

PC-NeRF: Parent-Child Neural Radiance Fields Using Sparse LiDAR Frames in Autonomous Driving Environments

Xiuzhong Hu, Guangming Xiong, Zheng Zang, Peng Jia, Yuxuan Han, Junyi Ma*

Abstract—Large-scale 3D scene reconstruction and novel view synthesis are vital for autonomous vehicles, especially utilizing temporally sparse LiDAR frames. However, conventional explicit representations remain a significant bottleneck towards representing the reconstructed and synthetic scenes at unlimited resolution. Although the recently developed neural radiance fields (NeRF) have shown compelling results in implicit representations, the problem of large-scale 3D scene reconstruction and novel view synthesis using sparse LiDAR frames remains unexplored. To bridge this gap, we propose a 3D scene reconstruction and novel view synthesis framework called parent-child neural radiance field (PC-NeRF). Based on its two modules, parent NeRF and child NeRF, the framework implements hierarchical spatial partitioning and multi-level scene representation, including scene, segment, and point levels. The multi-level scene representation enhances the efficient utilization of sparse LiDAR point cloud data and enables the rapid acquisition of an approximate volumetric scene representation. With extensive experiments, PC-NeRF is proven to achieve high-precision novel LiDAR view synthesis and 3D reconstruction in large-scale scenes. Moreover, PC-NeRF can effectively handle situations with sparse LiDAR frames and demonstrate high deployment efficiency with limited training epochs. Our approach implementation and the pre-trained models are available at <https://github.com/biter0088/pc-nerf>.

Index Terms—Neural Radiance Fields, 3D Scene Reconstruction, Autonomous Driving.

I. INTRODUCTION

LARGE-SCALE 3D scene reconstruction and novel view synthesis are essential for autonomous vehicles to conduct environmental exploration, motion planning, and closed-loop simulation [1]–[8], especially when the available sensor data is temporally sparse due to various practical factors [9]–[15]. Although conventional explicit representations can depict the reconstructed scene and synthetic view visually [16]–[18], they remain a significant bottleneck towards representing the scene at unlimited resolution [19], as the explicit representation is discrete. As a trendy method of implicit representation [20], [21], neural radiance fields (NeRF) [22] attract significant research interest in computer vision, robotics, autonomous driving, and augmented reality communities [6], [7], [23]–[41]. NeRF typically represents a scene using a fully connected deep network, which maps a single continuous 5D coordinate

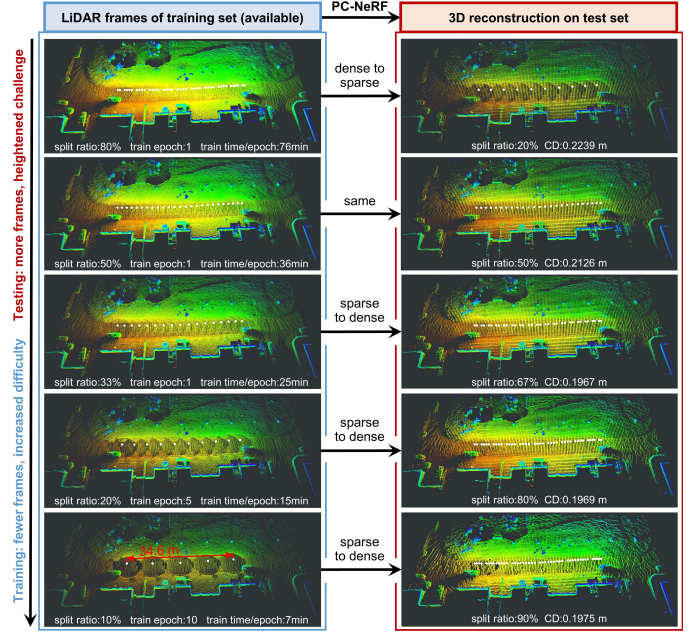


Fig. 1. PC-NeRF excels in 3D scene reconstruction and novel view synthesis, showcasing robustness to increased LiDAR frame sparsity with minimal training. Each subfigure depicts 3D scene reconstruction achieved by stitching real or synthetic LiDAR views with their corresponding poses. This scene involves frames 1151-1200 from the KITTI 00 sequence, encompassing diverse elements like the ground, grass, walls, and vehicles. White dots in each subfigure depict the LiDAR positions of each frame, and CD gauges 3D reconstruction accuracy, with smaller values indicating superior performance. As the proportion of LiDAR frames for training decreases, signifying increased sparsity, PC-NeRF achieves sparse-to-dense 3D reconstruction, as evident in the last three rows of subfigures. Moreover, utilizing only 33 % of LiDAR frames during training demonstrates advantages in both reconstruction quality and time consumption compared to using 50 % and 80 % of frames, as depicted in the first three rows of subfigures. More details are in Sec. IV-C.

(spatial location and viewing direction) to the volume density and view-dependent emitted radiance at that spatial location [22]. Hence, NeRF can construct a smooth, continuous, and differentiable scene representation [24], which helps utilize as much available sparse sensor data as possible. Most NeRF-related works have been carried out based on camera image data or indoor laser scan data [7], [19], [21], [23], [25]–[27], [32], [35]–[41], with only a few works using LiDAR point cloud data in large-scale outdoor environments [6], [29]–[31]. Unlike cameras, LiDAR has the capability to directly capture accurate distance information [42], [43], making it particularly

The research is funded by the National Natural Science Foundation of China under Grant 52372404. (Corresponding author: Junyi Ma.)

The authors are with the School of Mechanical Engineering, Beijing Institute of Technology, Beijing, 100081, China. (e-mail: 3120210302@bit.edu.cn; xionguangming@bit.edu.cn; zhengzang-biter@gmail.com; xtjp960722@163.com; yx_han_work@foxmail.com; junyi.ma@bit.edu.cn).

valuable for high-precision outdoor mapping in large-scale environments and dealing with complex scene geometry in NeRF [6], [28]. However, to leverage the successes of image-based NeRF methods, many LiDAR-based NeRF approaches project 3D LiDAR point clouds onto 2D range pseudo-images, resulting in substantial information loss [29], [30], [44]. Besides, in real autonomous vehicle applications, some unfavorable conditions, such as hardware failures and unstable communication in remote control tasks, may aggravate the temporal sparsity of LiDAR frames due to missing observation [9]–[15]. Therefore, exploring NeRF-based methods for effective utilization of sparse LiDAR frames is crucial to enhance vehicles’ autonomous capabilities, particularly evident in 3D scene reconstruction and novel view synthesis as depicted in Fig. 1.

This paper presents a parent-child neural radiance fields (PC-NeRF) framework for large-scale 3D scene reconstruction and novel LiDAR view synthesis optimized for efficiently utilizing sparse LiDAR frames in outdoor autonomous driving. PC-NeRF incorporates a hierarchical spatial partitioning approach and a multi-level scene representation. The hierarchical spatial partitioning approach first divides the driving environment into multiple large blocks, labelled as parent NeRFs, and then further partitions each parent NeRF by extracting child NeRFs. Unlike the collected LiDAR point clouds, which are often sparse and do not completely cover the object surfaces, each bounding-box-wise child NeRF represents a point cloud segment, encompassing a collection of closely located laser points and its surrounding area. The parent NeRF shares the network with child NeRFs within it for implicitly unified spatial representations. Based on the hierarchical spatial partitioning approach, we propose a multi-level scene representation, including scene, segment, and point levels. Compared to scene-level and point-level representations that represent the scene in its entirety and detail, respectively, segment-level representations try to represent the individual objects within the scene. Recognizing the inherent limitation of LiDAR point clouds in providing discrete samples of actual object surfaces, we choose the segment-level representation over the ideal object-level one. This choice facilitates the swift capture of the approximate object distribution in the environment, even in the presence of sparse LiDAR frames.

