

Optimized Feature Points and Keyframe Methods for VSLAM in High-Dynamic Indoor Environments

Zhuhua Hu, *Senior Member, IEEE*, Wenlu Qi^{id}, Kunkun Ding, Hao Qi, Yaochi Zhao^{id},
Xuebo Zhang^{id}, *Senior Member, IEEE*, and Mingfeng Wang^{id}, *Member, IEEE*

Abstract—VSLAM is one of the key technologies for indoor mobile robots, used to perceive the surrounding environment, achieve accurate positioning and mapping. However, traditional VSLAM algorithms based on the assumption of a static environment still face certain challenges. The movement, occlusion, and appearance changes of dynamic objects can lead to feature point-matching errors, making data association difficult and causing biases in motion estimation. In order to address this challenge, this paper proposes a dynamic feature point removal method and a closed-loop detection method for high dynamic scenes, aiming to effectively improve the robustness and positioning accuracy in dynamic environments. First, the YOLOv7-tiny object detection network and LK optical flow algorithm are combined to detect the dynamic area, and the adaptive threshold keyframe selection method is adopted to solve the problem of poor quality of keyframe caused by the existing heuristic threshold selection method. Then, this paper proposes a dynamic keyframe sequence creation method based on the angle difference between keyframes, which reduces the workload of loop back detection and accelerates the efficiency of loop back detection in the system. Next, the ParC_NetVLAD image matching algorithm is proposed. In this paper, ConvNeXt-Tiny network is used for feature extraction of images, and ParC-Net network and CBAM attention mechanism are added to the feature extraction network. Finally, NetVLAD is used to cluster the extracted local features to obtain global features that can represent images. Experiments are conducted on public TUM RGB-D datasets and in real-world situations. The proposed algorithm reduces the ATE (Absolute Trajectory Error) by 96.4% and the RPE (Relative Trajectory Error) by 82.8% on average in highly dynamic scenarios. In the Pittsburgh30k dataset, the average accuracy of loop closure detection has been improved by 2.6%.

Index Terms—Visual SLAM, dynamic scene, ORB-SLAM3, YOLOv7-tiny, keyframe sequences, ParC_NetVLAD.

I. INTRODUCTION

NOWADAYS, the technologies of Simultaneous Localization and Mapping (SLAM) have been widely studied because it allows for the instant creation of pose estimation of maps and sensors in unknown environments. Visual SLAM technology plays an important role in this area, for which the camera sensor has the advantages of low cost and rich image information. It has become a research hotspot [1]. However, the assumption that traditional visual SLAM is based on static environments limits its application scenarios. Although the RANSAC (Random Sample Consensus) algorithm can identify the feature points on dynamic objects as outliers and filter them, it is limited to low-dynamic scenes with fewer dynamic elements [2]. In highly dynamic scenes, when dynamic objects occupy a large area of the image, the feature points extracted by traditional visual SLAM (VSLAM) may be distributed on dynamic objects, resulting in a significant decrease in accuracy. The estimated trajectory will no longer be available. To solve this problem, the study of VSLAM in dynamic scenes has attracted widespread attention and has become a frontier topic in current research [3].

II. RELATED WORK

A. Research Motivation

In 1986, SmithSelf and Cheeseman first proposed SLAM technology, which has been developed for more than 30 years. In 2017, ORB-SLAM2 [4] was proposed as a comprehensive SLAM solution that supports multiple types of cameras and includes features such as map fusion, loop closing, and repositioning. In 2021, Campos et al. [5] proposed ORB-SLAM3, which is a visual SLAM algorithm that supports multiple cameras. This algorithm optimized many aspects, such as map initialization, repositioning, loop detection, keyframe selection, map construction, and sensor support, and has excellent performance in terms of running speed, tracking results, and mapping accuracy.

The above SLAM algorithms assume that the environment where the robot is located is stationary. Therefore, they cannot recognize complex and changing real scenes, especially when moving objects appear in the scene, and the accuracy and

Received 6 June 2023; revised 16 July 2024 and 16 October 2024; accepted 12 December 2024. Date of publication 13 January 2025; date of current version 3 March 2025. This work was supported in part by the National Natural Science Foundation of China under Grant 62161010 and Grant 62361024, in part by the Key Research and Development Project of Hainan Province under Grant ZDYF2022GXJS348 and Grant ZDYF2024GXJS021, and in part by Hainan Province Natural Science Foundation under Grant 623RC446. The Associate Editor for this article was S. S. Nedeveschi. (Corresponding author: Yaochi Zhao.)

Zhuhua Hu, Wenlu Qi, Kunkun Ding, Hao Qi, and Yaochi Zhao are with the School of Information and Communication Engineering, School of Cyberspace Security (School of Cryptology), Hainan University, Haikou 570228, China (e-mail: ealger_hu@hainanu.edu.cn; 1152908275@qq.com; 784239944@qq.com; 1282160577@qq.com; zhyc@hainanu.edu.cn).

Xuebo Zhang is with the College of Artificial Intelligence, Institute of Robotics and Automatic Information System, Nankai University, Tianjin 300350, China (e-mail: zhangxuebo@nankai.edu.cn).

Mingfeng Wang is with the Department of Mechanical and Aerospace Engineering, Brunel University London, UB8 3PH London, U.K. (e-mail: mingfeng.wang@brunel.ac.uk).

Digital Object Identifier 10.1109/TITS.2024.3520177

precision of positioning and mapping will be greatly reduced. First, motion segmentation methods based on the assumption of stationary image foreground mainly rely on the hypothesis that most of the features in the image are static and differentiate between static and moving feature points in the image sequence based on the geometric constraints of feature points. Kim et al. [6] constructed a static background based on the initial depth change between images and filtered out moving objects that were inconsistent with the static background in subsequent image sequences, but it was incomplete in filtering out pedestrians. Then, Sun et al. [7] proposed a motion removal method based on RGB-D data as a preprocessing module for filtering moving objects, which is one of the pioneering works in the field of dynamic simultaneous localization and mapping. The GMSK-SLAM system proposed by Wei et al. [8] combines grid-based motion statistical feature matching method and K-means clustering method to detect dynamic regions.

With the significant progress in deep learning for image processing in recent years, motion segmentation methods based on prior semantic information have been applied in dynamic SLAM systems. First, DS-SLAM [9] combines the SegNet [10] semantic segmentation network with conventional SLAM to remove the influence of individual movement in the environment. However, this framework significantly increases the computational pressure of the system by performing semantic segmentation for each frame of the image and cannot achieve real-time performance. Then, the DynaSLAM [11] algorithm combines the MaskRCNN [12] semantic segmentation network with multi-view geometry and uses the consistency of depth information to remove dynamic feature points caused by moving objects. However, this system does not optimize segmentation efficiency. Shao et al. [13] proposed a convolutional network with an improved Faster R-CNN to act as a semantic filter, which filters out patches with low-level semantic labels. The remaining high-level semantic information is then used for exact matching. To address the poor real-time performance of semantic methods in dynamic scenes, Singhet al. [14] proposed an approach that utilizes only keyframes to extract semantics, thereby reducing computational overhead. Specifically, the authors extracted semantics only on keyframes where the image content changed significantly. Liang et al. [15] proposed DIG-SLAM, which eliminates dynamic interference features by combining YOLOv7 instance segmentation, line segment detection, and K-means clustering to build a semantic map. However, the redundancy of line feature optimization limits the system's flexibility. The evolution of SLAM algorithms based on dynamic scenes is shown in Figure 1.

In summary, most existing visual SLAM systems generally suffer from poor robustness or real-time performance in highly dynamic environments. Previous research has focused on improving the accuracy of localization in visual SLAM algorithms, and the processing speed of these algorithms is constrained by the computational demands of their models. There has been little research on improving the real-time performance of visual SLAM systems in depth. Therefore, simultaneously improving the robustness and real-time performance

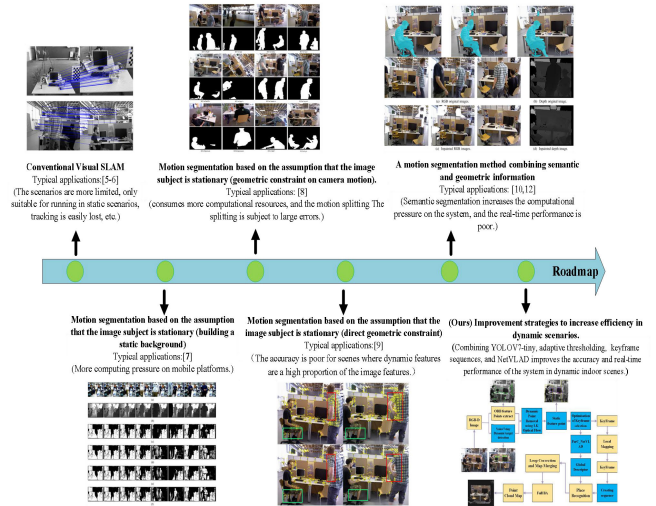


Fig. 1. Roadmap of dynamic SLAM.

of visual SLAM in indoor high-dynamic scenes is currently a hot issue.

B. Key Contributions

The main contributions of this paper are summarized below:

1. In the front-end part, the innovative use of lightweight target detection network YOLOv7-tiny combined with LK optical flow algorithm to remove the dynamic feature points in the detection frame, and based on the camera observation model, an adaptive threshold algorithm is added to solve the problems of poor keyframe quality of the original heuristic threshold selection and ensure that the number of keyframe feature points selected by the algorithm is sufficient and evenly distributed.

