

MF-MOS: A Motion-Focused Model for Moving Object Segmentation

Jintao Cheng, Kang Zeng, Zhuoxu Huang, Xiaoyu Tang, Jin Wu, Chengxi Zhang, Xieyuanli Chen, Rui Fan

Abstract—Moving object segmentation (MOS) provides a reliable solution for detecting traffic participants and thus is of great interest in the autonomous driving field. Dynamic capture is always critical in the MOS problem. Previous methods capture motion features from the range images directly. Differently, we argue that the residual maps provide greater potential for motion information, while range images contain rich semantic guidance. Based on this intuition, we propose MF-MOS, a novel motion-focused model with a dual-branch structure for LiDAR moving object segmentation. Novelty, we decouple the spatial-temporal information by capturing the motion from residual maps and generating semantic features from range images, which are used as movable object guidance for the motion branch. Our straightforward yet distinctive solution can make the most use of both range images and residual maps, thus greatly improving the performance of the LiDAR-based MOS task. Remarkably, our MF-MOS achieved a leading IoU of 76.7% on the MOS leaderboard of the SemanticKITTI dataset upon submission, demonstrating the current state-of-the-art performance. The implementation of our MF-MOS has been released at <https://github.com/SCNU-RISLAB/MF-MOS>.

I. INTRODUCTION

A key challenge for safe autonomous driving systems is the precise perception of moving objects *e.g.* pedestrians and other vehicles that share the traffic environment [1]. The LiDAR-based MOS task tackles the uncertainty perception in the traffic environment by segmenting the current moving objects such as pedestrians and cyclists while distinguishing a stationary vehicle from one that is moving [2]–[4]. Therefore, it helps develop such uncertainty perception and is essential in autonomous driving. This paper follows the mainstream

This research was supported by the National Natural Science Foundation of China under Grants 62001173 and 62233013, the Project of Special Funds for the Cultivation of Guangdong College Students’ Scientific and Technological Innovation (“Climbing Program” Special Funds) under Grants pdjh2022a0131 and pdjh2023b0141, and the Science and Technology Commission of Shanghai Municipal under Grant 22511104500 (*Corresponding author: Xiaoyu Tang*).

Jintao Cheng, Kang Zeng, Xiaoyu Tang are with the School of Electronic and Information Engineering, South China Normal University, Foshan 528225, China. tangxy@scnu.edu.cn

Zhuoxu Huang is with the Department of Computer Science, Aberystwyth University, Aberystwyth SY23 3DB, U.K. zh6@aber.ac.uk

Jin Wu is with the Department of Electronic and Computer Engineering, Hong Kong University of Science and Technology, Hong Kong, China. jin.wu.uestc@hotmail.com

Chengxi Zhang is with the School of Internet of Things Engineering, Jiangnan University, Wuxi, China. dongfangxy@163.com

Xieyuanli Chen is with the College of Intelligence Science and Technology, National University of Defense Technology, Changsha, China. xieyuanli.chen@nudt.edu.cn

Rui Fan is with the College of Electronics & Information Engineering, Shanghai Research Institute for Intelligent Autonomous Systems, the State Key Laboratory of Intelligent Autonomous Systems, and Frontiers Science Center for Intelligent Autonomous Systems, Tongji University, Shanghai 201804, China. rui.fan@ieee.org

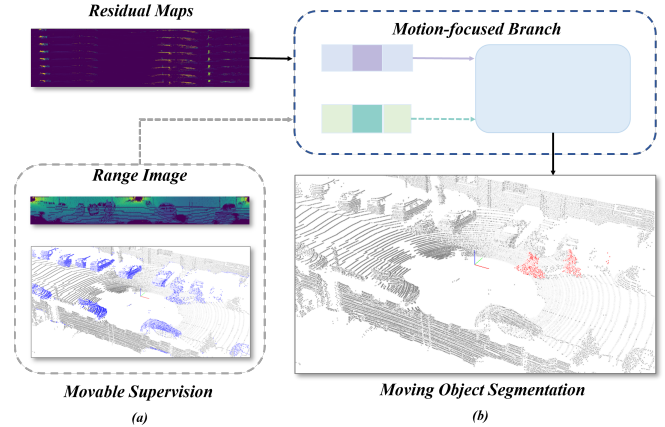


Fig. 1. Core idea of the proposed motion-focused model. The blue parts in (a) represent the point cloud of movable objects and the red parts in (b) represent the point cloud of moving objects. The moving objects are usually a subset of movable objects. Our MF-MOS emphasizes motion information (via residual maps) and utilizes movable features (via the range image) to provide semantic enhancement.

setting of MOS that segments the currently moving objects in a point cloud frame using range projection and residual maps from LiDAR data. Previous methods mainly tackle the dynamic capture with the range view images from point cloud scenes [1], [2], [5], while some of them [2], [6] apply residual maps as auxiliary guidance during dynamic capture.