To sum up, the primary contributions of this paper include:

- To our knowledge, our proposed PC-NeRF is the first NeRF-based large-scale 3D scene reconstruction and novel LiDAR view synthesis method using sparse LiDAR frames, even though NeRF is a dense volumetric representation typically constructed using large amounts of sensor data.
- To represent outdoor large-scale autonomous driving environments, PC-NeRF introduces a hierarchical spatial partitioning approach, progressively dividing the driving environment into parent and child NeRFs.
- Based on the hierarchical spatial partitioning approach, we propose a multi-level scene representation to optimize scene-level, segment-level, and point-level representations concurrently. The multi-level scene representation is capable of efficiently utilizing sparse LiDAR frames, along with achieving high-precision 3D scene reconstruction and novel LiDAR view

synthesis with minimal training epochs.

II. RELATED WORKS

With NeRF’s inherent advantages of continuous dense volumetric representation, NeRF-based techniques in novel view synthesis [29]–[31], [37]–[39], scene reconstruction [7], [34]–[40], and localization systems [6], [25]–[27], [32], [33], [45] have rapidly developed and are highly referential and informative. PC-NeRF utilizes NeRF’s capability for continuous scene representation modelling, employs LiDAR data as inputs, and hierarchically divides the scene spaces. Therefore, this section reviews the literature on LiDAR-based NeRF and space-division-based NeRF.

A. LiDAR-Based NeRF

Motivated by NeRF’s capability to render photo-realistic novel image views, several studies have investigated its potential application to LiDAR point cloud data for novel view synthesis [29]–[31] and robot navigation [6], [32], [33].

The goal of novel view synthesis is to generate a view of a 3D scene from a viewpoint where no real sensor image has been captured, providing the opportunity to observe real scenes from a virtual perspective. Neural Fields for LiDAR (NFL) [31] combines the rendering power of neural fields with a physically motivated model of the LiDAR sensing process, thus enabling it to accurately reproduce key sensor behaviors like beam divergence, secondary returns, and ray drop. Our work is inspired by the modeling of LiDAR sensing processes. LiDAR-NeRF [29] converts the 3D point cloud into the range pseudo image in 2D coordinates and then optimizes three losses, including absolute geometric error, point distribution similarity, and realism of point attributes. Similar to LiDAR-NeRF, NeRF-LiDAR [30] has also employed the spherical projection strategy and consists of three key components: NeRF reconstruction of the driving scenes, realistic LiDAR point clouds generation, and point-wise semantic label generation. However, when multiple laser points project onto the same pseudo-pixel, only the one with the smallest distance is retained [29]. This effect becomes particularly pronounced with small resolution range pseudo-images, leading to significant information loss in cases of spherical projections [44].

In robotics, LiDAR-based NeRF is usually proposed for localization and mapping. IR-MCL [32] focuses on the problem of estimating the robot’s pose in an indoor environment using 2D LiDAR data. With the pre-trained network, IR-MCL can synthesize 2D LiDAR scans for an arbitrary robot pose through volume rendering. However, the error between the synthesized and real scans is relatively large. NeRF-LOAM [6] presents a novel approach for simultaneous odometry and mapping using neural implicit representation with 3D LiDAR data. NeRF-LOAM employs sparse octree-based voxels combined with neural implicit embeddings, decoded into a continuous signed distance function (SDF) by a neural implicit decoder. However, NeRF-LOAM cannot currently operate in real-time with its unoptimized Python implementation. LocNDF [33] utilizes neural distance fields (NDFs) for robot localization, demonstrating the direct learning of NDFs from range sensor

observations. LocNDF has raised the challenge of addressing real-time constraints, and our work endeavors to investigate this challenge.

In contrast to projecting the LiDAR point cloud onto a range pseudo-image, our proposed PC-NeRF handles 3D LiDAR point cloud data directly. Besides learning the LiDAR beam emitting process, our proposed PC-NeRF explores the deployment performance of NeRF-based methods.

B. Space-Division-Based NeRF

When large-scale scenes such as where the autonomous vehicles drive need to be represented with high precision, the model capacity of a single NeRF is limited in capturing local details with acceptable computational complexity [7], [35], [36]. For large-scale 3D scene reconstruction tasks, Mega-NeRF [7], Block-NeRF [35], and Switch-NeRF [36] have adopted the multiple NeRF solution, with each NeRF responsible for different scene areas. Mega-NeRF decomposes a scene into cells with centroids and initializes a corresponding set of model weights. At query time, Mega-NeRF produces opacity and color for a given position and direction using the model weights closest to the query point. Like Mega-NeRF, Block-NeRF [35] proposes dividing large environments into individually trained Block-NeRFs, which are then rendered and combined dynamically at inference time. For rendering a target view, a subset of the Block-NeRFs are rendered and then composited based on their geographic location compared to the camera. Switch-NeRF [36] proposes a novel end-to-end large-scale NeRF with learning-based scene decomposition and designs a gating network to dispatch 3D points to different NeRF sub-networks. The gating network can be optimized with the NeRF sub-networks for different scene partitions by design with the Sparsely Gated Mixture of Experts. Besides the multiple NeRF solution, NeRF-LOAM and Shine-mapping employ the octree structure to recursively divide the scene into leaf nodes with basic scene unit voxels, simplifying the description of large-scale scenes [4], [6], [23]. These basic scene unit voxels attach an N-dimensional encoding at each vertex and share it with neighboring voxels. Thus, the attributes of any 3D location in the scene can be inferred from the vertex encoding values output by the neural network, which in turn achieves 3D reconstruction.

Regarding space-division-based NeRF on a smaller scale, further space division is employed for faster and higher-quality rendering [37]–[39]. DeRF [37] and KiloNeRF [38] adopt the multiple NeRF solution to represent scene details and speed up rendering. DeRF [37] decomposes the scene space and represents each decomposed part using a separate neural network (also named a decomposition head network), where each decomposition head network is defined over the entire scene space. Using Voronoi learnable decompositions, only one of the decomposition head networks works at any given spatial location. Thus, only one decomposition head network needs to be evaluated for each spatial location, resulting in an accelerated inference process. KiloNeRF [38] demonstrates that real-time rendering is possible using thousands of tiny MLPs instead of one extensive multilayer perceptron network

(MLP). Rather than representing the entire scene with a single, high-capacity MLP, KiloNeRF represents the scene with thousands of small MLPs. Similar to the octree scene partition on a large-scale scene, Neural Sparse Voxel Fields (NSVF) [39] also partitions the scene space with the sparse voxel octree and assigns the voxel embedding at each vertex.

