

MotionBEV: Attention-Aware Online LiDAR Moving Object Segmentation with Bird’s Eye View based Appearance and Motion Features

Bo Zhou*, Jiapeng Xie*, Yan Pan, Jiajie Wu, and Chuanzhao Lu

Abstract—Identifying moving objects is an essential capability for autonomous systems, as it provides critical information for pose estimation, navigation, collision avoidance, and static map construction. In this paper, we present MotionBEV, a fast and accurate framework for LiDAR moving object segmentation, which segments moving objects with appearance and motion features in the bird’s eye view (BEV) domain. Our approach converts 3D LiDAR scans into a 2D polar BEV representation to improve computational efficiency. Specifically, we learn appearance features with a simplified PointNet and compute motion features through the height differences of consecutive frames of point clouds projected onto vertical columns in the polar BEV coordinate system. We employ a dual-branch network bridged by the Appearance-Motion Co-attention Module (AMCM) to adaptively fuse the spatio-temporal information from appearance and motion features. Our approach achieves state-of-the-art performance on the SemanticKITTI-MOS benchmark. Furthermore, to demonstrate the practical effectiveness of our method, we provide a LiDAR-MOS dataset recorded by a solid-state LiDAR, which features non-repetitive scanning patterns and a small field of view.

Index Terms—Semantic Scene Understanding, Deep Learning Methods, LiDAR Moving Object Segmentation

I. INTRODUCTION

DYNAMIC objects such as pedestrians, cyclists, and moving vehicles are frequently present in traffic scenes. They can cause errors in localization and mapping [33, 6], and hinder downstream tasks like obstacle avoidance [14] and path planning [22]. In this context, online moving object segmentation (MOS) plays a critical role in enabling autonomous systems to obtain a more accurate perception of the environment and make reliable decisions.

This paper focuses on the task of LiDAR moving object segmentation (LiDAR-MOS), which involves identifying objects in a scene that are currently in motion using LiDAR (Light Detection and Ranging) sensors. This task is challenging due to the sparsity, uncertain distribution, and the noise of LiDAR point cloud, as well as the complexity and diversity of dynamic scenes. To fully understand the dynamics of the environment, it is necessary to exploit 4D spatio-temporal information from sequential LiDAR frames.

This work is partly supported by the National Natural Science Foundation (NNSF) of China under the Grants No. 62073075. (*Bo Zhou and Jiapeng Xie contributed equally to this work.) (Corresponding author: Bo Zhou.)

All the authors are with the School of Automation, Southeast University, Nanjing 210096, P. R. China (emails: zhoubo@seu.edu.cn; xiejiaopeng@seu.edu.cn; yanpan@seu.edu.cn; jiajiawu@seu.edu.cn; luz@seu.edu.cn)

Some map cleaning methods [16, 21] reject moving points with geometry information, but they often rely on global maps and can only run offline. Recent approaches leverage learning-based methods for map-free moving object segmentation. Some point-based methods [23, 29] extract features directly from the sequence of point clouds, but these methods typically require a significant amount of computational resources and are not suitable for real-time applications. Chen *et al.* [5] proposed LMNet, a method that leverages residual range images to capture temporal information and aligns features through direct concatenation. However, range view (RV) based methods may suffer from boundary-blurring issues resulting from back-projection, and the motion cues derived from residual range images may not be optimal for representing temporal information due to their sensitivity to changes in distance. Moreover, the direct feature alignment of multi-modality features can be limited in accuracy due to redundant and invalid information in motion features. Some subsequent studies [28, 17] have attempted to enhance the fusion of spatial-temporal information by introducing attention mechanisms and incorporating semantic information. However, these approaches still heavily rely on the quality of the motion features, resulting in decreased accuracy when motion features are of poor quality.

Motivated by the aforementioned challenges, we propose a Bird’s Eye View (BEV) based moving object segmentation method called MotionBEV. We convert 3D LiDAR point clouds into a 2D BEV representation to improve computational efficiency. The BEV representation maintains a consistent object size regardless of distance and has better spatial consistency compared to the RV representation. Such a view can provide a rich spatial context that is easy to interpret and process.

The contributions of this paper are summarized as follows:

- We propose a bird’s eye view-based method that exploits high-quality spatio-temporal information for LiDAR-MOS from appearance and motion features. Specifically, the method learns appearance features for each grid cell in the polar BEV images using a simplified PointNet [26], while extracting motion features through the height differences of vertical columns. Temporal information captured from such BEV-based motion features is robust to distance changes.
- We design a dual-branch network bridged by the Appearance-Motion Co-attention Module (AMCM) to adaptively fuse appearance and motion features. To avoid

excessive reliance on motion features, the AMCM dynamically assigns importance weights to appearance and motion features to balance their contributions. Furthermore, the AMCM enhances appearance features using motion features in an attention mechanism, ensuring that the fusion of appearance and motion features is effective and mutually reinforcing.

- The proposed method outperforms existing baselines on the SemanticKITTI-MOS benchmark [5, 2], with 69.7% IoU for moving class and an average inference time of 23ms (on an RTX 3090 GPU). In addition, we evaluate our method on a dataset collected by a solid-state LiDAR and demonstrate its practical effectiveness on LiDARs with non-repetitive scanning patterns and small field of view.

Our code will be publicly available ¹.

II. RELATED WORKS

While significant progress has been made in the field of image and video-based moving object segmentation [32, 19, 13, 35], challenges still remain for LiDAR-MOS task due to the sparse and uneven distribution of point clouds. This chapter focuses on LiDAR-based moving object segmentation, which can be broadly categorized into two main approaches: geometry-based methods and learning-based methods.

A. Geometry-based Methods

Geometry-based methods do not require training data but are typically offline and rely on pre-built maps. Lim *et al.* [21] removes dynamic objects by checking the pseudo occupancy ratio of each sector in the LiDAR scan and recovers the ground plane through region growing. Schauer and Nüchter [27] propose a ray casting-based method that estimates the occupancy probability in the grid space to distinguish static and dynamic points. To improve computational efficiency, the visibility-based method [25] projects the point cloud onto a range image, and clears dynamic points by comparing the visibility difference between frames and maps. Kim *et al.* [16] propose a remove-then-revert mechanism that iteratively retains static points from mistakenly removed points using a multi-resolution range image. Fan *et al.* [11] fuse ray casting-based and visibility-based methods to achieve high-precision dynamic object segmentation. Chen *et al.* [4] propose an automatic annotation method, using clustering-based object segmentation and tracking-based association to provide training labels for learning-based MOS. In general, these methods are ideal for static map construction but are not suitable for online MOS.