2. Based on the bag-of-words model in the back-end part, we dynamically create image sequences according to the angle differences between keyframes, replace the previous method of calculating a single BoW vector with the method of calculating the entire sequence's BoW vector, and then match the sequence. This reduces the workload and time consumption of loop detection in SLAM.

3. The ParC_NetVLAD image matching algorithm is proposed, which uses the ConvNeXt-Tiny network to extract features from images, and adds the ParC-Net network to the feature extraction network. Then, in order to extract the channel and spatial information of the feature map, the CBAM attention mechanism is added to the feature extraction network. Finally, NetVLAD is used to cluster the extracted local features to obtain the global features that can represent the image.

4. By verifying the innovative idea, the experimental results show that in the TUM dynamic scene dataset, compared with ORB-SLAM3, the positioning accuracy of our system is improved. The ATE of the proposed algorithm is reduced by 96.4% on average, and the RPE is reduced by 82.8% on average. The Pittsburgh dataset is tested, and the comparison experiments are carried out with NetVLAD, ResNet50, and other neural networks. The experimental results show that the

image-matching algorithm proposed in this paper has a higher accuracy. At the same time, the depth camera Astra Pro is tested on the actual dynamic scene, and the running trajectory of our system is more accurate compared with ORB-SLAM3. The algorithm in this paper has been improved in terms of positioning accuracy and real-time performance.

C. Paper Organization

The arrangement of this article is as follows. The opening section of this paper outlines the significance of the research work, related work, and challenges, as well as the main contributions of this paper. Section II introduces a visual SLAM system based on YOLOv7-tiny and NetVLAD. Section III provides experimental verification of the proposed technical scheme and analyses and discusses the experimental results. Section V concludes the whole paper and further points out future directions.

III. PROPOSED METHODS

In dynamic environments, the ORB-SLAM3 algorithm is affected by moving objects, leading to a decrease in positioning accuracy and poor robustness [16], [17]. To address this issue, we introduce the YOLOv7-tiny object detection algorithm into the front-end of the ORB-SLAM3 system to simultaneously detect targets in the input image and extract feature points [18]. After obtaining the semantic information in the image, the LK optical flow algorithm is used to determine the dynamic objects present in the image. Based on the target detection results, the dynamic ORB feature points are removed, and only static feature points are retained for pose calculation, thus improving the positioning accuracy of the visual SLAM system [19]. At the same time, an adaptive threshold method is used to optimize keyframe selection for higher quality keyframes. Then, we introduce a serialization process for keyframes before loop detection, combining the BoWs of several consecutive frames into one complete BoW. Based on this sequence (several consecutive frames), we perform loop closure detection to reduce the time required for SLAM loop detection. Next, the improved ConvNeXt network is used to extract features from the image, and the ParC-Net network and CBAM attention mechanism are added to the feature extraction network. Then, NetVLAD is used to cluster the extracted local features to obtain global features that can represent the image. The overall flowchart of the system is shown in Fig. 2, where the blue boxes are improved in this paper.

A. Dynamic Feature Point Culling Based on YOLOv7-tiny

YOLOv7 [20] is a deep learning-based object detection algorithm with high speed, accuracy, and efficiency, making it suitable for various object detection tasks in different scenarios. In this paper, YOLOv7-tiny is used as the detection model, which mainly consists of four parts: Input, Backbone, Neck, and Prediction. Figure 3 shows the front-end feature point optimization process.

Elimination of dynamic feature points using YOLOv7-tiny. In the dynamic object prediction box generated based on

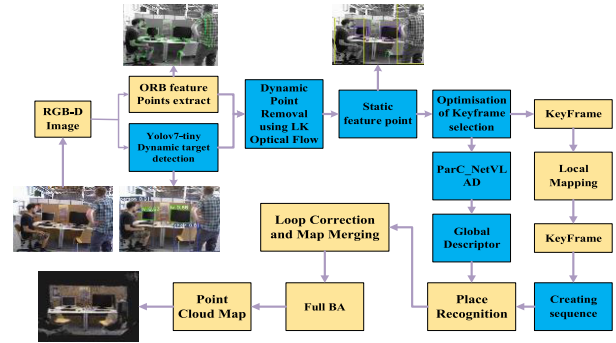


Fig. 2. Overview of our algorithm.

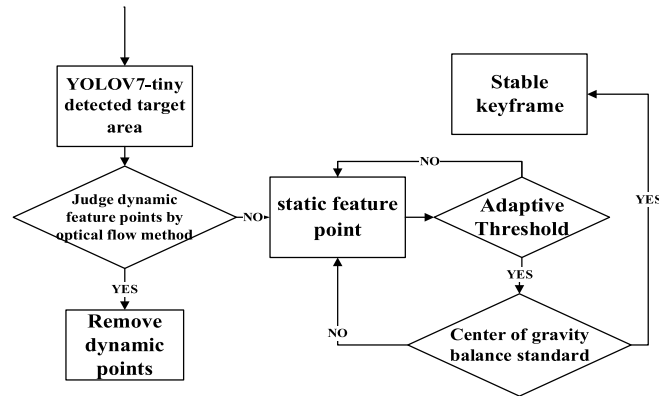


Fig. 3. Flowchart of front-end feature point optimization.

object detection, the dynamic feature points are determined by a priori knowledge, and the elimination process is described as follows:

In the tracking thread of ORB-SLAM3, the feature points are denoted as F_k , and k denotes the k -th frame processed by the system. After being processed by YOLOv7-tiny and the dynamic recognition algorithm, the recognized dynamic objects are denoted as D_k . Removing the dynamic feature points from all the feature points extracted by the ORB-SLAM3 algorithm gives the static feature points P_k .

To cull dynamic objects, this study introduces the YOLOv7-tiny target detection algorithm at the front end of the system to detect targets in the input image. The LK optical flow algorithm is then utilized to determine the dynamic objects in the target detection frame [21], [22], [23], [24]. The ORB feature point located at (x, y) in the image at time t can be expressed as $I(x, y, t)$. Based on the grayscale invariance assumption of the optical flow method, [25] derives that:

$$[I_x I_y]_k \begin{bmatrix} u \\ v \end{bmatrix} = -I_t, k = 1, \dots, w^2 \quad (1)$$

where u and I_x are the velocity and gradient of the motion of the feature point on the X-axis, respectively. v and I_y are the velocity and gradient on the Y-axis, respectively. I_t is the amount of variation of the feature point's grayscale against time. The symbol w represents the size of the window being set. If the optical flow of a feature point is greater than the threshold ε_{th} , the point is considered a dynamic point.

The static feature points are optimized for keyframes using the camera observation model proposed by Azimi et al. [26].

The adaptive threshold is determined based on the number of points where the line of sight changes, allowing for an accurate modeling and selection of the current frame as a key frame. In addition to satisfying adaptive thresholding, the stability of the frame must also be considered, with the center of gravity threshold used to determine if the points are evenly distributed throughout the image. After the image enters the SLAM system, it will be divided into 3×3 grids, and a 3×3 matrix will be generated. In each grid, the points with a conical area change greater than 30 degrees are counted as effective points, and this angle has been determined by previous experiments [27].

Suppose that the coordinates of each pixel of the image are x_i , the pixel value is p_i , and the direction coordinate of the centre of gravity is denoted as x .

$$\sum_{i=1}^n p_i (x_i - x) = 0 \quad (2)$$

Then the center of gravity in the x,y direction is:

$$\begin{cases} x = \frac{\sum_{i=1}^n p_i x_i}{\sum_{i=1}^n p_i} \\ y = \frac{\sum_{i=1}^n p_i y_i}{\sum_{i=1}^n p_i} \end{cases} \quad (3)$$

The centre of gravity of the matrix serves as the threshold T_1 ,

$$T_1 = \sqrt{x^2 + y^2} \quad (4)$$

T_2 Defined as:

$$T_2 = \frac{\sqrt{(c_{\max_1} - c_{\max_2})^2 + (r_{\max_1} - r_{\max_2})^2} * M_2}{2\sqrt{2}M_1} \quad (5)$$

where M_1 and M_2 are the two largest numbers within the matrix, and c and r are the numbers of rows and columns, respectively. If $T_2 > T_1$, the feature point distribution of the current frame is considered more appropriate, and the frame has sufficient stability.

B. The Dynamic Creation Method of Keyframe Sequences for Optimizing Closed-Loop Detection

The selection of keyframes is a crucial step in visual SLAM systems, significantly reducing information redundancy, improving frame quality, and minimizing error accumulation, thus providing a better foundation for back-end optimization. Particularly in highly dynamic environments with frequent scene changes, each frame's pose estimation contains some degree of error, and dynamic objects may lead to local mismatches and error accumulation. Keyframes can effectively confine the optimization problem within a smaller region, improving both the efficiency and effectiveness of optimization. By storing and comparing information from keyframes, the system can more accurately recognize loop

closures, thereby conducting global optimization to eliminate accumulated errors.

Based on the research of Bampis et al. [28], we create a sequence of keyframes in ORB-SLAM3 to recognize keyframes generated by the local mapping thread as a sequence with common features. R_{cw} is the rotation matrix that describes the rotation from the world coordinate system to the camera coordinate system. From ORB-SLAM3, we can derive the rotation matrix R_{cw2} of the current keyframe k_2 and the rotation matrix R_{cw1} of the previous keyframe k_1 and then calculate the rotation matrix R_{cw}^1 between the two frames.

$$R_{cw}^1 = R_{cw2} \times (R_{cw1})^{-1} \quad (6)$$

From the rotation matrix, we can find the angle of rotation between the two θ_1 :

$$\theta_1 = \arccos \frac{\text{tr}(R_{cw}^1) - 1}{2} \quad (7)$$

where $\text{tr}(X)$ is the sum of the diagonal elements of matrix X .