Our work introduces a hierarchical spatial partitioning approach that incorporates spatial division concepts at both large and small scales for representing autonomous driving scenes. We first partition the driving environment into multiple large blocks and then, within each block, extract the spatial extents of point cloud segments. We also employ the multiple NeRF solution, where the volumetric representation of a large block and the point cloud segments inside it are represented by a shared NeRF network. Unlike the octree scene partition and voxel vertex embedding, we concentrate on constructing the volumetric representation around and within individual point cloud segments, as the point cloud segments can represent the approximate spatial extent of the actual object surfaces and thus alleviate the sparsity of the LiDAR frames. Therefore, we opt for a segment-level representation instead of the ideal object-level representation to rapidly access to the detailed structural features of a scene.

III. PARENT-CHILD NEURAL RADIANCE FIELDS (PC-NeRF) FRAMEWORK

To explore the large-scale 3D scene reconstruction and novel LiDAR view synthesis based on LiDAR and NeRF, especially using sparse LiDAR frames, we propose a parent-child neural radiance field (PC-NeRF) framework, as shown in Fig. 2. We introduce a hierarchical spatial partitioning approach in Sec. III-A and Sec. III-B, dividing the entire autonomous vehicle driving environment into large blocks, i.e., parent NeRFs, and further dividing a block into geometric segments, represented by child NeRFs. A parent NeRF shares a network with child NeRFs within it. Based on the hierarchical spatial partitioning approach, Sec. III-C introduces a comprehensive multi-level scene representation aimed at collectively optimizing scene-level, segment-level, and point-level representations while efficiently utilizing sparse LiDAR frames. Moreover, we propose a two-step depth inference method in Sec. III-D to realize segment-to-point inference.

A. Spatial Partition of Parent NeRF

To efficiently represent the autonomous vehicle’s driving environment, our proposed hierarchical spatial partitioning approach constructs multiple rectangular-shaped parent NeRFs (as Fig. 2(a) illustrates) one after the other along the trajectory. A new parent NeRF is created when the autonomous vehicle orientation variation exceeds a given threshold. Besides, parent NeRFs with highly overlapping spatial areas should be merged. Constructing parent NeRF requires three considerations: the effective LiDAR point cloud selection, the driving environment representation, and the NeRF near/far bounds calculation. Since the point clouds become sparser the further away from the LiDAR origin and the NeRF itself is a dense volumetric representation, LiDAR points within a

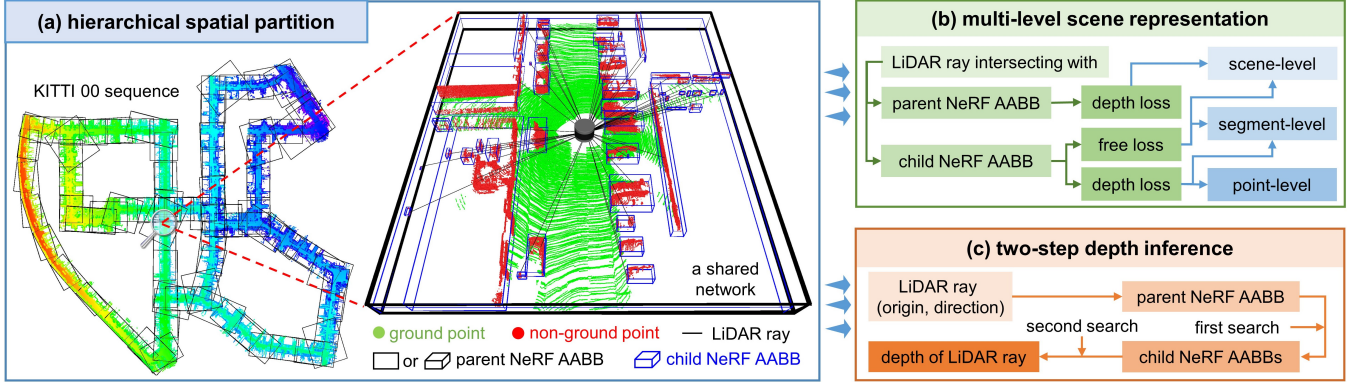


Fig. 2. Our PC-NeRF framework: (a) The hierarchical spatial partition divides the entire large-scale scene into large blocks, referred to as parent NeRFs. After multi-frame point cloud fusing, ground filtering, and non-ground point cloud clustering, a large block is further divided into point cloud geometric segments represented by a child NeRF. The parent NeRF shares a network with the child NeRFs within it. (b) In the multi-level scene representation, the surface intersections of the LiDAR ray with the parent and child NeRF AABBs and the LiDAR origin are used to divide the entire LiDAR ray into different line segments. The three losses on these line segments concurrently optimize the scene representation at the scene level, segment level, and point level, effectively leveraging sparse LiDAR frames. (c) For depth inference of each LiDAR ray, PC-NeRF searches in the parent NeRF AABBs to locate corresponding child NeRF AABBs and then refines its inference in the child NeRF AABBs for higher precision.

certain distance from the LiDAR origin rather than the holistic point clouds are chosen to train our proposed PC-NeRF model. Moreover, the driving environments featured with roads, walls, and vehicles can be tightly enclosed by bounding boxes. Therefore, each parent NeRF’s space is represented as a large Axis-Aligned Bounding Box (AABB) in our work, making it easier to calculate the related near and far bounds for rendering.

B. Spatial Partition of Child NeRF

To efficiently capture the approximate environmental distributions using sparse LiDAR frames, our proposed hierarchical spatial partitioning approach distributes the point clouds in a parent NeRF into multiple child NeRFs (as Fig. 2(a) illustrates). In contrast to the collected LiDAR point clouds, which are often sparse and do not completely cover the object surfaces, each child NeRF represents a point cloud segment, encompassing a collection of laser points close to each other and its surrounding area. Since the geometric point cloud segment quantity is limited and the child NeRFs’ space can also be represented by AABBs, constructing child NeRFs from the raw point cloud is a fast way to generate detailed environmental representations. With less total space volume, child NeRFs can represent a larger environment with the same model capacity. Moreover, by understanding the spatial distribution of point cloud segments, child NeRFs can address environmental representation inadequacies due to the sparsity of LiDAR frames.

The point cloud allocation for child NeRF can be divided into the following three steps, and the results are shown in Fig. 2(a). Step 1: Extract ground point clouds, such as road surfaces and sidewalks along roads, from the fused point cloud data. The ground plane is one significant geometric segment in the driving environment and is spatially adjacent to the other individual geometric segments. Therefore, distinguishing it helps to extract other geometric segments further accurately. Step 2: Cluster the remaining non-ground point clouds into

various segments using region-growing clustering. The advantages of the utilized clustering method are excellent scalability for operation on large-scale point clouds and the adaptive ability to cluster different shaped and sized objects. Step 3: Construct child NeRF AABBs. The AABBs of the segmented point clouds obtained from step 1 and step 2 can be used as child NeRFs’ spatial extent. After the division process above, the fused point clouds are distributed into a limited quantity of child NeRFs. Note that when the driving scene is complex, there are many artifacts existing in the results of ground extraction and point cloud clustering. However, our proposed method mitigates the negative effects of the inherent inaccuracy for these two tasks in complicated driving scenes, which is further introduced in Sec. III-C.

