z)$  to the pixel  $(h, w)$  in the camera coordinate system by using perspective projection. Then, we initialize all pixels in  $\mathbf{Y}$  by 0 and compute the projected labels in  $\mathbf{Y}$  by

$$\mathbf{Y}_{h,w} := \mathbf{y}_i. \quad (10)$$

The multi-class focal loss w.r.t. the LiDAR stream is defined as

$$\tilde{\mathcal{L}}_{foc} = \frac{1}{K} \sum_{s=1}^S \sum_{h=1}^H \sum_{w=1}^W \alpha_s \mathbb{1}\{\mathbf{Y}_{h,w} = s\} FL(\tilde{\mathbf{O}}_{s,h,w}), \quad (11)$$

where  $FL(p) = -(1-p)^2 \log(p)$  denotes the focal-loss function.  $\alpha_s$  denotes the weights w.r.t. the  $s$ -th class.  $K = \sum_{s=1}^S \sum_{h=1}^H \sum_{w=1}^W \mathbb{1}\{\mathbf{Y}_{h,w} = s\}$  indicates the number of available labels.  $\mathbb{1}\{\cdot\}$  indicates the indicator function. Then, the multi-class focal loss w.r.t. the camera stream is

$$\mathcal{L}_{foc} = \frac{1}{K} \sum_{s=1}^S \sum_{h=1}^H \sum_{w=1}^W \mathbb{1}\{\mathbf{Y}_{h,w} = s\} FL(\mathbf{O}_{s,h,w}), \quad (12)$$

The Lovász-softmax loss w.r.t. the LiDAR stream is

$$\tilde{\mathcal{L}}_{lov} = \frac{1}{S} \sum_{s=1}^S \overline{\Delta_{J_s}}(\tilde{\mathbf{m}}(s)), \quad (13)$$

where

$$\tilde{\mathbf{m}}_i(s) = \begin{cases} 1 - \tilde{\mathbf{O}}_{s,h,w} & \text{if } s = \mathbf{Y}_{h,w}, \\ \tilde{\mathbf{O}}_{s,h,w} & \text{otherwise.} \end{cases} \quad (14)$$

$\overline{\Delta_{J_s}}$  indicates the Lovász extension of the Jaccard index for class  $s$ . Here,  $(h, w)$  is obtained from the 3D coordinates  $(x, y, z)$  of  $\mathbf{P}_i$  by using perspective projection.  $\tilde{\mathbf{m}}(s) \in [0, 1]^N$  indicates the vector of errors. The Lovász-softmax loss w.r.t. the camera stream is defined as

$$\mathcal{L}_{lov} = \frac{1}{S} \sum_{s=1}^S \overline{\Delta_{J_s}}(\mathbf{m}(s)), \quad (15)$$

where

$$\mathbf{m}_i(s) = \begin{cases} 1 - \mathbf{O}_{s,h,w} & \text{if } s = \mathbf{Y}_{h,w}, \\ \mathbf{O}_{s,h,w} & \text{otherwise.} \end{cases} \quad (16)$$

By considering the objective functions of both LiDAR stream and camera stream, we formulate the objective function of the proposed two-stream network as

$$\mathcal{L} = \tilde{\mathcal{L}}_{foc} + \mathcal{L}_{foc} + \tilde{\lambda} \tilde{\mathcal{L}}_{lov} + \lambda \mathcal{L}_{lov} + \tilde{\gamma} \tilde{\mathcal{L}}_{per} + \gamma \mathcal{L}_{per}, \quad (17)$$

where  $\lambda$ ,  $\tilde{\lambda}$ ,  $\gamma$  and  $\tilde{\gamma}$  indicate the hyper-parameters that balance different losses.

### 3.5 Pipeline of post-processing

With the proposed perception-aware losses, PMF generates dense segmentation results with information from RGB images and point clouds. We then obtain the sparse prediction from the dense results. Let  $\tilde{\mathbf{O}} \in \mathbb{R}^{S \times H \times W}$  be the output probabilities of the LiDAR stream.  $S$  indicates the number of classes.  $H$  and  $W$  indicate the height and width of the predictions, respectively. Let  $\hat{\mathbf{Y}} \in \mathbb{R}^{H \times W}$  be the dense predictions from the LiDAR stream. Then, the dense predictions are computed by

$$\hat{\mathbf{Y}}_{h,w} = \arg \max_s \tilde{\mathbf{O}}_{s,h,w}. \quad (18)$$

Let  $\hat{\mathbf{y}} \in \mathbb{R}^N$  be the sparse predictions of point cloud  $\mathbf{P}$ . As shown in Figure 6, for each point  $\mathbf{P}_i$ , we first project the 3D coordinates  $(x, y, z)$  to the camera coordinate system by using perspective projection and compute the corresponding pixel  $(h, w)$  in the projected image. Then the semantic prediction  $\hat{\mathbf{y}}_i$  w.r.t. the point  $\mathbf{P}_i$  is computed by

$$\hat{\mathbf{y}}_i := \hat{\mathbf{Y}}_{h,w}. \quad (19)$$

Note that for the point cloud with multi-camera views, e.g., nuScenes, there are overlaps between different camera views. To address this issue, we conduct inference for each camera view and merge the results by assigning the predictions with the highest confidence scores to the points in the overlaps of different views.

### 3.6 Techniques to improve efficiency and effectiveness

On top of the two-stream network designed in [82], we further explore the techniques to improve the efficiency and effectiveness of the fusion network<sup>4</sup>.

**Dropping decoder of camera stream.** By introducing the perception-aware loss, the knowledge of the camera stream is distilled into the LiDAR stream. Therefore, we only use the predictions of LiDAR stream while dropping the results of the camera stream during inference. In this sense, we can also drop the decoder of the camera stream to speed up inference.

**Improved contextual module.** In [15], Cortinhal *et al.* have designed the contextual module that improves the ability of SalsaNet [1] for comprehending global contextual information. However, the proposed contextual module is explored under spherical projection and may be ineffective with perspective projection. In our experiments, we find that the high resolution of the projected point cloud feature is unnecessary due to the sparse nature of point clouds. Therefore, we insert extra down-sampling operation into the contextual module of LiDAR stream to reduce the resolution of point cloud features and improve the efficiency

4. We study the effect of the proposed techniques in Section 5.3

TABLE 1  
Comparisons of the number of points of SemanticKITTI and SemanticKITTI-FV.

Dataset splits	SemanticKITTI	SemanticKITTI-FV	Percentage
Training set	$2.35 \times 10^9$	$3.78 \times 10^8$	16.03%
Validation set	$4.99 \times 10^8$	$8.00 \times 10^7$	16.07%

TABLE 2  
Comparisons of the number of valid labels of SemanticKITTI and SemanticKITTI-FV.

Dataset splits	SemanticKITTI	SemanticKITTI-FV	Percentage
Training set	$2.28 \times 10^9$	$3.63 \times 10^8$	15.93%
Validation set	$4.77 \times 10^8$	$7.57 \times 10^7$	15.88%

of LiDAR stream. Moreover, to reduce the impact of the sparse issue of projected point clouds, we replace the convolutional layers in the contextual modules of LiDAR stream with sparse invariant convolutional layers [63].

**Fusing high-level LiDAR features.** In the two-stream network, the camera stream generates predictions with RGB images only, which, however, results in the unsatisfactory segmentation performance of the camera stream and may limit in performance of the fusion network. To address this issue, we further fuse the high-level features from the last stage of the LiDAR stream backbone into the camera stream by a concatenate operation. In this way, we improve the performance of the camera stream without introducing extra computational costs and thus boost the effectiveness of the perception-aware loss.

## 4 EXPERIMENTS

In this section, we first compare EPMF with the state-of-the-art methods on the benchmark data sets. Then, we provide distance-based evaluation and qualitative results of our methods. Last, we conduct experiments to evaluate the efficiency of our method.

We organize the experiments as follows. 1) We introduce the benchmark data sets and evaluation metrics in Section 4.1. 2) We provide the implementation details of our method in Section 4.2. 3) We evaluate the performance of our method on several benchmark data sets in Section 4.3. 4) We investigate the performance of our method under different distances in Section 4.4. 5) We discuss the efficiency of the proposed method in Section 4.5. 6) We show the qualitative results in Section 4.6. 7) We study the robustness of the proposed method on adversarial samples in Section 4.7.

### 4.1 Data sets and Evaluation Metrics

#### 4.1.1 Data sets

We empirically evaluate our method on several benchmark data sets, including SemanticKITTI-FV [3], nuScenes [7], and A2D2 [23].

SemanticKITTI is a large-scale data set based on the KITTI Odometry Benchmark [20], providing 43,000 scans with point-wise semantic annotation, where 21,000 scans (sequence 00-10) are available for training and validation. The data set has 19 semantic classes for the evaluation of semantic benchmarks. The point clouds are collected using a Velodyne HDL-64E sensor, which has 64 beams vertically. Since SemanticKITTI provides only the images of the front-view camera, we project the point clouds to a perspective view and keep only the available points on the images to build a subset of SemanticKITTI, namely, SemanticKITTI-FV. The comparisons between SemanticKITTI and SemanticKITTI-FV

are shown in Table 1 and Table 2. SemanticKITTI-FV has only 15.93% data for training compared with the full SemanticKITTI data set.

nuScenes contains 1,000 driving scenes with different weather and light conditions. The scenes are split into 28,130 training frames and 6,019 validation frames, which are collected with a Velodyne HDL-32E sensor. Unlike SemanticKITTI, which provides only the images of the front-view camera, nuScenes has 6 cameras for different views of LiDAR.

