

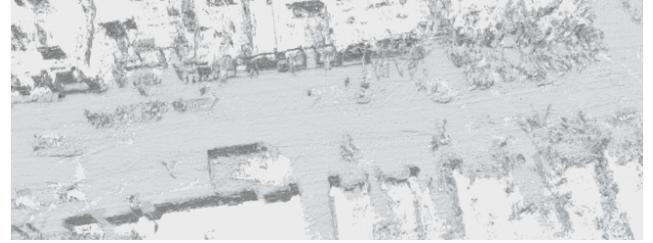
Neural Implicit Representation for Highly Dynamic LiDAR Mapping and Odometry

Qi Zhang^{1,*}, He Wang^{2,*}, Ru Li², and Wenbin Li¹

Abstract—Recent advancements in Simultaneous Localization and Mapping (SLAM) have increasingly highlighted the robustness of LiDAR-based techniques. At the same time, Neural Radiance Fields (NeRF) have introduced new possibilities for 3D scene reconstruction, exemplified by SLAM systems. Among these, NeRF-LOAM has shown notable performance in NeRF-based SLAM applications. However, despite its strengths, these systems often encounter difficulties in dynamic outdoor environments due to their inherent static assumptions. To address these limitations, this paper proposes a novel method designed to improve reconstruction in highly dynamic outdoor scenes. Based on NeRF-LOAM, the proposed approach consists of two primary components. First, we separate the scene into static background and dynamic foreground. By identifying and excluding dynamic elements from the mapping process, this segmentation enables the creation of a dense 3D map that accurately represents the static background only. The second component extends the octree structure to support multi-resolution representation. This extension not only enhances reconstruction quality but also aids in the removal of dynamic objects identified by the first module. Additionally, Fourier feature encoding is applied to the sampled points, capturing high-frequency information and leading to more complete reconstruction results. Evaluations on various datasets demonstrate that our method achieves more competitive results compared to current state-of-the-art approaches.

I. INTRODUCTION

Recently, Neural Radiance Fields (NeRF) [1] have emerged as a powerful tool for 3D scene reconstruction and novel view synthesis, largely due to their ability to generate highly detailed and realistic representations of complex scenes from sparse observations. NeRF's strength lies in its capacity to model scenes as continuous volumetric fields, enabling high-quality novel view synthesis and accurate 3D reconstructions from multiple input images. This capability makes NeRF particularly promising for Simultaneous Localization and Mapping (SLAM) applications, both indoors and outdoors. By capturing intricate scene details and providing dense volumetric representations, NeRF enhances SLAM systems, improving spatial understanding and object recognition. In indoor environments, several NeRF-based SLAM systems have been developed, including iMap [2], NICE-SLAM [3], NeRF-SLAM [4], Co-SLAM [5], GO-SLAM [6], and NGEL-SLAM [7]. These systems leverage camera sensors to deliver precise reconstructions. For outdoor scenarios, NeRF-based SLAM has been extended to incorporate LiDAR



(a) Ours in KITTI MOT 19



(b) NeRF-LOAM [8] in KITTI MOT 19

Fig. 1: The 3D reconstruction of KITTI MOT 19 using our proposed method and NeRF-LOAM [8].

sensors, as demonstrated in NeRF-LOAM [8], LONER [9], and PinSLAM [10]. Additionally, some systems combine visual and LiDAR data, such as CLONeR [11] and SiLVR [12]. However, most NeRF-based SLAM systems operate under the assumption that the environment is static or only minimally dynamic. This assumption presents significant challenges when applying these methods to the real world with highly dynamic objects in outdoor scenarios, where the scene reconstruction becomes considerably inaccurate.

In this paper, we aim to construct a dense 3D map of outdoor scenes in highly dynamic scenarios. To achieve this, we enhance NeRF-LOAM [8] by integrating an additional thread that detects the moving objects, generating their 3D bounding boxes. This allows us to categorize the LiDAR points into background and foreground segments, under the assumption that the background remains static while the foreground is dynamic. Building upon the foundation of NeRF-LOAM [8], we focus exclusively on calculating the Signed Distance Function (SDF) values for the background, ensuring that dynamic elements are accurately separated from the static environment.

Moreover, we propose a hybrid feature representation method that significantly benefits the generation of highly dynamic objects. We extended the octree structure in NeRF-LOAM [8] to support multiple resolutions, similar to the

*These authors contributed equally and are considered as co-first authors

¹Department of Computer Science, University of Bath, UK, {qz727, w.li}@bath.ac.uk

²School of Computer and Information Technology, ShanXi University, ShanXi, China, wh_sxu@foxmail.com, liru@sxu.edu.cn

approach used in NGLOD [13]. By considering multiple resolutions, our method can better capture the fine details and motion of dynamic objects across different scales. Inspired by Co-SLAM [5], we recognized that although parametric encoding alone can improve reconstruction results, it has limitations in filling holes and ensuring smooth transitions. To overcome these challenges, we combine multiple learnable features with Fourier feature positional encoding [14].