Similarly, you can find the rotation matrix R_{cw}^2 between a keyframe k_2 and its next keyframe k_3 :

$$R_{cw}^2 = R_{cw3} \times (R_{cw2})^{-1} \quad (8)$$

From the rotation matrix, we can find the angle of rotation between the two θ_2 :

$$\theta_2 = \arccos \frac{\text{tr}(R_{cw}^2) - 1}{2} \quad (9)$$

Meanwhile, we set an angle threshold $\theta_T = 0.1$. If $\theta_1 > \theta_T$ and $\theta_2 < \theta_T$, then no new sequence will be created, and new keyframes will continue to be added to the current sequence. If $\theta_1 < \theta_T$ and $\theta_2 > \theta_T$, then the current sequence ends, and a new sequence will be created. No new sequence will be created in all other cases. Additionally, to prevent sequences from being too long or too short, we set a length threshold $L_T = 15$ and a sequence length L_S . No new sequence will be created when the sequence length L_S is less than 2, and the current sequence ends and a new sequence is created when the sequence length L_S exceeds the length threshold L_T . The algorithm for creating new sequences is shown in Algorithm 1.

C. Improved NetVLAD-Based Loop Closure Detection Algorithm

Although bag-of-words model-based closed-loop detection can achieve certain results in some cases, this method is mainly based on artificially designed features, which are difficult to extract in complex environments [30]. Therefore, the bag-of-words model-based closed-loop detection algorithm is not effective in complex environments, which seriously limits the development of visual SLAM.

The closed-loop detection based on deep learning is essentially image re-identification, and the key technology is image matching technology. Based on this core, this paper proposes the ParC_NetVLAD image-matching algorithm. Firstly, to extract more representative features of the image, the ConvNeXt-Tiny [31] network (hereinafter referred to as ConvNeXt) [29] is used for feature extraction of the image. Secondly, to extract the location information features of the

Algorithm 1 Creating Sequence Algorithm

input: KeyFrame K_F , Length of the keyframe list L_K , Length of current sequence L_S
output: Results after verification of new sequences
 1 Initialization $\theta_T = 0.1$, $L_T = 15$, $i = 0$
 2 **for** $i \leftarrow 0$ to L_K **do**
 3 $\theta_{i+1} = \text{CalAngle}(K_{F(i)}, K_{F(i+1)})$
 4 //Calculate the angle of rotation between two keyframes
 5 **if** $L_S < 2$ **then**
 6 newsequence = 0
 7 **end if**
 8 **if** $L_S > L_T$ **then**
 9 newsequence = 1
 10 **end if**
 11 **if** $\theta_j > \theta_T$ and $\theta_{j+1} < \theta_T$ **then**
 12 newsequence = 0
 13 **end if**
 14 **if** $\theta_j < \theta_T$ and $\theta_{j+1} > \theta_T$ **then**
 15 newsequence = 1
 16 **end if**
 17 // For all other remaining cases, newsequence = 0
 18 **end for**

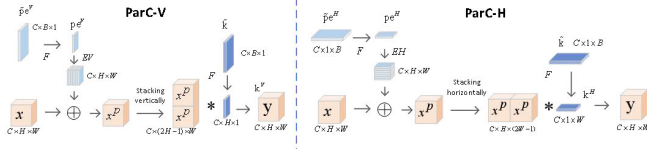


Fig. 4. ParC: (a) ParC-V module; (b) ParC-H module.

image, the ParC-Net [32] network is added to the feature extraction network. Then, to extract the channel and spatial information of the feature map, the CBAM [33] attention mechanism is added to the feature extraction network. Finally, NetVLAD [34] is used to cluster the local features extracted to obtain the global features that can represent the image.

IV. IMPROVED CONVNEXT NETWORK

A. ParC-Net Network

ParC-Net is a lightweight convolutional backbone network, whose core module, ParC (Position Information Aware Circular Convolution), extracts global features in both horizontal and vertical directions by associating the convolutional kernel weights with their positional information. By combining ParC-V and ParC-H convolution operations, ParC-Net is able to extract global features from all input pixels and achieve global information exchange at a relatively low computational cost, as illustrated in Figure 4.

B. CBAM Attention Mechanism

CBAM (Convolutional Block Attention Module) is a lightweight attention mechanism that enhances feature representation capability by combining information from both channel and spatial dimensions. The CBAM module first uses the Channel Attention Module (CAM) to extract the importance weights between channels, performing global average pooling and max pooling for each channel, and then applying these weights to optimize the feature map. Next, the Spatial Attention Module (SAM) further extracts spatial features by

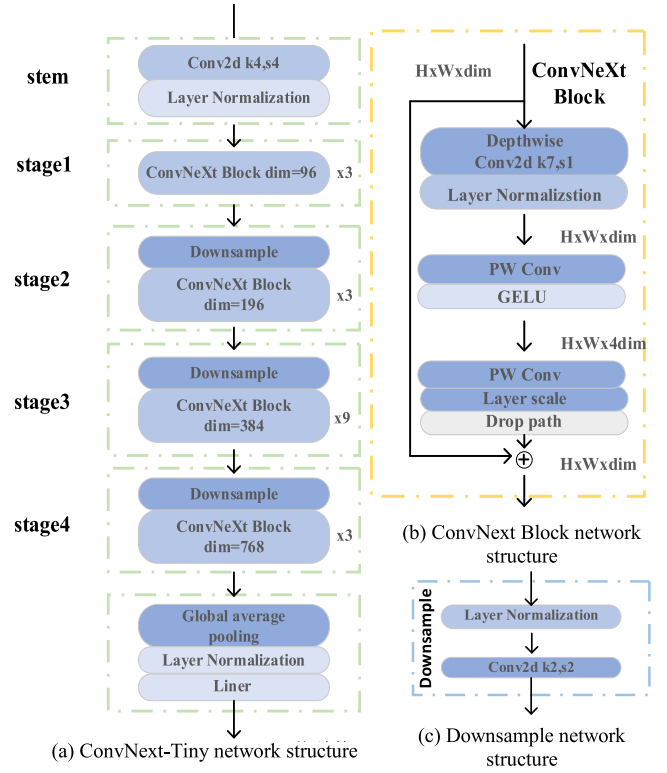


Fig. 5. ConvNeXt-Tiny.

performing global pooling, channel concatenation, convolution, and applying a Sigmoid function to generate a spatial attention map, which is multiplied by the channel-optimized feature map to produce the final output features. The CBAM module has a simple design, is easy to plug in, and can significantly improve network performance without altering the dimensions of the feature map.

C. Improved Convnext Network

ConvNeXt network was first used in image classification, and its accuracy has exceeded Swin Transformer [35] under the same FLOPs. Its structure is shown in Figure 5.

In order to better apply the ConvNeXt network in closed-loop detection, this paper removes the final global average pooling layer, regularization layer, and Linear layer. In addition, in order to extract the location information in image features, the ParC-Net network is introduced. In order to ensure the advantages of the ConvNeXt Block module and better utilize the ability of the ParC-Net module to focus on location information, ParC-Net is added to stage 3 and stage 4 of the ConvNeXt network. ParC is added to the last three ConvNeXt Block modules in stage 3 and the last ConvNeXt Block module in stage 4. The structure of the ParC-ConvNeXt Block module obtained after replacement is shown in Figure 6.

At this time, the information on the features extracted by the improved ConvNeXt network still cannot meet the needs. In order to better extract features, this paper adds the CBAM attention mechanism after the improved ParC_ConvNeXt network. CBAM attention mechanism can adjust the receptive field of the convolutional neural network according to the image content so that the model can better capture the

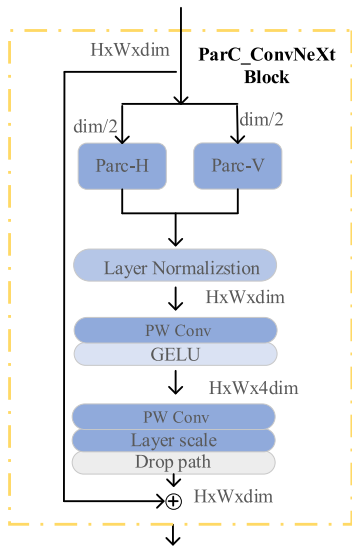


Fig. 6. ParC_ConvNeXt Block.

important information in the image. It can adaptively adjust the feature map at different levels so that the model can better distinguish between important features and noise. In addition, the CBAM attention mechanism can also reduce overfitting and improve the generalization ability of the model.

1) *Closed-Loop Detection Algorithm Based on Improved Netvlad Network*: Deep neural networks can learn deeper features of the image by convolution operation, and can accurately reflect image information. It is guaranteed that accurate image matching can still be achieved in complex environments so that closed-loop detection can have better effects. This paper studies the NetVLAD network, and on this basis proposes an improved ParC_NetVLAD, which is used for closed-loop detection. The NetVLAD network is improved on the basis of VLAD (Vector of Locally Aggregated Descriptors) [36]. The idea of VLAD is to gather all the local descriptors (such as SIFT, SURF, ORB, etc.) in an image and form a vector, which is used to represent the global descriptor of the whole image.

