

Neural Surface Reconstruction and Rendering for LiDAR-Visual Systems

Jianheng Liu, Chunran Zheng, Yunfei Wan, Bowen Wang, Yixi Cai and Fu Zhang

Abstract—This paper presents a unified surface reconstruction and rendering framework for LiDAR-visual systems, integrating Neural Radiance Fields (NeRF) and Neural Distance Fields (NDF) to recover both appearance and structural information from posed images and point clouds. We address the structural visible gap between NeRF and NDF by utilizing a visible-aware occupancy map to classify space into the free, occupied, visible unknown, and background regions. This classification facilitates the recovery of a complete appearance and structure of the scene. We unify the training of the NDF and NeRF using a spatial-varying scale SDF-to-density transformation for levels of detail for both structure and appearance. The proposed method leverages the learned NDF for structure-aware NeRF training by an adaptive sphere tracing sampling strategy for accurate structure rendering. In return, NeRF further refines structural in recovering missing or fuzzy structures in the NDF. Extensive experiments demonstrate the superior quality and versatility of the proposed method across various scenarios. To benefit the community, the codes will be released at <https://github.com/hku-mars/M2Mapping>.

I. INTRODUCTION

3D reconstruction and novel-view synthesis (NVS) are fundamental tasks in computer vision and robotics [1], [2]. The widespread availability of cameras and low-cost LiDAR sensors on robots provides rich multimodal data for 3D reconstruction and NVS. LiDAR-visual SLAM technology are widely used in 3D reconstruction tasks like LVI-SAM [3], R3LIVE [4], and FAST-LIVO [5]. However, these methods can only obtain colorized raw point maps, so the map resolution, density and accuracy are fundamentally limited by the LiDAR sensors. In practical digital-twin applications [6], high-quality watertight surface reconstructions along with photorealistic rendering are crucially important.

Recovering high-quality surface reconstructions or photorealistic rendering in the wild presents significant challenges. Explicit surface reconstruction methods, such as direct meshing [7], often perform poorly for noisy or misaligned LiDAR point cloud data and lead to artifacts. In contrast, implicit surface reconstruction methods estimate implicit functions, such as Poisson functions [8] or signed distance fields (SDFs) [9], [10], which define watertight manifold surfaces at their zero-level sets, offering a more reliable approach to surface reconstruction. The neural distance field (NDF) [11]–[13] enhanced by neural networks, further improves the continuity of distance fields through gradient regularization [14], demonstrating significant potential for capturing high-granularity details in surface reconstruction.

For scenes with accurate geometry, surface rendering can efficiently render solid objects using only the intersection points of pixel rays and surfaces. However, obtaining high-

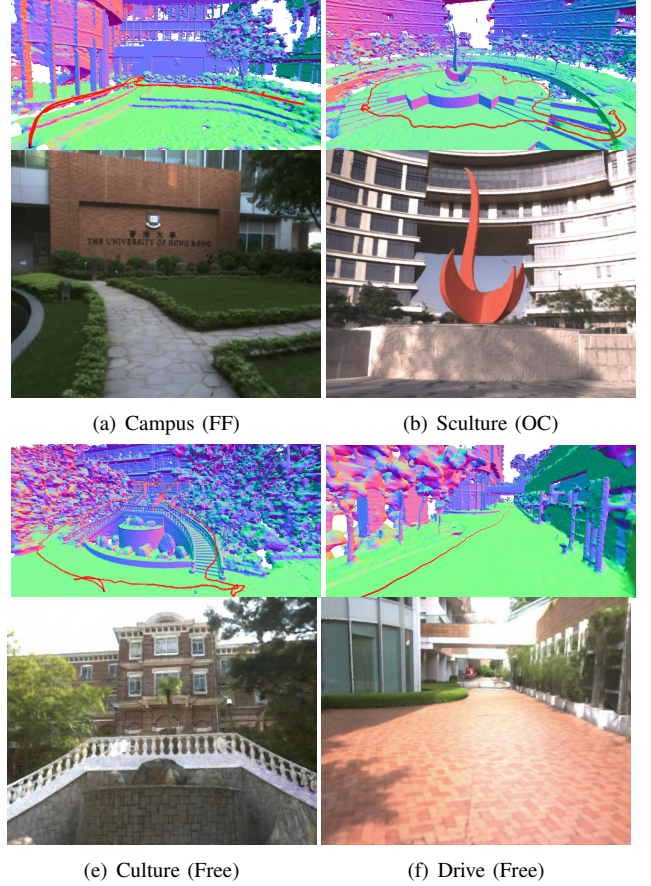


Fig. 1. We show our surface reconstruction results (color indicates the normal direction) and rendering results collected by real-world LiDAR-visual sensor systems under different trajectories (red lines): forward-facing (FF), object-centric (OC), and free-view (Free).

quality geometry is often challenging in practice, and the rendering quality is heavily influenced by the accuracy of the underlying geometry [7]. In cases where the geometry is incomplete or imprecise (e.g., due to the sparse and noisy LiDAR point measurements), surface rendering methods could degrade considerably, making volumetric rendering methods more appropriate [15]. Neural Radiance Fields (NeRFs) have emerged as a powerful 3D representation method for photorealistic novel-view synthesis and geometry reconstruction. NeRFs encode a scene’s volumetric density and appearance as a function of 3D spatial coordinates and viewing direction, generating high-quality images via volume rendering. However, the ambiguous density field of NeRF may exhibit inaccuracies and visual artifacts when extrapolating views. Recent works [16], [17] demonstrate that detailed implicit surfaces can be learned via NeRF by

converting NDFs to density fields. However, this approach requires multi-view images and extensive training for each location in construction, making it less suitable for free-view trajectories, where the region of interest around the trajectory often suffers from limited multi-view images. A low-cost LiDAR, which provides direct sample points on the surface, offers a promising solution to reduce the reliance on extensive multi-view images and enable free-view trajectories, as shown in Fig. 1.

This work aims to recover both the appearance and structural information of unbounded scenes using posed images and low-cost LiDAR data from any casual trajectories, which are typical situations in practical 3D reconstruction data collection. We propose a unified surface reconstruction and rendering framework for LiDAR-visual systems, capable of recovering both the NeRF and NDF. This approach combines NDF and NeRF to provide better structure and appearance for scenes. We directly apply signed ray distance in 3D space to efficiently approximate an NDF. Photorealistic images are rendered by integrating the density field from the NDF with the radiance field. The NDF enforces geometric consistency within NeRF and guides NeRF's samples to concentrate on surfaces using sphere tracing. NeRF, in turn, leverages photometric errors from rendering to refine and complete fuzzy or incomplete structures within the NDF. The main contributions of this paper are as follows:

- 1) An open-source neural rendering framework for LiDAR-visual systems achieving high-granularity surface reconstructions and photorealistic rendering using a neural signed distance field and neural radiance field.
- 2) A visible-aware prior occupancy map constructed from point clouds for efficient and complete rendering of unbounded scenes.
- 3) A neural signed distance field with spatial-varying scale leverages sphere tracing to enable structure-aware volume rendering for neural radiance field.

We perform extensive experiments to validate our insights and demonstrate the superiority of the proposed method across various scenarios for general types of trajectories.

II. RELATED WORKS

Neural Radiance Fields (NeRF) [18] is a groundbreaking method for novel-view synthesis that employs a neural network to model a scene's volumetric density and radiance. However, NeRF encounters significant challenges when representing surfaces, as it lacks a defined density level set for surface representations. Recent works [16], [17], [19] employ SDF as the structure field in volume rendering, enabling surface reconstructions from multi-view images. Recovering surfaces from images is resource-intensive while doing so from point clouds is much easier. Surface reconstructions from point clouds are predominantly achieved using implicit representations like Poisson functions [20] or signed distance fields [9], [10]. Recent advancements in neural implicit representations [12], [13], [21] model SDFs using neural networks, providing high-granularity and complete surface representations.

Point clouds, which provide direct structural information, have been effectively used to constrain the structure field in NeRF by aligning the rendering depth with depth measurements [22]. RGBD cameras that provide aligned color and depth information are well practiced in neural radiance fields [23], [24] but limited in indoor scene settings. For outdoor scenes, LiDAR providing more accurate measurements is more favorable than RGBD cameras and extends NeRF's applicability from room-scale to urban-scale scenes [?]. To address the sparsity of LiDAR point clouds, existing works [25], [26] use networks to densify LiDAR into depth images for NeRF training. While these above works still limit LiDAR usage for deep supervision in 2D image space to regularize the structure field. In this work, we directly apply structure constraints in 3D space by forming a continuous NDF from point clouds, and conduct a structure-aware NeRF rendering based on the NDF to avoid the limitations posed by LiDAR sparsity.