Our experimental results demonstrate that our methods are particularly effective for highly dynamic scenes, as it not only remove the dynamic objects and fills holes, but also enhances smoothness, which is crucial for maintaining the integrity and continuity of fast-moving objects. The key contributions are summarized as follows:

- We enhance NeRF-LOAM [8] by integrating a method to segment scenarios into the dynamic foreground and static background, removing dynamic LiDAR points and rebuilding ground points to facilitate accurate 3D mapping in highly dynamic outdoor environments.
- We extended the single-layer learnable features in the octree of NeRF-LOAM [8] to multiple layers and applying Fourier feature encoding to the sampled points, allowing us to achieve better reconstruction results.
- To improve the accuracy of SDF values in highly dynamic scenarios, we incorporate additional loss functions into the optimization process.

II. RELATED WORKS

Simultaneous Localization and Mapping (SLAM) is a foundational technology in robotics and autonomous driving, crucial for enabling machines to navigate and understand their surroundings. SLAM systems can be categorized based on the sensors they utilize. LiDAR-based SLAM methods, such as Lo-net [15], Deppco [16], and Pwclo-net [17], have gained prominence due to their robustness and illumination invariance. Relying on RGB cameras, SLAM systems can capture color and texture information, including ORB-SLAM2 [18], ORB-SLAM3 [19], and DynaSLAM [20], [21]. Combining LiDAR and RGB data, like Gaussian-LIC [22], leverages the strengths of both modalities.

In parallel with advancements in SLAM, Neural Implicit Representations (NeRF) [1] have significantly influenced 3D scene reconstruction techniques. Early approaches mapped point coordinates to signed distance functions (SDFs) [23] and occupancy fields [24], but these methods required access to ground truth 3D geometry, which limited their applicability. VDBFusion [25] stores TSDF values in sparse voxels, resulting in a highly accurate but incomplete reconstructed map, Puma [26] represents the map as a triangle mesh through Poisson reconstruction, enabling it to capture more detailed maps compared to common mapping methods. More recent work introduced differentiable rendering functions to represent scenes using only 2D images [27], [28]. While these approaches improved accessibility to 3D scene reconstruction, they often resulted in overly smooth representations for complex shapes [29].

In the realm of NeRF-based SLAM, several approaches have been proposed to integrate NeRF with SLAM systems. IMap [2] was an early effort in this direction but struggled with network-induced forgetting and capacity limitations in large-scale environments. NICE-SLAM [3] improved on this by subdividing the scene into uniform grids and using a pre-trained geometry decoder for better reconstruction. ESSLAM [30] employed axis-aligned feature planes to manage memory requirements and used an implicit Truncated Signed Distance Field (TSDF) [31] for improved geometry representation and faster convergence. Vox-Fusion [32] introduced an incremental scene representation using octrees, which leveraged embeddings at octree leaf nodes for interpolation.

For large-scale environments, NeRF-LOAM [8] represents the first attempt to combine neural implicit representations with LiDAR data for odometry and mapping. Despite its advancements, NeRF-LOAM [8] faces challenges with dynamic object removal and hole filling. Pin-SLAM [10] uses optimizable neural points to construct an implicit map, but the improvement in reconstruction quality is limited by the fixed resolution of the neural points.

To address these issues, our approach involves removing points from dynamic regions, performing static filling to close large gaps, and reconstructing the scene using multi-resolution and Fourier feature encoding, thereby enhancing the overall robustness and accuracy of the reconstruction.

III. METHOD

A. System Overview

The Figure. 2 is the overview of the proposed system. It can be separated into two sections. The left part of the image shows the background and foreground separation, which is illustrated in Section III-B. In each frame, we detect the moving objects, then combine them with the moving objects from previous frames to form the foreground. For the points in the foreground, if they are not ground points, we remove all of them and generate the ground points within the space.

The right part of the image illustrates the scene representation and the training process of the SDF values. For the same query point, the interpolated embedding of its corresponding node is obtained from different levels of the octree. This embedding is then concatenated with the coordinates of the sampled point after Fourier encoding, and the combined data is fed into an MLP network to predict the SDF value, as detailed in Section III-C. Furthermore, for the non-ground points in the foreground, we propose a dynamic region-based SDF loss function to minimize the loss of the SDF value for these points. Section III-D describes the final loss function.

B. Background and Foreground Separation

To accurately reconstruct dynamic outdoor scenes, we implement a background and foreground separation strategy that effectively distinguishes between static and dynamic points. This process involves detecting the dynamic objects, marking their trajectories of the dynamic occupations as dynamic foreground zones, and ensuring consistent ground surface reconstruction.