B. Learning-based Methods

Learning-based methods learn effective information directly from data with trainable deep neural networks, enabling the detection of dynamic objects without pre-built maps.

LiDAR semantic segmentation [34, 7, 15, 36] is a closely related task to LiDAR-MOS. However, most semantic segmentation methods can only distinguish between movable

and immovable objects, such as vehicles and buildings, but cannot differentiate between moving and non-moving objects, such as moving cars and parked cars. Scene-flow methods [8, 20, 9] estimate 3D scene flow between consecutive point clouds, enabling the identification of moving points. However, most scene-flow methods consider only two subsequent frames, resulting in limited accuracy for slowly moving objects. Some point-based methods [23, 29, 18, 30] extract temporal information directly from sequences of point clouds to improve accuracy and generalization. However, due to the high dimensionality of point clouds, these methods often suffer from computational inefficiencies and limited real-time performance.

Recently, Chen *et al.* [5] release a moving object segmentation dataset and benchmark based on SemanticKITTI [2]. At the same time, they propose LMNet, which extracts temporal features from the residual range images for online MOS. To better exploit spatio-temporal information, Sun *et al.* [28] presents a dual-branch structure to separately process spatial and temporal information and fuse them with motion-guided attention modules. Kim *et al.* [17] design a network that incorporates semantic information to further improve the MOS performance. These methods all exploit spatio-temporal information of range images, which can suffer from boundary-blurring issues during the back-projection process and be sensitive to changes in distance.

Compare to the range view projection, bird’s eye view (BEV) projection provides a more intuitive representation of object motion and spatial relationships in the scene. Mohapatra *et al.* [24] introduced a BEV-based method that exploits the disparity between two successive frames with a residual computation layer. Although this method runs very fast, it exhibits relatively low accuracy compared to RV-based methods. To improve the MOS performance, our approach leverages grid height differences between BEV images to extract effective temporal information and uses a lightweight PointNet [26] to learn appearance features for each vertical grid. The spatial and temporal features are deeply fused with the appearance-motion co-attention module. To handle complex scene changes, we add co-attention gates [32] before the motion-guided attention modules [19] to suppress the interference from redundant and invalid information. Benefitting from the BEV representation, our method achieves state-of-the-art performance without any post-processing, while also offering faster computation compared to range view-based methods.

III. PROPOSED APPROACH

In this section, we introduce our bird’s eye view-based framework for moving object segmentation, whose overall pipeline is illustrated in Fig. 1. First, in Section III-A, we describe the projection of 3D LiDAR point clouds to the polar bird’s eye view image. Then, we extract temporal information by computing the height differences of the projected local maps in Section III-B. The structure and modules of our network are presented in Section III-C, and the details of our network training are described in Section III-D.

¹<https://github.com/xiekkki/motionbev>

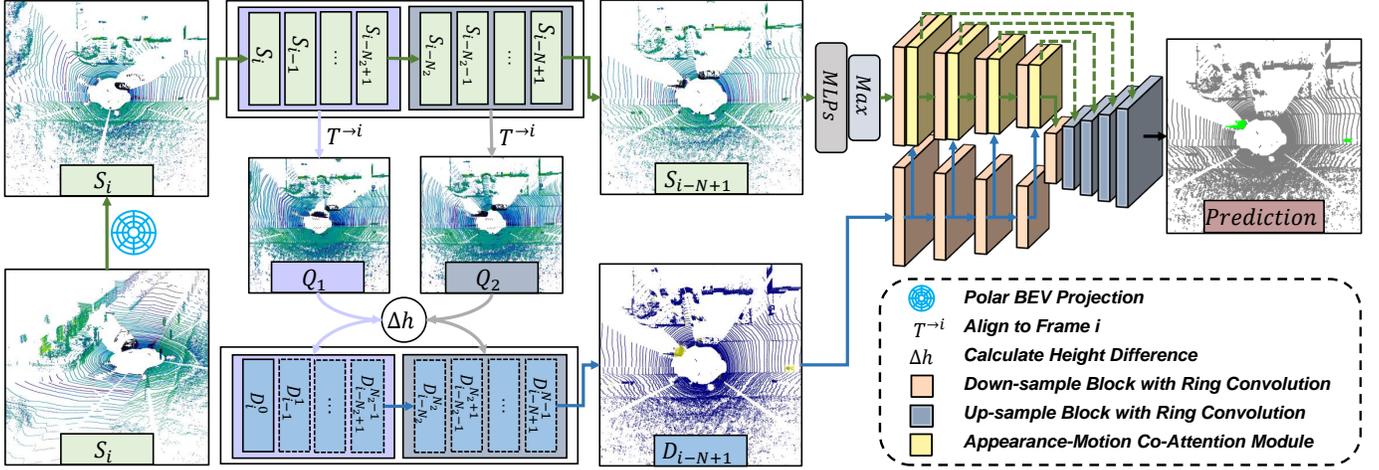


Fig. 1. Overview of MotionBEV. For each LiDAR scan, we first partition the points into grids using the polar BEV coordinates. We push the scan into the temporal window and transform all scans of the window to the current viewpoint. After N iterations of movement, we obtain the motion features with N channels. We then pop the scan out and employ PointNet to learn the appearance features. The appearance and motion features are fed into a dual-branch network bridged by AMCM to explore spatio-temporal information. Finally, the CNN outputs a quantized prediction and we decode it to the point domain.