Unlike NeRF, 3D Gaussian Splatting (3DGS) [27] explicitly stores appearance within Gaussian points initialized by structure-from-motion points. Its high efficiency and high-quality rendering establish it as a new benchmark in novel-view synthesis. However, the discrete representation of 3DGS makes it difficult to represent manifold surfaces, which fails to apply direct structure regularization in 3D space. Therefore, depth supervision [28] is more commonly used to regularize 3DGS' geometry. LiDAR point clouds are well-suit for initializing 3DGS [29], but 3DGS lacks a viable pipeline to refine the structure in reverse. Therein, we adhere to the volume rendering of SDF pipeline [16] to enable surface and appearance reconstructions within a unified framework and leverage both point clouds and images to supervise the structure field.

III. METHODOLOGY

Given a series of posed images and LiDAR point clouds, we aim to recover both the surface and the appearance of the scene. The overall pipeline is shown in Fig. 2. We first leverage ray casting to construct a visible-aware occupancy map from the LiDAR point cloud and images to classify the space where is need to be encoded (Sec. III-A). We utilize LiDAR point clouds to supervise an NDF to recover the structure of the scene (Sec. III-B). The NDF serves as the structure field for NeRF rendering, and the NeRF reversely refines the NDF through multi-view photometric errors (Sec. III-C). In this section, we elaborate on how our approach enables both the NDF and the NeRF to be aware of scenes' structures (Sec. III-C.1 and Sec. I-A) and the training strategies to deal with common issues shown in real-world applications (Sec. I-C).

A. Visible-aware Occupancy Map

We would like to enable NeRF training for generalized LiDAR-visual systems. LiDAR point clouds provide strong and accurate information of the environment, such as occupied, free, and unknown spaces. An occupancy map can be easily attained through ray casting to represent such

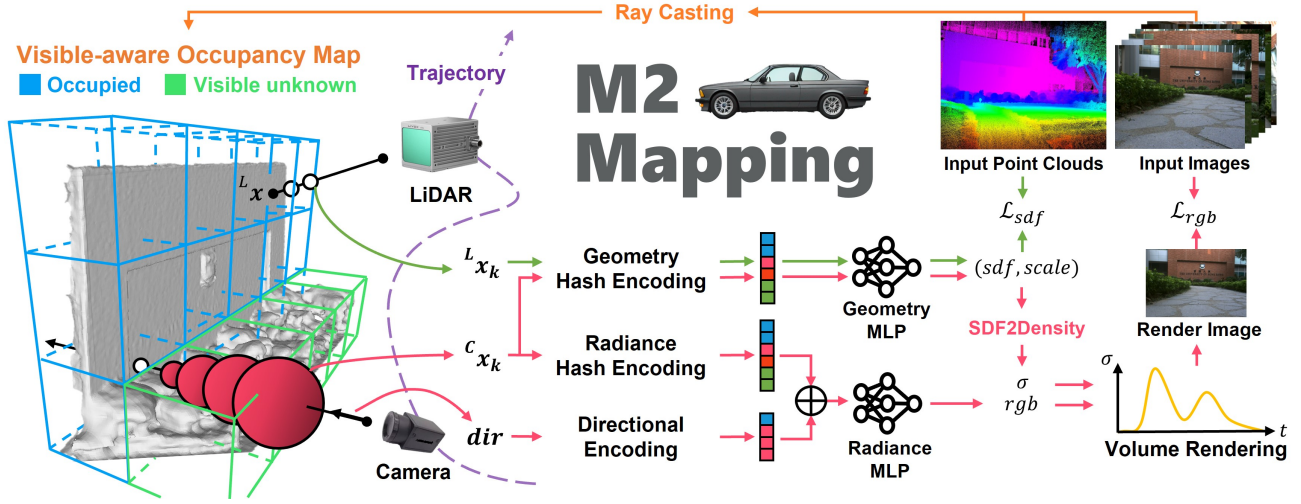


Fig. 2. The overall pipeline of the proposed multimodal mapping framework, named M2Mapping. Given a series of posed images and LiDAR point clouds, we first construct the visible-aware occupancy map via ray casting. The neural distance field is trained using the ray distance value from the point cloud. The neural radiance field guides the structure-aware sampling process of the neural radiance field and predicts the density of each point. The sample points and direction are encoded as features, and the MLP forwards the concatenated features to infer color. The volume rendering accumulates densities and color for novel-view synthesis. (\oplus denotes the concatenation operation)

information [30], which can serve as an auxiliary acceleration data structure for volume rendering to skip free space. However, due to the different field of views between LiDAR and cameras and the sparsity of the LiDAR point cloud, much space is visible to the camera but not to the LiDAR. Only using occupied information from LiDAR for NeRF training may lead to incomplete rendering results, therefore, we need to identify the space in need to be encoded for efficiency and complete rendering. To address this, we adopt an occupancy map and categorize the occupancy grid states as free, occupied, visible unknown, and invisible unknown, where the visible unknown and invisible unknown grids are separated from the unknown grids based on whether they are visible in the input image. As illustrated in Fig. 3, we cast images pixels’ rays in the built occupancy map and mark the traversed unknown grids before the first arrived occupied grids or map boundary as the visible unknown grids. We encode the neural fields within a visible-aware occupancy map that only includes the occupied and visible unknown grids to reduce the training workload.

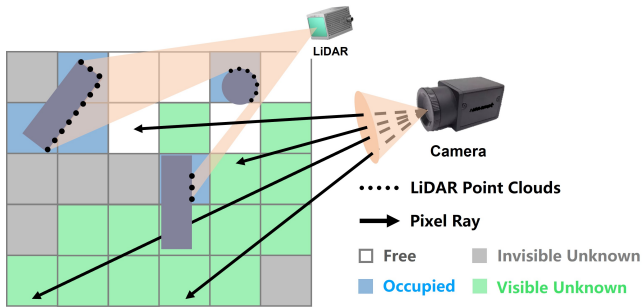


Fig. 3. Illustrations of our visible-aware occupancy map. The LiDAR observations derive a standard occupancy grid map (OGM) with states: free, occupied, and unknown. We further classify the unknown states into visible unknown and invisible unknown via images’ pixel ray casting. Note that the visible unknown grids arise from unknown grids which is visible from the camera without any occupied grid in between.

B. Geometry supervision

The neural signed distance field represents a scene as a distance field that maps 3D points $\mathbf{x} \in \mathbb{R}^3$ to signed distance value $s \in \mathbb{R}$. We leverage hash encoding [31] and tiny multi-layer perceptions (MLPs) to form a geometry neural network to encode a scalable signed distance field: $(\hat{s}, \hat{\beta}) = f_S(\mathbf{x})$, where $\hat{\beta} \in \mathbb{R}$ is our proposed spatial-varying scale (Sec. III-C.1). Most existing volume rendering methods for SDFs [16], [17] train the NDF via pure images. In contrast, we supervise the NDF using direct 3D LiDAR points.

Given a LiDAR point $\mathbf{L}\mathbf{x} = \mathbf{L}\mathbf{o} + t \mathbf{L}\mathbf{d}$, where $\mathbf{L}\mathbf{o}$ is the LiDAR origin and $\mathbf{L}\mathbf{d}$ is the LiDAR point’s ray direction, we uniformly sample points from the LiDAR origin to the point along the ray. For one sample point $\mathbf{L}\mathbf{x}_k = \mathbf{L}\mathbf{o} + t_k \mathbf{L}\mathbf{d}$, its ray distance value $s_k = t - t_k$ as the supervision ground truth value is transformed to the occupancy $o_k = \Phi(-s_k, \hat{\beta}_k)$, where $\Phi(v, h) = (1 + \exp(-v/h))^{-1}$ is the sigmoid function, and $\hat{\beta}_k$ is the predicted scale factor, which indicates the standard deviation in the Logistic function controlling the sharpness of the transition. The binary cross-entropy loss [13] is employed to approximate a signed distance field:

$$\mathcal{L}_{sdf} = -\frac{1}{N_L} \sum_k o_i \log(\hat{o}_k) + (1 - o_k) \log(1 - \hat{o}_k), \quad (1)$$

where $\hat{o}_k = \Phi(-\hat{s}_k, \hat{\beta}_k)$ is the geometry neural network predicted occupancy value, $(\hat{s}_k, \hat{\beta}_k) = f_S(\mathbf{L}\mathbf{x}_k)$, and N_L is total number of points sampled on a ray.