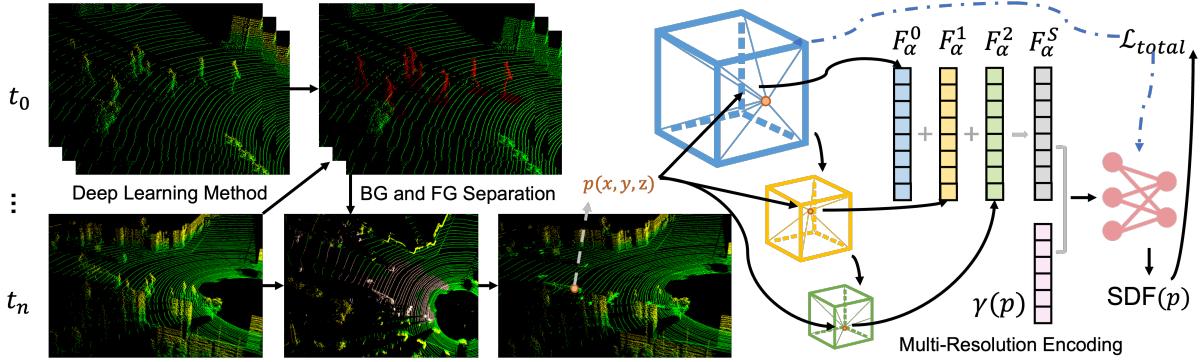


Fig. 2: The overview of the system. The left part of the image illustrates the process of background and foreground separation. We remove the dynamic points from the foreground, and generate a foreground mask (pink), resulting in a purely static scene. The right part of the image shows the training process of the neural SDF module. We interpolate the query point at different resolution levels to obtain the corresponding features, and finally combine the Fourier feature positional encoding and feed them into the MLP to predict the SDF value.

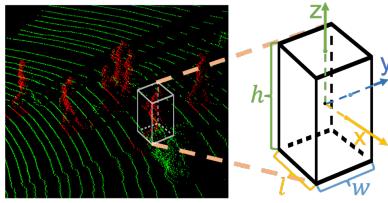


Fig. 3: Definition of the 3D box B .

Moving Object 3D Box Recognition:

In a separate thread, we use an existed deep-learning-based method to identify moving objects $\{O_N, \dots, O_M\} = \{O_i\}_{i=N}^M$ within the scene. These objects are tracked and enclosed within 3D bounding boxes, defined as $\{B_N, \dots, B_M\} = \{B_i\}_{i=N}^M$. As shown in Figure. 3, each bounding box $B_i = \{\mathbf{P}_{ci}, h_i, w_i, l_i\}$, where $\mathbf{P}_{ci} = \{x_{ci}, y_{ci}, z_{ci}\}^T$ represent the coordinates of the center of the 3D box B_i , and h (here the subscript c denotes the center), w and l denote its height, width, and, length respectively.

Dynamic Foreground Marking:

In the outdoor environment, we assume that if moving objects are detected, the space within their corresponding 3D bounding boxes is defined as foreground. During the 3D reconstruction, no structures should be built within the foreground area, except for the ground floor. This separation allows us to focus on reconstructing the static background accurately while treating the moving objects as outliers.

To achieve this, we maintain a global dynamic list L , where each element in the list corresponds to the dynamic mask D_i of a moving object O_i . The list is defined as $L = \{D_i\}_{i=N}^M$, representing the masks for all detected moving objects. The $D_i = \{\mathbf{p}_{li}, \mathbf{p}_{ri}\}$ is represented as:

$$\begin{aligned} \mathbf{p}_{li} &= \left((\mathbf{R}\mathbf{P}_{ci})_x + T_x - \frac{w_i}{2}, (\mathbf{R}\mathbf{P}_{ci})_y + T_y - \frac{l_i}{2} \right) \\ \mathbf{p}_{ri} &= \left((\mathbf{R}\mathbf{P}_{ci})_x + T_x + \frac{w_i}{2}, (\mathbf{R}\mathbf{P}_{ci})_y + T_y + \frac{l_i}{2} \right) \end{aligned} \quad (1)$$

where $(\mathbf{R}|T)$ is the camera pose with rotation matrix \mathbf{R} and translation vector T .

To update the list L , whenever a new frame is input, we also remove elements by checking the y_{ci} . If $y_{ci} < T_y$, the corresponding element is removed from the list.

Thus, when a new frame is input, if the LiDAR points are not ground points and fall within any mask D_i in the list $L = \{D_i\}_{i=N}^M$, these points are considered as part of the dynamic foreground and not be included in the static background reconstruction. (To check if the LiDAR points are ground points, we use the same method in [8].)

Ground Surface Generation in Foreground:

When dealing with dynamic foreground zones, our method assumes that only the ground surface G needs to be generated within these regions. However, ignoring dynamic objects entirely could lead to gaps in the ground's reconstruction. To address this, we first estimate an average height \bar{z}_G for the ground surface in the z-axis within a radius r around the dynamic mask D_i . The average height is calculated as:

$$\bar{z}_G = \frac{1}{|P_G|} \sum_{p \in P_G} z_p \quad (2)$$

where P_G is a set of ground points within the distance r around D_i , and z_p is the z-coordinate of each point in P_G .

To ensure a consistent ground surface representation, we then randomly generate LiDAR points $\hat{P}_G = \{\hat{\mathbf{p}}_1, \hat{\mathbf{p}}_2, \dots, \hat{\mathbf{p}}_n\}$ at the estimated height \bar{z}_G . These points $\hat{\mathbf{p}}_j$ are positioned such that:

$$\hat{\mathbf{p}}_j = (x_j, y_j, \bar{z}_G), \quad j = 1, 2, \dots, n \quad (3)$$

where (x_j, y_j) are randomly sampled within D_i .