A2D2 provides 38-class semantic segmentation images and point cloud labels for 41,277 non-sequential frames, which are collected with six cameras and five Velodyne VLP-16 sensors. Following [77], we split 22,408 scans for training, 2,274 for validation, and 13,264 for testing, respectively. Similar to nuScenes, the sensor suite of A2D2 provides full 360° coverage. However, the camera and LiDAR sensor orientations are optimized manually to minimize the blind spot around the vehicle and maximize camera and LiDAR field of view overlap, which makes A2D2 lack horizontal scan lines.

#### 4.1.2 Evaluation Metrics

To evaluate the performance of our method, we use the mean Intersection over Union (mIoU) as the evaluation metric following the official guidance in [3], [7]. To class  $s$ , the respective intersection over union  $\text{IoU}_s$  is defined by

$$\text{IoU}_s = \frac{|\mathcal{P}_s \cap \mathcal{G}_s|}{|\mathcal{P}_s \cup \mathcal{G}_s|}, \quad (20)$$

where  $\mathcal{P}_s$  is the set of point with a class prediction  $s$ ,  $\mathcal{G}_s$  is the set of label for class  $s$ , and  $|\cdot|$  represents the cardinality of the set. Then, the mIoU is formulated as  $\text{mIoU} = \frac{1}{S} \sum_{s=1}^S \text{IoU}_s$ .

### 4.2 Implementation details

We implement the proposed method in PyTorch [53], and use ResNet-34 [28] and SalsaNext [15] as the backbones of the camera stream and LiDAR stream, respectively. Because we process the point clouds in the camera coordinates, we incorporate ASPP [9] into the LiDAR stream network to adjust the receptive field adaptively. To leverage the benefits of existing image classification models, we initialize the parameters of ResNet-34 with the pre-trained ImageNet models from [53]. We also adopt hybrid optimization methods [75] to train the networks w.r.t. different modalities, *i.e.*, SGD with Nesterov [51] for the camera stream and Adam [35] for the LiDAR stream. We train the networks for 50 epochs with a batch size of 8 on SemanticKITTI. On nuScenes and A2D2, the batch size is 24 with 150 epochs for training as the point clouds of these data sets are more sparse. The learning rate starts at 0.001 and decays to 0 with a cosine policy [46]. We tune the hyper-parameters of the objective function in Eq. (17) using the weighting strategy proposed in [34]. We set the confidence threshold  $\tau$  to 0.7 following [82].<sup>5</sup> To prevent overfitting, a series of data augmentation strategies are used, including random horizontal flipping, random scaling, color jitter, 2D random rotation, and random cropping.

### 4.3 Comparisons on benchmark data sets

#### 4.3.1 Results on SemanticKITTI-FV

To evaluate our method on SemanticKITTI-FV, we compare EPMF with several state-of-the-art LiDAR-only methods including

5. The ablation study of  $\tau$  is in the supplementary materials of [82].

TABLE 3

Comparisons on SemanticKITTI-FV validation set. **L** indicates LiDAR-only methods. **L+C** indicates fusion-based methods. \* indicates the results based on our implementation. The **bold** numbers indicate the best results.

Method	Modality	Car	Bicycle	Motorcycle	Truck	Other-vehicle	Person	Bicyclist	Motorcyclist	Road	Parking	Sidewalk	Other-ground	Building	Fence	Vegetation	Trunk	Terrain	Pole	Traffic-sign	mIoU (%)
RandLANet [31]	L	92.0	8.0	12.8	74.8	46.7	52.3	46.0	0.0	93.4	32.7	73.4	0.1	84.0	43.5	83.7	57.3	73.1	48.0	27.3	50.0
RangeNet++ [50]	L	89.4	26.5	48.4	33.9	26.7	54.8	69.4	0.0	92.9	37.0	69.9	0.0	83.4	51.0	83.3	54.0	68.1	49.8	34.0	51.2
SequeezeSegV2 [67]	L	82.7	15.1	22.7	25.6	26.9	22.9	44.5	0.0	92.7	39.7	70.7	0.1	71.6	37.0	74.6	35.8	68.1	21.8	22.2	40.8
SequeezeSegV3 [69]	L	87.1	34.3	48.6	47.5	47.1	58.1	53.8	0.0	95.3	43.1	78.2	0.3	78.9	53.2	82.3	55.5	70.4	46.3	33.2	53.3
SalsaNext [15]	L	90.5	44.6	49.6	86.3	54.6	74.0	81.4	0.0	93.4	40.6	69.1	0.0	84.6	53.0	83.6	64.3	64.2	54.4	39.8	59.4
MinkowskiNet [13]	L	95.0	23.9	50.4	55.3	45.9	65.6	82.2	0.0	94.3	43.7	76.4	0.0	87.9	57.6	87.4	67.7	71.5	63.5	43.6	58.5
SPVNas [62]	L	96.5	44.8	63.1	59.9	64.3	72.0	86.0	0.0	93.9	42.4	75.9	0.0	88.8	59.1	88.0	67.5	73.0	63.5	44.3	62.3
Cylinder3D [81]	L	96.4	61.5	78.2	66.3	69.8	80.8	93.3	0.0	94.9	41.5	78.0	1.4	87.5	50.0	86.7	72.2	68.8	63.0	42.1	64.0
2DPASS [72]	L+C	96.6	60.8	71.6	82.0	77.8	78.2	92.0	0.2	93.9	44.4	75.7	3.1	89.5	60.9	88.7	72.6	74.0	61.5	45.5	<b>66.8</b>
2DPASS <sup>†</sup> [72]	L+C	93.6	54.9	71.6	81.6	43.5	71.5	82.2	0.2	94.0	33.8	75.1	0.2	88.7	57.3	88.5	67.1	75.1	56.0	40.2	61.8
PointPainting* [64]	L+C	94.7	17.7	35.0	28.8	55.0	59.4	63.6	0.0	95.3	39.9	77.6	0.4	87.5	55.1	87.7	67.0	72.9	61.8	36.5	54.5
RGBAL* [47]	L+C	87.3	36.1	26.4	64.6	54.6	58.1	72.7	0.0	95.1	45.6	77.5	0.8	78.9	53.4	84.3	61.7	72.9	56.1	41.5	56.2
PMF [82]	L+C	95.4	47.8	62.9	68.4	75.2	78.9	71.6	0.0	96.4	43.5	80.5	0.1	88.7	60.1	88.6	72.7	75.3	65.5	43.0	63.9
Uni-modal baseline	L	93.2	35.9	33.9	78.2	54.4	63.1	71.9	0.0	95.0	42.4	80.0	0.7	86.0	53.7	86.1	63.5	73.7	59.4	41.9	58.6
EPMF (Ours)	L+C	95.4	52.5	66.3	82.4	80.5	77.2	85.1	0.0	95.8	48.1	80.5	0.8	89.2	66.1	86.0	70.0	68.9	62.6	43.6	65.9

<sup>†</sup> Model trained under our settings.

TABLE 4

Comparisons on the nuScenes validation set. **L** indicates LiDAR-only methods. **L+C** indicates fusion-based methods. <sup>†</sup> indicates the results with test-time augmentation. The **bold** numbers indicate the best results.