C. Multi-level Scene Representation for Sparse LiDAR Frames

By observing the intersections of the LiDAR rays with parent and child NeRF AABBs (Fig. 2(a) and Fig. 3), we propose a multi-level scene representation (Fig. 2(b)), including scene-level, segment-level, and point-level representations. The multi-level scene representation enables the simultaneous extraction of global, local, and detailed information from the environment, enhancing the utilization of effective information in sparse LiDAR frames. When the LiDAR ray intersects with parent and child NeRF AABBs, we can obtain three intersecting positions of the LiDAR ray with child NeRF AABBs’ inner surface, child NeRF AABBs’ outer surface, and parent NeRF AABBs’ outer surface, as shown in and Fig. 3. The LiDAR origin and these three intersections divide the LiDAR ray into multiple line segments. By examining the distribution of object surface points on these line segments, we can derive the positions of individual object surface points (point-level representation), the spatial distribution of point segments formed by these object surface points (segment-level representation), and the overall distribution of object surface points in the whole space (scene-level representation).

In the global coordinate system, the LiDAR origin position of the i -th frame is represented as $\mathbf{o}_i = (x_i, y_i, z_i)$. The depth,

direction, and the corresponding child NeRF order number of the j -th laser point $\mathbf{p}_{ij} = (x_{ij}, y_{ij}, z_{ij})$ in the i -th frame point cloud are $d_{ij} = \|(\mathbf{p}_{ij} - \mathbf{o}_i)\|_2$, $\mathbf{d}_{ij} = (\mathbf{p}_{ij} - \mathbf{o}_i)/d_{ij}$ and k_{ij} respectively. Then, we can get one LiDAR ray $\mathbf{r}_{ij} = [\mathbf{o}_i, \mathbf{p}_{ij}, d_{ij}, \mathbf{d}_{ij}, k_{ij}]$. We apply LiDAR point cloud data to train our proposed PC-NeRF with model parameters θ . Besides, we further design three LiDAR losses and child NeRF segmented sampling, illustrated in Fig. 3.

In NeRF volume rendering, the depth value of a ray \mathbf{r} is synthesized from the weighted sampling depth values t between the near and far bounds (t_n and t_f) by:

$$d(\mathbf{r}) = \int_{t_n}^{t_f} w(t) \cdot t dt \quad (1)$$

where $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$ represents a ray with camera or LiDAR origin \mathbf{o} oriented as \mathbf{d} , and the volume rendering integration weights $w(t)$ is calculated by:

$$w(t) = \exp\left(-\int_{t_n}^t \sigma(s) ds\right) \cdot \sigma(t) \quad (2)$$

where $\exp(-\int_{t_n}^t \sigma(s) ds)$ is the visibility of $\mathbf{r}(t)$ from origin \mathbf{o} , and $\sigma(t)$ is the volumetric density at $\mathbf{r}(t)$.

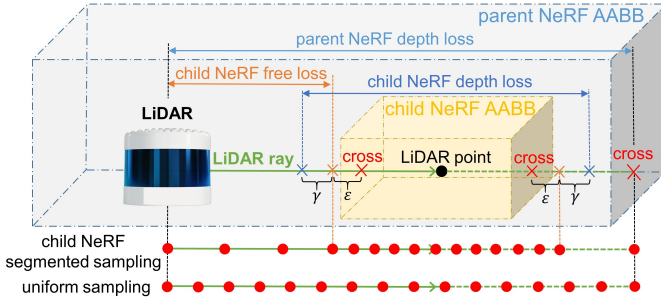


Fig. 3. Three LiDAR losses and child NeRF segmented sampling. The three LiDAR losses include parent NeRF depth loss, child NeRF depth loss, and child NeRF free loss. Using different sampling densities, the Child NeRF segmented sampling uniformly samples both inside and outside the intersection of the LiDAR ray with the Child NeRF.

Parent NeRF depth loss: In Sec. III-A, the autonomous vehicle driving environment is represented efficiently using parent NeRF. We use the intersection of the LiDAR ray \mathbf{r}_{ij} with its corresponding parent NeRF surface as the far bound f_{ij}^p , as seen in Fig. 3. To adequately represent the volumetric distribution of the whole driving environment, we set parent NeRF depth loss as:

$$\mathcal{L}_{ij}^{\text{pd}}(\theta) = \mathcal{L}'_{\text{L1}} \left(\int_{t_0}^{f_{ij}^p} w(t) \cdot t dt, d_{ij} \right) \quad (3)$$

where the integration lower limit t_0 is set to 0. Considering the space occupied by the LiDAR or the autonomous vehicle, setting t_0 to 0.5 m or 1 m is also recommended. Moreover, $\mathcal{L}'_{\text{L1}}(x, y) = 0.1 \cdot \text{SmoothL1Loss}(10 \cdot x, 10 \cdot y)$ is an extension of SmoothL1Loss that shifts the $|x - y|$ turning point from 1 m to 0.1 m to improve model sensitivity.

The commonly used depth inference approach uses the synthetic depth between the scene's near and far bounds as

the inference depth [22], [25], [27], [32], [34], which is used in Eq. 3 and calculated as follows:

$$\hat{d}_{ij}^{\text{p}} = \int_{t_0}^{f_{ij}^p} w(t) \cdot t dt \quad (4)$$

Segment-to-point hierarchical representation strategy: To find object surface points effectively and address the sparsity of LiDAR frames, we further propose a segment-to-point hierarchical representation strategy. The depth difference $f_{ij}^p - t_0$ between parent NeRF near and far bounds in large-scale outdoor scenes can be about ten times greater than that in indoor scenes, making it more difficult for the NeRF model to find the object surface through sampling. Given that it is much easier to find segment-level child NeRF space than to find point-level object surface points between the parent NeRF near and far bounds, we propose a segment-to-point hierarchical representation strategy, which first finds the child NeRFs that intersect the LiDAR ray and then finds object surface points inside the child NeRFs that intersect the LiDAR ray.

We use the intersections of the LiDAR ray \mathbf{r}_{ij} with its corresponding child NeRF surface as the child NeRF near and far bounds $[n_{ij}^c, f_{ij}^c]$, as seen in Fig. 3. It is pretty evident that the depth difference between child NeRF near and far bounds ($f_{ij}^c - n_{ij}^c$) is much smaller than that between parent NeRF near and far bounds. Considering that the object surface has a certain thickness and the object surface may appear at the child NeRF bounds, we slightly inflate the child NeRF near and far bounds as $[n_{ij}^c - \varepsilon, f_{ij}^c + \varepsilon]$, where ε is a small inflation coefficient, as shown in Fig. 3.