These artificially generated points \hat{P}_G are then incorporated into the reconstruction process to generate the signed distance function (SDF) for the background.

Foreground SDF Loss Function:

The proposed foreground SDF Loss Function is designed to improve the accuracy and consistency of the Signed

Distance Function (SDF) values in the NeRF-LOAM [8] in highly dynamic environments where objects are moving.

For each LiDAR point \mathbf{p}_i , consider a dynamic region D_i defined as a certain radius r around the point. The SDF values within such dynamic regions are regularized to ensure that the geometric representation is smooth and accurate.

$$\mathcal{L}_d = \frac{1}{|D_i|} \sum_{\mathbf{p}_j \in D_i} \left(\Psi(\mathbf{p}_j) - \frac{1}{|D_i|} \sum_{\mathbf{p}_k \in D_i} \Psi(\mathbf{p}_k) \right)^2 \quad (4)$$

Here, $\Psi(\mathbf{p}_j)$ is the SDF value at point \mathbf{p}_j within the dynamic region D_i , and the loss encourages the SDF values within the region to be consistent.

C. Implicit neural scene representation

Multi-Resolution Encoding:

Inspired by NGLOD [13], we extend the octree structure in NeRF-LOAM [8]. Specifically, we embed features only at the deepest H levels of the octree to balance reconstruction quality and training speed. In other words, features are embedded only at levels $d \in \{D_{max} - H + 1, \dots, D_{max}\}$, where we set $H = 3$, which was empirically found to provide a good balance, D_{max} is the maximum depth of the octree. For a query point \mathbf{p} , its corresponding multi-resolution encoding is represented as:

$$F_\alpha^s(\mathbf{p}) = \sum_{j=D_{max}-H+1}^{D_{max}} F_\alpha^j(\mathbf{p}) \quad (5)$$

The embedding at each layer is obtained through trilinear interpolation of the embeddings at the eight vertices of the nodes in the current layer:

$$F_\alpha^j(\mathbf{p}) = TriInpo(\mathbf{p}, \mathbf{e}_1^j, \dots, \mathbf{e}_8^j) \quad (6)$$

where $TriInpo$ refers to trilinear interpolation, the parameters include the sample point \mathbf{p} and the embeddings of the eight vertices \mathbf{e}^i of the voxel containing the sample point.

Fourier Features Positional Encoding:

Instead of using the frequency encoding adopted in NeRF [1] to encode the sampled points into a higher dimension, following the suggestion by Sun et al. [33], we input the coordinates of the sampled point, encoded and combined with the embedding of the point, into the neural SDF module to predict its SDF value. Given a query point \mathbf{p} , its corresponding Fourier positional encoding is represented:

$$\gamma(\mathbf{p}) = [\sin(2\pi B_1 \mathbf{p}), \cos(2\pi B_1 \mathbf{p}), \dots, \sin(2\pi B_k \mathbf{p}), \cos(2\pi B_k \mathbf{p})]^\top \quad (7)$$

where B_i are coefficients ($i \in 1, 2, \dots, k$) sampled from an isotropic Gaussian distribution. $B_i \sim \mathcal{N}(0, \sigma^2)$, and σ is chosen for each task and dataset with a hyperparameter sweep. k serves as a hyperparameter that determines the length of the positional encoding feature.

Finally, we concatenate the embedding after trilinear interpolation with the query point coordinates after Fourier

positional encoding, and input it into the neural SDF module to predict the SDF value. The SDF value is represented as:

$$\Psi(\mathbf{p}) = f(\gamma(\mathbf{p}), F_\alpha^s) \quad (8)$$

where Ψ is the predicted SDF value of a certain point.

D. Optimization

The foreground loss can be integrated with the existing SDF loss \mathcal{L}_s , free space loss \mathcal{L}_f , and Eikonal loss \mathcal{L}_e in NeRF-LOAM [8]. The final loss function could be a weighted sum of these losses:

$$\mathcal{L}_{total} = \lambda_s \mathcal{L}_s + \lambda_f \mathcal{L}_f + \lambda_e \mathcal{L}_e + \lambda_d \mathcal{L}_d \quad (9)$$

where λ_s , λ_f , λ_e , and λ_d are the weights that balance the contributions of each loss term.

IV. EXPERIMENTS

In this section, we first elaborate on the detailed settings of the experiments (Section IV-A). As shown in Section IV-B and Section IV-C, the reconstruction results of our method are more competitive both qualitatively and quantitatively. Finally, in the ablation study, we demonstrated the effectiveness of each component in our proposed method. (Section IV-D).

A. Experiment Setup

1) **Baseline:** Our method is compared with the current state-of-the-art methods NeRF-LOAM [8] and Pin-SLAM[10] and Puma [26] for the 3D reconstruction results and odometry. We also compare our method with the current state-of-the-art explicit reconstruction method VDBFusion [25] and the neural implicit reconstruction method SHINE-Mapping [29].