A. Input Representation

Given the LiDAR scan of the i^{th} frame in the point cloud sequence, denoted as $S_i = \{p_j \in \mathbb{R}^4\}_{j=0}^{M-1}$, $p_j = [x_j, y_j, z_j, 1]^T$, which contains M points represented in homogeneous coordinates, and the past $N-1$ consecutive frames $S_{i-1}, S_{i-2}, \dots, S_{i-N+1}$, the goal of MOS is to determine which points in S_i are actually moving. The current frame S_i provides appearance features for prediction, while the consecutive N frames $S_i, S_{i-1}, \dots, S_{i-N+1}$ provide motion features. First, we align the past frames to the viewpoint of S_i using the relative pose transformations $T \in \mathbb{R}^{4 \times 4}$ estimated by LiDAR odometry as:

$$S_{i-n}^{i} = \{T_{i-n}^i p_j \mid p_j \in S_{i-n}\}, \text{ where} \quad (1)$$

$$T_{i-n}^i = \prod_{k=1}^n T_{i-n}^{i-k+1}, n \in \{1, 2, \dots, N-1\}$$

To achieve online MOS, we project the 3D LiDAR points onto 2D image coordinates. As suggested in [34], we adopt the polar bird's eye view for point cloud coordinate partitioning to balance the uneven distribution of LiDAR points in space. For each point $p_j = (x_j, y_j, z_j) \in S_i$ represented in Cartesian coordinates, we use the following equations to convert it into a representation in polar coordinates:

$$\begin{pmatrix} \rho_j \\ \theta_j \\ z_j \end{pmatrix} = \begin{pmatrix} \sqrt{x_j^2 + y_j^2} \\ \arctan(y_j, x_j) \\ z_j \end{pmatrix} \quad (2)$$

Each point p_j is assigned to the corresponding grid of the polar BEV image according to its polar coordinates:

$$S_{(u,v),i} = \{p_j \mid p_j \in S_i, \frac{\rho_{max} - \rho_{min}}{w} \cdot (u-1) \leq \rho_j < \frac{\rho_{max} - \rho_{min}}{w} \cdot u, \frac{\theta_{max} - \theta_{min}}{h} \cdot (v-1) \leq \theta_j < \frac{\theta_{max} - \theta_{min}}{h} \cdot v\} \quad (3)$$

where $S_{(u,v),i}$ denotes all points of S_i contained in the $(u, v)^{\text{th}}$ grid, (h, w) are the height and width of the polar BEV image, $(\theta_{max}, \theta_{min})$ are the maximum and minimum limits of the

angle, and (ρ_{max}, ρ_{min}) are the maximum and minimum values of the distance. In the experiment section, we will study the impact of the size of the BEV image on MOS performance. The limits of angle and distance can be specified based on the characteristics of the LiDAR sensor and environment.

B. Motion Features Generation

The generation of motion features aims to extract temporal information from consecutive LiDAR frames to provide clues for moving object segmentation. LMNet [5] achieves this by calculating residual range images, but we found that these images are not efficient in representing temporal information due to their sensitivity to changes in distance. To address this issue, we propose a method that involves calculating the height difference between bird's eye view images to generate motion features.

To mitigate the effects of sparse point clouds, we refrain from directly comparing differences between frames. Instead, we maintain two adjacent temporal windows, Q_1 and Q_2 , with equal lengths of $N_2 = N/2$, and generate motion features by examining the height differences of corresponding grid cells in Q_1 and Q_2 . We first compensate for ego-motion by relative pose transformation and align Q_1 and Q_2 to the current local coordinate system. The polar BEV images are calculated according to the partition results of Eq. (2) and Eq. (3), where each pixel value $I_{(u,v),i}$ represents the height occupied by the $(u, v)^{\text{th}}$ grid in Q_i :

$$I_{(u,v),i} = \text{Max}\{Z_{(u,v),i}\} - \text{Min}\{Z_{(u,v),i}\}, \text{ where} \quad (4)$$

$$Z_{(u,v),i} = \{z_j \in p_j \mid p_j \in Q_{(u,v),i}, z_{min} < z_j < z_{max}\}$$

We obtain the region of interest by limiting the range of z to the interval (z_{min}, z_{max}) , which is the specific location where moving objects like vehicles and pedestrians usually appear. For the SemanticKITTI dataset, we take $(z_{min}, z_{max}) = (-4, 2)$. We subtract the projected BEV images I_1 and I_2 to obtain the residuals:

$$D_{(u,v),i}^0, D_{(u,v),i-1}^1, \dots, D_{(u,v),i-N_2+1}^{N_2-1} = I_{(u,v),1} - I_{(u,v),2}, \quad (5)$$

$$D_{(u,v),i-N_2}^{N_2}, D_{(u,v),i-N_2-1}^{N_2+1}, \dots, D_{(u,v),i-N+1}^{N-1} = I_{(u,v),2} - I_{(u,v),1}$$

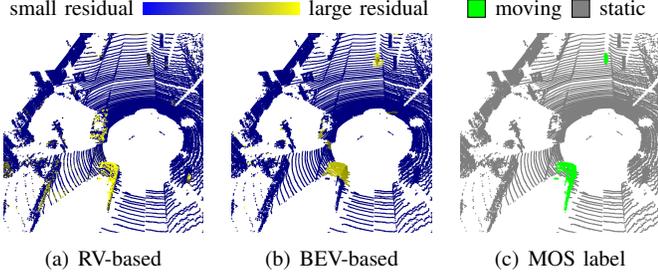


Fig. 2. Visualization and comparison of range view-based motion features (RV-based) and bird’s eye view-based motion features (BEV-based).

where $D_{(u,v),i}^k$ represents the motion feature in the $(u, v)^{th}$ grid of the k^{th} channel for the i^{th} frame. The value of k is determined based on the frame’s position in the temporal window. We shift both temporal windows to the next position as each new frame arrives. Each frame stays within the temporal window for N cycles, thus ensuring that the motion features retain a sequential characteristic with a length of N . To reduce the interference of noise and occlusion, we discard regions with residual values that are too large ($D_{(u,v),i} > 4$) or too small ($D_{(u,v),i} < 0.4$), as well as regions that might be blocked by nearby objects and therefore have too few points (less than 5).

It is worth noting that generating complete motion features requires the scan to pass through a temporal window of length N , which introduces a fixed delay of $(N - 1) \times T$, where T is the measurement period of the sensor. To offer a delay-free solution, one can opt to utilize only the motion features from the first channel, as D_i^0 can be computed immediately when S_i enters the temporal window. Nevertheless, adopting this approach may lead to a trade-off in segmentation accuracy.

For comparison, we apply back-projection to both the range view-based motion features and our bird’s eye view-based motion features to 3D space, as shown in Fig. 2. Our method attributes more salient motion features to moving cars and bicycles, whereas the range view-based method tends to assign larger residuals to the nearby parked car compared to the bicycle driving far away, making it difficult to distinguish whether distant objects are stationary or moving. Despite the distinctiveness of our motion features for moving object segmentation, noise and occlusion could result in erroneous residuals, mainly for points hidden behind objects or accidentally appearing. Therefore, it is not recommended to use motion features as MOS results directly. Nevertheless, most of these erroneous residuals can be distinguished by the subsequent CNN network with the aid of appearance features as they tend to appear on immovable objects.