For instance, MotionSeg3D [2] utilizes a dual-branch framework based on SalsaNext [7]. It simultaneously encodes spatial-temporal information from range images and incorporates a residual branch to enhance motion features. However, these approaches usually prioritize the semantic information of object appearance and detect whether an object can move while relegating the actual motion state of objects to the status of an auxiliary feature.

Drawing from the most intuitive observations that dynamic capture forms the foundational component in addressing the MOS problem, we put forth MF-MOS, a dual-branch structure for the LiDAR moving object segmentation task. The core idea of the MF-MOS is to focus on the dynamic information from the residual maps as a fundamental component of the network (see Fig. 1). Specifically, we design a primary motion branch to capture the dynamic from residual maps. Additionally, a semantic branch is used to integrate the semantic information of object appearance from range images into the motion branch. We have meticulously designed a kind of pooling layer which is more suitable for the two branches. We name it Strip Average Pooling Layer (SAPL).

Our method shares a similar envision with RVMOS [6] that segments movable objects from the range images. However, it primarily emphasizes the motion potential of objects rather than directly addressing their moving status, which is the core aspect of the MOS problem. Different from what, our MF-MOS goes a step further with the direct capture of the dynamic information, thus performing remarkably well in completing the task. More differently, we develop a distribution-based data augmentation to address the influence of different frame sampling of the residual maps to build a robust network. Furthermore, we introduce a 3D Spatial-Guided Information Enhancement Module (SIEM) that provides additional spatial guidance to both the primary motion branch and the semantic branch, thereby alleviating the potential loss of information.

Extensive experiments have demonstrated the superiority of our design. Leveraging the exceptional dynamic perceptual ability of the proposed MF-MOS, we substantially improve the performance of the MOS task on the SemanticKITTI dataset [1] and achieve the top spot on the leaderboard. In summary, our contributions can be summarized as follows:

- We target the direct capture of the dynamic information in the MOS task and propose a motion-focused network with a dual-branch structure named MF-MOS: (i) a primary motion branch to capture the motion feature from residual maps; (ii) a semantic branch to compute the semantic information of object appearance from range images.
- We propose a novel distribution-based data augmentation method that improves the network robustness. We also propose SIEM to refine both the motion branch and alleviate the loss of information.
- The proposed method attains the highest ranking on the SemanticKITTI-MOS benchmark for both the test and validation datasets. We also tested our method in different benchmarks to validate its robustness and superior performance.

II. RELATED WORK

A. MOS Based on Occupancy Map and Visibility

Previous methods address the MOS task by applying occupancy maps or adopting visibility-based methods. They both own the unique advantage of data-free learning. Firstly, inspired by Octomap [8], occupancy grids are often utilized in MOS tasks such as moving obstacles removal. These methods compute the motion information by comparing the occupancy maps between continuous frames and locating the dynamic points within the occupancy grid. For instance, J. Schauer et al. [9] proposed to remove pedestrians based on the differences in volumetric occupancy between different temporal scans. Likewise, S. Pagad et al. [10] targeted to remove dynamic objects on wide urban roads with 3D occupancy maps. Besides, H. Lim et al. [11] further proposed a pseudo occupancy map based on the height threshold that is robust to motion ambiguity. Secondly, the other type of

method adopts the visibility-based theory and applies visual projection data, *e.g.* range images for the MOS task. G. Kim et al. [3] proposed a motion points removal and reverting method based on multi-resolution range images. P. Chi et al. [12] proposed a static points construction algorithm via LiDAR and images. However, both the occupancy maps and the visibility-based methods rely on the previously obtained maps and pose information from Simultaneous Localization and Mapping (SLAM) systems, thus being limited in real-time applications.

B. MOS Based on Deep Learning

Recent approaches tend to apply popular deep learning models and capture spatial-temporal features from data directly. These methods usually adopt different view projections such as range-view projection, voxelization, and bird’s eye view projection. For instance, X. Chen et al. [1] presented a novel LiDAR moving segmentation method based on range-view images. Compared to previous approaches, this method prevents static points from being mistakenly removed by capturing canny features. After that, Motionseg3D [2] was proposed as an improved version. It proposed a dual-branch network with a refined module to optimize the MOS results. RVMOS [6] also illustrated a multi-branch segmentation framework to fuse semantic and motion information and further improve the MOS performance. Other methods [4], [13]–[15] adopt different view projections to address the LiDAR-MOS task. B. Mersch et al. [4] applied the 4D voxelization on LiDAR point clouds for efficiency. Similarly, [15] adopted the bird’s eye view projection on LiDAR point clouds and proposed a real-time network for the MOS task. We also use the range-view projection in the proposed MF-MOS. Distinct from previous networks, we design a motion-focused network that mainly captures the motion feature from the residual maps and generates semantic features from range images.

III. METHODOLOGY

We present our MF-MOS in detail in the following sections. Firstly, we start with the basic data projection from the LiDAR inputs to the range view and residual inputs. Then, we elaborate on the proposed MF-MOS and the SIEM. Finally, we describe our distribution-based data augmentation for the MOS task in detail.