C. Appearance supervision

The neural radiance field maps 3D points $\mathbf{x} \in \mathbb{R}^3$ and viewing directions $\mathbf{d} \in \mathbb{R}^3$ to view-dependent colors $\mathbf{c} \in \mathbb{R}^3$. We leverage hash encoding [31] for multi-resolution feature encoding and spherical harmonics encoding [32] for directional encoding, and these two encodings are concatenated and fed into a radiance MLP to predict view-dependent color: $\hat{\mathbf{c}}(\mathbf{x}, \mathbf{d}) = f_C(\mathbf{x}, \mathbf{d})$.

Given a pixel ray direction ${}^C\mathbf{d}$, we sample points $\{\mathbf{x}_k = {}^C\mathbf{o} + t_k {}^C\mathbf{d}\}$ from the camera origin ${}^C\mathbf{o}$ along the ray with our proposed structure-aware sampling (Sec. I-A) and the volume rendering gets a pixel color $\hat{\mathbf{C}}$ via blending the samples' colors and densities:

$$\hat{\mathbf{C}} = \sum_k \exp\left(-\sum_{j < k} \hat{\sigma}_j \delta_j\right) (1 - \exp(-\hat{\sigma}_k \delta_k)) \hat{\mathbf{c}}_k, \quad (2)$$

where $\delta_k = t_{k+1} - t_k$ is the distance between two adjacent samples along the ray, and the density $\hat{\sigma}_k$ is transformed from the signed distance value $(\hat{s}_k, \hat{\beta}_k) = f_S(\mathbf{x}_k)$ to enable volume rendering of SDF [33]:

$$\hat{\sigma}_k = \max\left(-\frac{\Phi(-\hat{s}_k, \hat{\beta}_k)}{\hat{\beta}_k} \frac{\partial \hat{s}_k}{\partial {}^C\mathbf{x}_k} \cdot {}^C\mathbf{d}, 0\right), \quad (3)$$

where $\hat{M}_k = \frac{\partial \hat{s}_k}{\partial {}^C\mathbf{x}_k} \cdot {}^C\mathbf{d}$ is the gradient of the signed distance value along the ray direction, so-called slope. The volume rendering formulation (Eq. 2) can be further abbreviated to $\hat{\mathbf{C}} = \sum_k T_k \alpha_k \hat{\mathbf{c}}_k$, where $\alpha_k = 1 - \exp(-\hat{\sigma}_k \delta_k)$ is the opacity and $T_k = \prod_{j < k} (\exp(-\hat{\sigma}_j \delta_j))$ is the transmittance.

To encode the scene's appearance, a photometric loss between the source pixel color \mathbf{C}_i and rendered pixel color $\hat{\mathbf{C}}_i$ is applied:

$$\mathcal{L}_{rgb} = \frac{1}{N_C} \sum_i \|\mathbf{C}_i - \hat{\mathbf{C}}_i\|^2, \quad (4)$$

N_C is total number of pixels sampled during the training. The photometric error indirectly adjusts the neural distance field by backpropagation via volume rendering (Eq. 2) and SDF-to-density transformation (Eq. 3).

1) *Spatial-varying Scale*: The sigmoid function $\Phi(\cdot)$ introduces a scale factor β to control both the sharpness of the NDF and NeRF, where the smaller scale indicates the sharper transition. In the past works, the scale factor is treated as a global scale either using truncated distance [24], a single learnable scale [17] or calculated scale [16], which is not well adaptive to various granularity structures, such as tiny or fuzzy structures. We extend the NDF to predict both the SDF values and scales at a location: $(\hat{s}, \hat{\beta}) = f_S(\mathbf{x})$. The spatial-varying learnable scale turns the neural signed distance field into a stochastic signed distance field, where the scale shows the confidence of the neural networks' prediction. For geometry supervision, a smaller scale makes the structure loss (Eq. 1) more sensitive to noise, which enables the networks to adaptively adjust the sharpness for different granularities geometry. For appearance supervision, a smaller scale returns a higher density (Eq. 3), which emphasizes the surface points to act like surface rendering. In reverse, a larger scale considers more points, which is more suitable for fuzzy structures with high uncertainty in geometry.

2) *Structure-aware sampling*: The volume rendering sampling focuses on identifying the most critical areas for final rendering. While the photometric-oriented NeRF training often renders artifacts in the air that overfit the images, we would like to conduct an accurate structure rendering based

on our learned NDF. Sphere tracing [34] is designed for finding the surface of a signed distance field. This method naturally distributes more samples near a surface and fewer samples when is away. Therefore, we propose to treat each step point in the sphere tracing as a sample that makes the sampling aware to structure information from the NDF for a geometry-consistent rendering, as shown in Fig. 3. The detailed algorithm of structure-aware sampling could be found in the supplementary materials¹.

D. Training

We further regularize the NDF with Eikonal [14] and curvature [35] loss, and the details can be found in supplementary materials. The overall training loss is defined as follows:

$$\mathcal{L} = \mathcal{L}_{sdf} + \lambda_{rgb} \mathcal{L}_{rgb} + \lambda_{eik} \mathcal{L}_{eik} + \lambda_{curv} \mathcal{L}_{curv}, \quad (5)$$

where λ_{rgb} , λ_{eik} and λ_{curv} are the weights for the photometric, Eikonal, and curvature losses, respectively.

IV. EXPERIMENTS

In this section, we conduct extensive experiments to evaluate our proposed system. Due to space limitations, the implementation details and part of the experiments including ablation study are illustrated in the supplementary materials.

A. Experiments Settings

1) *Baselines*: We compare our method with other state-of-the-art approaches. For structure-oriented methods, we include the voxel-based method VDBFusion [10], and the neural-network-based methods iSDF [12] and SHINE-Mapping [13]. For appearance-oriented methods, we include the NeRF-based method InstantNGP [31]. For depth-aided appearance-oriented methods, we include 3D Gaussian Splatting [27] initialized with dense point clouds, denoted as 3DGS[†], and RGBD-based methods H2Mapping [24] and MonoGS [28].

2) *Metrics*: To evaluate the quality of the geometry, we recover the implicit model to a triangular mesh using marching cubes [36] and calculate the Chamfer Distance (C-L1, cm) and F-Score (< 2 cm, %) with the ground truth mesh. For rendering quality evaluation, we use the Structural Similarity Index (SSIM), Peak Signal-to-Noise Ratio (PSNR), and Learned Perceptual Image Patch Similarity (LPIPS) to compare the rendered image with the ground truth image.

B. Replica Dataset

The Replica dataset [6] provides room-scale indoor simulation RGBD sensor data. We emphasize the issues of extrapolation rendering consistency by uniformly sampling positions and orientations in each scene to generate the extrapolation dataset from Replica. The compared surface reconstruction, interpolation, and extrapolation rendering results are shown in Tab. I. As illustrated in Fig. 4 (Left), we capture more precise geometric details in slim objects while

¹The supplementary materials could be found in <https://github.com/hku-mars/M2Mapping/blob/main/supplement.pdf>

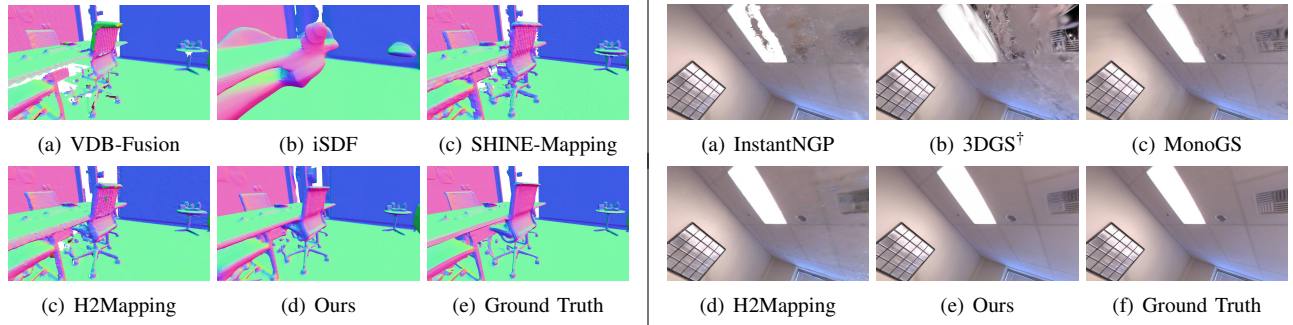


Fig. 4. (Left) Reconstructed mesh on the Replica dataset’s office-2 and colors indicate the direction of the surface normal. (Right) Extrapolation rendering results on the Replica dataset’s room-0. Our method can capture more precise geometric details in slim objects and retain both the structure and texture in extrapolation rendering.