Method	Modality	Barrier	Bicycle	Bus	Car	Construction	Motorcycle	Pedestrian	Traffic-cone	Trailer	Truck	Driveable	Other-flat	Sidewalk	Terrain	Mannade	Vegetation	mIoU (%)
RangeNet++ [50]	L	66.0	21.3	77.2	80.9	30.2	66.8	69.6	52.1	54.2	72.3	94.1	66.6	63.5	70.1	83.1	79.8	65.5
PolarNet [77]	L	74.7	28.2	85.3	90.9	35.1	77.5	71.3	58.8	57.4	76.1	96.5	71.1	74.7	74.0	87.3	85.7	71.0
Salsanext [15]	L	74.8	34.1	85.9	88.4	42.2	72.4	72.2	63.1	61.3	76.5	96.0	70.8	71.2	71.5	86.7	84.4	72.2
SVASeg [79]	L	73.1	44.5	88.4	86.6	48.2	80.5	77.7	65.6	57.5	82.1	96.5	70.5	74.7	74.6	87.3	86.9	74.7
PVKD [30]	L	76.2	40.0	90.2	94.0	50.9	77.4	78.8	64.7	62.0	84.1	96.6	71.4	76.4	76.3	90.3	86.9	76.0
AMVNet [42]	L	79.8	32.4	82.2	86.4	62.5	81.9	75.3	72.3	83.5	65.1	97.4	67.0	78.8	74.6	90.8	87.9	76.1
Cylinder3D [81]	L	76.4	40.3	91.3	93.8	51.3	78.0	78.9	64.9	62.1	84.4	96.8	71.6	76.4	75.4	90.5	87.4	76.1
RPVNet [70]	L	78.2	43.4	92.7	93.2	49.0	85.7	80.5	66.0	66.9	84.0	96.9	73.5	75.9	76.0	90.6	88.9	77.6
SDSeg3D [39]	L	77.5	49.4	93.9	92.5	54.9	86.7	80.1	67.8	65.7	86.0	96.4	74.0	74.9	74.5	86.0	82.8	77.7
SDSeg3D <sup>†</sup> [39]	L	78.2	52.8	94.5	93.1	54.5	88.1	82.2	69.4	67.3	86.6	96.4	74.5	75.2	75.3	87.1	84.1	78.7
WaffleIron [54]	L	78.7	51.3	93.6	88.2	47.2	86.5	81.7	68.9	69.3	83.1	96.9	74.3	75.6	74.2	87.2	85.2	77.6
WaffleIron <sup>†</sup> [54]	L	79.8	53.8	94.3	87.6	49.6	89.1	83.8	70.6	72.7	84.9	97.1	75.8	76.5	75.9	87.8	86.3	79.1
RangeFormer [36]	L	78.0	45.2	94.0	92.9	58.7	83.9	77.9	69.1	63.7	85.6	96.7	74.5	75.1	75.3	89.1	87.5	78.1
SphereFormer [38]	L	77.7	43.8	94.5	93.1	52.4	86.9	81.2	65.4	73.4	85.3	97.0	73.4	75.4	75.0	91.0	89.2	78.4
SphereFormer <sup>†</sup> [38]	L	78.7	46.7	95.2	93.7	54.0	88.9	81.1	68.0	74.2	86.2	97.2	74.3	76.3	75.8	91.4	89.7	79.5
2DPASS [72]	L+C	74.4	44.3	93.6	92.0	54.0	79.7	78.9	57.2	72.5	85.7	96.2	72.7	74.1	74.5	87.5	85.4	76.4
2DPASS <sup>†</sup> [72]	L+C	77.0	50.4	95.9	94.2	56.0	86.0	81.9	64.4	76.9	88.6	96.8	75.4	76.7	76.4	88.9	86.6	79.5
2D3DNet [22]	L+C	78.3	55.1	95.4	87.7	59.4	79.3	80.7	70.2	68.2	86.6	96.1	74.9	75.7	75.1	91.4	89.9	79.0
PMF [82]	L+C	74.1	46.6	89.8	92.1	57.0	77.7	80.9	70.9	64.6	82.9	95.5	73.3	73.6	74.8	89.4	87.7	76.9
PMF-R50 [82]	L+C	74.9	55.4	91.0	93.0	60.5	80.3	83.2	73.6	67.2	84.5	95.9	75.1	74.6	75.5	90.3	89.0	79.0
Uni-modal baseline	L	74.9	19.2	73.5	89.3	36.9	63.1	68.2	52.0	64.5	74.4	96.8	73.0	75.6	75.1	87.8	86.0	69.4
EPMF (Ours)	L+C	79.7	55.8	96.0	92.4	65.6	86.4	80.9	74.3	68.1	87.0	97.0	75.6	76.2	76.3	90.2	88.1	<b>80.6</b>

SalsaNext [15], Cylinder3D [81], 2DPASS [72], etc. Following [15], [33], [81], we use sequence 08 for validation. The remaining sequences (00-07 and 09-10) are used as the training set. We evaluate the release models of the state-of-the-art LiDAR-only methods on our data set. Because SPVNas [62] did not release its best model, we report the result of the best-released model (with 65G MACs). In addition, we re-implement two fusion-based methods, *i.e.*, RGBAL [47] and PointPainting [64] on our data set. To further understand the effectiveness of our proposed fusion strategies, we construct a *uni-modal baseline*, which uses the same network architecture of LiDAR stream of the proposed two-stream network. Note that we do not use test time augmentation (TTA) during evaluation as this technique is time-consuming and cannot be adopted to real applications on the car.

From Table 3, EPMF outperforms the uni-modal baseline by 7.3% in mIoU. Compared with PMF, our EPMF also achieves 2.0%

improvements in mIoU. However, EPMF performs slightly worse than the pre-trained model of 2DPASS [72] on SemanticKITTI-FV. Note that our EPMF is trained on SemanticKITTI-FV with only 16.03% points of SemanticKITTI. For fair comparisons, we also train 2DPASS using the officially released code and evaluate the model on SemanticKITTI-FV. In this case, our EPMF outperforms 2DPASS by 4.1% in mIoU.

#### 4.3.2 Results on nuScenes

Following [81], to evaluate our method on more complex scenes, we compare EPMF with the state-of-the-art methods on the nuScenes LiDAR-seg validation set. The experimental results are shown in Table 4. Note that the point clouds of nuScenes are sparser than those of SemanticKITTI (35k points/frame vs. 125k points/frame). Thus, it is more challenging for 3D segmentation tasks. In this case, EPMF achieves the best performance on nuScenes validation set. Specifically, EPMF outperforms the best LiDAR-only method, *i.e.*,

TABLE 5

Comparisons on the nuScenes test set. **L** indicates LiDAR-only methods. **L+C** indicates fusion-based methods. \* represents the results reported in the leader board of nuScenes. † indicates the results with test-time augmentation. ‡ denotes the results with additional fine-tuning with class re-sampling before the leaderboard submission. The **bold** numbers indicate the best results.

Method	Modality	Barrier	Bicycle	Bus	Car	Construction	Motorcycle	Pedestrian	Traffic-cone	Trailer	Truck	Driveable	Other-flat	Sidewalk	Terrain	Made	Vegetation	mIoU (%)
PolarNet [77]	L	72.2	16.8	77.0	86.5	51.1	69.7	64.8	54.1	69.7	63.4	96.6	67.1	77.7	72.1	87.1	84.4	69.4
AMVNet [42]	L	79.8	32.4	82.2	86.4	62.5	81.9	75.3	72.3	83.5	65.1	97.4	67.0	78.8	74.6	90.8	87.9	76.1
Cylinder3D* [81]	L	82.8	29.8	84.3	89.4	63.0	79.3	77.2	73.4	84.6	69.1	97.7	70.2	80.3	75.5	90.4	87.6	77.2
SPVNAS* [62]	L	80.0	30.0	91.9	90.8	64.7	79.0	75.6	70.9	81.0	74.6	97.4	69.2	80.0	76.1	89.3	87.1	77.4
AF2S3Net* [12]	L	78.9	52.2	89.9	84.2	77.4	74.3	77.3	72.0	83.9	73.8	97.1	66.5	77.5	74.0	87.7	86.8	78.3
RangeFormer [36]	L	83.9	46.1	89.4	89.2	70.3	83.3	75.4	72.5	81.4	71.1	95.6	68.5	77.3	73.4	89.3	86.9	78.3
RangeFormer† [36]	L	85.6	47.4	91.2	90.9	70.7	84.7	77.1	74.1	83.2	72.6	97.5	70.7	79.2	75.4	91.3	88.9	80.1
SphereFormer [38]	L	81.5	39.7	93.4	87.5	66.4	75.7	77.2	70.6	85.6	73.6	97.6	64.8	79.8	75.0	92.2	89.0	78.1
SphereFormer‡ [38]	L	83.3	39.2	94.7	92.5	77.5	84.2	84.4	79.1	88.4	78.3	97.9	69.0	81.5	77.2	93.4	90.2	<b>81.9</b>
2DPASS†‡ [72]	L+C	81.7	55.3	92.0	91.8	73.3	86.5	78.5	72.5	84.7	75.5	97.6	69.1	79.9	75.5	90.2	88.0	80.8
2D3DNet‡ [22]	L+C	83.0	59.4	88.0	85.1	63.7	84.4	82.0	76.0	84.8	71.9	96.9	67.4	79.8	76.0	92.1	89.2	80.0
PMF [82]	L+C	80.1	35.7	79.7	86.0	62.4	76.3	76.9	73.6	78.5	66.9	97.1	65.3	77.6	74.4	89.5	87.7	75.5
PMF-R50 [82]	L+C	82.1	40.3	80.9	86.4	63.7	79.2	79.8	75.9	81.2	67.1	97.3	67.7	78.1	74.5	89.9	88.5	77.0
EPMF (Ours)	L+C	76.9	39.8	90.3	87.8	72.0	86.4	79.6	76.6	84.1	74.9	97.7	66.4	79.5	76.4	91.1	87.9	79.2

TABLE 6  
Comparisons on A2D2 test set - Part1. The **bold** numbers indicate the best results.