2) **Datasets:** We evaluate our method on three public Lidar datasets, including KITTI [34], MaiCity [26], and Newer College [35] datasets. We use the MOT challenge in the KITTI dataset to evaluate our system with a highly dynamic outdoor dataset. MaiCity [26] is a synthetic LiDAR dataset of outdoor urban environments. Newer College [35] is a real radar dataset on a university campus. KITTI [34] does not provide ground truth maps, therefore, we only provide the qualitative results of the reconstruction. The other two datasets provide registered dense point clouds as ground truth references for quantitative evaluation.

3) **Evaluation Metric:** To ensure a fair comparison, for odometry, we use the Root Mean Square Error (RMSE) of the Absolute Trajectory Error (ATE) to evaluate. For mapping, we use the evaluation criteria employed in most methods [8], [26], [29], including accuracy, completion, Chamfer-L1 distance, and F-score. These metrics are calculated by comparing the ground truth and the predicted mesh.

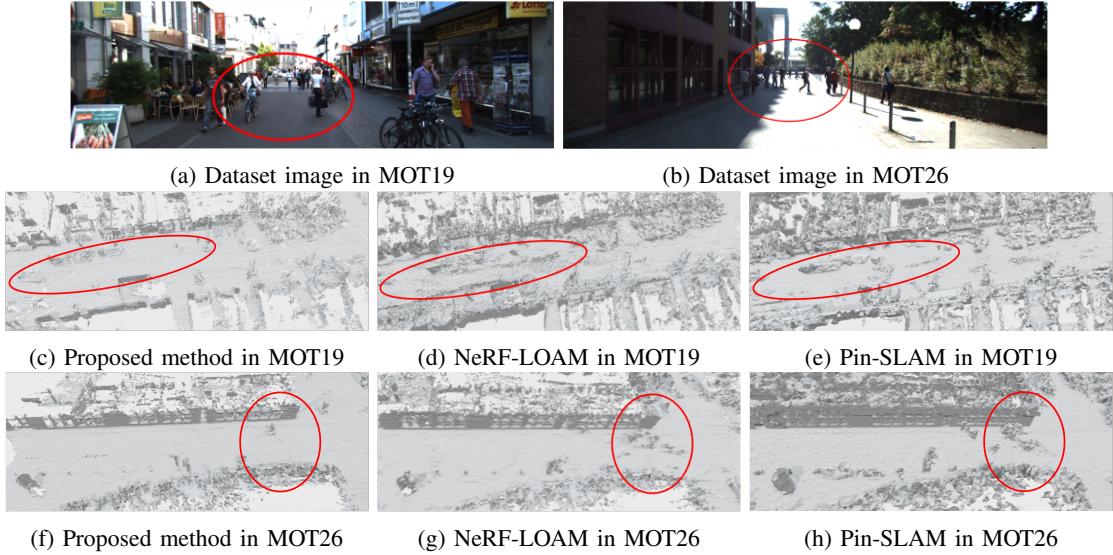


Fig. 4: 3D reconstruction results among Proposed method, NeRF-LOAM [8], Pin-SLAM [10]. (a) and (b): Original images from the datasets MOT19 and MOT26. The red ovals highlight dynamic pedestrians of the scene. (c) to (h): These show how the proposed Method, NeRF-LOAM [8], and Pin-SLAM [10] perform on the MOT19 dataset. The grayscale images are reconstructions, and the red ovals indicate the areas of focus for comparison.

TABLE I: Quantitative evaluation of the reconstruction quality on the MaiCity and NewerCollege. The term ‘w/o GT pose’ refers to reconstructions using odometry-estimated poses, while ‘w GT pose’ refers to that using ground truth poses.

Dataset	Method	Comp.[cm] \downarrow	Acc.[cm] \downarrow	C-II[cm] \downarrow	F-score[%] \uparrow
MaiCity	Puma	9.14	7.89	8.51	68.04
	Pin-SLAM	<u>6.27</u>	<u>5.11</u>	<u>5.69</u>	<u>85.30</u>
	NeRF-LOAM	9.00	5.33	7.17	81.99
	Ours	5.71	4.12	4.91	88.06
	VDBFusion	17.23	2.88	10.06	91.40
	SHINE	4.24	3.88	4.06	89.73
NewerCollege	NeRF-LOAM	<u>4.11</u>	3.48	<u>3.79</u>	<u>93.30</u>
	Ours	3.68	<u>3.41</u>	3.54	94.11
	Puma	71.91	15.30	43.60	57.27
	Pin-SLAM	<u>15.25</u>	<u>11.55</u>	<u>13.40</u>	82.08
	NeRF-LOAM	16.60	11.70	14.14	77.67
	Ours	15.05	10.91	12.98	81.93
w GT pose	VDBFusion	22.72	4.73	13.72	91.13
	SHINE	<u>14.36</u>	8.32	11.34	90.65
	NeRF-LOAM	15.59	<u>6.86</u>	<u>11.24</u>	<u>91.83</u>
	Ours	11.14	7.63	9.26	93.12

4) **Implementation Details:** All experiments were run on an RTX 4090 GPU. The machine-learning-based moving object detection method employed was QD-3DT [36]. The distance parameter r was set to 0.3 meters, meaning that points within a 0.3-meter radius around the moving object were considered reference ground points. For Fourier positional encoding, the variance parameter σ^2 was configured to 50, a value that proved effective across all datasets. The loss function parameter λ_L was set to 50. All other parameters were configured according to the settings specified in [8].