A. Data Preprocessing

Range images serve as a lightweight 2D representation of point cloud data. We project the LiDAR point cloud into the range image and residual map following the stander setting of previous work [1], [2], [6]. After getting the range images of different frames of point clouds, we obtain past k -frame residual maps \mathbf{I}_{res} by the pixel level variance calculation between frames as follows:

$$\mathbf{I}_{res}^k(\mathbf{u}, \mathbf{v}) = \left| \frac{\mathbf{I}_{RV}^k(\mathbf{u}, \mathbf{v}) - \mathbf{I}_{RV}^0(\mathbf{u}, \mathbf{v})}{\mathbf{I}_{RV}^0(\mathbf{u}, \mathbf{v})} \right|, \quad (1)$$

where \mathbf{I}_{RV} represents the range image, u and v are the transformed pixel coordinates in the 2D image space.

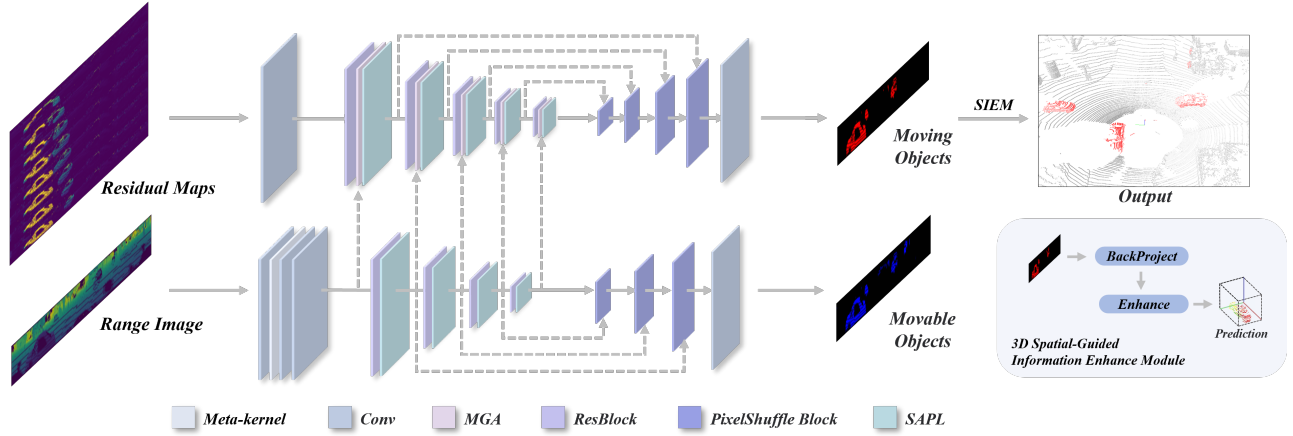


Fig. 2. The overall of MF-MOS is a dual-input-dual-output branching structure. The **semantic branch** (the bottom one) which takes the range image as input is used to predict movable objects in the current frame, and the **motion branch** (the upper one) takes the residual maps as input to predict the moving objects. The intermediate feature maps obtained from the encoder of the semantic branch are fused into the motion branch through the MGA module. To obtain further refined segmentation results, we use the output of the the motion branch as the input of the SIEM to obtain the final point cloud segmentation results.

B. Network Structure

1) Dual-Branch Motion-focused Framework with SAPL:

We design a novel dual-branch network that focuses on residual maps. Fig. 2 illustrates the overall architecture of the proposed method. Inspired by [1], [2], SalsaNext has demonstrated powerful performance in the MOS task, hence we employ it as the backbone in motion (via residual maps) and semantic (via the range image) branches.

The proposed semantic branch utilizes the range image to effectively extract features of movable objects and adds meta kernel [16] for feature-level enhancement. In contrast, the feature of the motion branch emphasizes the dynamism of the current features. To enhance the fusion of features from two distinct inputs, we adopt a fusion strategy [17] in the motion branch. The encoded feature outputs from the residual maps are combined with those from each layer of the range image, serving as inputs to the subsequent layers of the coding module in the motion branch. This fusion process facilitates the integration of complementary information from both inputs and promotes the effective utilization of features for subsequent processing in the motion branch. The dual-branch framework can be illustrated as:

$$\mathbf{F}_s = \text{sigmoid}(\text{Conv}_{1 \times 1}(\mathbf{F}_{\text{semantic}})) \otimes \mathbf{F}_{\text{motion}}, \quad (2)$$

where $\mathbf{F}_{\text{motion}}$ represents the feature map of the motion branch in the input fusion module, while $\mathbf{F}_{\text{semantic}}$ represents the feature map of the semantic branch in the input fusion module. Followed by (2), we can contact two branch features. \mathbf{F}_s presents the fusion result after sigmoid processing.

$$\mathbf{F}_f = \text{softmax}(\text{Conv}_{1 \times 1}(\text{Pool}(\mathbf{F}_s))) \times C, \quad (3)$$