Child NeRF free loss: To make the process of finding the segment-level child NeRF's corresponding space faster, we propose the child NeRF free loss in Eq. 5. As no opaque objects exist in $(t_0, n_{ij}^c - \varepsilon)$ and no objects can be observed in $(f_{ij}^c + \varepsilon, f_{ij}^p)$ in numerous cases, the loss is calculated by:

$$\mathcal{L}_{ij}^{\text{cf}}(\theta) = \int_{t_0}^{n_{ij}^c - \varepsilon} w(t)^2 dt + \int_{f_{ij}^c + \varepsilon}^{f_{ij}^p} w(t)^2 dt \quad (5)$$

Child NeRF depth loss: Based on child NeRF free loss, we propose child NeRF depth loss to find object surface points inside the child NeRF, which is calculated by:

$$\mathcal{L}_{ij}^{\text{cd}}(\theta) = \mathcal{L}'_{\text{L1}} \left(\int_{n_{ij}^c - \varepsilon}^{f_{ij}^c + \varepsilon} w(t) \cdot t dt, d_{ij} \right) \quad (6)$$

To let the child NeRF free loss and the child NeRF depth loss have a smooth transition at the child NeRF bounds, $\mathcal{L}_{ij}^{\text{cd}}(\theta)$ is further modified to:

$$\mathcal{L}_{ij}^{\text{cd}}(\theta) = \mathcal{L}'_{\text{L1}} \left(\int_{n_{ij}^c - \varepsilon - \gamma}^{f_{ij}^c + \varepsilon + \gamma} w(t) \cdot t dt, d_{ij} \right) \quad (7)$$

where γ is a constant designed to represent the smooth transition interval on a LiDAR ray between the child NeRF free loss and the child NeRF depth loss, as seen in Fig. 3. Child NeRF free loss and child NeRF depth loss are employed separately to supervise the space from the LiDAR origin to the vicinity of the child NeRF AABB and the spatial extent of the child NeRF AABB after it expands γ (in m). For the transition from

free space to the object’s surface, the child NeRF free loss is stringent, while the child NeRF depth loss is comparatively lenient. Both functions collaboratively contribute to accurately capturing the real surface of the object. Meanwhile, they help mitigate the impact of incomplete ground extraction and insufficient point cloud clustering discussed in Sec. III-B.

To sum up, the total training loss from one LiDAR ray \mathbf{r}_{ij} contains all the losses mentioned above:

$$\mathcal{L}_{ij}(\boldsymbol{\theta}) = \lambda_{\text{pd}} \mathcal{L}_{ij}^{\text{pd}}(\boldsymbol{\theta}) + \lambda_{\text{cf}} \mathcal{L}_{ij}^{\text{cf}}(\boldsymbol{\theta}) + \lambda_{\text{cd}} \mathcal{L}_{ij}^{\text{cd}}(\boldsymbol{\theta}) \quad (8)$$

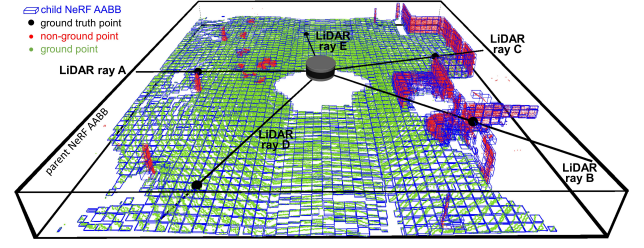
where λ_{pd} , λ_{cf} , and λ_{cd} are the parameters to jointly optimize different losses.

Child NeRF segmented sampling: To efficiently find the objects in large-scale scenes, even using sparse LiDAR frames, we propose a child NeRF segmented sampling method for objects more likely to be found in and around the child NeRFs. Assuming that N points are sampled uniformly along the LiDAR rays, the child NeRF segmented sampling is sampling $\lambda_{\text{in}} \cdot N$ points in $[n_{ij}^c - \varepsilon, f_{ij}^c + \varepsilon]$ and sampling $(1 - \lambda_{\text{in}}) \cdot N$ in $[t_0, f_{ij}^p]$, as shown in Fig. 3. Therefore, child NeRF segmented sampling guarantees that at least $\lambda_{\text{in}} \cdot N$ points are sampled inside child NeRF, which means that at least $\lambda_{\text{in}} \cdot N$ sampling points are sampled near the real object.

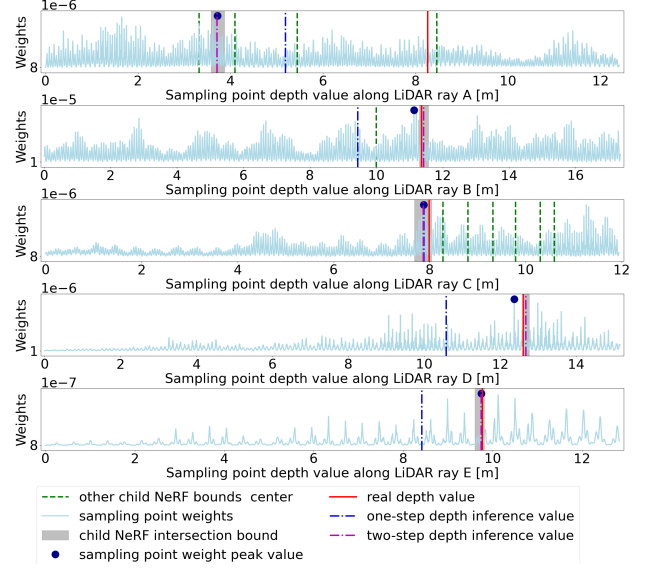
D. Two-step Depth Inference

Most current depth inference methods are one-step depth inference methods, similar to Eq. 4. In contrast, we provide a two-step depth inference method (as Fig. 2(c) illustrates) to infer more accurately, especially training with sparse LiDAR frames. We first search in the parent NeRF AABBS to acquire the child NeRF AABBS potentially intersecting the LiDAR ray \mathbf{r}_{ij} and then conduct further inference in the child NeRF AABBS’s near and far bounds $[n_{ij}^c, f_{ij}^c]$ intersecting with the LiDAR ray, as shown in Fig. 4(a). To this end, we first select the child NeRF whose AABBS outer sphere intersects the ray and then use the Axis Aligned Bounding Box intersection test proposed by Haines [46], which can readily process millions of voxels or AABBS in real time [39]. If the LiDAR ray does not intersect any child NeRF AABBS, the space extent of all child NeRF AABBS for this LiDAR ray should be slightly inflated in incremental steps, and the above inference should be performed again.

To mitigate the risk of misinterpreting free space as object space, we constrain the inferred depth values within the near and far bounds of a single child NeRF, as depicted in the first subfigure of Fig. 4(b). In complex scenes with LiDAR rays intersecting multiple child NeRF AABBS, we select the AABBS containing the peak weight value (calculated by Eq. 2). This choice is illustrated in the first and third subfigures of Fig. 4(b), considering the LiDAR ray’s inability to penetrate opaque objects during emission. Besides, when the peak weight value is not in any child NeRF AABBS, we opt for the child NeRF with the maximum weight integration W (from Eq. 9), i.e., the most probable existence range of an object on the LiDAR ray, as shown in the second subfigure of Fig. 4(b). When the LiDAR ray only intersects one child NeRF AABBS, we can directly select it as the object area, as shown in the last two



(a) LiDAR rays intersect with different AABBS.



(b) Weight distribution and depth inference along 5 LiDAR rays.

Fig. 4. Parent-child NeRF’s two-step depth inference effect illustration. The five subfigures in Fig. 4(b) represent depth value inference results for the five LiDAR rays in Fig. 4(a), where the weight distribution data comes from our proposed PC-NeRF model trained on the KITTI 00 sequence 1151-1200 frame scene in Sec. IV-B.

subfigures of Fig. 4(b). Note that if the child NeRF weight integration W is minimal, the LiDAR ray doesn’t intersect any child NeRF, and no depth value needs to be inferred. The inferred depth value of each LiDAR ray is calculated by Eq. 10.

$$W = \int_{n_{ij}^c}^{f_{ij}^c} w(t) dt \quad (9)$$

$$\hat{d}_{ij} = \frac{\int_{n_{ij}^c}^{f_{ij}^c} w(t) \cdot t dt}{\int_{n_{ij}^c}^{f_{ij}^c} w(t) dt} \quad (10)$$

IV. EXPERIMENTS

A. Experiment Setups

Datasets: We evaluate our proposed PC-NeRF using 13 data sequences from two publicly available outdoor datasets, including the MaiCity [18] and KITTI odometry datasets [47]. The 00 and 01 sequences of the MaiCity dataset contain 64-beam noise-free synthetic LiDAR frames in virtual urban-like environments. The KITTI odometry dataset also contains 64-beam LiDAR data collected by vehicles driving in real-world

environments and provides a localization benchmark with ground-truth vehicle poses. The semantic labels of point clouds in KITTI are from SemanticKITTI [48], which further filters out movable objects. This filtering enhances the accuracy and stability of 3D reconstruction and novel view synthesis, simplifies processing, and improves computational efficiency.