C. Network Structure

Our network structure is built upon PolarNet [34], a polar bird’s eye view-based method which achieves end-to-end LiDAR semantic segmentation through an encoder-decoder architecture. In order to explore spatio-temporal information for MOS, we modify PolarNet into a dual-branch structure and adaptively fuse appearance and motion features with the appearance-motion co-attention module.

1) *Appearance Features Encoding*: In contrast to the manual extraction of appearance features, we adopt a simplified

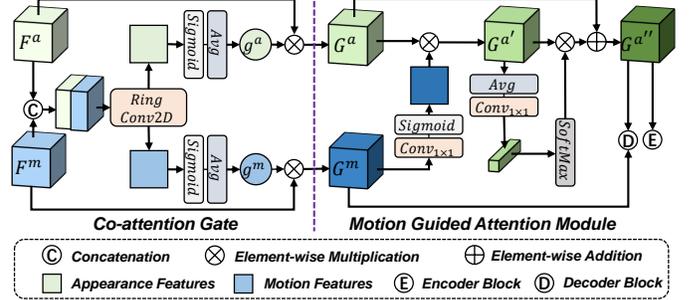


Fig. 3. Structure of the Appearance-Motion Co-attention Module (AMCM).

PointNet [26] followed by max-pooling to learn the point cloud distribution across the vertical column of a grid. For the grid cell of (u, v) coordinate in the polar BEV image, the appearance features can be represented as follows:

$$F_{(u,v),i}^a = \text{MaxPool}\{\text{MLPs}(p_j) \mid p_j \in S_{(u,v),i}\} \quad (6)$$

Since the representation is learned in the polar coordinate system, the left and right ends of the bird’s eye view image should be connected in physical space. We use the same ring convolution kernel as [34] to extract complete spatial information and replace all the conventional convolution kernels in the network.

2) *Appearance-Motion Co-attention Module*: We utilize the Appearance-Motion Co-attention Module (AMCM) to facilitate interaction between motion-appearance features from two different branches. Inspired by visual segmentation methods, our AMCM consists of two parts: the co-attention gate [32] and the motion guided attention module [19], as shown in Fig. 3.

First, the co-attention gate is used to adaptively compute the importance of appearance features and motion features, in order to balance their contributions and suppress the redundant and misleading information. Given the appearance features F_i^a and motion features F_i^m from the i^{th} layer, we capture the relative relationship between multi-modal features using cross-channel concatenation and ring convolution operations, resulting in a fused feature $H \in \mathbb{R}^{h \times w \times 2}$. Next, we split feature H into two sub-branches, and apply the sigmoid function and global average pooling on each channel to obtain a pair of co-attention scores g_i^a and g_i^m , reflecting the importance of each modality feature. We combine the appearance features and motion features to construct the gating function:

$$g_i = \text{Avg}(\text{Sigmoid}(\text{Conv}(\text{Cat}(F_i^a, F_i^m)))) \quad (7)$$

where g_i contains the pair of co-attention scores g_i^a and g_i^m . Higher co-attention scores indicate that the corresponding modality feature contains more effective information for accurate segmentation. We apply the co-attention scores to the features, generating gated appearance features G_i^a and gated motion features G_i^m , where

$$G_i^a = F_i^a \otimes g_i^a, G_i^m = F_i^m \otimes g_i^m \quad (8)$$

The motion guided attention module follows the co-attention gate, and utilizes motion features to emphasize important positions or elements in the appearance features. Specifically, we first fuse the gated motion features G_i^m and the gated

appearance features G_i^a by spatial attention, generating motion salient feature $G_i^{a'}$:

$$G_i^{a'} = G_i^a \otimes \text{Sigmoid}(\text{Conv}_{1 \times 1}(G_i^m)) \quad (9)$$

Then, we use channel attention to enhance the response of key attributes, generating the final spatio-temporal fused feature $G_i^{a''}$ of size $C \times h \times w$, with the following formula:

$$G_i^{a''} = G_i^{a'} \otimes [\text{Softmax}(\text{Conv}_{1 \times 1}(\text{Avg}(G_i^{a'}))) \cdot C] + G_i^a \quad (10)$$

By combining the co-attention gate and motion guided attention module, the AMCM is able to perform adaptive multi-modality features fusion.

D. Implementation Details

We use PyTorch to build our network and train it on an NVIDIA RTX 3090 GPU with a batch size of 8. The input sequences for training are randomly shuffled before each epoch. Widely-used data augmentation techniques such as random flipping, rotation, and small-scale translation are applied to enhance the training data. The loss function of the network is a linear combination of weighted cross-entropy (\mathcal{L}_{wce}) and Lovász-Softmax (\mathcal{L}_{ls}) [3] losses:

$$\mathcal{L} = \mathcal{L}_{wce} + \mathcal{L}_{ls}, \quad (11)$$

where the weighted cross-entropy loss is defined as:

$$\mathcal{L}_{wce}(y, \hat{y}) = - \sum \alpha_i p(y_i) \log(p(\hat{y}_i)), \alpha_i = 1/\sqrt{f_i}, \quad (12)$$

and the Lovász-Softmax loss is given by:

$$\mathcal{L}_{ls} = \frac{1}{|C|} \sum_{c \in C} \Delta_{J_c}^-(m(c)), m_i(c) = \begin{cases} 1 - x_i(c) & \text{if } c = y_i(c) \\ x_i(c) & \text{otherwise} \end{cases} \quad (13)$$

In the equations, y_i and \hat{y}_i represent the true and predicted labels, and f_i is the frequency of the i^{th} class. $|C|$ is the number of classes, $\Delta_{J_c}^-$ represents the Lovász extension of the Jaccard index. Additionally, $x_i(c) \in [0, 1]$ and $y_i(c) \in \{1, 1\}$ represent the predicted probability and ground truth label of pixel i for class c , respectively.

We use stochastic gradient descent (SGD) to minimize \mathcal{L}_{wce} and \mathcal{L}_{ls} , with a momentum of 0.9 and weight decay of 0.0001. The initial learning rate is set to 0.005, and it is decayed by a factor of 0.99 after each epoch. No pre-trained weights are used, and the network is trained from scratch until the validation loss converges.