Method	Car	Bicycle	Pedestrian	Truck	Small-vehi	Traffic-signal	Traffic-signal	Utility-vehi	Sidebars	Bumper	Curbstone	Solid-line	Irrelevant signs	Road-blocks	Tractor	Non-driveable	Zebra-crossing	Obstacles	Poles	mIoU (%)
SqueezeSeg [66]	9.7	0.0	0.0	15.8	0.0	0.7	64.4	0.0	0.4	0.0	2.2	15.6	0.5	15.9	0.0	0.0	0.0	0.0	0.3	8.9
SqueezeSegV2 [67]	15.4	0.2	8.6	63.8	0.0	16.8	61.7	0.6	0.1	0.0	14.8	24.7	12.7	33.2	0.0	5.8	0.0	0.2	5.2	16.4
DarkNet53 [3]	15.2	0.8	6.1	68.5	0.0	15.5	63.8	0.4	0.3	0.0	17.3	23.8	13.3	35.6	0.0	6.3	0.0	3.9	7.6	17.2
PolarNet [77]	23.8	10.1	18.2	69.7	9.6	49.1	58.5	0.0	11.3	0.0	28.3	37.6	24.8	42.8	0.0	14.8	0.0	8.0	11.0	23.9
Cylinder3D [80]	32.6	6.0	14.9	74.7	20.9	51.2	65.9	3.6	6.8	0.0	33.2	42.8	24.1	40.7	0.1	16.7	0.0	10.8	5.4	24.2
2DPASS [72]	26.2	1.1	16.9	78.7	12.4	39.0	66.2	2.8	4.8	0.0	19.9	20.0	18.8	31.4	0.0	6.0	0.0	9.9	3.4	18.9
PMF [82]	84.5	36.0	47.5	89.4	38.1	75.0	82.8	4.6	64.9	0.0	54.8	72.8	41.2	77.1	0.0	27.1	0.0	30.8	29.6	41.8
Uni-modal baseline	42.3	5.6	13.7	75.2	0.6	32.2	65.6	0.2	13.3	0.0	30.6	50.5	12.3	53.9	0.1	6.9	0.0	6.1	9.2	21.8
EPMF(Ours)	79.4	42.3	49.4	90.7	79.8	78.4	84.4	11.8	69.1	0.0	56.8	73.9	47.6	78.3	0.0	26.0	0.0	35.5	33.0	<b>45.0</b>

SphereFormer, by 2.2% in mIoU. Compared with PMF, our EPMF achieves 3.7% improvements in mIoU.

In addition, we further provide the results on nuScenes test set. Different from the results on the nuScenes validation set, existing methods always use additional techniques, including test-time augmentation or fine-tuning with class re-sampling that improves performance on rare classes, to pursue better results on the leaderboard. Unlike existing methods, we directly evaluate our model on the test set without additional techniques. For fair comparisons, we also evaluate the performance of the released models of SphereFormer without TTA or fine-tuning. As shown in Table 5, our EPMF outperforms SphereFormer [38] by 1.1% in mIoU under the same evaluation settings. Compared with PMF, our EPMF consistently achieves 3.7% improvements in mIoU on nuScenes test set. These results are consistent with our expectations. Since EPMF incorporates RGB images, our fusion strategy is capable of addressing such challenging segmentation under extremely sparse point clouds.

#### 4.3.3 Results on A2D2

To further evaluate the performance of our method on more sparse and irregular point clouds, we conduct experiments on A2D2 and compare our EPMF to existing methods. For a fair comparison, we select the model with the best validation performance and report the evaluation results on the test set. Note that both Cylinder3D and

2DPASS did not report the results on A2D2 test set, we also train Cylinder3D and 2DPASS using the officially released code and report the results on the test set. As shown in Table 6 and Table 7, both our PMF and EPMF outperform the LiDAR-only method by a large margin. Specifically, our EPMF outperforms Cylinder3D by 20.8% in mIoU. Compared to uni-modal baseline, EPMF achieves 23.2% improvements in mIoU, which indicates the effectiveness of our proposed fusion strategies. Moreover, our EPMF outperforms PMF consistently, with 3.2% improvement in mIoU.

#### 4.4 Distance-based Evaluation

In 3D LiDAR perception, the point cloud becomes sparser with the increase of perception distance. As long-range perception is important to the safety of autonomous cars, we further conduct a distance evaluation on the benchmark datasets and investigate the performance of our method under different distances. As shown in Figure 7, both PMF and EPMF outperform the uni-modal baseline by a large margin under different distances on nuScenes and A2D2, which indicates that our fusion strategy can incorporate the information from RGB images effectively.

We also notice that EPMF cannot consistently outperform PMF or the uni-modal baseline on SemanticKITTI-FV. We argue that this phenomenon is mainly caused by the following reasons. First, we compare the number of classes at different distances w.r.t. the validation set of SemanticKITTI-FV and nuScenes, as well as

TABLE 7  
Comparisons on A2D2 test set - Part2. The **bold** numbers indicate the best results.

Method	RD restricted area	Animals	Grid structure	Signal corpus	Drivable cobbleston	Electronic traffic	Slow drive area	Nature object	Parking area	Sidewalk	Ego car	Painted driv instr	Traffic guide obj	Dashed line	RD normal street	Sky	Buildings	Blurred area	Rain dirt	mIoU (%)
SqueezeSeg [66]	0.0	0.0	0.0	0.0	0.0	0.0	64.5	0.0	13.7	0.0	0.0	0.1	0.2	77.7	10.4	27.7	0.0	0.0	8.9	
SqueezeSegV2 [67]	29.5	0.0	10.3	5.5	2.7	0.0	1.9	76.4	3.8	29.2	0.0	6.4	12.4	17.1	85.8	12.1	50.9	0.0	0.0	16.4
DarkNet53 [3]	38.7	0.0	10.8	4.4	3.3	0.0	0.0	77.9	3.1	31.5	0.0	9.4	7.3	15.7	86.4	12.9	55.2	0.0	0.0	17.2
PolarNet [77]	55.6	0.0	14.8	11.9	7.0	0.0	4.4	81.6	12.8	42.5	0.0	12.7	11.5	31.8	90.3	9.2	57.0	0.0	0.0	23.9
Cylinder3D [80]	57.7	0.0	20.5	10.6	10.4	1.5	2.5	82.9	9.1	47.4	0.0	19.1	8.0	33.8	89.7	13.0	63.5	0.0	0.0	24.2
2DPASS [72]	37.0	0.0	13.8	6.7	5.0	0.0	0.3	74.8	5.9	35.1	0.0	8.6	10.0	11.9	85.6	6.9	57.5	0.0	0.0	18.9
PMF [82]	65.4	0.0	46.6	22.6	26.3	0.0	3.7	91.7	17.8	48.6	0.0	48.9	64.1	69.2	94.1	53.5	78.4	0.0	0.0	41.8
Uni-modal baseline	45.7	0.0	14.0	2.7	9.1	0.0	0.1	74.6	3.1	38.1	0.0	12.3	11.5	35.5	90.2	14.2	59.0	0.0	0.0	21.8
EPMF(Ours)	70.3	0.0	52.1	33.3	29.4	0.0	1.5	92.4	15.3	57.4	0.0	55.0	66.3	70.0	94.5	54.2	80.5	0.0	0.0	<b>45.0</b>

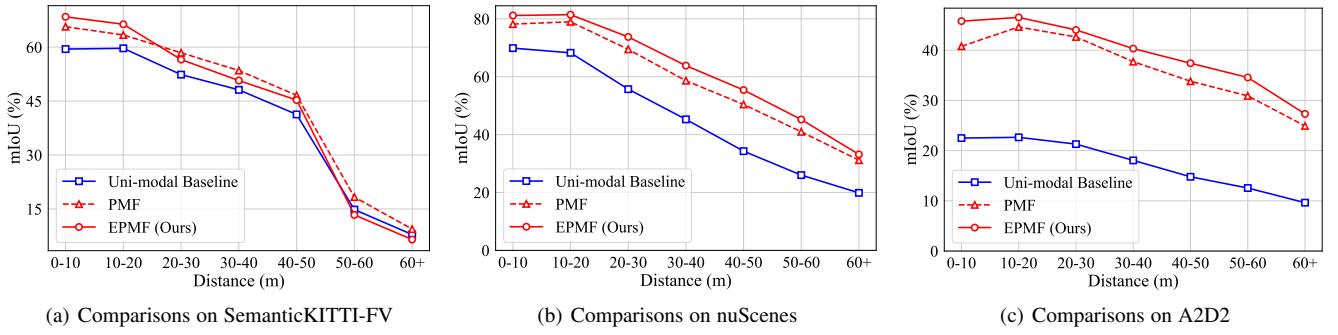


Fig. 7. Distance-based evaluation on SemanticKITTI-FV, nuScenes and A2D2. As the distance increases, the point cloud becomes sparser.

TABLE 8  
Comparisons of the number of classes at different distances.

Dataset	Distance (m)						
	0-10	10-20	20-30	30-40	40-50	50-60	60+
SemanticKITTI-FV	18	19	19	19	19	11	8
nuScenes	16	16	16	16	16	16	16
A2D2	33	35	35	34	34	33	34

TABLE 9  
Percentage of point cloud distribution (%) at different distances.

Dataset	Distance (m)						
	0-10	10-20	20-30	30-40	40-50	50-60	60+
SemanticKITTI-FV	34.7	41.5	12.9	5.2	2.7	1.5	1.5
nuScenes	70.8	15.7	6.3	3.2	1.8	1.0	1.2
A2D2	27.4	36.9	19.5	9.4	4.1	1.6	1.1

the test set of A2D2. From Table 8, the number of classes on SemanticKITTI-FV decreases at long distances. Since it is difficult for our method to beat the baseline on all the semantic classes of SemanticKITTI-FV, EPMF performs worse than the uni-modal baseline when the distance is larger than 50m on SemanticKITTI-FV, which covers only 3% of point clouds (See Table 9). Second, to make a trade-off between efficiency and effectiveness, we insert a down-sampling operation into the contextual module of LiDAR stream to improve its efficiency. As the point clouds of SemanticKITTI are dense, reducing the resolution of point cloud features also leads to higher classification errors at long distances. Nevertheless, EPMF achieves better performance when the distance is less than 20m which covers 76.2% of the points. Overall, our EPMF outperforms PMF by 2.0% in mIoU with 2.06× acceleration

TABLE 10  
Inference time of different methods on SemanticKITTI using TensorRT.  
“-” indicates the results that are not available. For a fair comparison, Cylinder3D is accelerated by sparse convolution.