Frame sparsity: Frame sparsity represents the proportion of the test set (unavailable during training) when dividing the LiDAR dataset into training and test sets. Increased frame sparsity implies fewer LiDAR frames for training and more for model testing, posing heightened challenges across various tasks. Following the dataset splitting method employed in previous works like NeRF-LOAM [6] and NFL [31], we initially take one test frame from every five frames, setting the frame sparsity to 20%. Following that, we systematically vary frame sparsity by selecting test frames as follows: one from every four frames (25%), one from every three frames (33%), one from every two frames (50%), two from every three frames (67%), three from every four frames (75%), and four from every five frames (80%). Therefore, the test frame selection for {20%, 25%, 33%, 50%} frame sparsities shows a regular increase, as does the test frame selection for {50%, 67%, 75%, 80%} frame sparsities. In addition, it is reasonable to use 50% (the frame sparsity threshold, i.e., the same number of LiDAR frames for model training and model performance testing) as a transition from {20%, 25%, 33%} to {67%, 75%, 80%}. We further explore extreme dataset splitting by selecting nine test frames from every ten frames, resulting in a frame sparsity of 90%. In summary, selecting these eight frame sparsities of {20%, 25%, 33%, 50%, 67%, 75%, 80%, 90%} is meaningful and represents the range of 0% to 100% well.

Metrics: We evaluate our method’s performance in novel LiDAR view synthesis, single-frame 3D reconstruction, and 3D scene reconstruction. We use the error metrics presented in IR-MCL [32] and UrbanNeRF [34] to evaluate novel LiDAR view synthesis and single-frame 3D reconstruction performance. By comparing each synthesized LiDAR frame with its corresponding real LiDAR frame on the test set and averaging the metrics across all test frames, we report the average absolute error of LiDAR ray depth (Dep. Err. [m]), the average accuracy of LiDAR ray depth at 0.2 m threshold (Dep. Acc@0.2m [%]), chamfer distance (CD [m]), and F-score at 0.2 m threshold (F@0.2m). For 3D scene reconstruction, we concatenate individual synthesized and real LiDAR frames using the test set poses to generate the reconstructed and real LiDAR point cloud maps. We use the reconstruction metrics commonly used in most reconstruction methods [4], [6], [18], [49], i.e., accuracy (Acc. [m]), completion (Comp. [m]), chamfer distance (Map CD [m]), and F-score at 0.2 m threshold (Map F@0.2m), to evaluate the 3D scene reconstruction results between the reconstructed and real LiDAR point cloud maps. Given the inherent interdependence between 3D scene reconstruction and single-frame 3D reconstruction, the experiments in Sec. IV-B, Sec. IV-C, and Sec. IV-D focus primarily on evaluating metrics related to single-frame 3D Reconstruction.

Baselines: A standard pipeline for generating new LiDAR

views is constructing a 3D point cloud map and then using the ray-casting approach [50] to query new point clouds from the map. We implement this pipeline, which voxelizes the 3D point cloud map (voxel size: 0.05 m) into a 3D voxel map to speed up the query but may slightly reduce accuracy, as a baseline method named **MapRayCasting**. In addition, we also extend the original NeRF model proposed by Mildenhall [22] by replacing a camera ray with a LiDAR ray as IR-MCL [32], named **OriginalNeRF**. For the memory consumption, OriginalNeRF and PC-NeRF include the model file and the child NeRFs bounds file, while MapRayCasting includes a voxel map. Furthermore, we also employ ground truth poses in **NeRF-LOAM** [6] to reconstruct the mesh map of the environments. In the whole execution of NeRF-LOAM, a shared network of just two fully connected layers with 256 units per layer is used. Therefore, it is difficult for NeRF-LOAM to perform the novel view synthesis task using the network weights obtained from training. Accordingly, we use the mesh obtained from NeRF-LOAM for the 3D scene reconstruction task.

Training details: We train our proposed PC-NeRF and all the baselines with an NVIDIA GeForce RTX 3090 and an Intel i9-11900K CPU and use Adam [51] as the training optimizer. Since reconstructing and storing large-scale outdoor scenes as NeRF models require a relatively lengthy training period [7], [34], [35], [41], we aim to achieve optimal results with minimal training epochs, particularly emphasizing one epoch for swift deployment. After training with the current number of epochs, if invalid inferences occur on certain LiDAR rays during the two-step depth inference method, we will increment the number of epochs, e.g., 2, 5, or 10. Besides, considering the relatively time-consuming nature of each epoch’s training, as indicated in Tab. I and IV, we have chosen not to include experimental results for the maximum possible epochs in the paper. Based on extensive experimental testing, the initial learning rate is set to 4×10^{-5} and adjusted using Pytorch’s MultiStepLR strategy with milestones at [5, 120] and an adjustment factor of 0.1. The experimental results on the KITTI and MaiCity datasets, including Sec. IV-B, Sec. IV-C, and Sec. IV-D, demonstrate that our proposed PC-NeRF achieves satisfactory results with only one training epoch in most scenes. In Sec. IV-B and Sec. IV-C, λ_{pd} , λ_{cf} , λ_{cd} , λ_{in} , and γ of our proposed PC-NeRF are set to 1, 10^6 , 10^5 , 0.1, and 2.0 m, respectively. This group of parameters is not specifically tuned for a single scene but is valid for all experimental scenes, and the corresponding ablation studies are shown in Sec. IV-D. Both our PC-NeRF and OriginalNeRF use the hierarchical volume sampling strategy along the LiDAR ray proposed in NeRF [22], where the points number N_c and N_f for coarse and fine sampling along the ray are 768 and 1536, respectively. However, for coarse sampling, our PC-NeRF uses the Child NeRF segmented sampling proposed in Sec. III-C, while OriginalNeRF samples uniformly along the LiDAR rays.

B. Evaluation for Novel LiDAR View Synthesis and 3D Reconstruction in Different Scales

We qualitatively and quantitatively evaluate the inference accuracy and deployment potential of our PC-NeRF across

different scales. For small-scale evaluation, we use 50 consecutive LiDAR frames from the MaiCity and KITTI datasets as a single scene, training and evaluating each model on this scene. The frame sparsity is set to 20% and the experimental results are presented in Fig. 5, Tab. I, and Tab. II.