Pool denotes the adaptive average pooling layer, where C represents the number of channels in the feature map. \mathbf{F}_f is the output after normalization softmax process.

$$\mathbf{F}_o = \mathbf{F}_c + \mathbf{F}_{\text{res}}, \quad (4)$$

where \mathbf{F}_o denotes the connection of \mathbf{F}_c and \mathbf{F}_{res} . We can note that SalsaNext [7] includes a 2×2 pooling layer with a standard square pooling kernel during the feature encoding phase. It is obvious that range and residual maps are generally not aligned in terms of height and width, and the use of a square pooling kernel for feature extraction can easily lead to partial feature loss. Therefore, in the proposed two-branch input network, we modify the pooling layer down-sampled by the SalsaNext encoder with the PixelShuffle layer up-sampled by the decoder, which is called the Strip Average Pooling Layer (SAPL), expressed as follows:

$$\mathbf{F}'_{x_0, y_0} = \frac{1}{h \times w} \sum_{i=0}^w \sum_{j=0}^h \mathbf{F}_{w \times x_0 + i, h \times y_0 + j}, \quad (5)$$

where \mathbf{F} represents the input feature map for stripe pooling, while \mathbf{F}' represents the output feature map after stripe pooling. h and w denote the height and width of the stripe pooling kernel, respectively. x_0 and y_0 represent the coordinates of the output feature map, corresponding to the x and y dimensions.

In the encoder, we modify the original 2×2 pooling kernel of SalsaNext [7] by replacing it with a pooling kernel size of 2×4 . Additionally, in the PixelShuffle operation, we adjust the ratio of channel-to-image aspect conversion to match the modified rectangular pooling operation.

2) 3D Spatial-Guiding Information Enhancement Module:

To compensate for the information loss caused by the data dimension reduction during the conversion from point clouds to range images, we propose the SIEM to refine the segmentation results of the dual-branch network. SIEM transforms the last layer feature map of the first stage decoder through a back-projection process into the point cloud space, resulting in the initial feature point cloud. After voxelization, the initial point cloud is fed into our proposed 3D Spatial-Guided Block (SGB) for further processing. As shown in Fig. 3, the input of the SGB module first goes through three different 3D spatial convolution processes to decompose the features into

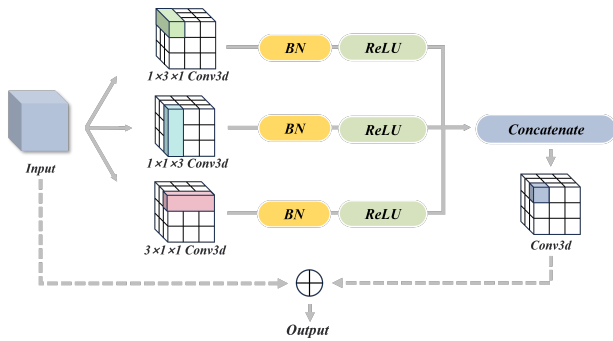


Fig. 3. Enhancing 3D Spatial Information with the SGB. The SGB partitions and enriches features across dimensions before fusion, aiming to distill insights from sparse point clouds.

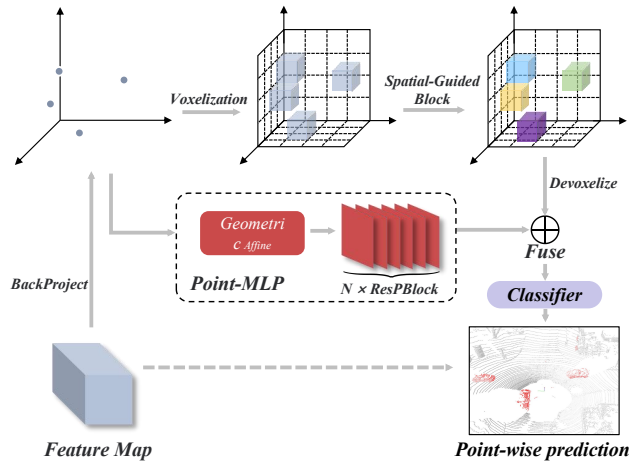


Fig. 4. Illustration of the SIEM. The process involves voxelization of the initial feature map, followed by SGB and Devoxelization. The resulting output is fused with the Point-MLP output and classified.

multiple dimensions and enhance the information in each dimension. And then, the outputs from last step are fused by concatenation and 3D convolution in order to capture effective information from the sparse point cloud to a great extent. The final output of the SGB is obtained by skip connection with the original module’s input to avoid gradient dispersion.

Finally, we apply a de-voxelization operation to the processed point cloud and fuse it with the point cloud data after MLP feature extraction. Afterward, a classifier is employed to output the refined per-point segmentation results. Overall, The SIEM can be shown in Fig. 4.