Method	#FLOPs	#Params.	Inference time	mIoU
PointPainting [64]	51.0 G	28.1 M	2.3 ms	54.5%
RGBAL [47]	55.0 G	13.2 M	2.7 ms	56.2%
SalsaNext [15]	31.4 G	6.7 M	1.6 ms	59.4%
Cylinder3D [81]	-	55.9 M	62.5 ms	64.9%
PMF [82]	854.7 G	36.3 M	22.3 ms	63.9%
EPMF (Ours)	418.0 G	34.2 M	14.2 ms	<b>65.9%</b>

on SemanticKITTI-FV.

#### 4.5 Efficiency analysis

In this section, we evaluate the efficiency of EPMF on GeForce RTX 3090. Note that we consider the efficiency of PMF in two aspects. First, since predictions of the camera stream are fused into the LiDAR stream, we remove the decoder of the camera stream to speed up the inference. Second, both our PMF and EPMF are built on 2D convolutions and can be easily optimized by existing inference toolkits, e.g., TensorRT. In contrast, Cylinder3D is built on 3D sparse convolutions [24] and is difficult to be accelerated by TensorRT. We report the inference time of different models optimized by TensorRT in Table 10. From the results, our PMF achieves the best performance on nuScenes and is 2.8× faster than Cylinder3D (22.3 ms vs. 62.5 ms) with fewer parameters. While compared to PMF, our EPMF achieves 1.6× acceleration (14.2 ms vs. 22.3 ms) with 2.0% improvements in mIoU.

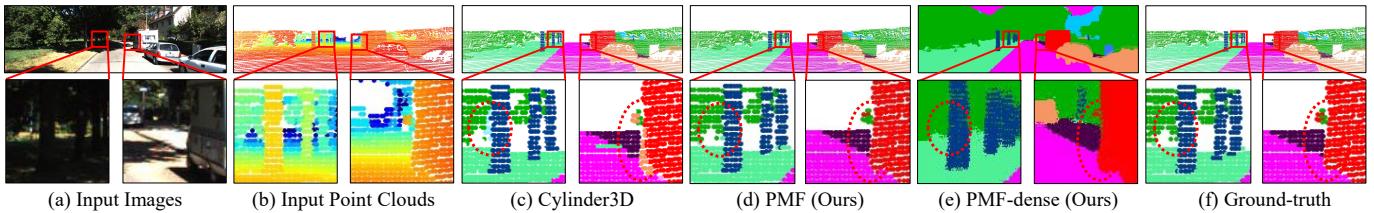


Fig. 8. Qualitative results on SemanticKITTI-FV. The red dashed circle indicates the difference between the results of PMF and the baseline.

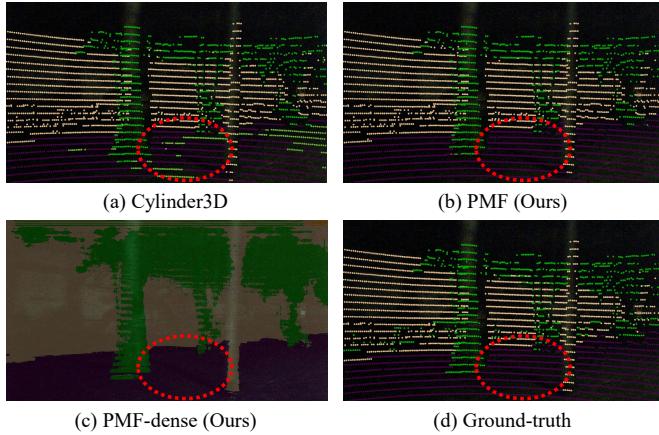


Fig. 9. Qualitative results on nuScenes. We use the corresponding images (night) as the background of both the predictions and labels. We highlight the difference between the results of PMF and the baseline with the red dashed circle.

#### 4.6 Qualitative evaluation

To better understand the benefits of PMF, we visualize the predictions of PMF on the benchmark data sets. From Figure 8, compared with Cylinder3D, PMF achieves better performance at the boundary of objects. For example, as shown in Figure 8 (d), the truck segmented by PMF has a more complete shape. More critically, PMF is robust in different lighting conditions. Specifically, as illustrated in Figure 9, PMF outperforms the baselines on more challenging scenes (*e.g.*, night). In addition, as demonstrated in Figure 8 (e) and Figure 9 (c), PMF generates dense segmentation results that combine the benefits of both the camera and LiDAR, which are significantly different from existing LiDAR-only and fusion-based methods.

In Figure 10, we show the error maps of both EPMF and the unimodal baseline on three benchmark data sets under different scenes. From the results, our EPMF can fully utilize both information from point clouds and RGB images and thus performs well in some challenging scenes with extremely sparse point clouds or poor lighting conditions.

#### 4.7 Robustness analysis

To investigate the robustness of EPMF on adversarial samples, we first insert extra objects (*e.g.*, a traffic sign) into the images and keep the point clouds unchanged. In addition, we implement a camera-only method, *i.e.*, FCN [45], on SemanticKITTI as the baseline. Note that we do not use any adversarial training technique during training. As demonstrated in Figure 11, the camera-only

TABLE 11  
Comparisons to SegFormer on nuScenes at day/night time. **C** and **L** indicate the model with camera or LiDAR as input, respectively. **L+C** indicates the model with both camera and LiDAR inputs.

Method	Input	Day	Night
SegFormer-B5 [68]	C	67.6%	41.6%
SegFormer-B5 [68]	L	56.6%	46.9%
EPMF (Ours)	L+C	81.4%	62.0%

TABLE 12  
Ablation study for the network components on the SemanticKITTI-FV validation set. **PP** denotes perspective projection. **RF** denotes the residual-based fusion module. **PL** denotes perception-aware loss. The **bold** number is the best result.

	Baseline	PP	ASPP	RF	PL	mIoU (%)
1	✓					57.2
2	✓	✓				57.6
3	✓	✓	✓			59.7
4	✓			✓	✓	55.8
5	✓	✓	✓	✓		61.7
6	✓	✓	✓	✓	✓	<b>63.9</b>

methods are easily affected by changes in the input images. In contrast, because EPMF integrates reliable point cloud information, the noise in the images is reduced during feature fusion and imposes only a slight effect on the model performance.

To investigate the performances of EPMF under different lighting conditions, we evaluate the performance of EPMF on nuScenes validation set at day/night time. We initialize SegFormer-B5 by the officially released weights on CityScape and train the model SegFormer-B5 with the camera or LiDAR as inputs. Note that some of the classes may not appear at night, we only report the mIoU of the available classes. From Table 11, as the camera only provides limited information at night, SegFormer with camera-only inputs performs worse than the LiDAR-only counterpart. By fusing both information from the camera and LiDAR, EPMF achieves better performance at night, with 15.1% improvements in mIoU compared to SegFormer-B5.

## 5 ABLATION STUDY

### 5.1 Effect of network components

We study the effect of the network components of PMF, *i.e.*, perspective projection, ASPP, residual-based fusion modules, and perception-aware loss. The experimental results are shown in Table 12. Since we use only the front-view point clouds of SemanticKITTI, we train SalsaNext as the baseline on our data set using the officially released code. Comparing the first and second lines in Table 12, perspective projection achieves only a 0.4% mIoU