IV. EXPERIMENTS

This section presents the experimental evaluation of our method for moving object segmentation. First, we introduce the datasets and evaluation metrics in Section IV-A. In Section IV-B, we provide quantitative and qualitative comparisons between our method and other state-of-the-art methods. In Section IV-C, we conduct an ablation study to demonstrate the effectiveness of our BEV-based motion features generation method, as well as the contributions of different components in the network to the overall performance. In Section IV-D, we show the practical effect of our method on a solid-state LiDAR. Lastly, in Section IV-E, we showcase the runtime performance of our method.

A. Experiment Setups

Datasets-SemanticKITTI. We evaluate our method by comparing its performance with state-of-the-art methods on the SemanticKITTI-MOS [5] dataset. The original SemanticKITTI dataset [2, 12] comprises 22 labeled point cloud sequences collected by a single Velodyne HDL-64E LiDAR. SemanticKITTI-MOS maps all semantic classes to two categories: moving and static. The dataset is split into train sequences 00-07, 09-10 (19,130 frames), validation sequence 08 (4,071 frames), and test sequences 11-21 (20,351 frames). The odometry estimation is obtained from Suma [1].

Datasets-SipailouCampus. To further assess the effectiveness of our method across different LiDAR sensors, we collected a dataset with a solid-state LiDAR Livox Avia mounted on an unmanned vehicle, at different areas of Southeast University Sipailou Campus. The Livox Avia has a narrower field of view (70.4° horizontally and 77.2° vertically) and is equipped with non-repetitive scanning mode, which presents new challenges for MOS. We manually annotate all sequences with moving labels, resulting in a dataset of 26,279 frames, of which 15,585 frames contain dynamic objects, and each frame has 24,000 points. We use 5 sequences (16,887 frames) for training, 1 sequence (3,191 frames) for validation, and 2 sequences (6,201 frames) for testing. The odometry estimation is obtained from FAST-LIO [31]. We will release this dataset to facilitate further research.

Evaluation Metrics. Following the previous work [5], we use the Jaccard Index or Intersection-over-Union (IoU) metric [10] over moving objects to quantify MOS performance:

$$\text{IoU} = \frac{\text{TP}}{\text{TP} + \text{FP} + \text{FN}} \quad (14)$$

where TP, FP, and FN represent the number of true positives, false positives, and false negatives in the prediction of the moving class, respectively.

B. Evaluation and Comparisons on SemanticKITTI

We evaluate our method by submitting moving object segmentation predictions on the SemanticKITTI-MOS benchmark. We consider two alternatives for our method: one utilizes only D_i^0 (referred to as "without delay"), while the other incorporates complete motion features with a temporal window size of 8. Our method is compared against various baseline methods, including a) range view-based methods: LMNet [5], MotionSeg3D v1 (with kNN) [28], MotionSeg3D v2 (with the point refine module), and RVMOS [17]; b) an offline method: AutoMOS [4]; c) point-based methods: 4DMOS (with the binary Bayes filter) [23] and InsMOS [30]; d) and another BEV-based method: LiMoSeg [24]. To ensure a comprehensive comparison, we also provide all the evaluation results on the validation set (sequence 08). The experimental outcomes for each method are sourced from the SemanticKITTI-MOS benchmark and their respective original papers. Notably, all reported results are derived from training on the original SemanticKITTI dataset, and any results based on additional data [28, 4, 30] are omitted to maintain fairness.

As depicted in Table I, RVMOS demonstrates superior performance on the hidden test set, which can be primarily

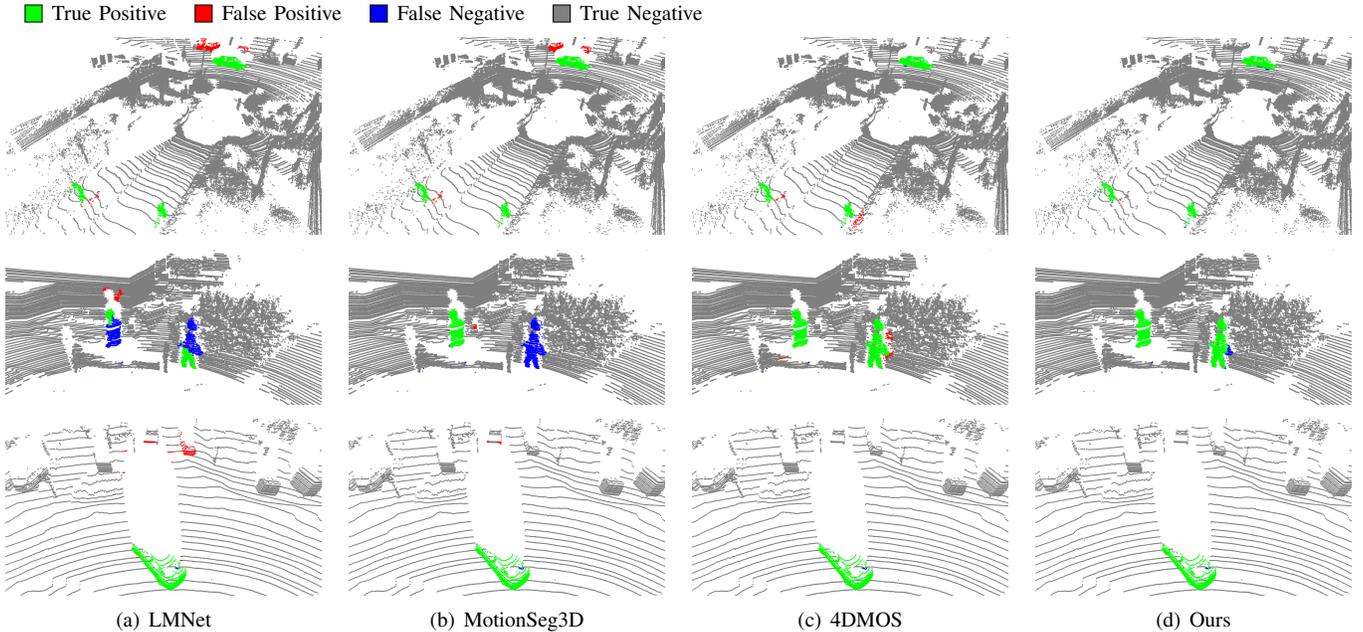


Fig. 4. Qualitative results of different methods for LiDAR-MOS on SemanticKITTI validation set. Best viewed in color.

attributed to its utilization of both semantic labels and moving object labels during training. Our method achieves the best result among approaches that solely relied on moving object labels, with an IoU of 69.7%.