C. Data Augmentation

We propose a streamlined yet effective distribution-based data augmentation to improve our MF-MOS. We maintain the idea of motion focus and enhance the learning process in the temporal domain. As shown in Fig. 5, different residual maps with different frame strides usually represent different ranges of temporal information. To better enhance the motion features, we propose to generate the residual maps using multiple frame strides instead of a fixed stride. Given a frame stride $\Delta t \in [1, 2, 3]$, the correspondence residual maps is represented as $I_{\Delta t}$. To prevent data redundancy

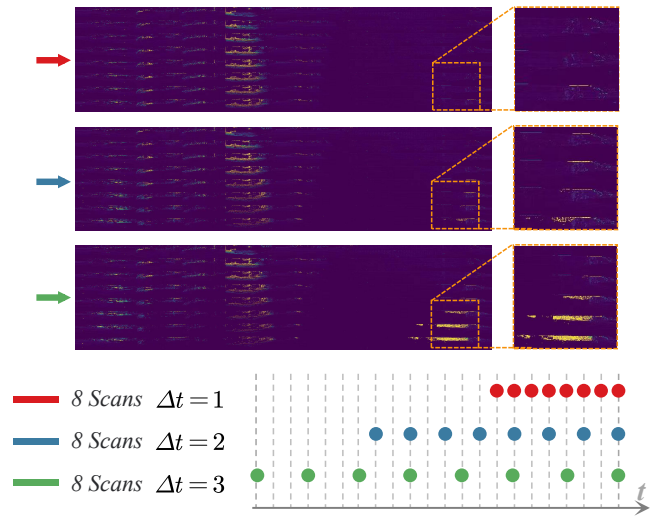


Fig. 5. K -frames residual maps using different frame stride Δt . The red-boxed region shows residual feature responses correspondence to the different moving speeds of objects. A larger Δt corresponds to slower-moving objects. Here we show results from eight-frame residual maps.

during training, we augment the data based on the given distributions of the frame strides. In other words, instead of feeding all $I_{\Delta t}$ into the network, we choose one Δt based on a designed distribution probability in every training iteration. We evaluate different distributions and different ranges of the Δt and report ablation results in Sec. IV-C.

D. Loss Function

During the training process, the total loss function of the proposed algorithm includes the motion-branch losses and the range-branch losses. The sum of the total loss $\mathcal{L}_{\text{Total}}$ can be followed as:

$$\mathcal{L}_{\text{Total}} = \mathcal{L}_{\text{Semantic}} + \mathcal{L}_{\text{Motion}}, \quad (6)$$

where $\mathcal{L}_{\text{Semantic}}$ represents the losses of the semantic branch with the range image, and $\mathcal{L}_{\text{Motion}}$ is the losses of motion branch with residual maps.

Both of semantic and motion branch used the weighted cross-entropy \mathcal{L}_{wce} and Lov’asz-Softmax losses \mathcal{L}_{ls} . The loss function for each individual branch is as follows:

$$\mathcal{L} = \mathcal{L}_{\text{wce}} + \mathcal{L}_{\text{ls}}, \quad (7)$$

IV. EXPERIMENTS

We conduct extensive experiments to comprehensively evaluate our MF-MOS. In the following sections, we first illustrate the experimental setup for the MOS task, then report the basic validation/test results on the two widely used MOS datasets SemanticKITTI-MOS [1] and Apollo [18] to demonstrate the generalization ability of our approach. Following these results, we design our ablation studies fastidiously on the SemanticKITTI-MOS dataset to evaluate the rationality of our MF-MOS.

A. Experiment Setups

We utilize the SemanticKITTI-MOS dataset [1] as the main training and evaluating benchmark in our experiment.

TABLE I
COMPARISONS RESULT ON SEMANTICKITTI-MOS DATASET.

Methods	Publication	Validation (%)	Test (%)
SpSequenceNet [19]	CVPR 2020	-	43.2
KPConv [20]	ICCV 2019	-	60.9
Cylinder3D [21]	CVPR 2021	66.3	61.2
LMNet [1]	ICRA 2021	63.8	60.5
4DMOS [4]	RAL 2022	71.9	65.2
MotionSeg [2]	IROS 2022	71.4	70.2
RVMOS [6]	RAL 2022	71.2	74.7
InsMOS [14]	IROS 2023	73.2	75.6
MF-MOS(Ours)	-	76.1	76.7

TABLE II
COMPARISONS RESULT ON APOLLO DATASET.

Methods	IoU (%)
MotionSeg3D (<i>cross-val</i>) [2]	7.5
LMNet (<i>cross-val</i>) [1]	16.9
LMNet (<i>fine-tune</i>) [1]	65.9
MF-MOS (<i>cross-val</i>)	49.9
MF-MOS (<i>fine-tune</i>)	70.7

SemanticKITTI-MOS is the most popular and authoritative dataset for the MOS task with richly labeled moving objects. We use the standard data splits for training, validating, and testing following previous works [1], [2], [6]. As illustrated in Sec. III-B, we use the semantic labels for movable objects during training to provide additional supervision for the primary motion branch. Additionally, we also perform validation experiments on the Apollo dataset [18], accompanied by some quantitative analysis following the standard experiment setting in [5].