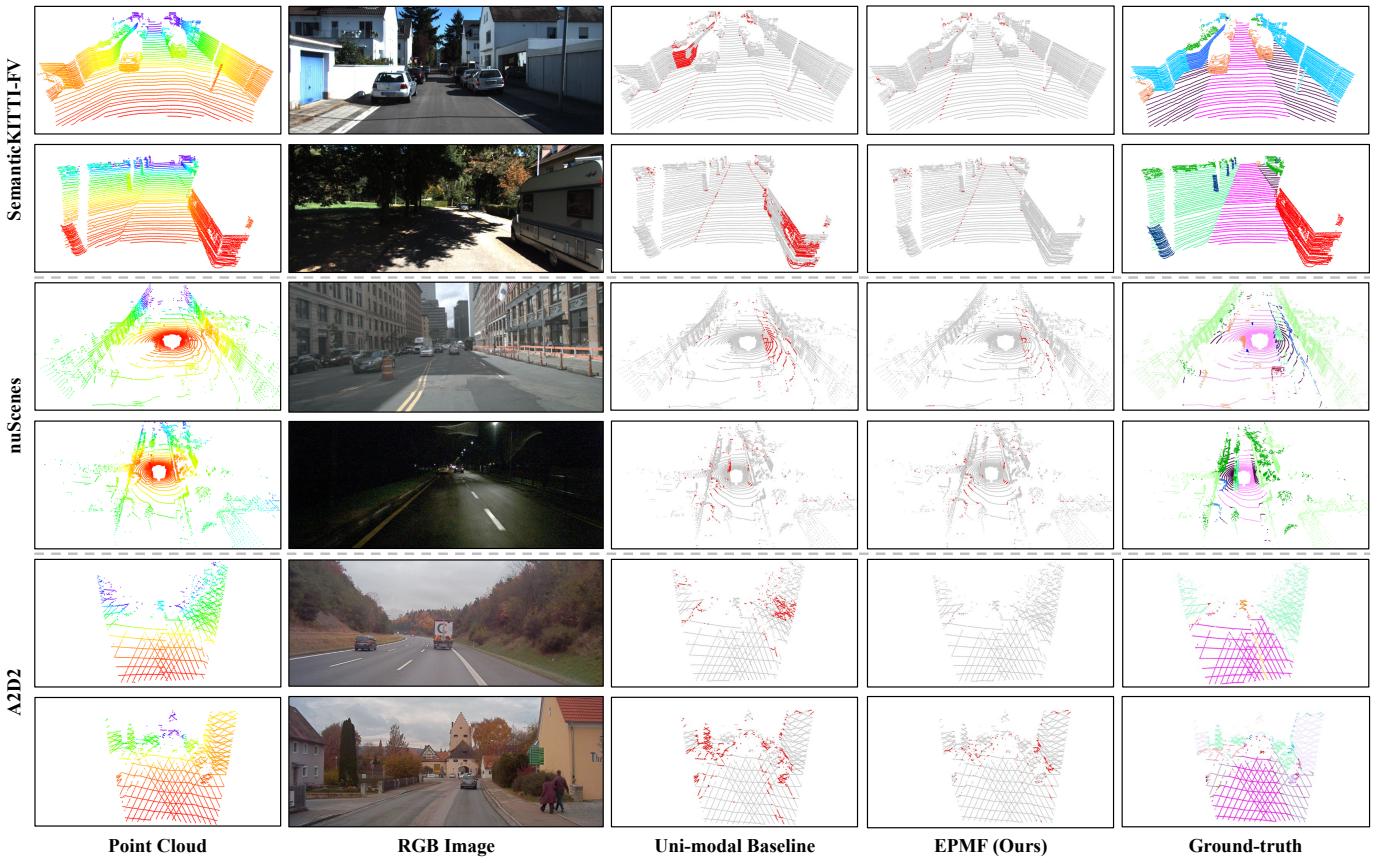


Fig. 10. Qualitative results on three benchmark data sets. We show the error maps of both our method and the uni-modal baseline, in which the red points indicate mispredictions. Zoom in for more details.

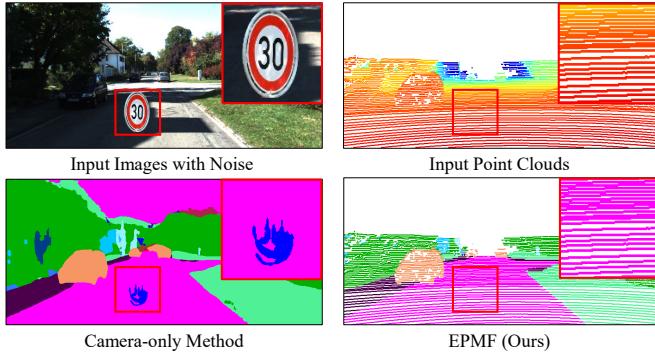


Fig. 11. Comparisons of EPMF and camera-only methods on adversarial samples. The camera-only methods use only RGB images as inputs, while PMF uses both images and point clouds as inputs. We highlight the inserted traffic sign with a red box.

improvement over spherical projection with LiDAR-only input. In contrast, comparing the fourth and fifth lines, perspective projection brings a 5.9% mIoU improvement over spherical projection with multimodal data inputs. From the third and fifth lines, our fusion modules bring 2.0% mIoU improvement to the fusion network. Moreover, comparing the fifth and sixth lines, the perception-aware losses improve the performance of the network by 2.2% in mIoU.

TABLE 13  
Ablation study for the proposed improved techniques on SemanticKITTI-FV validation set.

Proposed strategies	mIoU	#FLOPs	#Params.
PMF [82]	63.9%	859.7 G	36.4 M
+dropping decoder of camera stream	63.9%	854.7 G	36.3 M
+cross-modal alignment and crop	64.4%	739.9 G	36.3 M
+improved contextual module	64.8%	418.0 G	34.2 M
+fusing high-level LiDAR feature	65.9%	418.0 G	34.2 M

TABLE 14  
Model performance with different image masked ratios on SemanticKITTI-FV. We report the mIoU (%) with/without fine-tuning.

Masked Ratio (%)	0	10	20	30	40	50
without fine-tuning	65.9	65.5	64.8	63.0	60.0	54.8
with fine-tuning	65.7	65.4	65.0	64.2	62.9	60.6

## 5.2 Effect of perception-aware loss

To investigate the effect of perception-aware loss, we visualize the predictions of the LiDAR stream networks with and without perception-aware loss in Figure 12. From the results, perception-aware loss helps the LiDAR stream capture the perceptual information from the images. For example, the model trained with perception-aware loss learns the complete shape of cars, while the baseline model focuses only on the local features of points. As the perception-aware loss introduces the perceptual difference between

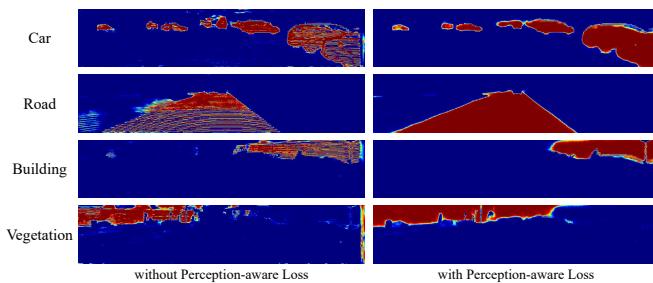


Fig. 12. Comparisons of the predictions w.r.t. the networks trained with and without perception-aware loss. **PL** denotes the perception-aware loss. Red indicates predictions with higher confidence scores. We only show the predictions of Car for the sake of clarity.

the RGB images and the point clouds, it enables an effective fusion of the perceptual information from the data of both modalities. As a result, our PMF generates dense predictions that combine the benefits of both the images and point clouds.

### 5.3 Effect of the improved techniques

In this section, we explore the effectiveness of the proposed improved techniques on SemanticKITTI-FV. As shown in Table 13, on top of PMF, dropping the decoder of the camera stream saves 5.0 GFLOPs in computation budgets without performance degradation. The proposed CAC further reduces 114.8 GFLOPs in computation costs by removing the useless area of RGB images. With the proposed improved contextual module, we achieve  $2.04\times$  acceleration compared to PMF (418.0 GFLOPs vs. 854.7 GFLOPs) with 0.9% improvement in mIoU. By fusing the high-level features of LiDAR stream into the camera stream, we further achieve 1.1% improvements in mIoU.

When conducting CAC, if the LiDAR has a larger vertical FOV than cameras, keeping point clouds outside the image also results in partial blank image inputs. To investigate the impact of these areas without image information, we partially mask the RGB image from bottom to top with ratios from 10% to 50%, and evaluate the model performance on SemanticKITTI-FV. As shown in Table 14, masking 10% of RGB image only results in 0.4% degradation of mIoU. With the increasing of the masked ratio, the model performance suffers from significant degradation. Nevertheless, the impact of masked images can be mitigated by introducing the image mask into training or fine-tuning.

## 6 CONCLUSION

In this work, we have proposed a perception-aware multi-sensor fusion scheme for 3D LiDAR semantic segmentation. Unlike existing methods that conduct feature fusion in the LiDAR coordinate system, we project the point clouds to the camera coordinate system to enable a collaborative fusion of the perceptual features from the two modalities. By fusing complementary information from both cameras and LiDAR, PMF is robust in complex outdoor scenes with extremely sparse point clouds or poor lighting conditions. Moreover, we propose EPMF which improves the efficiency and effectiveness of PMF. Specifically, we introduce cross-modal alignment and cropping to reduce unnecessary computation. Besides, we also adjust the architecture of the fusion network by fusing high-level features into the camera stream and exploring the design of contextual modules under perspective projection. The experimental

results on three benchmarks show the superiority of our method. In the future, we will extend EPMF to other challenging tasks, e.g., object detection and semantic scene completion.