The VGG19 convolution network is used in the NetVLAD network to extract image features. Although the extracted image features can contain a lot of information about the image, the image features extracted by the VGG19 network lack position information, and the deeper features in the image cannot be better used. In addition, the VGG19 network has a large number of parameters in the training process, which is easy to overfit the training data, resulting in the decline of model performance. In this paper, aiming at the shortcomings of VGG19 in image feature extraction, the improved ConvNeXt network is proposed to extract image features. Finally, the extracted features are processed by the NetVLAD module to generate the required $K \times D$ as features. The improved network is shown in Figure 7.

In the whole improved network, the input image with dimension of $3 \times 480 \times 640$ is input. The image is output with feature dimension of $738 \times 15 \times 20$ after ParC_ConvNeXt module. The feature map is sent to CBAM module to extract channel and space information and keep the dimension of the feature map unchanged. Then the feature map is sent to

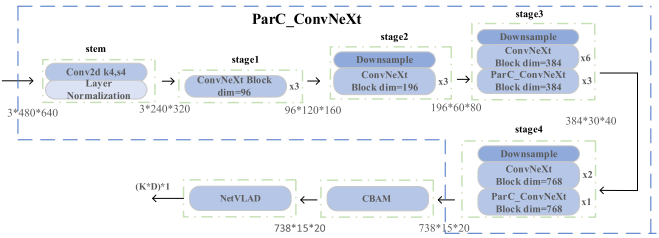


Fig. 7. Network structure of ParC_ConvNeXt_T_NetVLAD.

TABLE I
EXPERIMENTAL ENVIRONMENT

Hardware/Software component	Specification
Operating system	Ubuntu18.04
GPU	RTX 3090(24GB)
CPU	12 vCPU Intel(R) Xeon(R) Platinum 8255C CPU @ 2.50GHz
PyTorch	1.8.0
Python	3.7

NetVLAD module for dimension reduction, so that the feature map becomes the required $K \times D$ dimension feature, i.e. 64×738 . Since the dataset used in this study contains only similar image pairs and dissimilar image pairs, the triplet loss function [37] is selected for training.

V. EXPERIMENTAL RESULTS AND ANALYSIS

A. Dataset and Experimental Environment

In this experiment, we mainly used the TUM [38] dataset, Pittsburgh30K [39] dataset, and Tokyo247 [40] dataset. The TUM dataset is obtained by collecting RGB-D and monocular image sequences and corresponding motion capture live trajectories in indoor and outdoor environments. The dataset is challenging due to the presence of dynamic objects, changing lighting conditions, and large camera motion, making it a good test platform for evaluating the robustness and accuracy of SLAM algorithms. Pittsburgh30K is a subset dataset created by the Department of Computer Science at Carnegie Mellon University in 2012, containing more than 30,000 high-resolution street-view images. Tokyo247 is a large-scale street-view image dataset created by Waseda University and the University of Tokyo, containing more than 240,000 images. These datasets are widely used in computer vision tasks and urban planning research and are valuable in human activity recognition, autonomous driving, and intelligent transportation systems.

The experimental environment is shown in Table I:

B. Evaluation Metrics

In this paper, the trajectory accuracy is evaluated to verify the optimization degree of our method. The pose estimation error analysis experiment uses the evo tool, as well as the evaluate_ate and evaluate_rpe tools on the TUM website to test and compare the camera pose estimated by the ORB-SLAM3 system CameraTrajectory.txt with the real pose given by

TABLE II
COMPARISON OF ABSOLUTE TRAJECTORY ERROR BETWEEN ORB-SLAM3 AND OUR ALGORITHM (ATE/UNIT:M)

Datasets	ORB-SLAM3				Ours				Improvements%			
	Mean	Median	RMSE	STD	Mean	Median	RMSE	STD	Mean	Median	RMSE	STD
walking_static	0.3256	0.3173	0.3598	0.1524	0.0057	0.0051	0.0066	0.0032	98.25%	98.39%	98.17%	97.90%
walking_xyz	0.5534	0.4824	0.6426	0.3251	0.0127	0.0109	0.0151	0.0081	97.71%	97.74%	97.65%	97.51%
walking_rpy	0.7378	0.7168	0.821	0.3558	0.0262	0.0194	0.0405	0.0309	96.45%	97.29%	95.07%	91.32%
walking_half	0.5086	0.4889	0.5466	0.194	0.0169	0.0142	0.0223	0.0145	96.68%	97.10%	95.92%	92.53%
sitting_static	0.0067	0.0060	0.0077	0.0038	0.0054	0.0046	0.0062	0.0031	19.40%	23.33%	19.48%	18.42%

TABLE III
RESULTS OF TRANSLATIONAL RELATIVE POSE ERROR (RPE/UNIT:M)

Datasets	ORB-SLAM3				Ours				Improvements%			
	Mean	Median	RMSE	STD	Mean	Median	RMSE	STD	Mean	Median	RMSE	STD
walking_static	0.1021	0.0213	0.2225	0.1977	0.0099	0.0081	0.0124	0.0074	90.30%	61.97%	94.43%	96.26%
walking_xyz	0.2621	0.1643	0.3578	0.2435	0.0183	0.0158	0.0212	0.0107	93.02%	90.38%	94.07%	95.61%
walking_rpy	0.2701	0.1427	0.3903	0.2817	0.0456	0.0283	0.0828	0.0691	83.12%	80.17%	78.79%	75.47%
walking_half	0.1444	0.0556	0.2373	0.1882	0.0253	0.0215	0.0321	0.0197	82.48%	61.33%	86.47%	89.53%
sitting_static	0.0085	0.0076	0.0097	0.0046	0.0068	0.0060	0.0077	0.0037	20.00%	21.05%	20.62%	19.57%

the dataset groundtruth.txt. The evaluation parameters mainly adopt relative pose error (RPE) and absolute trajectory error (ATE) [38].

The comparison with the error of ORB-SLAM3 and the relative lift rate η is calculated as follows:

$$\eta = \frac{\text{orbislam3} - \text{ours}}{\text{orbislam3}} \times 100\% \quad (10)$$

In the loop detection experiment, the features of each query image are first clustered, and then the L2 distance is compared with the features within the class. The L2 distance calculation formula is shown as follows.

$$d_2 = \sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2} \quad (11)$$

Then select the 20 images with the highest similarity and sort them, and finally compare whether the 20 images contain the query image scene in order. When comparing the similarity, we set $N = [1, 5, 10, 20]$. If the first image is a correct match, then add 1 to the counter under $N = 1$. If the first five images have a correct match, then add 1 to the counter under $N = 5$. Similarly, for $N = 10$ and $N = 20$, the same operation is carried out. After all the query images are matched, the value of the counter under each N is counted. Then divide it with the total number of query images to find the accuracy.

C. Experimental Results Analysis

TUM is a dataset for indoor dynamic environments. The results will be compared with ORB-SLAM3 [5], Dyna

SLAM [11], DS-SLAM [9], SOF-SLAM [41], and OFM-SLAM [42]. Table II compares the ATE with ORB-SLAM3, and Table III and Table IV compare the RPE. Four metrics were measured for the performance of the VSLAM system, namely Standard Deviation (STD), Root Mean Square Error (RMSE), Median Error (Median), and Mean Error (Mean). In dynamic scenarios, our performance is somewhat improved compared to most algorithms. Visual interference can be better suppressed when the object is moving faster, resulting in more accurate localization and mapping. However, in the case of slowly moving objects, the optimization strategy used in the algorithm has only a minor impact on system accuracy, resulting in a slight improvement over the ORB-SLAM3 algorithm. Therefore, the proposed algorithm can be considered as a promising solution for SLAM in dynamic environments with fast-moving objects.

To ensure the effectiveness of each module, we conducted ablation experiments on YOLOv7-tiny and the adaptive threshold module. The experimental results are shown in Table V (\checkmark denotes the existence of this module, and \times denotes the removal of this module), and the evaluation criterion is the RMSE value. Table V clearly shows the impact of each module. From the data results, it can be concluded that using both modules simultaneously achieve the minimum root mean square error for the algorithm.

In Table VI, our algorithm is compared with other SLAM methods in dynamic environments. From Table VI, it can be seen that the highest positioning accuracy is achieved

TABLE IV
RESULTS OF ROTATIONAL RELATIVE POSE ERROR (RPE/UNIT:DEG)

Datasets	ORB-SLAM3				Ours				Improvements%			
	Mean	Median	RMSE	STD	Mean	Median	RMSE	STD	Mean	Median	RMSE	STD
walking_static	2.6563	2.9369	3.0119	1.4053	0.2571	0.2278	0.2986	0.1518	90.32%	92.24%	90.09%	89.20%
walking_xyz	4.9684			4.5852	0.5081	0.4164	0.644	0.3956	89.77%	86.98%	90.48%	91.37%
walking_rpy	5.2563	3.1978	6.7628	5.4278	0.6988	0.3264	4.2352	4.1771	86.71%	87.87%	43.96%	23.04%
walking_half	2.8593	2.6908	7.5570	5.4278	0.6397	0.3264	4.2352	4.1771	86.71%	87.87%	43.96%	23.04%
sitting_static	2.8593	1.3501	4.5654	3.5575	0.6397	0.5782	0.7276	0.3465	77.63%	57.17%	84.06%	90.26%
sitting_static	0.2646	0.2528	0.2929	0.1255	0.2390	0.2232	0.2653	0.1151	9.67%	11.71%	9.42%	8.29%

TABLE V
ABLATION EXPERIMENT ATE

YOLOv7-tiny	Adaptive threshold	walking_xyz	walking_half	walking_s_tatic	walking_rpy
✓	✓	0.0151	0.0223	0.0066	0.0405
✓	✗	0.0158	0.0239	0.0113	0.0535
✗	✓	0.2988	0.1930	0.0274	0.3276
✗	✗	0.6426	0.5466	0.3598	0.821