TABLE I
EVALUATION AND COMPARISON ON THE SEMANTICKITTI VALIDATION SET AND THE SEMANTICKITTI-MOS BENCHMARK.

Methods	IoU (Validation)	IoU (Test)
LMNet	66.4	58.3
MotionSeg3D, v1	68.1	62.5
MotionSeg3D, v2	71.4	64.9
AutoMOS	-	54.3
4DMOS	77.2	65.2
LiMoSeg	52.6	-
Ours, without delay	68.1	63.9
Ours	76.5	69.7
LMNet*	67.1	62.5
RVMOS*	71.2	74.7
InsMOS*	73.2	70.6

* indicates the method exploiting semantic labels.

* indicates the method exploiting instance labels.

- indicates that the result is not available yet.

Best results in bold.

Fig. 4 shows the qualitative comparison of LMNet, MotionSeg3D, 4DMOS, and our method on the SemanticKITTI validation set. Range view-based methods like LMNet are prone to boundary-blurring issues caused by back-projection, and MotionSeg3D partially improved this issue through the point refine module. However, due to the sensitivity of residual range image to distance, MotionSeg3D sometimes makes mistakes in classifying distant objects. 4DMOS performs well in segmenting dynamic objects at various distances but exhibits less distinct boundaries between objects. In contrast, our method effectively fuses appearance and motion features based on BEV representation and shows superior performance.

C. Ablation Study

In this section, we perform several ablation experiments to assess the contribution of our motion features and network components to the MOS performance. All experiments are conducted on the SemanticKITTI validation set (sequence 08).

As shown in Table II, we use vanilla PolarNet as the baseline, and vertically compare three different ways of fusing appearance and motion features: a) directly concatenating appearance and motion features (DC) following LMNet, b) exploring feature interactions with dual-branch structure bridged by motion guided attention module (MGA), and c) incorporating co-attention gate (CAG) on top of b) to learn appearance-motion co-attention. In addition, we also horizontally evaluate the performance of these network setups using different motion features as input for cross-comparison: range view-based motion features (F_{RV}^m), and our bird's eye view-based motion features (F_{BEV}^m). Both motion features are generated with consecutive 8 frames.

In general, we observe improvements in IoU for all setups using the proposed BEV-based motion features. When employing a dual-branch structure with motion guided attention module for deep multi-modality interaction, the setup with BEV-based motion features gets a stronger performance boost compared to the one with RV-based motion features. This indicates that our BEV-based motion features provide better temporal information compared to LMNet's range residual features. By combining the co-attention gate and motion-guided attention module, our method achieves the best performance and obtains a 16.6 percentage points improvement compared to the baseline. Overall, using BEV-based motion features and fusing appearance and motion features with the appearance-motion co-attention module (CAG + MGA) results in the best performance.

Another ablation study is about the settings of motion features generation, as presented in Fig. 5. Firstly, we compare the MOS performance of using different BEV image sizes

TABLE II
ABLATION STUDY OF NETWORK COMPONENTS AND MOTION FEATURES
ON THE SEMANTICKITTI VALIDATION SET.

Baseline and components	F_{RV}^m		F_{BEV}^m	
	IoU [%]	Δ	IoU [%]	Δ
PolarNet (without F^m)	59.99	-	59.99	-
PolarNet + DC	66.31	+6.32	70.78	+10.79
PolarNet + MGA	67.21	+7.22	75.74	+15.75
PolarNet + CAG + MGA	69.40	+9.41	76.54	+16.55

Δ means the improvement compared to the baseline.

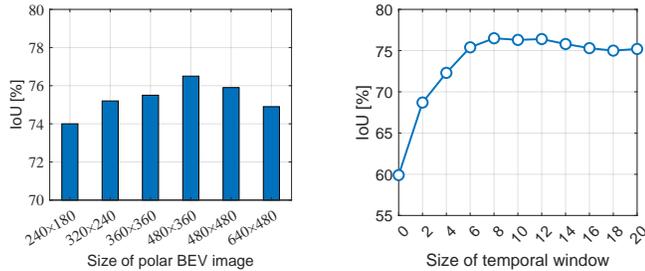


Fig. 5. Ablation studies. The left figure shows the ablation study on the MOS performance vs. the size of the BEV image. The right figure shows the ablation study on the MOS performance vs. the temporal window size N .

for motion features generation. All experiments are under the same settings except for the different representation sizes. The best-performing model employs polar BEV motion features with a representation size of 480×360 . Excessive enlargement of BEV images leads to the presence of numerous empty or sparsely populated grid cells, posing challenges in computing height differences and generating meaningful motion features. Furthermore, larger BEV image sizes render the motion features generation overly sensitive to abrupt or small objects in the environment, resulting in the emergence of unnecessary residuals.

We also examine the influence of the temporal window size N on MOS performance. Since two sub-maps of the same size need to be constructed, N must be set to an even number. It is noteworthy that the network achieves a high IOU even when N is 0, indicating the absence of motion features. One possible reason is that the BEV approach effectively captures the spatial relationships between objects in the scene, enabling the network to infer the motion status of objects based on their positions on the road. We observe that the MOS performance improves as N increases, reaching its peak in the range of 8 to 12. However, the further increase of N leads to a slight decline in accuracy. This could be attributed to the temporal window becoming excessively long, causing the first sub-map Q_1 to be too distant from the current frame, which may result in erroneous motion features. This issue is particularly evident in scenes with significant field-of-view changes, such as during vehicle turning.

D. Evaluation on SipailouCampus

In this section, we demonstrate the performance of our method on a solid-state LiDAR with a small field of view and non-repetitive scanning mode. First, in order to test the generalization ability of the network, we exclude the intensity channel of the point cloud, train all the models on the training set of SemanticKITTI, and evaluate them

TABLE III
EVALUATION AND COMPARISON ON SIPAILOUCAMPUS DATASET.