## REFERENCES

- [1] E. E. Aksoy, S. Baci, and S. Cavdar. Salsanet: Fast road and vehicle segmentation in lidar point clouds for autonomous driving. *IEEE Intelligent Vehicles Symposium*, pages 926–932, 2020.
- [2] V. Badrinarayanan, A. Kendall, and R. Cipolla. Segnet: A deep convolutional encoder-decoder architecture for image segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 39(12):2481–2495, 2017.
- [3] J. Behley, M. Garbade, A. Milioto, J. Quenzel, S. Behnke, C. Stachniss, and J. Gall. Semantickitti: A dataset for semantic scene understanding of lidar sequences. In *IEEE International Conference on Computer Vision*, pages 9297–9307, 2019.
- [4] M. Berman, A. R. Triki, and M. B. Blaschko. The lovász-softmax loss: A tractable surrogate for the optimization of the intersection-over-union measure in neural networks. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 4413–4421, 2018.
- [5] A. Bochkovskiy, C.-Y. Wang, and H.-Y. M. Liao. Yolov4: Optimal speed and accuracy of object detection. *arXiv preprint arXiv:2004.10934*, pages 1–17, 2020.
- [6] G. J. Brostow, J. Shotton, J. Fauqueur, and R. Cipolla. Segmentation and recognition using structure from motion point clouds. In *European Conference on Computer Vision*, pages 44–57. Springer, 2008.
- [7] H. Caesar, V. Bankiti, A. H. Lang, S. Vora, V. E. Lioung, Q. Xu, A. Krishnan, Y. Pan, G. Baldan, and O. Beijbom. nuscenes: A multimodal dataset for autonomous driving. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 11621–11631, 2020.
- [8] H. Cai, C. Gan, T. Wang, Z. Zhang, and S. Han. Once-for-all: Train one network and specialize it for efficient deployment. In *International Conference on Learning Representations*, pages 1–15, 2020.
- [9] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 40(4):834–848, 2017.
- [10] L.-C. Chen, G. Papandreou, F. Schroff, and H. Adam. Rethinking atrous convolution for semantic image segmentation. *arXiv preprint arXiv:1706.05587*, pages 1–14, 2017.
- [11] B. Cheng, I. Misra, A. G. Schwing, A. Kirillov, and R. Giridhar. Masked-attention mask transformer for universal image segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 1290–1299, 2022.
- [12] R. Cheng, R. Razani, E. Taghavi, E. Li, and B. Liu. 2-s3net: Attentive feature fusion with adaptive feature selection for sparse semantic segmentation network. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 12547–12556, 2021.
- [13] C. Choy, J. Gwak, and S. Savarese. 4d spatio-temporal convnets: Minkowski convolutional neural networks. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 3075–3084, 2019.
- [14] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, and B. Schiele. The cityscapes dataset for semantic urban scene understanding. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 3213–3223, 2016.
- [15] T. Cortinhal, G. Tzelepis, and E. Erdal Aksoy. Salsanext: Fast, uncertainty-aware semantic segmentation of lidar point clouds. In *International Symposium on Visual Computing*, pages 207–222. Springer, 2020.
- [16] P. Dollar, C. Wojek, B. Schiele, and P. Perona. Pedestrian detection: An evaluation of the state of the art. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 34(4):743–761, 2011.
- [17] C. Feichtenhofer, A. Pinz, and A. Zisserman. Convolutional two-stream network fusion for video action recognition. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 1933–1941, 2016.
- [18] C. Gan, H. Zhao, P. Chen, D. Cox, and A. Torralba. Self-supervised moving vehicle tracking with stereo sound. In *IEEE International Conference on Computer Vision*, pages 7053–7062, 2019.
- [19] A. Geiger, P. Lenz, C. Stiller, and R. Urtasun. Vision meets robotics: The kitti dataset. *The International Journal of Robotics Research*, 32(11):1231–1237, 2013.
- [20] A. Geiger, P. Lenz, and R. Urtasun. Are we ready for autonomous driving? the kitti vision benchmark suite. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 3354–3361. IEEE, 2012.
- [21] A. Geiger, F. Moosmann, Ö. Car, and B. Schuster. Automatic camera and range sensor calibration using a single shot. In *IEEE International Conference on Robotics and Automation*, pages 3936–3943. IEEE, 2012.

- [22] K. Genova, X. Yin, A. Kundu, C. Pantofaru, F. Cole, A. Sud, B. Breitung, B. Shucker, and T. Funkhouser. Learning 3d semantic segmentation with only 2d image supervision. In *International Conference on 3D Vision*, pages 361–372. IEEE, 2021.
- [23] J. Geyer, Y. Kassahun, M. Mahmudi, X. Ricou, R. Durgesh, A. S. Chung, L. Hauswald, V. H. Pham, M. Mühllegg, S. Dorn, et al. A2d2: Audi autonomous driving dataset. *arXiv preprint arXiv:2004.06320*, pages 1–10, 2020.
- [24] B. Graham, M. Engelcke, and L. Van Der Maaten. 3d semantic segmentation with submanifold sparse convolutional networks. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 9224–9232, 2018.
- [25] Y. Guo, Y. Zheng, M. Tan, Q. Chen, J. Chen, P. Zhao, and J. Huang. Nat: Neural architecture transformer for accurate and compact architectures. In *Advances in Neural Information Processing Systems*, pages 737–748, 2019.
- [26] Y. Guo, Y. Zheng, M. Tan, Q. Chen, Z. Li, J. Chen, P. Zhao, and J. Huang. Towards accurate and compact architectures via neural architecture transformer. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(10):6501–6516, 2021.
- [27] S. Han, H. Mao, and W. J. Dally. Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding. In *International Conference on Learning Representations*, pages 1–14, 2016.
- [28] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 770–778, 2016.
- [29] G. Hinton, O. Vinyals, and J. Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, pages 1–9, 2015.
- [30] Y. Hou, X. Zhu, Y. Ma, C. C. Loy, and Y. Li. Point-to-voxel knowledge distillation for lidar semantic segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 8479–8488, 2022.
- [31] Q. Hu, B. Yang, L. Xie, S. Rosa, Y. Guo, Z. Wang, N. Trigoni, and A. Markham. Randla-net: Efficient semantic segmentation of large-scale point clouds. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 11108–11117, 2020.
- [32] Z. Huang, X. Wang, L. Huang, C. Huang, Y. Wei, and W. Liu. Ccnet: Criss-cross attention for semantic segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 603–612, 2019.
- [33] M. Jaritz, T.-H. Vu, R. d. Charette, E. Wirbel, and P. Pérez. xmuda: Cross-modal unsupervised domain adaptation for 3d semantic segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 12605–12614, 2020.
- [34] A. Kendall, Y. Gal, and R. Cipolla. Multi-task learning using uncertainty to weigh losses for scene geometry and semantics. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 7482–7491, 2018.
- [35] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations*, pages 1–15, 2015.
- [36] L. Kong, Y. Liu, R. Chen, Y. Ma, X. Zhu, Y. Li, Y. Hou, Y. Qiao, and Z. Liu. Rethinking range view representation for lidar segmentation. In *IEEE International Conference on Computer Vision*, pages 228–240, 2023.
- [37] G. Krispel, M. Opitz, G. Waltner, H. Possegger, and H. Bischof. Fuseseg: Lidar point cloud segmentation fusing multi-modal data. In *The IEEE Winter Conference on Applications of Computer Vision*, pages 1874–1883, 2020.
- [38] X. Lai, Y. Chen, F. Lu, J. Liu, and J. Jia. Spherical transformer for lidar-based 3d recognition. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 17545–17555, 2023.
- [39] J. Li, H. Dai, and Y. Ding. Self-distillation for robust lidar semantic segmentation in autonomous driving. In *European Conference on Computer Vision*, pages 659–676. Springer, 2022.
- [40] G. Lin, C. Shen, A. Van Den Hengel, and I. Reid. Efficient piecewise training of deep structured models for semantic segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 3194–3203, 2016.
- [41] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár. Focal loss for dense object detection. In *IEEE International Conference on Computer Vision*, pages 2980–2988, 2017.
- [42] V. E. Liong, T. N. T. Nguyen, S. Widjaja, D. Sharma, and Z. J. Chong. Amvnet: Assertion-based multi-view fusion network for lidar semantic segmentation. *arXiv preprint arXiv:2012.04934*, pages 1–10, 2020.
- [43] J. Liu, B. Zhuang, Z. Zhuang, Y. Guo, J. Huang, J. Zhu, and M. Tan. Discrimination-aware network pruning for deep model compression. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(8):4035–4051, 2021.
- [44] Z. Liu, S. Zhou, C. Suo, P. Yin, W. Chen, H. Wang, H. Li, and Y. Liu. Lpd-net: 3d point cloud learning for large-scale place recognition and environment analysis. In *IEEE International Conference on Computer Vision*, pages 2831–2840, 2019.
- [45] J. Long, E. Shelhamer, and T. Darrell. Fully convolutional networks for semantic segmentation. *IEEE Conference on Computer Vision and Pattern Recognition*, pages 3431–3440, 2015.
- [46] I. Loshchilov and F. Hutter. Sgdr: Stochastic gradient descent with warm restarts. In *International Conference on Learning Representations*, pages 1–16, 2017.
- [47] K. E. Madawy, H. Rashed, A. E. Sallab, O. Nasr, H. Kamel, and S. Yogamani. Rgb and lidar fusion based 3d semantic segmentation for autonomous driving. *IEEE Intelligent Transportation Systems Conference*, pages 7–12, 2019.
- [48] S. Mehta, M. Rastegari, A. Caspi, L. Shapiro, and H. Hajishirzi. Espnet: Efficient spatial pyramid of dilated convolutions for semantic segmentation. In *European Conference on Computer Vision*, pages 552–568, 2018.
- [49] G. P. Meyer, J. Charland, D. Hegde, A. Laddha, and C. Vallespi-Gonzalez. Sensor fusion for joint 3d object detection and semantic segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 1–8, 2019.
- [50] A. Milioto, I. Vizzo, J. Behley, and C. Stachniss. Rangenet++: Fast and accurate lidar semantic segmentation. In *IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 4213–4220. IEEE, 2019.
- [51] Y. E. Nesterov. A method for solving the convex programming problem with convergence rate  $O(1/K^2)$ . In *Proceedings of the USSR Academy of Sciences*, volume 269, pages 543–547, 1983.
- [52] S. Niu, J. Wu, Y. Zhang, Y. Guo, P. Zhao, J. Huang, and M. Tan. Disturbance-immune weight sharing for neural architecture search. *Neural Networks*, 144:553–564, 2021.
- [53] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems*, pages 8026–8037, 2019.
- [54] G. Puy, A. Boulch, and R. Marlet. Using a waffle iron for automotive point cloud semantic segmentation. In *IEEE International Conference on Computer Vision*, pages 3379–3389, 2023.
- [55] C. R. Qi, H. Su, K. Mo, and L. J. Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 652–660, 2017.
- [56] C. R. Qi, L. Yi, H. Su, and L. J. Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. In *Advances in Neural Information Processing Systems*, pages 5099–5108, 2017.
- [57] A. Rényi. On measures of entropy and information. In *Proceedings of the Fourth Berkeley Symposium on Mathematical Statistics and Probability, Volume 1: Contributions to the Theory of Statistics*, volume 4, pages 547–562. University of California Press, 1961.
- [58] R. B. Rusu, Z. C. Marton, N. Blodow, M. Dolha, and M. Beetz. Towards 3d point cloud based object maps for household environments. *Robotics and Autonomous Systems*, 56(11):927–941, 2008.
- [59] T. Shan and B. Englot. Lego-loam: Lightweight and ground-optimized lidar odometry and mapping on variable terrain. In *IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 4758–4765. IEEE, 2018.
- [60] K. Simonyan and A. Zisserman. Two-stream convolutional networks for action recognition in videos. In *Advances in Neural Information Processing Systems*, pages 568–576, 2014.
- [61] C. Sitawarin, A. N. Bhagoji, A. Mosenia, M. Chiang, and P. Mittal. Darts: Deceiving autonomous cars with toxic signs. *arXiv preprint arXiv:1802.06430*, pages 1–18, 2018.
- [62] H. Tang, Z. Liu, S. Zhao, Y. Lin, J. Lin, H. Wang, and S. Han. Searching efficient 3d architectures with sparse point-voxel convolution. In *European Conference on Computer Vision*, pages 685–702. Springer, 2020.
- [63] J. Uhrig, N. Schneider, L. Schneider, U. Franke, T. Brox, and A. Geiger. Sparsity invariant cnns. In *International Conference on 3D Vision*, pages 11–20. IEEE, 2017.
- [64] S. Vora, A. H. Lang, B. Helou, and O. Beijbom. Pointpainting: Sequential fusion for 3d object detection. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 4604–4612, 2020.
- [65] P. Wang, P. Chen, Y. Yuan, D. Liu, Z. Huang, X. Hou, and G. Cottrell. Understanding convolution for semantic segmentation. In *IEEE Winter Conference on Applications of Computer Vision*, pages 1451–1460. IEEE, 2018.
- [66] B. Wu, A. Wan, X. Yue, and K. Keutzer. Squeezeseg: Convolutional neural nets with recurrent crf for real-time road-object segmentation from