B. Map Results

In this section, we first demonstrate the reconstruction results of our method compared to NeRF-LOAM [8] and Pin-

SLAM [10] on the MOT dataset. Subsequently, we present the quantitative performance of our method in comparison to Puma [26], Pin-SLAM [10] and NeRF-LOAM [8] on the MaiCity [26] and Newer College [35] datasets. None of these four methods rely on ground truth poses provided by the dataset as a reference. Finally, to provide a more comprehensive comparison of the reconstruction results, we compared our method with several reconstruction-focused methods, including VDBFusion [25] and SHINE [29].

The Fig.4 highlights the differences between our proposed method, NeRF-LOAM [8] and Pin-SLAM [10] in the context of dynamic object removal in outdoor street scenes. The proposed method excels in reconstructing scenes with moving

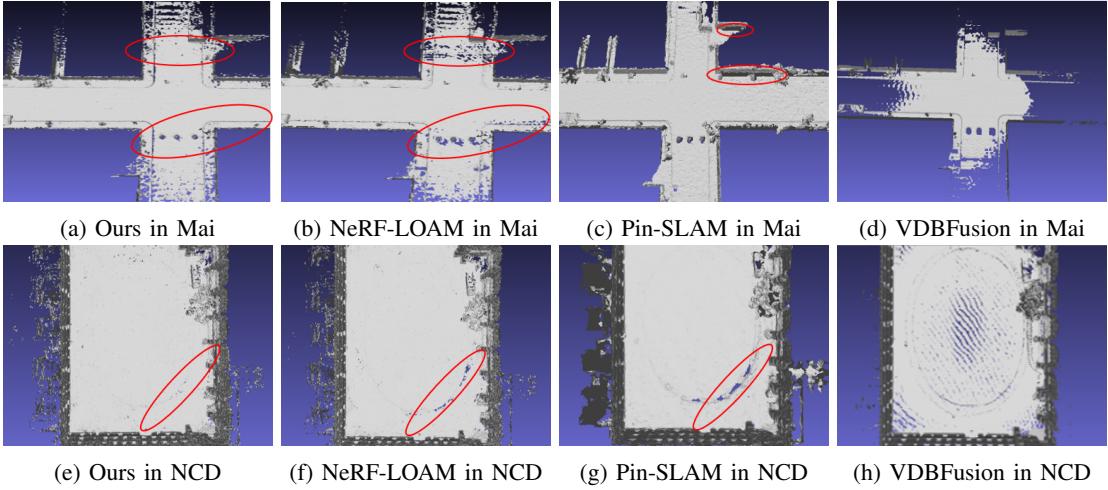


Fig. 5: Qualitative visualization of the map quality on the MaiCity dataset (Mai) (top) and Newer College dataset (NCD) (down). The areas highlighted by the red ellipses emphasize the contributions of our proposed method.

TABLE II: RMSE results of odometry. Mai for Maicity, NCD for NewerCollege. MOT for KITTI MOT chanlange dataset.

Method	Mai	NCD	MOT19	MOT20	MOT1-18
Pin-SLAM	0.07	0.09	3.41	18.72	9.53
NeRF-LOAM	0.19	0.14	<u>3.32</u>	<u>12.33</u>	<u>9.36</u>
Ours	<u>0.16</u>	0.01	3.29	10.81	8.83

objects, as evidenced by the clarity in the red-oval regions of the grayscale Fig.4 (c) and (f). While NeRF-LOAM [8] and Pin-SLAM [10] struggle to reconstruct the street that has moving objects, our proposed method maintains a more faithful representation of outdoor environment.

The quantitative evaluation results are shown in Table I. Although NeRF-LOAM [8] achieves high reconstruction accuracy by separating ground points from non-ground points, its completeness is relatively low. In comparison, our method achieves higher completeness without compromising accuracy, attributed to the multi-resolution octree implementation and Fourier positional encoding, which allow the network to learn higher-dimensional information. When using ground truth poses, VDBFusion [25] is reconstruction results are less complete compared to neural implicit methods due to its what-you-see-is-what-you-get storage characteristics. SHINE [29] improves reconstruction quality through regularization and optimization strategies, but our hybrid encoding method still achieves more competitive results. Fig.5 provides a qualitative analysis of the reconstructed maps. Due to the explicit map representation, VDBFusion [25] is unable to reasonably complete unobserved areas, resulting in low completeness. Pin-SLAM [10] are neural points that can theoretically predict the SDF value at any location, but they are prone to generating unrealistic artifacts.