As shown by the white dashed ellipses in the second column of Fig. 5(b), MapRayCasting successfully reconstructs the overall environment but introduces shadow artifacts. In the third to sixth column of Fig. 5(b), the one-step depth inference method fails to reconstruct the environment, while the two-step depth inference method outlined in Sec. III-D exhibits effective performance. Besides, employing the two-step depth inference method, our proposed PC-NeRF and OriginalNeRF additionally infer structures that do not belong to the test frames (as shown by the white dashed box in the third and fifth columns of Fig. 5(b)) but are genuinely present in the actual scene (as shown in Fig. 5(a)) which is more beneficial for 3D scene reconstruction. As demonstrated in the third column of the fourth row in Fig. 5(b), OriginalNeRF faces challenges in scene reconstruction with the two-step depth inference. Additionally, due to the invalid inferences on some LiDAR rays when utilizing the two-step depth inference method, we rely on the remaining valid inferences for quantitative evaluation in Tab. I. These challenges and invalid inferences occur because OriginalNeRF cannot quickly obtain the approximate environment distribution during training, leading to the failure of the first step in the two-step depth inference method. To sum up, implementing the two-step depth inference method with PC-NeRF yields the best qualitative results.

As shown in Tab. I, our proposed PC-NeRF outperforms MapRayCasting and OriginalNeRF, demonstrating superior accuracy in novel LiDAR view synthesis and single-frame 3D reconstruction. Besides, our proposed PC-NeRF demonstrates excellent deployment potential as it yields superior results with just one epoch of training, maintaining consistent training set across various scenes. Tab. I shows that by combining training time, memory consumption, novel LiDAR view synthesis accuracy, and single-frame 3D reconstruction accuracy, our proposed PC-NeRF is far superior to OriginalNeRF and MapRayCasting. Additionally, in Tab. II (rows with a frame sparsity of 20%), we present the 3D scene reconstruction results for the experimental groups (“PC-NeRF” + “two-step”) from Tab. I, showcasing the superior performance of our proposed PC-NeRF over NeRF-LOAM in 3D scene reconstruction accuracy on the KITTI and MaiCity datasets. This superior performance is because the network structure of NeRF-LOAM is lightweight, prioritizing the real-time implementation of the algorithm while slightly reducing the 3D scene reconstruction accuracy. In summary, applying the two-step depth inference method to PC-NeRF yields the highest quantitative results.

For large-scale evaluation, according to Sec. III-A, we divide the KITTI 03 sequence (800 consecutive frames, $555 \times 120 \times 7.2 \text{ m}^3$) into 32 sequential blocks. The spatial extent of each block is about $61.5 \times 42.5 \times 3 \text{ m}^3$, overlapping with that of neighboring blocks. In the 30th block, our proposed PC-NeRF trains three epochs because one or two epochs of training result in low inference accuracy with two-step depth inference.

Additionally, our proposed PC-NeRF trains only one epoch on all the remaining 31 blocks. As a comparison, OriginalNeRF trains three epochs on all 32 blocks. It is mentioned in the previous two paragraphs that the two-step depth inference method for OriginalNeRF may result in invalid inferences for some LiDAR rays, whereas for PC-NeRF, it performs well. Therefore, in Tab. III, we use the one-step depth inference method for OriginalNeRF and apply the two-step depth inference method for PC-NeRF. From Tab. III and Fig. 6, it is evident that our proposed PC-NeRF is trained to achieve high accuracy in novel LiDAR view synthesis and single-frame 3D reconstruction. This further underscores its promising potential for deployment in large-scale environments.

C. Evaluation for Novel LiDAR View Synthesis and 3D Reconstruction using Sparse LiDAR Frames

To assess the novel LiDAR view synthesis and 3D reconstruction capabilities of our proposed PC-NeRF using sparse LiDAR frames, we perform experiments on the MaiCity and KITTI datasets, treating 50 consecutive frames as a scene. Each PC-NeRF model is trained and evaluated based on a single scene. The sparse LiDAR frames are obtained by adjusting the frame sparsity, as described in the second paragraph of Sec. IV-A. As seen in Tab. IV, Fig. 1, and Fig. 7, with only one epoch training, our proposed PC-NeRF can reconstruct environments very well, even if the frame sparsity reaches 67%. In Tab. IV, specifically from rows 2 to 25, it is evident that as the proportion of the training set in the entire dataset gradually decreases from 4/5 to 1/10, indicating an increase in training difficulty, the three scenes start requiring more epochs at 75% (> 67%), 67%, and 67% of frame sparsity, respectively, to achieve valid inference by the two-step depth inference over all LiDAR rays. With no more than ten training epochs, our proposed PC-NeRF performs well even when the frame sparsity is 90%. Considering the training time, memory consumption, and inference results for these scenes, a frame sparsity of 67% emerges as a crucial threshold. Therefore, we subsequently perform additional tests on the KITTI and MaiCity datasets, as depicted in rows 26-31 of Tab. IV, to further verify PC-NeRF’s performance at this crucial threshold. By employing only one of every three LiDAR frames (frame sparsity at 67%) for training, our proposed PC-NeRF maintains robust reliability on these two datasets with just one epoch training. This robust reliability is because our PC-NeRF model concurrently optimizes scene-level, segment-level, and point-level scene representations, capable of utilizing the available point cloud information as efficiently as possible.

Additionally, we report the 3D scene reconstruction outcomes for the experimental groups of Tab. IV in Tab. II for the same frame sparsities. Tab. II shows that our proposed PC-NeRF achieves higher 3D scene reconstruction accuracy than NeRF-LOAM on the KITTI and MaiCity datasets at different frame sparsity. Therefore, our proposed PC-NeRF effectively tackles the challenge posed by sparse LiDAR frames, improving the accuracy of 3D reconstruction using sparse LiDAR frames.