- 3d lidar point cloud. In *IEEE International Conference on Robotics and Automation*, pages 1887–1893. IEEE, 2018.
- [67] B. Wu, X. Zhou, S. Zhao, X. Yue, and K. Keutzer. Squeezesegv2: Improved model structure and unsupervised domain adaptation for road-object segmentation from a lidar point cloud. In *IEEE International Conference on Robotics and Automation*, pages 4376–4382. IEEE, 2019.
- [68] E. Xie, W. Wang, Z. Yu, A. Anandkumar, J. M. Alvarez, and P. Luo. Segformer: Simple and efficient design for semantic segmentation with transformers. *Advances in Neural Information Processing Systems*, 34:12077–12090, 2021.
- [69] C. Xu, B. Wu, Z. Wang, W. Zhan, P. Vajda, K. Keutzer, and M. Tomizuka. Squeezesegv3: Spatially-adaptive convolution for efficient point-cloud segmentation. In *European Conference on Computer Vision*, pages 1–19. Springer, 2020.
- [70] J. Xu, R. Zhang, J. Dou, Y. Zhu, J. Sun, and S. Pu. Rpnet: A deep and efficient range-point-voxel fusion network for lidar point cloud segmentation. In *IEEE International Conference on Computer Vision*, pages 16024–16033, 2021.
- [71] S. Xu, H. Li, B. Zhuang, J. Liu, J. Cao, C. Liang, and M. Tan. Generative low-bitwidth data free quantization. In *European Conference on Computer Vision*, pages 1–17. Springer, 2020.
- [72] X. Yan, J. Gao, C. Zheng, C. Zheng, R. Zhang, S. Cui, and Z. Li. 2dpass: 2d priors assisted semantic segmentation on lidar point clouds. In *European Conference on Computer Vision*, pages 677–695. Springer, 2022.
- [73] Y. Yuan and J. Wang. Ocnet: Object context network for scene parsing. *arXiv preprint arXiv:1809.00916*, pages 1–22, 2018.
- [74] F. Zhang, J. Fang, B. Wah, and P. Torr. Deep fusionnet for point cloud semantic segmentation. In *European Conference on Computer Vision*, volume 2, pages 644–663, 2020.
- [75] W. Zhang, Z. Wang, and C. C. Loy. Exploring data augmentation for multi-modality 3d object detection. *arXiv preprint arXiv:2012.12741*, pages 1–12, 2020.
- [76] Y. Zhang, T. Xiang, T. M. Hospedales, and H. Lu. Deep mutual learning. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 4320–4328, 2018.
- [77] Y. Zhang, Z. Zhou, P. David, X. Yue, Z. Xi, B. Gong, and H. Foroosh. Polarinet: An improved grid representation for online lidar point clouds semantic segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 9601–9610, 2020.
- [78] H. Zhao, J. Shi, X. Qi, X. Wang, and J. Jia. Pyramid scene parsing network. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 2881–2890, 2017.
- [79] L. Zhao, S. Xu, L. Liu, D. Ming, and W. Tao. Svaseg: Sparse voxel-based attention for 3d lidar point cloud semantic segmentation. *Remote Sensing*, 14(18):4471, 2022.
- [80] X. Zhu, H. Zhou, T. Wang, F. Hong, W. Li, Y. Ma, H. Li, R. Yang, and D. Lin. Cylindrical and asymmetrical 3d convolution networks for lidar-based perception. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(10):6807–6822, 2021.
- [81] X. Zhu, H. Zhou, T. Wang, F. Hong, Y. Ma, W. Li, H. Li, and D. Lin. Cylindrical and asymmetrical 3d convolution networks for lidar segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 9939–9948, 2021.
- [82] Z. Zhuang, R. Li, K. Jia, Q. Wang, Y. Li, and M. Tan. Perception-aware multi-sensor fusion for 3d lidar semantic segmentation. In *IEEE International Conference on Computer Vision*, pages 16280–16290, 2021.



**Sitao Chen** is a Master student in the School of Software Engineering at South China University of Technology. He received his Bachelor Degree in the School of Mechanical & Automotive in 2023 from South China University of Technology in Guangzhou, China. His research interests include 3D scene understanding for autonomous driving.



**Rong Li** is currently an Associate Researcher at the Hong Kong University of Science and Technology (Guang Zhou). She received her Bachelor's Degree in 2019 and Master's Degree in 2022, both from the School of Software Engineering at South China University of Technology, China. Her research interests include LiDAR perception.



**Kui Jia** is currently a Professor at the School of Electronic and Information Engineering, South China University of Technology, Guangzhou, China. He received his Bachelor Degree in marine engineering from Northwestern Polytechnic University, Xi'an, China, in 2001, the Master Degree in electrical and computer engineering from National University of Singapore in 2003, and the Ph.D. degree in computer science from Queen Mary University of London, U.K., in 2007. His recent research focuses on theoretical deep learning and its applications in vision and robotic problems, including deep learning of 3D data and deep transfer learning.



**Qicheng Wang** is Ph.D. student in the Department of Mathematics at the Hong Kong University of Science and Technology. He received his Bachelor Degree in Electronic Engineering in 2007 from Tsinghua University, Beijing. His research interest lies in neural network and computer vision for autonomous driving.



**Yuanqing Li** is currently a Professor with the School of Automation and Engineering, South China University of Technology, Guangzhou, China. He received the Bachelor Degree in applied mathematics from Wuhan University, Wuhan, China, in 1988, the Master Degree in applied mathematics from South China Normal University, Guangzhou, China, in 1994, and the Ph.D. degree in control theory and applications from the South China University of Technology, Guangzhou, in 1997. His research interests include blind signal processing, sparse representation, machine learning, brain-computer interface, and EEG and functional magnetic resonance imaging data analysis.



**Mingkui Tan** is currently a Professor with the School of Software Engineering, South China University of Technology, Guangzhou, China. He received the Bachelor Degree in Environmental Science and Engineering in 2006 and the Master Degree in Control Science and Engineering in 2009, both from Hunan University in Changsha, China. He received the Ph.D. degree in Computer Science from Nanyang Technological University, Singapore, in 2014. From 2014–2016, he worked as a Senior Research Associate on computer vision in the School of Computer Science, University of Adelaide, Australia. His research interests include machine learning, sparse analysis, deep learning and large-scale optimization.



**Zhuangwei Zhuang** is a Ph.D. student in the School of Software Engineering at South China University of Technology, and is currently working as an intern at RoboSense, Shenzhen, China. He received his Bachelor Degree in Automation and Engineering in 2016 and Master Degree in Software Engineering in 2018, both from South China University of Technology in Guangzhou, China. His research interests include model compression and 3D scene understanding for autonomous driving.