Our code is implemented in PyTorch. The experiments are conducted on 2 NVIDIA Tesla A100 GPUs. We train our MF-MOS for 150 epochs with an initial learning rate of 0.01 and a decay factor of 0.99 in every epoch. The batch size is set to 8 for each GPU. We use the SGD optimizer with a momentum of 0.9 during training. Following the standard evaluation in MOS, we adopt the intersection over union (IoU) [22] of the moving objects to quantify our results in all experiments.

B. Comparison with SoTA Methods

We first report the validation and test results on the SemanticKITTI-MOS [1] dataset in Tab. I. We achieve state-of-the-art performance on both the validation set and the test set. Remarkably, we improve the validation IoU by 2.9% compared to the last SoTA [14]. For the test set measurement, we upload our moving object segmentation results to the benchmark server and report the IoU from the leaderboard. Results show that our proposal maintains consistent superiority on the test set. Our performance stays on top and suppresses all the other methods with an IoU of 76.7%.

We also report the validation result on the Apollo dataset [18] in Tab. II. Following the standard setting of previous approaches [5], [14], we adopt protocols including transfer

TABLE III
ABLATION EXPERIMENTS WITH PROPOSED MODULES.

Methods	Component			IoU (%)
	Dual-Branch	SIEM	Data Aug	
LMNet [1]				63.82
MF-MOS (<i>i</i>)	-	-	-	64.96
MF-MOS (<i>ii</i>)	✓	-	-	71.44
	-	✓	-	69.59
	-	-	✓	67.34
MF-MOS (<i>iii</i>)	✓	✓	-	74.13
	✓	-	✓	73.12
	-	✓	✓	70.47
	✓	✓	✓	76.12

TABLE IV
THE PROPOSED MODULES PERFORMANCE (%) ON OTHER METHODS.

Method	Baseline	w/ Aug	w/ SIEM
LMNet [1]	63.82	+1.15	+0.46
MotionSeg3D [2]	68.07	+0.46	+3.30

learning and end-to-end fine-tuning for validation experiments on Apollo. The *cross-val* setting refers to cross-validation and the *fine-tune* setting refers to end-to-end fine-tuning. Both pre-train weights are obtained from the SemanticKITTI-MOS training. Our MF-MOS shows significant improvements in both settings.

C. Ablation Studies

We conduct ablation experiments on the proposed MF-MOS and its different components. The results are shown in Tab. III. Firstly, without any refinement module, our motion-focused framework proves to be superior. With only the motion branch capturing motion information from the residual maps (setting *i*), MF-MOS exhibits a significant improvement (+1.16% IoU) compared to LMNet [1], which uses range images as the main inputs. Additionally, each of the proposed components consistently enhances the performance of our baseline to varying degrees (setting *ii*), with the dual-branch structure alone providing the most significant improvement. To further demonstrate the indispensability of each of our components, we design ablation experiments on different combinations of them in setting *iii*. The last row demonstrates that our full MF-MOS achieved the best performance.

To assess the effectiveness of the proposed SIEM and the distribution-based data augmentation at a deeper level, we apply them to other baseline models, including LMNet [1] and Motionseg3D [2]. The results present in Tab. IV highlight their versatility and efficacy in improving the performance of both algorithms. Remarkably, our SIEM module significantly improves the performance of MotionSeg3D, achieving a +3.3% increase in IoU.

As illustrated in Sec. III-C, we evaluate different distributions and various ranges of Δt in our distribution-based data augmentation. We select Δt based on the designed distribution probability in each training iteration and used