TABLE I

QUANTITATIVE SURFACE RECONSTRUCTION, INTERPOLATION(I) AND EXTRAPOLATION(E) RENDERING RESULTS ON THE REPLICA DATASET.

Metrics	Methods	Office-0	Office-1	Office-2	Office-3	Office-4	Room-0	Room-1	Room-2	Avg.
C-L1[cm]↓	VDBFusion	0.618	0.595	0.627	0.685	0.633	0.637	0.558	0.656	0.626
	iSDF	2.286	4.032	3.144	2.352	2.120	1.760	1.712	2.604	2.501
	SHINE-Mapping	0.753	0.663	0.851	0.965	0.770	0.650	0.844	0.826	0.790
	H2Mapping	0.557	0.529	0.585	0.644	0.616	0.568	0.523	0.616	0.580
	Ours	0.494	0.476	0.501	0.517	0.531	0.486	0.455	0.532	0.499
F-Score[%]↑	VDBFusion	97.129	97.107	97.007	96.221	96.993	97.871	98.311	96.117	97.095
	iSDF	75.829	55.429	71.232	77.274	81.281	78.253	84.801	76.815	75.114
	SHINE-Mapping	94.397	95.606	92.150	87.855	94.427	97.010	91.307	92.601	93.169
	H2Mapping	98.430	98.279	97.935	97.077	97.391	98.850	99.003	96.904	97.983
	Ours	98.970	98.771	98.657	98.415	98.158	99.238	99.496	97.677	98.673
SSIM(I)↑	InstantNGP	0.981	0.980	0.963	0.960	0.966	0.964	0.964	0.967	0.968
	3DGS [†]	0.987	0.980	0.980	0.976	0.980	0.977	0.982	0.980	0.980
	H2Mapping	0.963	0.960	0.931	0.929	0.941	0.914	0.929	0.935	0.938
	MonoGS	0.985	0.982	0.973	0.970	0.976	0.971	0.976	0.977	0.975
	Ours	0.985	0.983	0.973	0.970	0.977	0.968	0.975	0.977	0.976
PSNR(I)↑	InstantNGP	43.538	44.340	38.273	37.603	39.772	37.926	38.859	39.568	39.984
	3DGS [†]	43.932	43.237	39.182	38.564	41.366	39.081	41.288	41.431	41.010
	H2Mapping	38.307	38.705	32.748	33.021	34.308	31.660	33.466	32.809	34.378
	MonoGS	43.648	43.690	37.695	37.539	40.224	37.779	39.563	40.134	40.034
	Ours	44.369	44.935	39.652	38.874	41.318	38.541	40.775	40.705	41.146
LPIPS(I)↓	InstantNGP	0.046	0.069	0.096	0.087	0.089	0.079	0.094	0.094	0.082
	3DGS [†]	0.040	0.075	0.069	0.068	0.065	0.060	0.057	0.075	0.064
	H2Mapping	0.109	0.163	0.186	0.170	0.154	0.177	0.167	0.177	0.163
	MonoGS	0.031	0.048	0.061	0.057	0.052	0.055	0.050	0.056	0.051
	Ours	0.031	0.044	0.055	0.050	0.050	0.062	0.056	0.053	0.050
SSIM(E)↑	InstantNGP	0.972	0.961	0.934	0.938	0.952	0.918	0.936	0.941	0.944
	3DGS [†]	0.936	0.897	0.924	0.917	0.925	0.881	0.915	0.919	0.914
	H2Mapping	0.957	0.955	0.932	0.925	0.937	0.866	0.916	0.917	0.926
	MonoGS	0.974	0.961	0.945	0.942	0.950	0.912	0.942	0.946	0.947
	Ours	0.980	0.976	0.960	0.964	0.970	0.955	0.963	0.965	0.967
PSNR(E)↑	InstantNGP	39.874	39.120	31.274	32.135	34.458	32.587	33.024	32.266	34.341
	3DGS [†]	31.220	29.959	27.411	26.442	28.324	27.541	28.429	27.139	28.307
	H2Mapping	36.740	37.841	31.427	31.144	31.988	28.815	31.192	30.603	32.468
	MonoGS	39.197	38.818	29.740	29.664	31.632	29.949	31.126	30.621	32.593
	Ours	41.965	42.215	35.056	37.465	38.667	36.427	37.294	36.722	38.226
LPIPS(E)↓	InstantNGP	0.085	0.117	0.176	0.185	0.149	0.180	0.162	0.156	0.151
	3DGS [†]	0.117	0.177	0.167	0.181	0.155	0.204	0.170	0.174	0.168
	H2Mapping	0.119	0.163	0.195	0.221	0.166	0.250	0.188	0.199	0.188
	MonoGS	0.061	0.104	0.136	0.137	0.108	0.170	0.123	0.117	0.120
	Ours	0.046	0.062	0.084	0.081	0.066	0.114	0.084	0.087	0.078

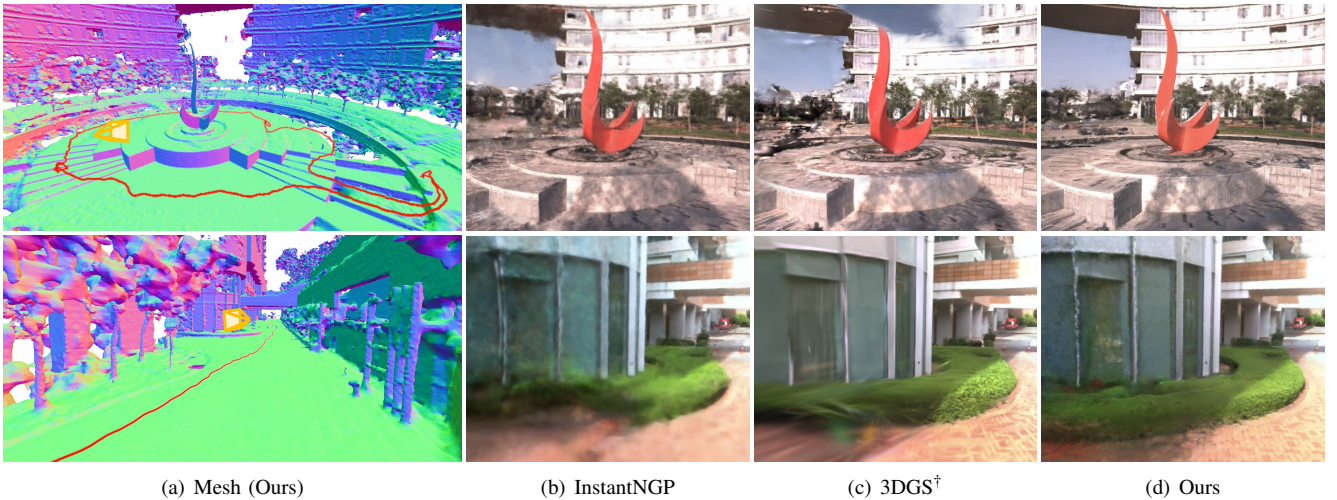


Fig. 5. The qualitative results of FAST-LIVO2 datasets (scenes from top to bottom are Sculpture and Drive). We show our surface reconstruction results on the left, and the red line indicates the training path and the orange cameras indicates the extrapolation views for the right side’s rendering results.

maintaining smoothness. For novel-view synthesis, the proposed structure-aware rendering method yields competitive results in interpolation rendering and achieves a significant improvement in extrapolation rendering compared to other methods. As depicted in Fig. 4 (Right), our method retains both the structure and texture in views of the ceiling, which were not seen with a vertical view in the training dataset.

TABLE II
QUANTITATIVE RESULTS ON THE FAST-LIVO2 DATASET.