TABLE VI

THE ABSOLUTE TRAJECTORY ERROR OF DIFFERENT ALGORITHMS(ATE)

	walking_static	walking_xyz	walking_rpy	walking_half	sitting_s_tatic
ORB-SLAM3	0.3598	0.6426	0.8210	0.5466	0.0077
DynaSLAM	0.0073	0.0131	0.0279	0.0174	0.0058
DS-SLAM	0.0081	0.0247	0.4442	0.0303	0.0065
SOF-SLAM	0.0070	0.0180	0.0270	0.0290	0.0100
OFM-SLAM	0.0081	0.0168	0.0326	0.0252	0.0065
Ours	0.0066	0.0151	0.0405	0.0223	0.0062

by DynaSLAM and our algorithm. Since DynaSLAM uses a MaskRCNN semantic segmentation network that performs per-pixel execution, the recognition precision is higher. Based on YOLOV7-tiny, our visual SLAM strategy can quickly detect and track moving objects. Although it has slightly lower accuracy than the DynaSLAM algorithm, it has a significant advantage in speed, making it suitable for real-time scenarios that require fast processing. We performed time tests on our algorithm and DynaSLAM algorithm on the w_xyz dataset, and Table VII lists the processing time per frame for both our algorithm and the DynaSLAM algorithm. Compared with the DynaSLAM algorithm, it shows superior performance in terms of real-time processing speed.

Table VIII presents a comparison of the testing time for loop closure detection between using the sequential keyframe method in ORB-SLAM3 and ORB-SLAM3 alone. It can be seen that in datasets with indoor loops, such as fr1_room and

TABLE VII

TRACKING TIME COMPARISON (UNIT MS)

Algorithm	Time
Dyna-SLAM[11]	2452
Ours	29
Improvements	98.78%

TABLE VIII

LOOP CLOSING TOTAL TIME COMPARISON (UNIT MS)

Datasets	ORB-SLAM2[4]	ORB-SLAM3	Ours	Improvements
fr1_room	1170	379	346	8.71%
fr2_desk	1584	571	544	4.73%

fr2_desk, the total time for loop closure detection is relatively reduced when using the sequential keyframe method.

Figure 8 shows the trajectory and error distribution estimated by ORB-SLAM3, DynaSLAM, and our algorithm. In the figure, the black, blue, and red lines denote the true trajectory, the estimated trajectory, and the error between them. The graphic clearly illustrates the high level of accuracy in the estimated pose by the algorithm presented in this article, which closely matches the ground truth trajectory.

Figure 9 shows the error distribution and trajectory estimation results of the proposed algorithm across four highly dynamic datasets. The results indicate that the proposed method achieves relatively low relative pose error, with the estimated trajectory closely aligning with the ground truth. The deviation between the estimated and actual trajectories is largely maintained within the range of 0 to 0.04, demonstrating the high localization accuracy of the proposed approach.

The improved ParC_NetVLAD algorithm for loop detection is compared with the original algorithm NetVLAD, ResNet50_NetVLAD and other algorithms, and the comparison results are shown in Figure 10. It can be seen from Figure 10 that in the Pittsburgh30k dataset, the improved model has an accuracy of 83.83% when $N = 1$, which is 79.55% higher than the original NetVLAD model. It is 2.15% and 2.46% higher than the ResNet50_NetVLAD model and the MobilenetV2_NetVLAD model, respectively. When

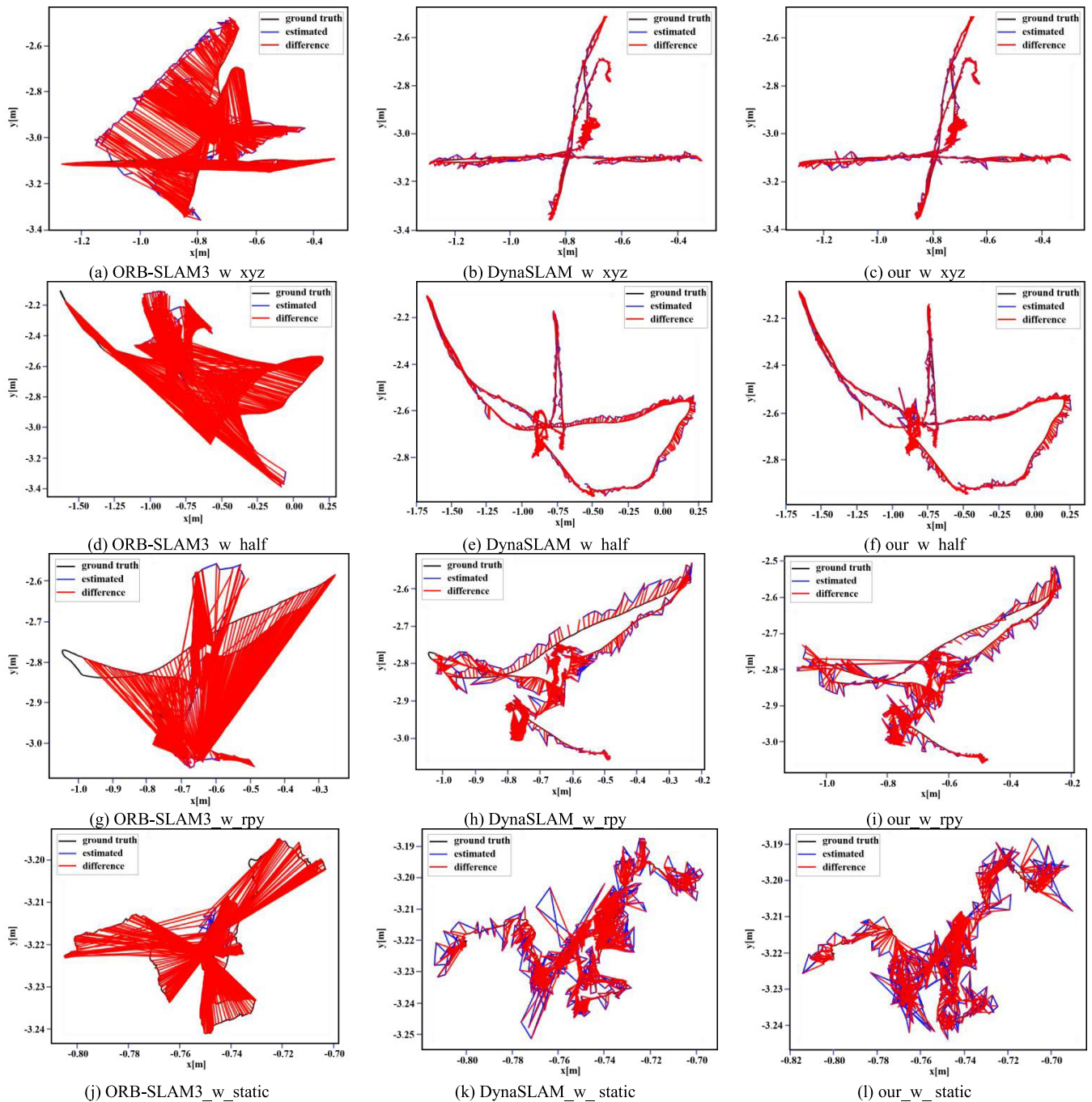


Fig. 8. Comparison of estimated trajectories and real trajectories in highly dynamic environments.

$N = 5$, it is 2.55% higher than the original model NetVLAD, 2.17% and 2.35% higher than the ResNet50_NetVLAD model and the MobilenetV2_NetVLAD model, respectively. When $N = 10$, it is 1.69% higher than the original model NetVLAD, 1.55% and 1.55% higher than the ResNet50_NetVLAD model and the MobilenetV2_NetVLAD model, respectively. When $N = 20$, it is 1.72% higher than the original model NetVLAD, 1.81% and 1.9% higher than the ResNet50_NetVLAD model and the MobilenetV2_NetVLAD model, respectively.

At the same time, in order to verify the effectiveness of the improved ParC_ConvNeXt model. This paper will use

CBAM_NetVLAD network and ParC_ConvNeXt_NetVLAD network to test on the Pittsburgh30k dataset, as shown in Figure 11(a). It can be seen from Figure 11(a) that the accuracy of the ParC_ConvNeXt module is higher than that without the use of the module. When $N = 1$, it is 1.86% higher; when $N = 5$, it is 2.11% higher; when $N = 10$, it is 1.53% higher; when $N = 20$, it is 1.55% higher. Figure 11(b) is a comparative test to verify the effectiveness of the CBAM attention mechanism module. It can be seen that the network with the CBAM attention mechanism module is 1.6% higher than that without the CBAM attention mechanism module when $N = 1$; when $N = 5$, it is 1.8% higher; when $N = 10$, it is