Methods	Retrain	IoU (Validation)	IoU (Test)
LMNet		5.37	6.88
MotionSeg3D		6.83	6.72
4DMOS		78.54	82.30
Ours		50.44	52.02
Ours-h		70.94	71.51
LMNet	✓	54.27	56.16
MotionSeg3D	✓	65.64	66.84
4DMOS	✓	87.30	88.89
Ours	✓	89.22	90.80

Best results in bold.

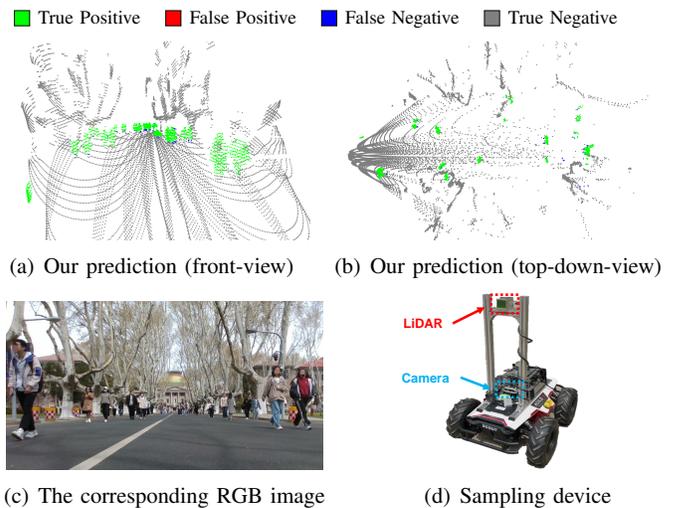


Fig. 6. Visualization of MOS result of a high dynamic scene on Sipailou-Campus validation set.

on the validation and test set of SipailouCampus dataset. The results are shown in Table III. Due to the significant differences in the range image projection between Velodyne LiDAR and Livox LiDAR, LMNet and MotionSeg3D (with the point refine module) fail to obtain reliable predictions. Our method does not perform satisfactorily without any parameter modification, as the BEV projection causes information loss in the vertical direction, making it difficult for the network to effectively segment objects within the same vertical grid (such as trees, pedestrians under trees, and the ground). By changing the overall height of the point cloud to match the ground position near the sensor with the SemanticKITTI dataset, we obtained a significant performance improvement without any additional training (corresponding to "Ours-h" in Table III). In comparison, the method 4DMOS (with the binary Bayes filter), which processes directly on the point cloud, demonstrates strong generalization ability.

Furthermore, we retrain all models on the training set of SipailouCampus to obtain the final performance of all methods on Livox LiDAR data. To effectively capture and process LiDAR data within a different field of view, we make adjustments to the image size of projection-based methods. Specifically, the range image size of LMNet and MotionSeg3D is set to 512×512 , and the BEV image size of our method is set to 480×64 . Due to the irregular scanning pattern of Livox LiDAR, range view-based methods (LMNet, MotionSeg3D)

still cannot achieve high accuracy. After retraining, our method outperforms 4DMOS, and reaches 89.3% IoU on the validation set and 90.8% IoU on the test set. Fig. 6 presents an example of the prediction results of our network.

E. Runtime

We assess the runtime of different approaches by measuring the average inference time (ms) on the SemanticKITTI validation set, as shown in Table IV. The evaluations for LMNet, MotionSeg3D, 4DMOS, and our method are conducted on a machine with an Intel Core i9-12900K CPU and a single NVIDIA RTX 3090 GPU. The results for RVMOS and InsMOS are obtained from their respective papers and are also evaluated using the RTX 3090 GPU. Our method achieves an inference time of only 23ms. However, incorporating the full set of motion features from all channels within the temporal window would lead to an extra fixed delay of $(N - 1) \times 100$ ms. In practical usage, one can balance accuracy and runtime by controlling the length of the temporal window or deciding whether to use the complete motion features or not.

TABLE IV
COMPARISON OF INFERENCE TIME (MS) ON SEMANTICKITTI
VALIDATION SET.

LMNet	MotionSeg3D	RVMOS	4DMOS	InsMOS	Ours
35	110	29	132	127	23

V. CONCLUSIONS

In this paper, we present a simple yet effective method for LiDAR-based online moving object segmentation. Our method exploits the spatio-temporal information from appearance and motion features in the bird’s eye view domain and performs multi-modality feature fusion with a dual-branch network bridged by the appearance-motion co-attention module. Evaluation on the SemanticKITTI-MOS dataset demonstrates the proposed MotionBEV is the most accurate method among approaches utilizing only MOS labels, while also offering a balanced compromise between accuracy and latency. Moreover, experimental results on a dataset collected by a Livox LiDAR show our method’s practical effectiveness on different types of LiDAR sensors.

REFERENCES

- [1] Jens Behley and Cyrill Stachniss. “Efficient Surfel-Based SLAM using 3D Laser Range Data in Urban Environments.” In: *Robot.: Sci. Syst.* Vol. 2018. 2018, p. 59.
- [2] Jens Behley et al. “Semantickitti: A dataset for semantic scene understanding of lidar sequences”. In: *Proc. IEEE Int. Conf. Comput. Vis.* 2019, pp. 9297–9307.
- [3] Maxim Berman, Amal Rannen Triki, and Matthew B Blaschko. “The lovasz-softmax loss: A tractable surrogate for the optimization of the intersection-over-union measure in neural networks”. In: *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.* 2018, pp. 4413–4421.
- [4] Xieyuanli Chen et al. “Automatic labeling to generate training data for online LiDAR-based moving object segmentation”. In: *IEEE Robot. Automat. Lett.* 7.3 (2022), pp. 6107–6114.
- [5] Xieyuanli Chen et al. “Moving object segmentation in 3D LiDAR data: A learning-based approach exploiting sequential data”. In: *IEEE Robot. Automat. Lett.* 6.4 (2021), pp. 6529–6536.
- [6] Xieyuanli Chen et al. “Suma++: Efficient lidar-based semantic slam”. In: *Proc. IEEE/RSJ Int. Conf. Intell. Robots Syst.* IEEE. 2019, pp. 4530–4537.
- [7] Tiago Cortinhal, George Tzelepis, and Eren Erdal Aksoy. “Salsanext: Fast, uncertainty-aware semantic segmentation of lidar point clouds”. In: *Proc. IEEE Veh. Symp. (IV)*. Springer. 2020, pp. 207–222.
- [8] Fangqiang Ding et al. “Self-Supervised Scene Flow Estimation With 4-D Automotive Radar”. In: *IEEE Robot. Automat. Lett.* 7.3 (2022), pp. 8233–8240.