Metrics	Methods	Campus	Sculpture	Culture	Drive	Avg.
SSIM \uparrow	InstantNGP	0.789	0.698	0.670	0.697	0.714
	3DGS †	0.849	0.769	0.726	0.778	0.780
	Ours	0.834	0.729	0.727	0.764	0.764
PSNR \uparrow	InstantNGP	28.880	22.356	21.563	24.145	24.236
	3DGS †	31.310	24.128	21.764	<u>25.837</u>	<u>25.760</u>
	Ours	<u>30.681</u>	23.453	24.695	25.941	26.193
LPIPS \downarrow	InstantNGP	0.255	0.376	0.428	0.416	0.369
	3DGS †	0.182	0.266	0.361	0.296	0.276
	Ours	<u>0.210</u>	<u>0.321</u>	0.350	0.293	<u>0.293</u>

C. FAST-LIVO2 Dataset

For real-world datasets, we evaluate three types of trajectory datasets collected with a camera and a LiDAR from the FAST-LIVO2 datasets [5]: forward-facing (Campus), object-centric (Sculpture), and free-view (Culture, Drive) trajectories, as shown in Fig. 1. We use the localization results from FAST-LIVO2 as ground truth poses and use a train-test split for dataset interpolation evaluation, where every 8th photo is used for the test set and the rest for the training set [27]. The quantitative results are shown in Tab. I. RGBD-based methods [24], [28] are not included as they are not compatible with LiDAR-visual systems. Our method demonstrates competitive rendering performance with the state-of-the-art method 3DGS [27] and outperforms others in complex free-view trajectory scenes. The interpolation qualitative results are shown in Fig. 3 (Left). Our method captures precise geometric details in fuzzy objects (leaf) and avoids overfitting occurring in InstantNGP and 3DGS †

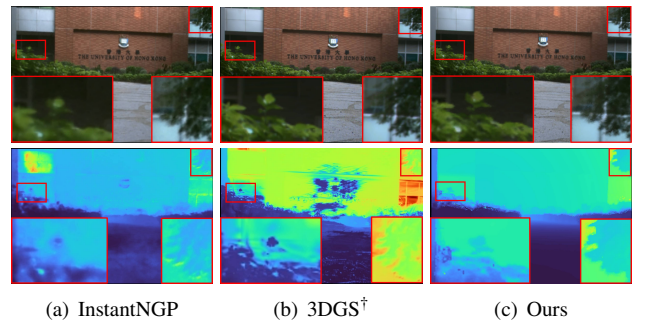


Fig. 6. Rendering results between different baselines on the FAST-LIVO2 dataset’s Campus scene.

(middle wall’s depth). We further present our complete and detailed surface reconstruction results and extrapolation rendering results on the FAST-LIVO2 datasets, as shown in Fig 1 and Fig. 2. In extrapolating views, all methods show some degree of degradation in rendering quality, where our method renders more consistent results than other methods, especially in the free-trajectory-type scene, where the camera moves freely, and lacks co-visible views.

V. CONCLUSION

This study presents a NeRF-based Surface Reconstruction and Rendering for LiDAR-visual systems. It bridges the gap between NDF and NeRF, harnessing the strengths of both to provide high-quality structure and appearance for a scene. The framework leverages structural information from LiDAR point clouds to build a visible-aware occupancy map prior for efficiency and a scalable NDF for high-granularity surface reconstruction. The NDF enables a structure-aware sampling for NeRF to conduct accurate structure rendering. The NeRF also utilizes photometric errors from rendering to refine structures. Extensive experiments validate the effectiveness of the proposed method across various scenarios, demonstrating its potential for real-world applications in computer vision and robotics.

REFERENCES

- [1] Y. Wang, Z. Su, N. Zhang, R. Xing, D. Liu, T. H. Luan, and X. Shen, "A survey on metaverse: Fundamentals, security, and privacy," *IEEE Communications Surveys & Tutorials*, vol. 25, no. 1, pp. 319–352, 2022.
- [2] V. Makoviychuk, L. Wawrzyniak, Y. Guo, M. Lu, K. Storey, M. Macklin, D. Hoeller, N. Rudin, A. Allshire, A. Handa, *et al.*, "Isaac gym: High performance gpu-based physics simulation for robot learning," *arXiv preprint arXiv:2108.10470*, 2021.
- [3] T. Shan, B. Englot, C. Ratti, and R. Daniela, "Lvi-sam: Tightly-coupled lidar-visual-inertial odometry via smoothing and mapping," in *IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2021, pp. 5692–5698.
- [4] J. Lin and F. Zhang, "R 3 live: A robust, real-time, rgb-colored, lidar-inertial-visual tightly-coupled state estimation and mapping package," in *2022 International Conference on Robotics and Automation (ICRA)*. IEEE, 2022, pp. 10672–10678.
- [5] C. Zheng, W. Xu, Z. Zou, T. Hua, C. Yuan, D. He, B. Zhou, Z. Liu, J. Lin, F. Zhu, *et al.*, "Fast-livo2: Fast, direct lidar-inertial-visual odometry," *arXiv preprint arXiv:2408.14035*, 2024.
- [6] J. Straub, T. Whelan, L. Ma, Y. Chen, E. Wijmans, S. Green, J. J. Engel, R. Mur-Artal, C. Ren, S. Verma, *et al.*, "The replica dataset: A digital replica of indoor spaces," *arXiv preprint arXiv:1906.05797*, 2019.
- [7] J. Lin, C. Yuan, Y. Cai, H. Li, Y. Ren, Y. Zou, X. Hong, and F. Zhang, "Immsh: An immediate lidar localization and meshing framework," *IEEE Transactions on Robotics*, 2023.
- [8] 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.
- [9] H. Oleynikova, Z. Taylor, M. Fehr, R. Siegwart, and J. Nieto, "Voxblox: Incremental 3d euclidean signed distance fields for on-board mav planning," in *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2017, pp. 1366–1373.
- [10] 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.
- [11] J. Chibane, G. Pons-Moll, *et al.*, "Neural unsigned distance fields for implicit function learning," *Advances in Neural Information Processing Systems*, vol. 33, pp. 21 638–21 652, 2020.
- [12] J. Ortiz, A. Clegg, J. Dong, E. Sucar, D. Novotny, M. Zollhoefer, and M. Mukadam, "tsdf: Real-time neural signed distance fields for robot perception," *arXiv preprint arXiv:2204.02296*, 2022.
- [13] X. Zhong, Y. Pan, J. Behley, and C. Stachniss, "Shine-mapping: Large-scale 3d mapping using sparse hierarchical implicit neural representations," *arXiv preprint arXiv:2210.02299*, 2022.
- [14] A. Gropp, L. Yariv, N. Haim, M. Atzmon, and Y. Lipman, "Implicit geometric regularization for learning shapes," in *Proceedings of the 37th International Conference on Machine Learning*, 2020, pp. 3789–3799.
- [15] Z. Wang, T. Shen, M. Nimier-David, N. Sharp, J. Gao, A. Keller, S. Fidler, T. Müller, and Z. Gojcic, "Adaptive shells for efficient neural radiance field rendering," *ACM Transactions on Graphics (TOG)*, vol. 42, no. 6, pp. 1–15, 2023.
- [16] L. Yariv, J. Gu, Y. Kasten, and Y. Lipman, "Volume rendering of neural implicit surfaces," *Advances in Neural Information Processing Systems*, vol. 34, pp. 4805–4815, 2021.
- [17] P. Wang, L. Liu, Y. Liu, C. Theobalt, T. Komura, and W. Wang, "Neus: Learning neural implicit surfaces by volume rendering for multi-view reconstruction," *Advances in Neural Information Processing Systems*, vol. 34, pp. 27 171–27 183, 2021.
- [18] 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.
- [19] M. Oechsle, S. Peng, and A. Geiger, "Unisurf: Unifying neural implicit surfaces and radiance fields for multi-view reconstruction," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 5589–5599.
- [20] M. Kazhdan, M. Bolitho, and H. Hoppe, "Poisson surface reconstruction," in *Proceedings of the fourth Eurographics symposium on Geometry processing*, vol. 7, no. 4, 2006.
- [21] J. Liu and H. Chen, "Towards real-time scalable dense mapping using robot-centric implicit representation," in *IEEE International Conference on Robotics and Automation (ICRA)*, 2024.
- [22] K. Deng, A. Liu, J.-Y. Zhu, and D. Ramanan, "Depth-supervised nerf: Fewer views and faster training for free," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 12 882–12 891.
- [23] 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.
- [24] C. Jiang, H. Zhang, P. Liu, Z. Yu, H. Cheng, B. Zhou, and S. Shen, "H_{2}-mapping: Real-time dense mapping using hierarchical hybrid representation," *IEEE Robotics and Automation Letters*, 2023.
- [25] Z. Xie, J. Zhang, W. Li, F. Zhang, and L. Zhang, "S-nerf: Neural radiance fields for street views," in *International Conference on Learning Representations (ICLR)*, 2023.
- [26] 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," in *IEEE International Conference on Robotics and Automation (ICRA)*, 2024.
- [27] B. Kerbl, G. Kopanas, T. Leimkühler, and G. Drettakis, "3d gaussian splatting for real-time radiance field rendering," *ACM Transactions on Graphics*, vol. 42, no. 4, pp. 1–14, 2023.
- [28] H. Matsuki, R. Murai, P. H. J. Kelly, and A. J. Davison, "Gaussian Splatting SLAM," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024.
- [29] S. Hong, J. He, X. Zheng, H. Wang, H. Fang, K. Liu, C. Zheng, and S. Shen, "Liv-gaussmap: Lidar-inertial-visual fusion for real-time 3d radiance field map rendering," *arXiv preprint arXiv:2401.14857*, 2024.
- [30] Y. Ren, Y. Cai, F. Zhu, S. Liang, and F. Zhang, "Rog-map: An efficient robocentric occupancy grid map for large-scene and high-resolution lidar-based motion planning," *arXiv preprint arXiv:2302.14819*, 2023.
- [31] T. Müller, A. Evans, C. Schied, and A. Keller, "Instant neural graphics primitives with a multiresolution hash encoding," *arXiv preprint arXiv:2201.05989*, 2022.
- [32] D. Verbin, P. Hedman, B. Mildenhall, T. Zickler, J. T. Barron, and P. P. Srinivasan, "Ref-nerf: Structured view-dependent appearance for neural radiance fields," in *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2022, pp. 5481–5490.
- [33] Y. Wang, I. Skorokhodov, and P. Wonka, "Hf-neus: Improved surface reconstruction using high-frequency details," *Advances in Neural Information Processing Systems*, vol. 35, pp. 1966–1978, 2022.
- [34] R. Bán and G. Valasek, "Automatic step size relaxation in sphere tracing," 2023.
- [35] H. Yang, Y. Sun, G. Sundaramoorthi, and A. Yezzi, "Steik: Stabilizing the optimization of neural signed distance functions and finer shape representation," in *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- [36] W. E. Lorensen and H. E. Cline, "Marching cubes: A high resolution 3d surface construction algorithm," *ACM Siggraph Computer Graphics*, vol. 21, no. 4, pp. 163–169, 1987.
- [37] Z. Li, T. Müller, A. Evans, R. H. Taylor, M. Unberath, M.-Y. Liu, and C.-H. Lin, "Neuralangelo: High-fidelity neural surface reconstruction," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 8456–8465.
- [38] J. T. Barron, B. Mildenhall, D. Verbin, P. P. Srinivasan, and P. Hedman, "Mip-nerf 360: Unbounded anti-aliased neural radiance fields," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 5470–5479.