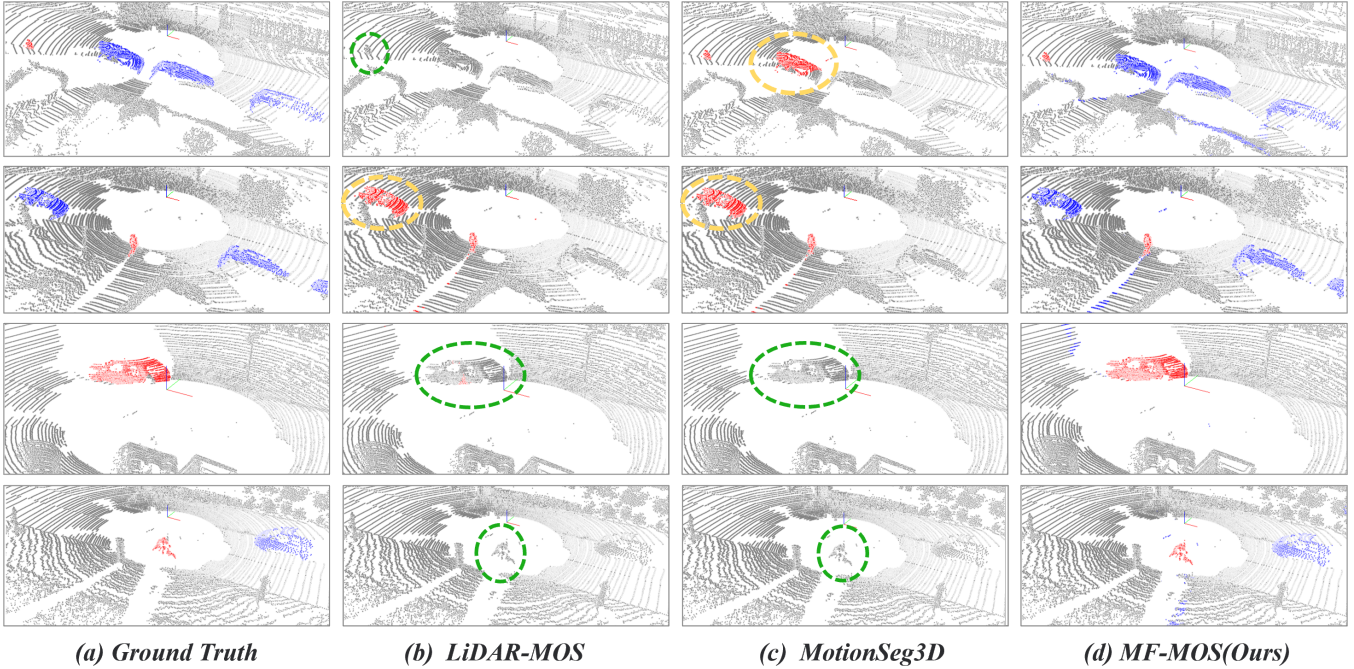


Fig. 6. The trinary model visualization comparison exhibits discernible distinctions, where the blue points correspond to movable objects while the red points correspond to moving objects. Those with green circles denote false negatives, and those with yellow circles indicate false positives.

TABLE V

ABLATION EXPERIMENTS OF DATA AUGMENTATION. THE SUM OF DISTRIBUTION PROBABILITY OF Δt IS EQUAL TO 1. $\Delta t=\text{max}$ MEANS USE THE LARGEST FRAME STRIDE IN TESTING.

distribution probability of Δt					IoU (%)	
1	2	3	4	5	$\Delta t=1$	$\Delta t=\text{max}$
0.33	0.33	0.33	-	-	71.41	71.91
0.25	0.25	0.25	0.25	-	71.12	72.13
0.2	0.2	0.2	0.2	0.2	68.69	70.13
0.4	0.3	0.3	-	-	70.28	71.02
0.5	0.25	0.25	-	-	71.62	73.12
0.6	0.2	0.2	-	-	70.52	70.51

$\Delta t = 1$ and $\Delta t = \text{max}$ during testing. The results are presented in Tab. V. Initially, we test different range values with an average distribution of Δt . When using a wider range of Δt , the performance deteriorates in the testing with $\Delta t = 1$, while the performance was consistently better in the testing with $\Delta t = \text{max}$, indicating non-robustness to different inputs. We further examine the effect of different distribution probabilities for Δt within the range $\Delta t \in [1, 2, 3]$, as it exhibits the highest robustness. Gradually increasing the proportion of $\Delta t = 1$ during training, we achieve the best performance with distribution probabilities $[0.5, 0.25, 0.25]$ for $\Delta t = [1, 2, 3]$, respectively.

D. Qualitative Analysis

In order to more intuitively compare our algorithm with other SoTA algorithms, we perform a visual qualitative anal-

TABLE VI

MODEL INFERENCE TIME (MS) RESULTS.

InsMOS	MotionSeg-v1	MF-MOS-v1	MotionSeg-v2	MF-MOS-v2
193.68	45.93	37.48	117.01	96.19

ysis on the SemanticKITTI dataset. As shown in Fig. 6, both LMNet and MotionSeg3D have misjudgments of movable objects and missed determinations of moving objects. Compared with the SoTA algorithms, we can effectively remove the influence of movable objects through a model based on residual maps, and accurately capture moving objects.

E. Runtime

All the comparative experiments are performed on a single V100 GPU for inference. As shown in Tab. VI, in the one-stage comparison, MF-MOS-v1 outperforms other models. In the two-stage comparison, MF-MOS-v2 achieves real-time processing speed and surpasses Motionseg3D-v2 in terms of performance.

V. CONCLUSIONS

This paper presents a dual-branch motion-based LiDAR moving object segmentation framework, a spatial-guided information enhancement module, and a distribution-based data augmentation method. Extensive experimental results demonstrate that 1) the framework MF-MOS in this study achieves the highest accuracy on both the validation and test sets, respectively, and 2) the proposed model demonstrates superior performance and generalization capabilities, making it applicable to other range-based methods.