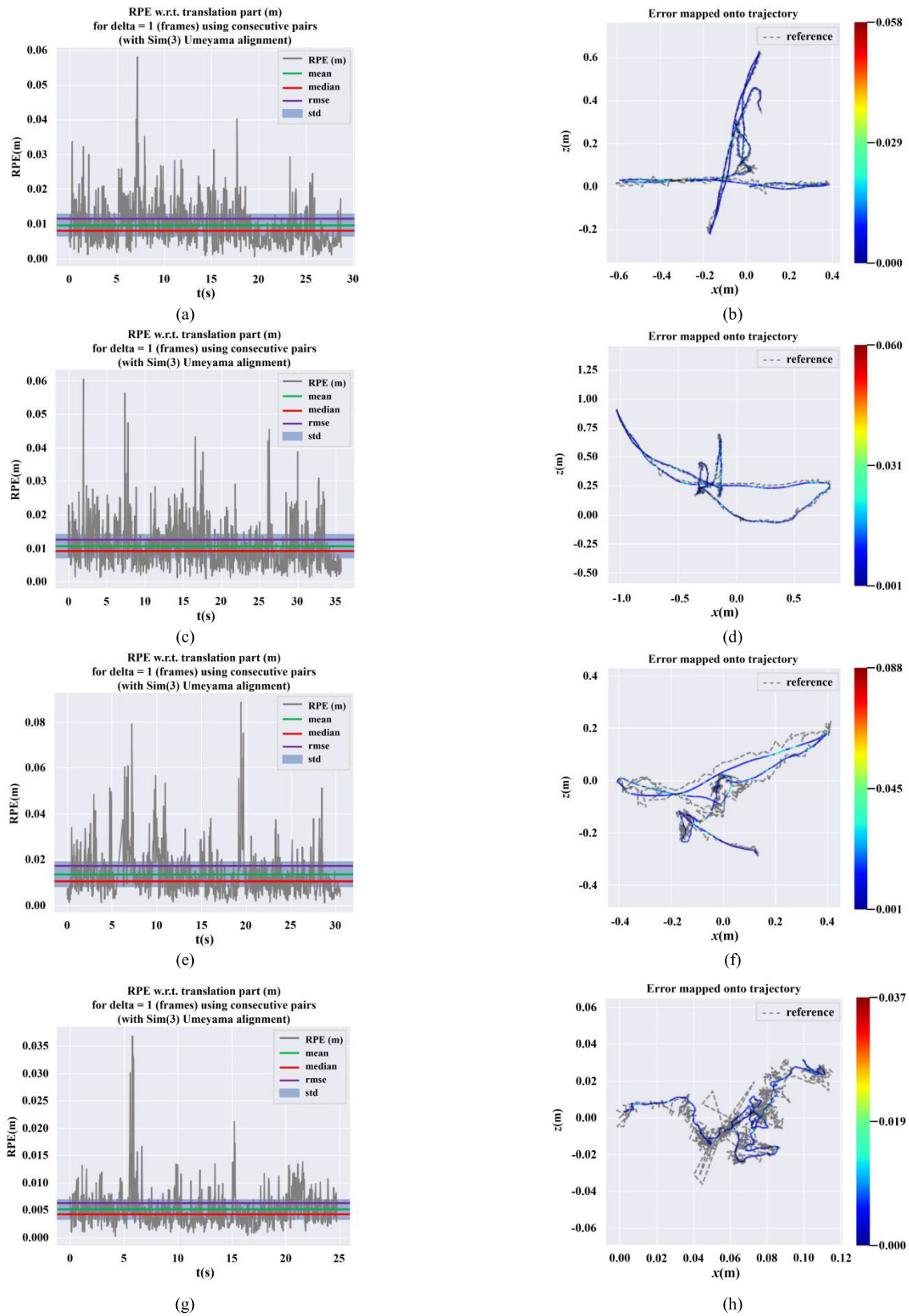
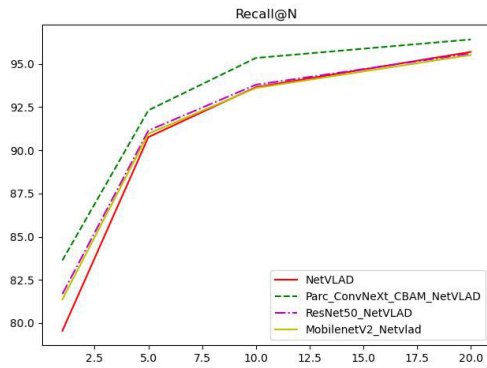


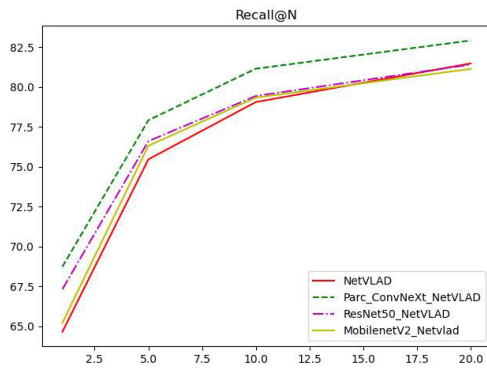
Fig. 9. Relative pose error. (a) (b) walking-xyz. (c) (d) walking-halfsphere. (e) (f) walking-rpy. (g) (h) walking-static.

1.12% higher; when $N = 20$, it is 1.24% higher. In summary, the closed-loop detection based on the improved NetVLAD algorithm studied in this paper has a certain improvement in

the accuracy of image matching, which can better improve the closed-loop detection performance of the image scene matching system.

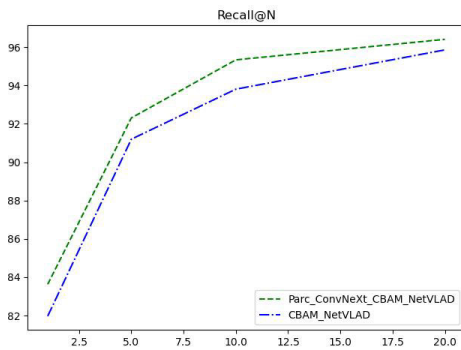


(a) Pittsburgh30k dataset

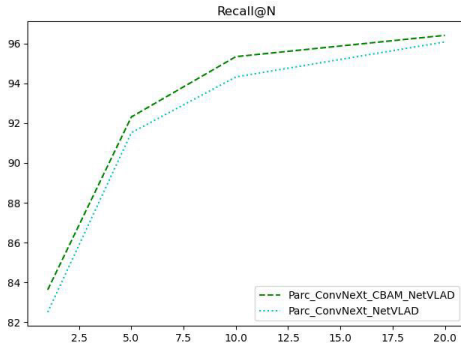


(b) Tokyo247 dataset

Fig. 10. Experimental results of the dataset.



(a) ParC ConvNeXt module test



(b) CBAM module test

Fig. 11. Comparison of ablation experiments.

D. Effects of the Actual Scene

To further validate the effectiveness of our proposed algorithm, experiments were conducted in real-world scenarios



Fig. 12. Mobile robot and Astra Pro camera.

TABLE IX
ASTRA PRO PARAMETERS

Astra Pro parameter	
Parameter	Values
Measuring range	0.6m - 8m
Resolution of color image	1280×720@30FPS
Resolution of depth map	1280×1024@7FPS
Precision	(±1-3)mm@1m
Colour FOV	H66.1°and V40.2°
Depth FOV	H58.4°and V45.5°
Delay (ms)	30-45



Fig. 13. Experimental site.

using a Jetson Nano mobile robot as the moving platform and an Astra Pro RGB-D camera produced by ORBBEC Corporation as the visual sensor. The appearance and parameters are shown in Figure 12 and Table IX, respectively. The depth measurement part of Astra Pro adopts structured light technology and performs depth ranging through the triangulation principle, showing excellent performance in indoor environments.

The experimental site is a rectangular conference room, as depicted in Figure 13. The intelligent robot departs from point A and moves in a clockwise direction of A, B, C, D, and back to A. At the same time, people will walk back and forth at points B and C during the movement, as shown in Figure 13, allowing for the recognition of dynamic objects in indoor scenes.

Finally, after multiple evaluations in real-world environments, the actual running trajectories of indoor dynamic scenes



Fig. 14. Comparison of ORB-SLAM3 with ORB feature extraction in our system. This method removes the dynamic features of the walking human body.

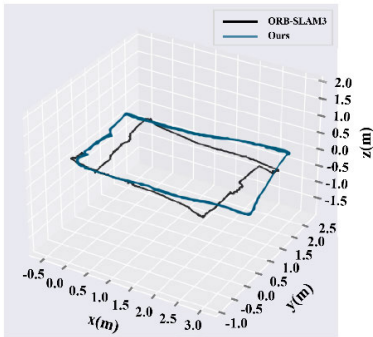


Fig. 15. ORB-SLAM3 compared to the actual operating trajectory of our system.

based on ORB-SLAM3 and our algorithm are displayed in Figure 15. It can be observed that the actual trajectory based on ORB-SLAM3 has a significant deviation under the interference of dynamic objects while the one based on our method is much more accurate and can better close the loop.

In the office experiment, we used the camera signal and the corresponding timestamp to determine the trajectory calculation value of the passing point. We set the starting point and the endpoint at the position of point A, and set one point at the position of B, C, and D. In addition, we also added a point in the middle of two points, forming a total of nine passing points.