Supplementary Material for Neural Surface Reconstruction and Rendering for Lidar-Visual Systems

I. SUPPLEMENTARY DETAILS OF METHODOLOGY

A. Structure-aware sampling

The overall algorithm is shown in Alg. 1 and Fig. 1. Given any desired rendering direction \mathbf{d} , the algorithm starts from the camera origin (Alg. 1-1), and iteratively steps forward (Alg. 1-4) along the ray direction according to its signed distance value. The gradient of the signed distance field along the ray direction M is estimated using linear interpolation (Alg.1-8). A filtered gradient m is updated with a relaxation coefficient $\gamma = 0.7$ (Alg.1-9), which adaptively determines the next step size δ_i (Alg. 1-3) for efficiency. At every step, we ensure that the space between two adjacent steps is intersected (Alg. 1-6), wherein we take predicted scale β_i into account to avoid triggering false revert steps in unconverged NDF. To avoid poor local minima, we keep marching even behind surfaces (Alg.1-4), and until a ray's transmittance (Alg. 1-7) falls below a threshold $\epsilon_T = 0.001$. The filtered slope m is verified by sphere tracing, and smaller step sizes near surfaces yield a higher sampling frequency for more accurate slope estimations. Therefore, we apply the estimated filtered slope M in the SDF-to-density transition to avoid expensive analytical or numerical gradient calculations. We further utilize the visible-aware occupancy map to skip the free space for efficiency, as shown in Fig. 1's skip step.

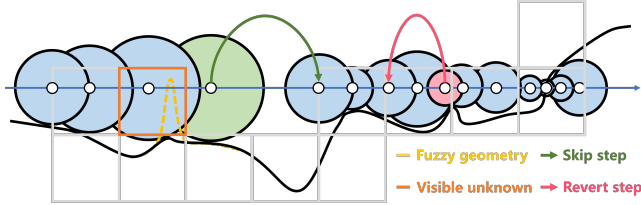


Fig. 1. Illustrations of visible-aware occupancy map and structure-aware sampling based on adaptive sphere tracing. Each circle has a radius equal to the SDF value at the sample point.

B. Losses

The Eikonal loss [14] and curvature loss [35] are further employed for both NeRF and NSDF samples to prevent the zero-everywhere and overfitting solutions for the SDF:

$$\begin{aligned}\mathcal{L}_{eik} &= \frac{1}{N_i} \sum_i (\|\nabla f_S(\mathbf{x}_i)\|_2 - 1)^2, \\ \mathcal{L}_{curv} &= \frac{1}{N_i} \sum_i |\nabla^2 f_S(\mathbf{x}_i)|,\end{aligned}\quad (1)$$

where $\nabla f_S(\mathbf{x}_i)$ and $\nabla^2 f_S(\mathbf{x}_i)$ are the gradient and Hessian of the SDF at \mathbf{x}_i calculated by numerical differentiation with a progressively smaller step size [37].

The overall training loss is defined as follows:

$$\mathcal{L} = \mathcal{L}_{sdf} + \lambda_{rgb} \mathcal{L}_{rgb} + \lambda_{eik} \mathcal{L}_{eik} + \lambda_{curv} \mathcal{L}_{curv}, \quad (2)$$

where $\lambda_{eik} = 0.1$ and $\lambda_{curv} = 5 \times 10^{-4}$ are the weights for the Eikonal, and curvature losses, respectively. We schedule the

λ_{rgb} linearly increases from 10^{-4} to 10 during the training to ensure that the NeRF learns on a well-structured NSDF to avoid local minima.

Algorithm 1: Structure-aware Sampling

Input: Neural distance field f_S , point on rendering ray $\mathbf{x}(t) = \mathbf{o} + t\mathbf{d}$, termination conditions ϵ_T , relaxation coefficient $\gamma \in (0, 1)$.

Output: ray samples $\{t_i, m_i\}$.

Notation: Slope M , filtered slope m , transmittance T .

```

1  $t_i := 0, (s_i, \beta_i) := f_S(\mathbf{x}(t_i)), m := -1, T := 1$ 
2 while  $s_i > 0 \cup T > \epsilon_T$  do
3    $\delta_i := |s_i| * \frac{2}{1-m}$   $\triangleright$  Adaptive step
4    $t_{i+1} := t_i + \delta_i$ 
5    $(s_{i+1}, \beta_{i+1}) := f_S(\mathbf{x}(t_{i+1}))$ 
6   if  $\delta_i \leq |s_i| + |s_{i+1}| + 3\beta_i$  then
7      $T := T * \exp(-\sigma(s_i, \beta_i, m) \delta_i)$ 
8      $M := (s_{i+1} - s_i) / \delta_i$ 
9      $m := \gamma m + (1 - \gamma) M$ 
10     $(t_i, m_i) := (t_{i+1}, m)$   $\triangleright$  Sample step
11  end
12 else
13    $m := -1$   $\triangleright$  Revert step
14 end
15 end
```

C. Training

1) *Outlier removal:* The zero-level set of the NDF defines the fitting surface, and the inferred signed distance values of the input LiDAR points naturally indicate the reconstruction error. In dynamic scenes, points on dynamic objects are supervised with an SDF value of zero. However, LiDAR points on the static background traversing through these dynamic points produce a supervised SDF value greater than zero. Given that dynamic points are typically sparse, this region is often dominated by the background points and has a SDF value greater than zero. Based on this observation, we propose an outlier removal strategy that periodically infers the signed distance values of training LiDAR points and eliminates points whose predicted signed distance values are more than ϵ away from 0. This enables the NDF remains a static structure field and also helps to erase dynamic objects in rendering, as shown in Fig. 6.