C. Odometry Results

Dynamic scenes pose significant challenges to odometry pose estimation. As shown in Table II, our proposed method

TABLE III: Ablation study of the proposed method. We present the quantitative results on the MaiCity01 sequence.

	w/o FFE	w/o MF	Full
Comp.[cm] \downarrow	9.59	6.00	5.71
Acc.[cm] \downarrow	4.82	4.49	4.12
C-11[cm] \downarrow	7.21	5.00	4.91
F-score[%] \uparrow	85.39	85.72	88.06

achieved competitive pose estimation accuracy especially in highly dynamic scenarios (MOT dataset) and NewerCollege.

D. Ablation Study

We conducted ablation experiments to demonstrate the relative performance of our proposed method.

As shown in Fig.4, our dynamic object removal module effectively removes dynamic points and fills the gaps in the scene after removing the dynamic points. Moreover, Table III presents the quantitative evaluation of our method with different encoding strategies on Maicity [26]. The completeness of the reconstruction results is poor without the Fourier feature positional encoding. The combination of multi-resolution octree structure and Fourier feature positional encoding achieves the overall best performance.

V. CONCLUSIONS

In this paper, we present a 3D scene reconstruction system for dynamic outdoor environments, extending NeRF-LOAM [8]. Our method excels in handling dynamic foregrounds by reconstructing only the ground surface while accurately modeling static backgrounds.

We use dynamic foreground masking and a novel ground height estimation method to ensure realistic reconstruction despite moving objects. The integration of a multi-resolution octree with Fourier feature positional encoding optimizes memory use and maintains scene quality.