REFERENCES

- [1] X. Chen, S. Li, B. Mersch, L. Wiesmann, J. Gall, J. Behley, and C. Stachniss, "Moving object segmentation in 3d lidar data: A learning-based approach exploiting sequential data," *IEEE Robotics and Automation Letters*, vol. 6, pp. 6529–6536, Oct 2021.
- [2] J. Sun, Y. Dai, X. Zhang, J. Xu, R. Ai, W. Gu, and X. Chen, "Efficient spatial-temporal information fusion for lidar-based 3d moving object segmentation," in *2022 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 11456–11463, Oct 2022.
- [3] G. Kim and A. Kim, "Remove, then revert: Static point cloud map construction using multiresolution range images," in *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 10758–10765, Oct 2020.
- [4] B. Mersch, X. Chen, I. Vizzo, L. Nunes, J. Behley, and C. Stachniss, "Receding moving object segmentation in 3d lidar data using sparse 4d convolutions," *IEEE Robotics and Automation Letters*, vol. 7, pp. 7503–7510, July 2022.
- [5] X. Chen, B. Mersch, L. Nunes, R. Marcuzzi, I. Vizzo, J. Behley, and C. Stachniss, "Automatic labeling to generate training data for online lidar-based moving object segmentation," *IEEE Robotics and Automation Letters*, vol. 7, pp. 6107–6114, July 2022.
- [6] J. Kim, J. Woo, and S. Im, "Rvmos: Range-view moving object segmentation leveraged by semantic and motion features," *IEEE Robotics and Automation Letters*, vol. 7, pp. 8044–8051, July 2022.
- [7] T. Cortinhal, G. Tzelepis, and E. Erdal Aksoy, "Salsanext: Fast, uncertainty-aware semantic segmentation of lidar point clouds," in *Advances in Visual Computing* (G. Bebis, Z. Yin, E. Kim, J. Bender, K. Subr, B. C. Kwon, J. Zhao, D. Kalkofen, and G. Baci, eds.), (Cham), pp. 207–222, Springer International Publishing, 2020.
- [8] A. Hornung, K. M. Wurm, M. Bennewitz, C. Stachniss, and W. Burgard, "Octomap: An efficient probabilistic 3d mapping framework based on octrees," *Autonomous robots*, vol. 34, pp. 189–206, 2013.
- [9] J. Schauer and A. Nüchter, "The peopleremover—removing dynamic objects from 3-d point cloud data by traversing a voxel occupancy grid," *IEEE Robotics and Automation Letters*, vol. 3, pp. 1679–1686, July 2018.
- [10] S. Pagad, D. Agarwal, S. Narayanan, K. Rangan, H. Kim, and G. Yalla, "Robust method for removing dynamic objects from point clouds," in *2020 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 10765–10771, May 2020.
- [11] H. Lim, S. Hwang, and H. Myung, "Eraser: Egocentric ratio of pseudo occupancy-based dynamic object removal for static 3d point cloud map building," *IEEE Robotics and Automation Letters*, vol. 6, pp. 2272–2279, April 2021.
- [12] P. Chi, H. Liao, Q. Zhang, X. Wu, J. Tian, and Z. Wang, "Online static point cloud map construction based on 3d point clouds and 2d images," *The Visual Computer*, pp. 1–16, 2023.
- [13] Y. Sun, W. Zuo, H. Huang, P. Cai, and M. Liu, "Pointmoseg: Sparse tensor-based end-to-end moving-obstacle segmentation in 3-d lidar point clouds for autonomous driving," *IEEE Robotics and Automation Letters*, vol. 6, pp. 510–517, April 2021.
- [14] N. Wang, C. Shi, R. Guo, H. Lu, Z. Zheng, and X. Chen, "Insmos: Instance-aware moving object segmentation in lidar data," *arXiv preprint arXiv:2303.03909*, 2023.
- [15] B. Zhou, J. Xie, Y. Pan, J. Wu, and C. Lu, "Motionbev: Attention-aware online lidar moving object segmentation with bird's eye view based appearance and motion features," 2023.
- [16] L. Fan, X. Xiong, F. Wang, N. Wang, and Z. Zhang, "Rangedet: In defense of range view for lidar-based 3d object detection," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 2918–2927, October 2021.
- [17] H. Li, G. Chen, G. Li, and Y. Yu, "Motion guided attention for video salient object detection," in *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 7274–7283, 2019.
- [18] W. Lu, Y. Zhou, G. Wan, S. Hou, and S. Song, "L3-net: Towards learning based lidar localization for autonomous driving," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 6389–6398, 2019.
- [19] H. Shi, G. Lin, H. Wang, T.-Y. Hung, and Z. Wang, "Spsequencenet: Semantic segmentation network on 4d point clouds," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 4574–4583, 2020.
- [20] H. Thomas, C. R. Qi, J.-E. Deschaud, B. Marcotegui, F. Goulette, and L. J. Guibas, "Kpconv: Flexible and deformable convolution for point clouds," in *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 6411–6420, 2019.
- [21] 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 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 9939–9948, June 2021.
- [22] M. Everingham, L. Van Gool, C. K. Williams, J. Winn, and A. Zisserman, "The pascal visual object classes (voc) challenge," *International journal of computer vision*, vol. 88, pp. 303–338, 2010.