- [9] Guanting Dong et al. “Exploiting rigidity constraints for lidar scene flow estimation”. In: *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.* 2022, pp. 12776–12785.
- [10] Mark Everingham et al. “The pascal visual object classes (voc) challenge”. In: *Int. J. Comput. Vision* 88 (2010), pp. 303–338.
- [11] Tingxiang Fan et al. “DynamicFilter: an Online Dynamic Objects Removal Framework for Highly Dynamic Environments”. In: *Proc. IEEE Int. Conf. Robot. Automat.* IEEE. 2022, pp. 7988–7994.
- [12] A. Geiger, P. Lenz, and R. Urtasun. “Are we ready for Autonomous Driving? The KITTI Vision Benchmark Suite”. In: *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.* 2012, pp. 3354–3361.
- [13] Jhony H Giraldo, Sajid Javed, and Thierry Bouwmans. “Graph moving object segmentation”. In: *IEEE Trans. Pattern Anal. Mach. Intell.* 44.5 (2020), pp. 2485–2503.
- [14] Binghua Guo, Nan Guo, and Zhisong Cen. “Obstacle avoidance with dynamic avoidance risk region for mobile robots in dynamic environments”. In: *IEEE Robot. Automat. Lett.* 7.3 (2022), pp. 5850–5857.
- [15] Qingyong Hu et al. “Randla-net: Efficient semantic segmentation of large-scale point clouds”. In: *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.* 2020, pp. 11108–11117.
- [16] Giseop Kim and Ayoung Kim. “Remove, then revert: Static point cloud map construction using multiresolution range images”. In: *Proc. IEEE/RSJ Int. Conf. Intell. Robots Syst.* IEEE. 2020, pp. 10758–10765.
- [17] Jaeyul Kim, Jungwan Woo, and Sunghoon Im. “RVMOS: Range-View Moving Object Segmentation Leveraged by Semantic and Motion Features”. In: *IEEE Robot. Automat. Lett.* 7.3 (2022), pp. 8044–8051.
- [18] Thomas Kreutz, Max Mühlhäuser, and Alejandro Sanchez Guinea. “Unsupervised 4D LiDAR Moving Object Segmentation in Stationary Settings with Multivariate Occupancy Time Series”. In: *Proc. IEEE Winter Conf. Appl. Comput. Vis.* 2023, pp. 1644–1653.
- [19] Haofeng Li et al. “Motion guided attention for video salient object detection”. In: *Proc. IEEE Int. Conf. Comput. Vis.* 2019, pp. 7274–7283.
- [20] Ruibo Li et al. “Rigidflow: Self-supervised scene flow learning on point clouds by local rigidity prior”. In: *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.* 2022, pp. 16959–16968.
- [21] Hyungtae Lim, Sungwon Hwang, and Hyun Myung. “ERASOR: Egocentric ratio of pseudo occupancy-based dynamic object removal for static 3D point cloud map building”. In: *IEEE Robot. Automat. Lett.* 6.2 (2021), pp. 2272–2279.
- [22] Yuanfu Luo et al. “Porca: Modeling and planning for autonomous driving among many pedestrians”. In: *IEEE Robot. Automat. Lett.* 3.4 (2018), pp. 3418–3425.
- [23] Benedikt Mersch et al. “Receding moving object segmentation in 3d lidar data using sparse 4d convolutions”. In: *IEEE Robot. Automat. Lett.* 7.3 (2022), pp. 7503–7510.
- [24] Sambit Mohapatra et al. “Limoseg: Real-time bird’s eye view based lidar motion segmentation”. In: *arXiv:2111.04875* (2021).
- [25] François Pomerleau et al. “Long-term 3D map maintenance in dynamic environments”. In: *Proc. IEEE Int. Conf. Robot. Automat.* IEEE. 2014, pp. 3712–3719.
- [26] Charles R Qi et al. “Pointnet: Deep learning on point sets for 3d classification and segmentation”. In: *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.* 2017, pp. 652–660.
- [27] Johannes Schauer and Andreas Nüchter. “The peopleremover-removing dynamic objects from 3-d point cloud data by traversing a voxel occupancy grid”. In: *IEEE Robot. Automat. Lett.* 3.3 (2018), pp. 1679–1686.
- [28] Jiadai Sun et al. “Efficient Spatial-Temporal Information Fusion for LiDAR-Based 3D Moving Object Segmentation”. In: *Proc. IEEE/RSJ Int. Conf. Intell. Robots Syst.* IEEE. 2022, pp. 11456–11463.
- [29] Yuxiang Sun et al. “PointMoSeg: Sparse tensor-based end-to-end moving-obstacle segmentation in 3-D lidar point clouds for autonomous driving”. In: *IEEE Robot. Automat. Lett.* 6.2 (2020), pp. 510–517.
- [30] Neng Wang et al. “InsMOS: Instance-Aware Moving Object Segmentation in LiDAR Data”. In: *arXiv preprint arXiv:2303.03909* (2023).
- [31] Wei Xu et al. “Fast-lid2: Fast direct lidar-inertial odometry”. In: *IEEE Trans. Robot.* 38.4 (2022), pp. 2053–2073.
- [32] Shu Yang et al. “Learning motion-appearance co-attention for zero-shot video object segmentation”. In: *Proc. IEEE Int. Conf. Comput. Vis.* 2021, pp. 1564–1573.
- [33] Ji Zhang and Sanjiv Singh. “LOAM: Lidar odometry and mapping in real-time.” In: *Robot.: Sci. Syst.* Vol. 2. 9. Berkeley, CA. 2014, pp. 1–9.
- [34] Yang Zhang et al. “Polarnet: An improved grid representation for online lidar point clouds semantic segmentation”. In: *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.* 2020, pp. 9601–9610.
- [35] Wangbo Zhao et al. “Modeling Motion with Multi-Modal Features for Text-Based Video Segmentation”. In: *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.* 2022, pp. 11737–11746.
- [36] Xinge Zhu et al. “Cylindrical and asymmetrical 3d convolution networks for lidar segmentation”. In: *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.* 2021, pp. 9939–9948.