REFERENCES

- [1] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng, “Nerf: Representing scenes as neural radiance fields for view synthesis,” *Communications of the ACM*, vol. 65, no. 1, pp. 99–106, 2021.
- [2] E. Sucar, S. Liu, J. Ortiz, and A. J. Davison, “imap: Implicit mapping and positioning in real-time,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 6229–6238.
- [3] Z. Zhu, S. Peng, V. Larsson, W. Xu, H. Bao, Z. Cui, M. R. Oswald, and M. Pollefeys, “Nice-slam: Neural implicit scalable encoding for slam,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 12 786–12 796.
- [4] A. Rosinol, J. J. Leonard, and L. Carlone, “Nerf-slam: Real-time dense monocular slam with neural radiance fields,” in *2023 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2023, pp. 3437–3444.
- [5] H. Wang, J. Wang, and L. Agapito, “Co-slam: Joint coordinate and sparse parametric encodings for neural real-time slam,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 13 293–13 302.
- [6] Y. Zhang, F. Tosi, S. Mattoccia, and M. Poggi, “Go-slam: Global optimization for consistent 3d instant reconstruction,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 3727–3737.
- [7] Y. Mao, X. Yu, K. Wang, Y. Wang, R. Xiong, and Y. Liao, “Ngel-slam: Neural implicit representation-based global consistent low-latency slam system,” *arXiv preprint arXiv:2311.09525*, 2023.
- [8] J. Deng, Q. Wu, X. Chen, S. Xia, Z. Sun, G. Liu, W. Yu, and L. Pei, “Nerf-loam: Neural implicit representation for large-scale incremental lidar odometry and mapping,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 8218–8227.
- [9] S. Isaacson, P.-C. Kung, M. Ramanagopal, R. Vasudevan, and K. A. Skinner, “Loner: Lidar only neural representations for real-time slam,” *IEEE Robotics and Automation Letters*, 2023.
- [10] Y. Pan, X. Zhong, L. Wiesmann, T. Posewsky, J. Behley, and C. Stachniss, “Pin-slam: Lidar slam using a point-based implicit neural representation for achieving global map consistency,” *arXiv preprint arXiv:2401.09101*, 2024.
- [11] A. Carlson, M. S. Ramanagopal, N. Tseng, M. Johnson-Roberson, R. Vasudevan, and K. A. Skinner, “Cloner: Camera-lidar fusion for occupancy grid-aided neural representations,” *IEEE Robotics and Automation Letters*, vol. 8, no. 5, pp. 2812–2819, 2023.
- [12] Y. Tao, Y. Bhalgat, L. F. T. Fu, M. Mattamala, N. Chebrolu, and M. Fallon, “Silvr: Scalable lidar-visual reconstruction with neural radiance fields for robotic inspection,” *arXiv preprint arXiv:2403.06877*, 2024.
- [13] T. Takikawa, J. Litalien, K. Yin, K. Kreis, C. Loop, D. Nowrouzezahrai, A. Jacobson, M. McGuire, and S. Fidler, “Neural geometric level of detail: Real-time rendering with implicit 3d shapes,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 11 358–11 367.
- [14] M. Tancik, P. Srinivasan, B. Mildenhall, S. Fridovich-Keil, N. Raghavan, U. Singhal, R. Ramamoorthi, J. Barron, and R. Ng, “Fourier features let networks learn high frequency functions in low dimensional domains,” *Advances in neural information processing systems*, vol. 33, pp. 7537–7547, 2020.
- [15] Q. Li, S. Chen, C. Wang, X. Li, C. Wen, M. Cheng, and J. Li, “Lo-net: Deep real-time lidar odometry,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 8473–8482.
- [16] W. Wang, M. R. U. Saputra, P. Zhao, P. Gusmao, B. Yang, C. Chen, A. Markham, and N. Trigoni, “Deeppco: End-to-end point cloud odometry through deep parallel neural network,” in *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2019, pp. 3248–3254.
- [17] G. Wang, X. Wu, Z. Liu, and H. Wang, “Pwclo-net: Deep lidar odometry in 3d point clouds using hierarchical embedding mask optimization,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 15 910–15 919.
- [18] R. Mur-Artal, J. M. M. Montiel, and J. D. Tardos, “Orb-slam: a versatile and accurate monocular slam system,” *IEEE transactions on robotics*, vol. 31, no. 5, pp. 1147–1163, 2015.
- [19] R. Mur-Artal and J. D. Tardós, “Orb-slam2: An open-source slam system for monocular, stereo, and rgbd cameras,” *IEEE transactions on robotics*, vol. 33, no. 5, pp. 1255–1262, 2017.
- [20] B. Bescos, J. M. Fácil, J. Civera, and J. Neira, “Dynaslam: Tracking, mapping, and inpainting in dynamic scenes,” *IEEE Robotics and Automation Letters*, vol. 3, no. 4, pp. 4076–4083, 2018.
- [21] B. Bescos, C. Campos, J. D. Tardós, and J. Neira, “Dynaslam ii: Tightly-coupled multi-object tracking and slam,” *IEEE robotics and automation letters*, vol. 6, no. 3, pp. 5191–5198, 2021.
- [22] X. Lang, L. Li, H. Zhang, F. Xiong, M. Xu, Y. Liu, X. Zuo, and J. Lv, “Gaussian-lic: Photo-realistic lidar-inertial-camera slam with 3d gaussian splatting,” *arXiv preprint arXiv:2404.06926*, 2024.
- [23] J. J. Park, P. Florence, J. Straub, R. Newcombe, and S. Lovegrove, “DeepSDF: Learning continuous signed distance functions for shape representation,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 165–174.
- [24] L. Mescheder, M. Oechsle, M. Niemeyer, S. Nowozin, and A. Geiger, “Occupancy networks: Learning 3d reconstruction in function space,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 4460–4470.
- [25] I. Vizzo, T. Guadagnino, J. Behley, and C. Stachniss, “Vdbfusion: Flexible and efficient tsdf integration of range sensor data,” *Sensors*, vol. 22, no. 3, p. 1296, 2022.
- [26] I. Vizzo, X. Chen, N. Chebrolu, J. Behley, and C. Stachniss, “Poisson surface reconstruction for lidar odometry and mapping,” in *2021 IEEE international conference on robotics and automation (ICRA)*. IEEE, 2021, pp. 5624–5630.
- [27] M. Niemeyer, L. Mescheder, M. Oechsle, and A. Geiger, “Differentiable volumetric rendering: Learning implicit 3d representations without 3d supervision,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 3504–3515.
- [28] V. Sitzmann, M. Zollhöfer, and G. Wetzstein, “Scene representation networks: Continuous 3d-structure-aware neural scene representations,” *Advances in Neural Information Processing Systems*, vol. 32, 2019.
- [29] X. Zhong, Y. Pan, J. Behley, and C. Stachniss, “Shine-mapping: Large-scale 3d mapping using sparse hierarchical implicit neural representations,” in *2023 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2023, pp. 8371–8377.
- [30] M. M. Johari, C. Carta, and F. Fleuret, “Eslam: Efficient dense slam system based on hybrid representation of signed distance fields,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 17 408–17 419.
- [31] D. Azinović, R. Martin-Brualla, D. B. Goldman, M. Nießner, and J. Thies, “Neural rgbd surface reconstruction,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 6290–6301.
- [32] X. Yang, H. Li, H. Zhai, Y. Ming, Y. Liu, and G. Zhang, “Vox-fusion: Dense tracking and mapping with voxel-based neural implicit representation,” in *2022 IEEE International Symposium on Mixed and Augmented Reality (ISMAR)*. IEEE, 2022, pp. 499–507.
- [33] S. Sun, M. Mielle, A. J. Lilienthal, and M. Magnusson, “3qfp: Efficient neural implicit surface reconstruction using tri-quadtrees and fourier feature positional encoding,” *arXiv preprint arXiv:2401.07164*, 2024.
- [34] A. Geiger, P. Lenz, and R. Urtasun, “Are we ready for autonomous driving? the kitti vision benchmark suite,” in *2012 IEEE conference on computer vision and pattern recognition*. IEEE, 2012, pp. 3354–3361.
- [35] M. Ramezani, Y. Wang, M. Camurri, D. Wisth, M. Mattamala, and M. Fallon, “The newer college dataset: Handheld lidar, inertial and vision with ground truth,” in *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2020, pp. 4353–4360.
- [36] H.-N. Hu, Y.-H. Yang, T. Fischer, T. Darrell, F. Yu, and M. Sun, “Monocular quasi-dense 3d object tracking,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 45, no. 2, pp. 1992–2008, 2022.