In this indoor flat experimental environment, the car moves horizontally and the coordinates change very little in the vertical direction, so this paper only focuses on the plane positioning accuracy. To obtain the true value of the trajectory, the passing points of the vehicle were marked with millimeter accuracy through manual measurement during the experiment. In this paper, the two-dimensional coordinates are used to represent the position of the camera and the conversion relationship between the coordinate system where

TABLE X
EXPERIMENTAL RESULTS OF TRAJECTORY LOCALIZATION IN THE OFFICE ENVIRONMENT

	Ours	ORB-SLAM3	Ground truth
Route 1	(-0.0003,-0.0003)	(-0.0004,-0.0003)	(-0.0001,-0.0005)
Route 2	(-0.0417,0.9185)	(-0.0472,0.9234)	(-0.0591,0.9296)
Route 3	(0.0281,2.0783)	(-0.0391,2.1288)	(0.0436,2.0787)
Route 4	(1.2409,2.0664)	(1.1703,2.1391)	(1.2366,2.0964)
Route 5	(3.0418,2.2329)	(2.9859,2.3148)	(2.9987,2.2875)
Route 6	(3.0636,1.3308)	(2.7756,1.3489)	(3.0198,1.3713)
Route 7	(2.8889,0.4121)	(2.8343,0.3916)	(2.9619,0.3988)
Route 8	(1.9701,0.3202)	(1.9048,0.3062)	(1.9807,0.3010)
Route 9	(0.0681,0.0071)	(0.0384,0.0236)	(0.0724,0.0060)
RMSE	0.0421	0.1056	—

the aligned trajectory is located and the office coordinate system is determined through the calculation of the starting position and direction. In this way, we can calculate the true coordinates of the passing points. The experimental results are shown in Table X. According to the RMSE of the ATE of the passing points, the positioning accuracy of the algorithm in this paper is significantly improved, which is 60.13% higher than that of ORB-SLAM3.

VI. CONCLUSION

Moving objects may introduce motion blur, occlusion, and environmental appearance changes, which make it difficult for VSLAM systems to maintain accurate and consistent estimations. In this paper, an efficient deep neural network object detection method YOLOv7-tiny algorithm is applied to ORB-SLAM3, and the adaptive threshold method is used to replace the heuristic threshold in ORB-SLAM3 to ensure the quality of key frames selected by the algorithm in dynamic scenes. By filtering dynamic objects in the environment and extracting feature points of static regions, the accuracy of the visual SLAM algorithm in dynamic scenes and the robustness of the system are improved.

Then, in loop detection, keyframes are processed based on the bag of binary features. According to the angle difference between keyframes, the key frame sequence is dynamically divided into key frame sequences to reduce the time consumption of the loop detection part of SLAM. At the same time, the NetVLAD is improved, the ParC_ConvNeXt network is designed to focus on the location information in the image feature map, and the CBAM attention mechanism is added to the network to focus on the feature channel and spatial information. Finally, the obtained image features are sent to the NetVLAD pooling network for dimensionality reduction.

Our algorithm reduces the average ATE by 96.4% and the average RPE by 82.8% in high-dynamic scenes, which is better than most dynamic visual SLAM. In addition, compared to the localization accuracy of DynaSLAM, the time consumption of the tracking thread in our algorithm is reduced by more than 98%. Therefore, our algorithm significantly improves the computational speed while enhancing the localization precision.

However, the performance of our algorithm on the TUM low-dynamic scenes is far inferior to the ORB-SLAM3 algorithm's localization accuracy. This is also a common drawback of many visual SLAM algorithms based on dynamic scenes. To improve the adaptability of algorithms in different environments, more flexible and intelligent methods for detecting different environmental changes need to be developed.

The experimental results on the Pittsburgh30 dataset and Tokyo247 dataset show that the accuracy of the improved algorithm is up to 4.25% higher than that of the other methods when $N = 1$, 2.55% higher when $N = 5$, 1.69% higher when $N = 10$, and 1.9% higher when $N = 20$. In summary, this paper improves the positioning accuracy and real-time performance of visual SLAM in indoor high-dynamic situations.

Future research efforts should focus on optimizing the algorithm model and improving the accuracy of indoor dynamic scene localization. Additionally, it would be worthwhile to explore leveraging semantic information extracted from object detection to build semantic maps, enabling the system to handle higher-level tasks.

ACKNOWLEDGMENT

The authors would like to thank the referees for their constructive suggestions.

REFERENCES

- [1] X. Shen et al., "A closed-loop detection algorithm for online updating of bag-of-words model," in *Proc. 9th Int. Conf. Comput. Data Eng.*, Haikou, China, Jan. 2023, pp. 34–40.
- [2] D. Cai, R. Li, Z. Hu, J. Lu, S. Li, and Y. Zhao, "A comprehensive overview of core modules in visual SLAM framework," *Neurocomputing*, vol. 590, Jul. 2024, Art. no. 127760.
- [3] Y. Fu, B. Han, Z. Hu, X. Shen, and Y. Zhao, "CBAM-SLAM: A semantic SLAM based on attention module in dynamic environment," in *Proc. 6th Asian Conf. Artif. Intell. Technol. (ACAIT)*, Changzhou, China, Dec. 2022, pp. 1–6.
- [4] R. Mur-Artal and J. D. Tardós, "ORB-SLAM2: An open-source SLAM system for monocular, stereo, and RGB-D cameras," *IEEE Trans. Robot.*, vol. 33, no. 5, pp. 1255–1262, Oct. 2017.
- [5] C. Campos, R. Elvira, J. J. G. Rodríguez, J. M. M. Montiel, and J. D. Tardós, "ORB-SLAM3: An accurate open-source library for visual, visual-inertial, and multimap SLAM," *IEEE Trans. Robot.*, vol. 37, no. 6, pp. 1874–1890, Dec. 2021.
- [6] D.-H. Kim and J.-H. Kim, "Effective background model-based RGB-D dense visual odometry in a dynamic environment," *IEEE Trans. Robot.*, vol. 32, no. 6, pp. 1565–1573, Dec. 2016.
- [7] Y. Sun, M. Liu, and M. Q.-H. Meng, "Improving RGB-D SLAM in dynamic environments: A motion removal approach," *Robot. Auto. Syst.*, vol. 89, pp. 110–122, Mar. 2017.
- [8] H. Wei, T. Zhang, and L. Zhang, "GMSK-SLAM: A new RGB-D SLAM method with dynamic areas detection towards dynamic environments," *Multimedia Tools Appl.*, vol. 80, nos. 21–23, pp. 31729–31751, Sep. 2021.
- [9] C. Yu et al., "DS-SLAM: A semantic visual SLAM towards dynamic environments," in *Proc. IEEE/RSJ Int. Conf. Intell. Robots Syst. (IROS)*, Oct. 2018, pp. 1168–1174.
- [10] V. Badrinarayanan, A. Kendall, and R. Cipolla, "SegNet: A deep convolutional encoder-decoder architecture for image segmentation," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 39, no. 12, pp. 2481–2495, Dec. 2017.
- [11] B. Bescos, J. M. Fácil, J. Civera, and J. Neira, "DynaSLAM: Tracking, mapping, and inpainting in dynamic scenes," *IEEE Robot. Autom. Lett.*, vol. 3, no. 4, pp. 4076–4083, Jul. 2018.
- [12] K. He, G. Gkioxari, P. Dollár, and R. Girshick, "Mask R-CNN," in *Proc. IEEE Int. Conf. Comput. Vis. (ICCV)*, Oct. 2017, pp. 2980–2988.
- [13] C. Shao, L. Zhang, and W. Pan, "Faster R-CNN learning-based semantic filter for geometry estimation and its application in vSLAM systems," *IEEE Trans. Intell. Transp. Syst.*, vol. 23, no. 6, pp. 5257–5266, Jun. 2022.
- [14] G. Singh, M. Wu, M. V. Do, and S.-K. Lam, "Fast semantic-aware motion state detection for visual SLAM in dynamic environment," *IEEE Trans. Intell. Transp. Syst.*, vol. 23, no. 12, pp. 23014–23030, Dec. 2022, doi: [10.1109/TITS.2022.3213694](https://doi.org/10.1109/TITS.2022.3213694).
- [15] R. Liang, J. Yuan, B. Kuang, Q. Liu, and Z. Guo, "DIG-SLAM: An accurate RGB-D SLAM based on instance segmentation and geometric clustering for dynamic indoor scenes," *Meas. Sci. Technol.*, vol. 35, no. 1, Sep. 2023, Art. no. 015401, doi: [10.1088/1361-6501/acfb2d](https://doi.org/10.1088/1361-6501/acfb2d).
- [16] Y. Chen et al., "BEVSOC: Self-supervised contrastive learning for calibration-free BEV 3-D object detection," *IEEE Internet Things J.*, vol. 11, no. 12, pp. 22167–22182, Jun. 2024.
- [17] H. Qi, Y. Fu, Z. Hu, J. Wu, and Y. Zhao, "A lightweight semantic vslam approach based on adaptive thresholding and speed optimization," *J. Beijing Univ. Aeronaut. Astronaut.*, 2024, doi: [10.13700/jb.1001-5965.2023.0552](https://doi.org/10.13700/jb.1001-5965.2023.0552).
- [18] R. Li, Y. Zhao, Z. Hu, W. Qi, and G. Liu, "Tohf: A feature extractor for resource -constrained indoor vslam," *J. Syst. Simul.*, pp. 1–12, doi: [10.16182/j.issn1004731x.joss.23-1334](https://doi.org/10.16182/j.issn1004731x.joss.23-1334).
- [19] S. Yang, A. Xu, P. Li, M. Chen, P. Du, and K. Shao, "Visual SLAM algorithm based on YOLOv5 in dynamic scenario," in *Proc. China Autom. Congr. (CAC)*, 2023, pp. 2640–2645.
- [20] C.-Y. Wang, A. Bochkovskiy, and H.-Y. Mark Liao, "YOLOv7: Trainable bag-of-freebies sets new state-of-the-art for real-time object detectors," 2022, *arXiv:2207.02696*.
- [21] A. Bochkovskiy, C.-Y. Wang, and H.-Y. M. Liao, "YOLOv4: Optimal speed and accuracy of object detection," 2020, *arXiv:2004.10934*.
- [22] Y. Ai, T. Rui, M. Lu, L. Fu, S. Liu, and S. Wang, "DDL-SLAM: A robust RGB-D SLAM in dynamic environments combined with deep learning," *IEEE Access*, vol. 8, pp. 162335–162342, 2020.
- [23] Y. Fan et al., "Semantic SLAM with more accurate point cloud map in dynamic environments," *IEEE Access*, vol. 8, pp. 112237–112252, 2020.
- [24] S. Han and Z. Xi, "Dynamic scene semantics SLAM based on semantic segmentation," *IEEE Access*, vol. 8, pp. 43563–43570, 2020.
- [25] Z. Chang, H. Wu, Y. Sun, and C. Li, "RGB-D visual SLAM based on Yolov4-tiny in indoor dynamic environment," *Micromachines*, vol. 13, no. 2, p. 230, Jan. 2022.
- [26] A. Azimi, A. Hosseinaveh Ahmadabadian, and F. Remondino, "PKS: A photogrammetric key-frame selection method for visual-inertial systems built on ORB-SLAM3," *ISPRS J. Photogramm. Remote Sens.*, vol. 191, pp. 18–32, Sep. 2022.
- [27] A. Hosseinaveh et al., "Automatic image selection in photogrammetric multi-view stereo methods," in *Proc. Eurograph. Assoc.*, 2012, pp. 1–8.
- [28] L. Bampis, A. Amanatiadis, and A. Gasteratos, "Fast loop-closure detection using visual-word-vectors from image sequences," *Int. J. Robot. Res.*, vol. 37, no. 1, pp. 62–82, Jan. 2018.
- [29] Z. Liu, H. Mao, C.-Y. Wu, C. Feichtenhofer, T. Darrell, and S. Xie, "A ConvNet for the 2020s," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2022, pp. 11966–11976.
- [30] Z. Hu, W. Qi, K. Ding, G. Liu, and Y. Zhao, "An adaptive lighting indoor vSLAM with limited on-device resources," *IEEE Internet Things J.*, vol. 11, no. 17, pp. 28863–28875, Sep. 2024, doi: [10.1109/JIOT.2024.3406816](https://doi.org/10.1109/JIOT.2024.3406816).
- [31] Y.-H. Choi and S.-C. Kee, "Monocular depth estimation using a Laplacian image pyramid with local planar guidance layers," *Sensors*, vol. 23, no. 2, p. 845, Jan. 2023.
- [32] H. Zhang, W. Hu, and X. Wang, "ParC-Net: Position aware circular convolution with merits from ConvNets and transformer," in *Proc. 17th Eur. Conf. Comput. Vis.*, Tel Aviv, Israel. Cham, Switzerland: Springer, 2022, pp. 613–630.
- [33] S. Woo, J. Park, J. Y. Lee, and I. S. Kweon, "CBAM: Convolutional block attention module," in *Proc. Eur. Conf. Comput. Vis. (ECCV)*, 2018, pp. 3–19.
- [34] R. Arandjelovic, P. Gronat, A. Torii, T. Pajdla, and J. Sivic, "NetVLAD: CNN architecture for weakly supervised place recognition," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2016, pp. 5297–5307.
- [35] Z. Liu et al., "Swin transformer: Hierarchical vision transformer using shifted windows," in *Proc. IEEE/CVF Int. Conf. Comput. Vis. (ICCV)*, Oct. 2021, pp. 9992–10002.