2) *Directional embedding scheduler:* To synthesize photorealistic novel-view images, the neural radiance field considers the view direction \mathbf{d} to output the view-dependent color \mathbf{c} at each position \mathbf{x} : $\mathbf{c} = f_C(\mathbf{x}, \mathbf{d})$, where the view direction is encoded using a 4-degree spherical harmonics encoding [32]. The tightly coupled position and view direction in training make the rendering quality in extrapolation views show a degree of color degeneration, especially at places that have only an image observations observing from similar directions, as shown in Fig. 2 (d). We consider that the surface's color is composed of view-independent diffuse

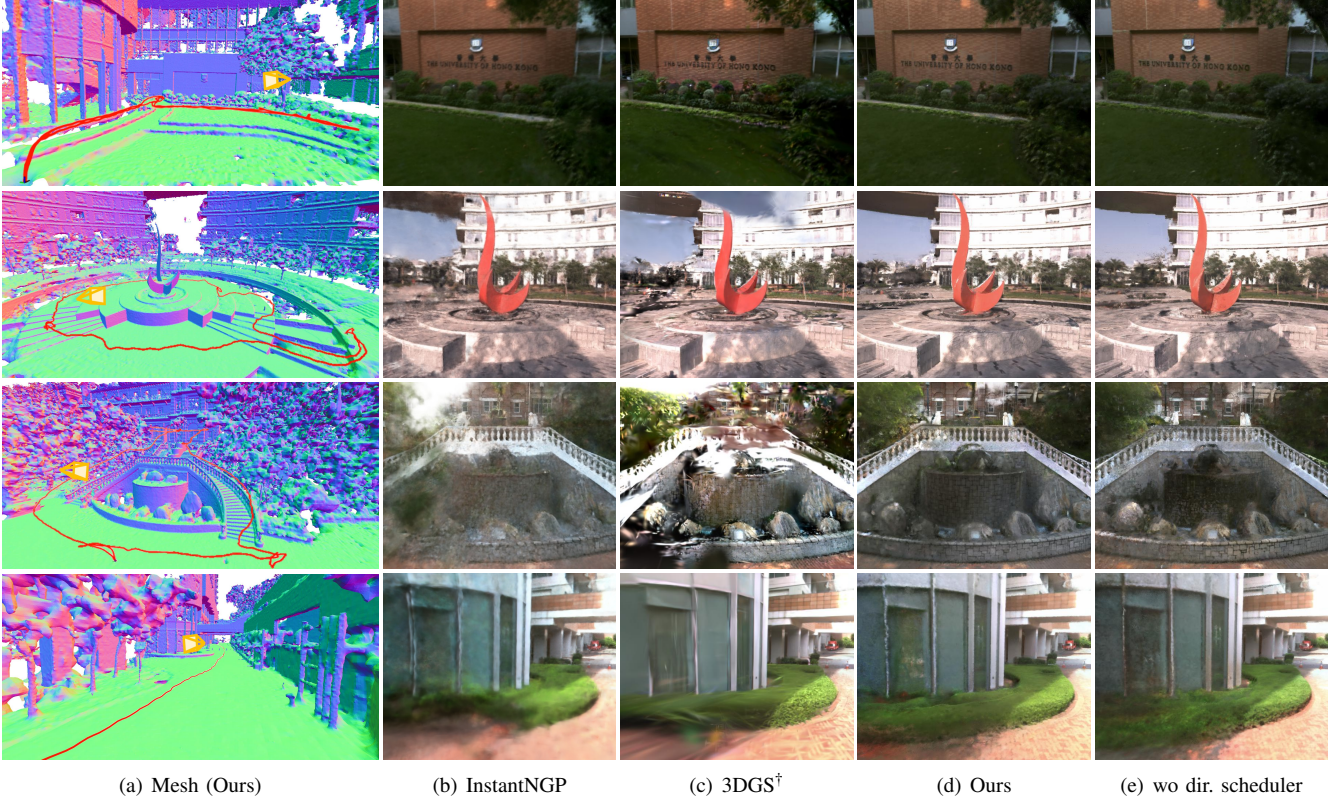


Fig. 2. The qualitative results of FAST-LIVO2 datasets (scenes from top to bottom are Campus, Sculpture, Culture, and Drive). We show our surface reconstruction results on the left, and the red line indicates the training path and the orange cameras indicates the extrapolation views for the right side’s rendering results.

color and view-dependent specular color. We schedule the degree of directional embedding to learn the surface diffuse color at the degree of 0 (i.e., view-independent feature) in the beginning and gradually increase from 0 to 4 during training to learn high-frequency view-dependent specular color.

3) *Scene contraction for background rendering*: To address the infinitely far background space, we define a boundary space that extends outward from the occupied map for background color rendering. For each ray, n_b points are uniformly sampled in the boundary space and contracted [38] as follows:

$$\text{contract}(\mathbf{x}) = \begin{cases} \mathbf{x} & , |\mathbf{x}| \leq 1 \\ \left(1 + B \left(1 - \frac{1}{|\mathbf{x}|}\right)\right) \left(\frac{\mathbf{x}}{|\mathbf{x}|}\right) & , |\mathbf{x}| > 1 \end{cases}, \quad (3)$$

where B is the size of the extending boundary space and is defined as the same voxel size as the occupancy map. The contracted samples are used to color the infinite space via volume rendering with false surfaces.

II. SUPPLEMENTARY EXPERIMENTS

A. Implementation Details

We represent our neural fields following a similar architecture to InstantNGP [31], utilizing a combination of multi-resolution hash encoding and a tiny MLP decoder. The hash encoding resolution spans from 2^5 to 2^{21} with 16 levels, and each level contains 2^{19} two-dimensional feature vectors.

Given any position \mathbf{x} , the hash encoding concatenates each level’s interpolation features to form a feature vector of size 32. One encoding feature is fed to a geometry MLP with 64-width and 3 hidden layers to obtain the SDF value and scale. Another encoding feature is concatenated with the spherical harmonics encoding of view directions and fed to an appearance MLP with 64-width and 3 hidden layers to obtain the view-dependent color. We sample 8192 rays for NDF training and follow InstantNGP [31] to fix the batch size of point samples as 256000 while the batch size of rays is dynamic per iteration. We take 20000 iterations for training, with outlier removal performed every 2000 iterations. We implement our method using LibTorch and CUDA. All experiments are conducted on a platform equipped with an Intel i7-13700K CPU and an NVIDIA RTX 4090 GPU.

B. FAST-LIVO2 Dataset

We present more surface reconstruction results and extrapolation rendering results on the FAST-LIVO2 datasets, as shown in Fig. 2. We further dive into the performance differences between volume-rendering-based methods (InstantNGP and Ours) and rasterization-based methods (3DGS), as shown in Fig. 4. InstantNGP and Our methods show more reasonable rendering results and fewer artifacts overall, such as the ground in the down-left image. While 3DGS shows powerful capability to capture high-frequency textures, such as the carved wall in the middle image.

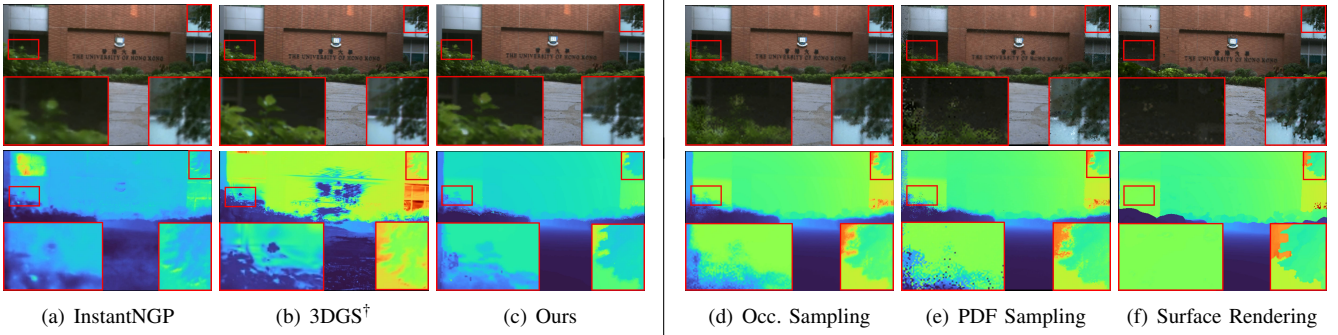


Fig. 3. (Left) Rendering results between different baselines on the FAST-LIVO dataset’s Campus scene. (Right) Rendering results with different sampling and rendering methods.

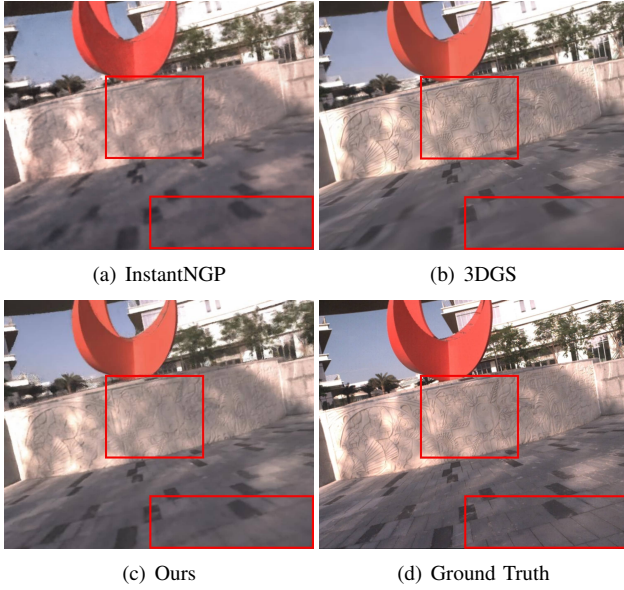


Fig. 4. We show the performance differences between volume-rendering-based methods (InstantNGP and Ours) and rasterization-based method (3DGS).

C. Ablation Study

1) *Spatial-varying scale*: To demonstrate the necessity of spatial-varying scale, we conducted experiments on the Replica dataset’s room-2 scene, where we applied $2cm$ normal noise to the ground truth depth. As shown in Fig. 5 (a-c), a large scale β results in smoother surface reconstructions with a loss of details, while a small scale leads to overfitting and noisy surface reconstructions. A spatial-varying scale is introduced to the neural distance field to adapt to the scene’s various granularity and preserve levels of detail for objects of different scales. In Fig. 5 (d-e), the rendering scale that accumulates the scale along the rays, shows that rays traverse fuzzy geometries return larger scale β indicating more uncertainty along these rays. Then the larger scale resulting in lower density (Eq. 3) makes the structure-aware sampling return more samples on these rays to meet the transmittance requirements.

2) *Structure-aware sampling*: To validate the proposed structure-aware sampling strategy, we compare the qualitative reconstruction results with occupancy sampling [31],

probabilistic density function (PDF) sampling [18] and surface rendering, as shown in Fig. 3 (Right). The occupancy sampling takes uniform samples in every occupied grid. The PDF sampling first takes uniform samples along the ray to obtains the cumulative distribution function, and then generates samples using inverse transform sampling. Surface rendering renders the scene with the first intersecting surfaces’ radiance. The proposed structure-aware sampling strategy better adapts to the SDF prior and focuses more on the surface than previous sampling methods, avoiding the skipping of fuzzy geometries as seen in surface rendering for more accurate structure rendering.

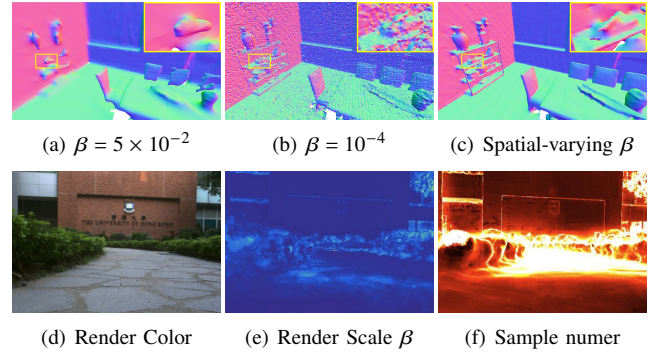


Fig. 5. In the first row, we show different scale settings in the Replica Room-2 scene. The average spatial-varying β in the scene is 1.6×10^{-4} . In the second row, we show the render scale (the lighter color corresponds to the bigger scale) and sample number per ray (the lighter color corresponds to the bigger number) in Campus scenes.

3) *Outlier removal*: To validate the necessity of outlier removal (Sec. I-C.1) for real-world applications, we conducted experiments on the FAST-LIVO2 dataset’s Drive scene, where dynamic objects (cart pusher and car) move in the scene. As shown in Fig. 6, the neural distance field with outlier removal can remove false surfaces caused by dynamic objects and regularize low density in the space traversed by the dynamic objects, and the more consistent background appearance filters out temporary dynamic objects in rendering.

4) *Directional embedding scheduler*: We propose to use a directional embedding scheduler (Sec. I-C.2) to enhance the generalization of the neural radiance field in extrapolating

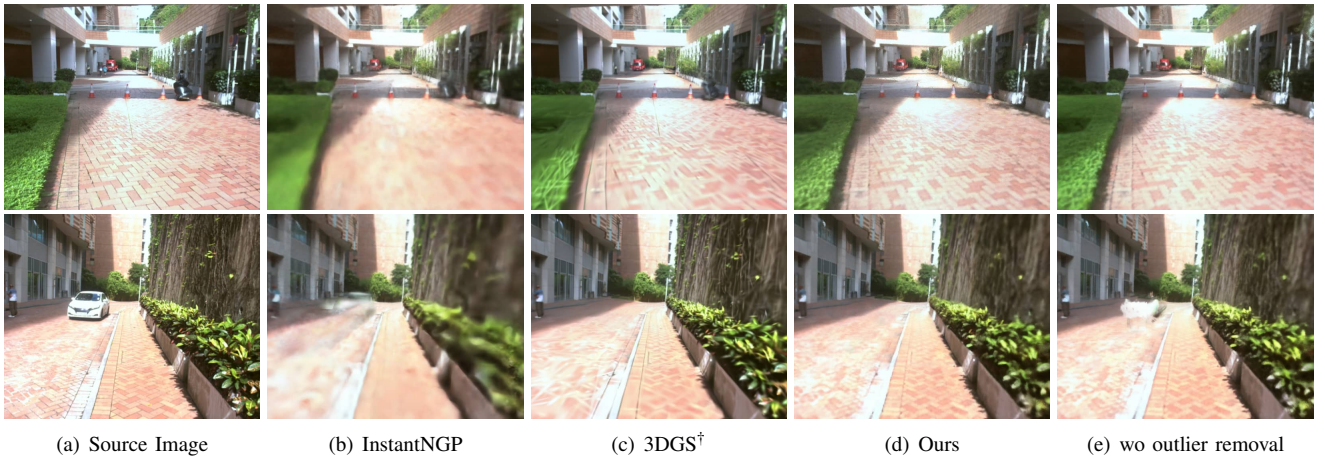


Fig. 6. Ablation study of outlier removal on the FAST-LIVO2 Drive dataset where dynamic objects (cart pusher and car) move in the scene.

TABLE I
QUANTITATIVE RESULTS ON THE FAST-LIVO2 DATASET.

Metrics	Methods	Campus	Sculpture	Culture	Drive	Avg.
SSIM↑	InstantNGP	0.789	0.698	0.670	0.697	0.714
	3DGS [†]	0.849	0.769	0.726	0.778	0.780
	Ours	<u>0.834</u>	<u>0.729</u>	0.727	<u>0.764</u>	<u>0.764</u>
	wo dir. sched.	0.833	0.726	0.729	0.767	0.764
PSNR↑	InstantNGP	28.880	22.356	21.563	24.145	24.236
	3DGS [†]	31.310	24.128	21.764	<u>25.837</u>	<u>25.760</u>
	Ours	<u>30.681</u>	<u>23.453</u>	24.695	25.941	26.193
	wo dir. sched.	30.572	23.534	24.797	26.072	26.244
LPIPS↓	InstantNGP	0.255	0.376	0.428	0.416	0.369
	3DGS [†]	0.182	0.266	<u>0.361</u>	<u>0.296</u>	0.276
	Ours	<u>0.210</u>	<u>0.321</u>	0.350	0.293	<u>0.293</u>
	wo dir. sched.	0.213	0.330	0.351	0.293	0.297

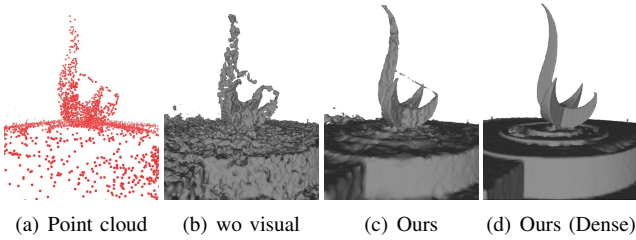


Fig. 7. We show the Sculpture scene’s reconstruction results from downsampled point clouds (a) without (b) and with (c) visual supervision. And (d) shows the reconstruction result from dense point clouds.

views, which can be verified in Fig. 2 (d). The rendering image with the directional embedding scheduler can better generalize the extrapolating views for more consistent color, especially in the free-view trajectory scenes, like Drive. Meanwhile, the directional embedding scheduler has little influence on the interpolation rendering results, as shown in Tab. I (wo dir. sched.).

5) *Visual-aided structure*: To validate the influence of image measurements on surface reconstructions, we first downsampled the point cloud in the Sculpture scene and compared the reconstruction results with and without visual supervision \mathcal{L}_{rgb} . As shown in Fig. 7, the neural radiance fields can complete the structure from sparse point clouds and avoid overfitting in scenes.

D. Training and Rendering Efficiency

So far the bottleneck of the efficiency is greatly impeded by the unbalanced sphere tracing’s steps, once there is a ray that does not converge to the surface, the other converged rays need to wait until it is finished or reach the maximum steps. The training time for a Replica scene takes about 20 minutes and the rendering time for a 1200x680 image takes about 0.07 seconds (13Hz). In the future, it could be solved by conducting a ray-oriented rendering for ray rendering instead of batch rendering.