- [36] H. Jégou, M. Douze, C. Schmid, and P. Pérez, "Aggregating local descriptors into a compact image representation," in *Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit.*, Jun. 2010, pp. 3304–3311.
- [37] F. Schroff, D. Kalenichenko, and J. Philbin, "FaceNet: A unified embedding for face recognition and clustering," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2015, pp. 815–823.
- [38] J. Sturm, N. Engelhard, F. Endres, W. Burgard, and D. Cremers, "A benchmark for the evaluation of RGB-D SLAM systems," in *Proc. IEEE/RSJ Int. Conf. Intell. Robots Syst.*, Oct. 2012, pp. 573–580.
- [39] A. Torii, J. Sivic, T. Pajdla, and M. Okutomi, "Visual place recognition with repetitive structures," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, Jun. 2013, pp. 883–890.
- [40] A. Torii, R. Arandjelovic, J. Sivic, M. Okutomi, and T. Pajdla, "24/7 place recognition by view synthesis," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2015, pp. 1808–1817.
- [41] L. Cui and C. Ma, "SOF-SLAM: A semantic visual SLAM for dynamic environments," *IEEE Access*, vol. 7, pp. 166528–166539, 2019.
- [42] X. Zhao, T. Zuo, and X. Hu, "OFM-SLAM: A visual semantic SLAM for dynamic indoor environments," *Math. Problems Eng.*, vol. 2021, pp. 1–16, Apr. 2021.



Zhuhua Hu (Senior Member, IEEE) received the B.Eng. and M.Eng. degrees from Jilin University in 2002 and 2005, respectively, and the Ph.D. degree from Hainan University in 2019. He was a Software Engineer at the Ningbo BIRD Research Institute of China from 2005 to 2006. He was a Software Engineer at Nanjing Research Institute, ZTE, from 2006 to 2007. He was a Minister of the Software Department at Shanghai Aoxun Information Technology Company Ltd., from 2007 to 2009. He has been a Professor and Doctorial Tutor at the

School of Information and Communication Engineering, Hainan University, since 2020. He is currently a high-level talent in Hainan Province. He led the "multimodal information intelligent processing and decision control" innovation team, and has published more than 110 academic articles in journals, such as *IEEE TRANSACTIONS ON MULTIMEDIA*, *IEEE TRANSACTIONS ON IMAGE PROCESSING*, *IEEE TRANSACTIONS ON VEHICULAR TECHNOLOGY*, *IEEE TRANSACTIONS ON AEROSPACE AND ELECTRONIC SYSTEMS*, *IEEE TRANSACTIONS ON GEOSCIENCE AND REMOTE SENSING*, *IEEE TRANSACTIONS ON INTELLIGENT TRANSPORTATION SYSTEMS*, *IEEE TRANSACTIONS ON COMPUTER-AIDED DESIGN OF INTEGRATED CIRCUITS AND SYSTEMS*, *OE*, and *COMPAG*, authorized 13 patents, and hosted more than ten large-scale commercial projects that have been successfully implemented. His current research interests include multimodal information processing and decision control. He is a CCF Senior Member. He acted as a reviewer of *IEEE INTERNET OF THINGS JOURNAL*, *Ocean Engineering*, *IEEE-ACM TRANSACTIONS ON NETWORKING*, *Engineering Applications of Artificial Intelligence*, *Image and Vision Computing*, *IEEE TRANSACTIONS ON EMERGING TOPICS IN COMPUTING*, *ICASSP2023&2024*, and *ICME 2024*.



Wenlu Qi received the B.S. degree in electronic information from Taiyuan University of Science and Technology, China, in 2021. She is currently pursuing the M.S. degree with the School of Information and Communication Engineering, Hainan University, China. Her research interests include is visual SLAM.



Kunkun Ding received the B.S. degree in electronic information science and technology from Henan University Minsheng College, China, in 2022. He is currently pursuing the M.S. degree with the School of Information and Communication Engineering, Hainan University, China. His research interests include is visual SLAM.



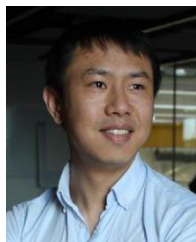
Hao Qi received the B.S. degree in communication engineering from the University of Jinan, China, in 2020. He is currently pursuing the master's degree in information and communication engineering with the School of Information and Communication Engineering, Hainan University. His current research interests include computer vision and SLAM.



Yaochi Zhao received the M.S. degree in pattern recognition and intelligent system from Central South University, Changsha, China, in 2005, and the Ph.D. degree from Tianjin University, Tianjin, China, in 2023. She worked at the Ningbo BIRD Research Institute and Shanghai Wingtech Communication Company Ltd., for three years. Later, she was engaged in teaching and research work at the College of Information Science and Technology, Hainan University, Haikou, China, where she is currently an Associate Professor at the School of Cyberspace Security. Her current research interests include image processing and computer vision and visual SLAM.



Xuebo Zhang (Senior Member, IEEE) received the B.Eng. degree in automation from Tianjin University, Tianjin, China, in 2002, and the Ph.D. degree in control theory and control engineering from Nankai University, Tianjin, in 2011, where he is currently a Professor with the Institute of Robotics and Automatic Information System, Tianjin Key Laboratory of Intelligent Robotics. His research interests include mobile robotics, motion planning, visual sensor network, localization, and mapping. He is a Technical Editor of *IEEE/ASME TRANSACTIONS ON MECHATRONICS* and an Associate Editor of the *ASME Journal of Dynamic Systems, Measurement, and Control*.



Mingfeng Wang (Member, IEEE) received the B.Eng. degree in mechanical design and automation and the M.Eng. degree in mechanical engineering from Central South University, Changsha, China, in 2008 and 2012, respectively, and the Ph.D. degree in mechanical engineering from the University of Cassino and South Latium, Italy, in 2016. He is currently a Senior Lecturer with the Robotics and Autonomous Systems, Brunel University London, London, U.K. He has published more than 30 articles that have been presented at peer-reviewed international conferences or journals. His research interests include cover novel design and development of humanoid robots, precision farming robots, continuum robots and hexapod robots, which including robotic technologies in terms of mechanical design, kinematic and dynamic analysis, motion planning, motor control, system design and integration, and fabricating and debugging.