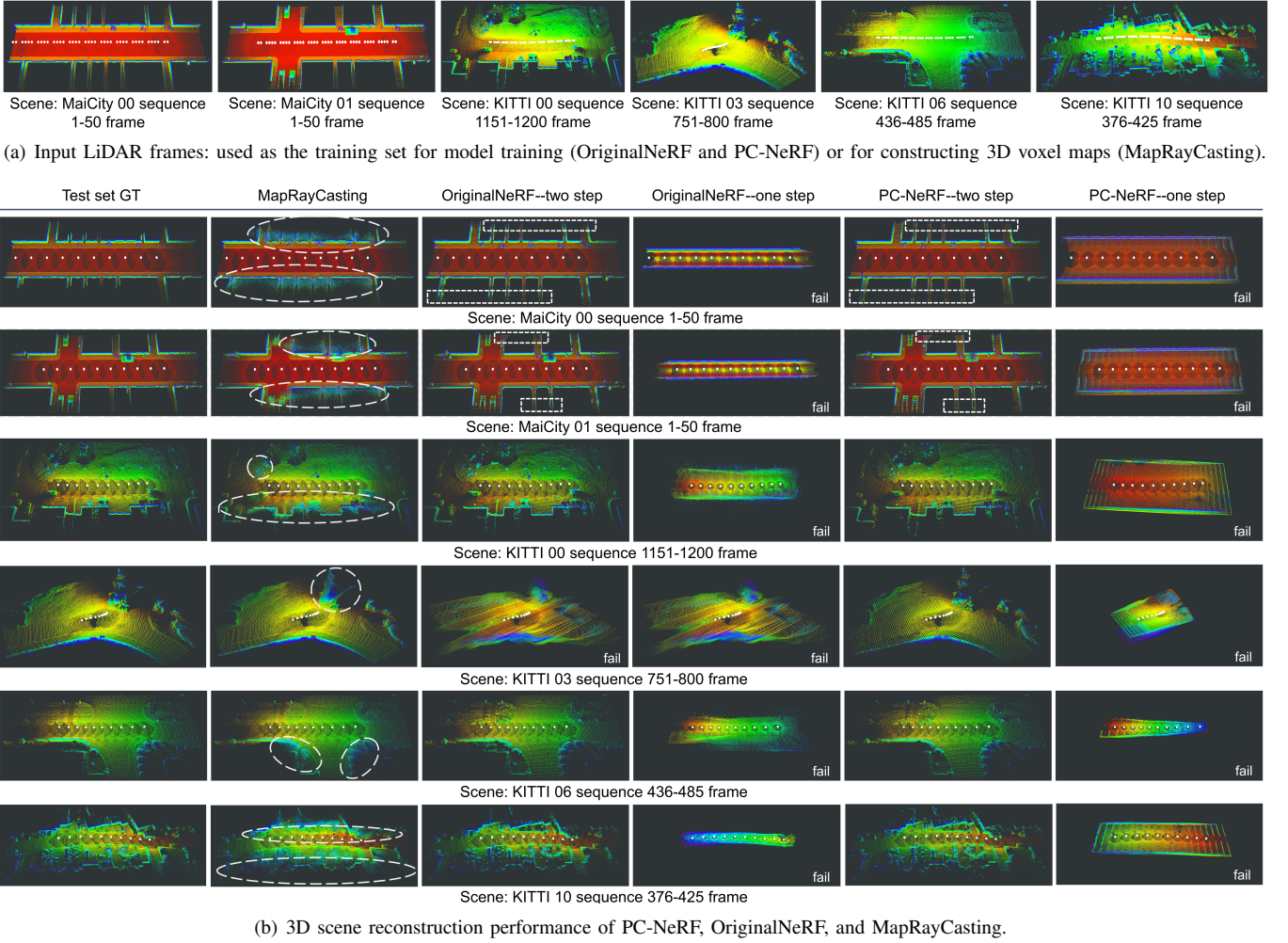


Fig. 5. 3D scene reconstruction on the MaiCity and KITTI datasets. Each subfigure represents the result of concatenating multiple LiDAR frames using real poses. The white dots in each subfigure represent the LiDAR positions of each frame. In Fig. (b): “one/two-step” denotes the one/two-step depth inference method. Fig. (b) illustrates the inference results corresponding to Tab. I. The “PC-NeRF--two step” column in Fig. (b) corresponds to the rows with 20 % frame sparsity in Tab. IV.

D. Ablation Study

To validate the effectiveness of our method’s components, we conduct ablation experiments probing various aspects of the proposed PC-NeRF. These experiments are tailored to scrutinize the individual components’ functionality rather than determining the optimal parameter set. We conduct extensive ablation tests and select the results on the KITTI 00 sequence 1151-1200 frame scene and the KITTI 01 sequence 1001-1050 frame scene for analysis. We train for one epoch and employ the two-step depth inference method for depth inference.

Effect of parent NeRF depth loss: As shown in Tab. V, on the KITTI 00 sequence 1151-1200 frame (frame sparsity = 20 %) scene, increasing λ_{pd} appropriately improves the accuracy of novel LiDAR view synthesis and single-frame 3D reconstruction. Additionally, compared to $\lambda_{cf} = 10^6$ and $\lambda_{cd} = 10^5$, the lower value of $\lambda_{pd} = 100$ is the ideal choice. When increasing the frame sparsity to 67 %, i.e., drastically reducing the number of laser points for training the PC-NeRF model, maintaining λ_{pd} to a small value (e.g., 0 and 1) achieves a relatively pleasing result on the KITTI

00 sequence 1151-1200 frame scene. This relatively pleasing result is because the parent NeRF free loss proposed in Sec. III-C is used to supervise the volumetric distribution over the entire scene space, relying on massive actual LiDAR points to obtain the volumetric distribution within the entire scene space. Therefore, to tackle the sparsity of the LiDAR frames, it is not appropriate to set too larger proportion of the parent NeRF free in the total loss. Besides, with $\lambda_{pd} = 1$, $\lambda_{cf} = 10^6$, and $\lambda_{cd} = 10^5$, we have gotten considerable results with one-epoch training on 13 sequences from MaiCity and KITTI datasets in Sec. IV-B and Sec. IV-C. In summary, in order to tackle the scene complexity and effectively utilize the sparse LiDAR frames, we set λ_{pd} to 1 to ensure that our proposed PC-NeRF achieves a relatively high novel LiDAR view synthesis and single-frame 3D reconstruction accuracy in different scenes through fast training.

Effect of child NeRF free loss and child NeRF depth loss: In our proposed PC-NeRF, the child NeRF free loss controlled by λ_{cf} optimizes the scene-level and segment-level environmental representation. As shown in rows 2 to 6 of Tab. VI, the child NeRF free loss helps improve novel LiDAR

TABLE I
NOVEL LiDAR VIEW SYNTHESIS AND SINGLE-FRAME 3D RECONSTRUCTION ON SMALL-SCALE SCENES (FRAME SPARSITY: 20 %, BOLD: BEST RESULTS FOR THE CORRESPONDING SCENE/DATASET)

Dataset	Method	Train epoch	Train time/ epoch [min]	Memory consumption	Inference method	Dep. Err. [m] ↓	Dep. Acc@ 0.2m [%] ↑	CD [m] ↓	F@ 0.2m ↑
MaiCity 00-01 sequence	MapRayCasting	-	-	5.8 MB	-	0.390	81.760	0.261	0.863
	OriginalNeRF	1	71	12.1 MB	one-step	3.581	0.154	3.336	0.000
				12.1 MB + 300.7 KB	two-step	0.505	85.008	0.265	0.921
	PC-NeRF	1	92	12.1 MB	one-step	1.890	2.382	1.367	0.063
				12.1 MB + 300.7 KB	two-step	0.347	87.877	0.179	0.945
KITTI 00-10 sequence	MapRayCasting	-	-	17.5 MB	-	0.461	68.754	0.265	0.837
	OriginalNeRF	1	58	12.1 MB	one-step	4.111	1.944	2.803	0.078
				12.1 MB + 917.6 KB	two-step	0.749	52.658	0.297	0.818
	PC-NeRF	1	76	12.1 MB	one-step	5.366	1.281	4.478	0.039
				12.1 MB + 917.6 KB	two-step	0.592	59.149	0.244	0.863

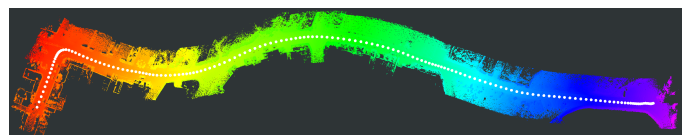
TABLE II
3D SCENE RECONSTRUCTION ON THE KITTI AND MAICITY DATASET (TYPE OF DEPTH INFERENCE FOR PC-NeRF: TWO-STEP)

Dataset	Method	Frame sparsity [%]	Acc. [m] ↓	Comp. [m] ↓	Map CD ↓ [m]	Map F@ ↑ 0.2m
MaiCity 00-01 sequence	NeRF-LOAM	20	0.029	0.032	0.030	99.253
		67	0.026	0.066	0.046	98.045
		90	0.027	0.122	0.074	95.777
	PC-NeRF	20	0.023	0.028	0.025	99.836
		67	0.023	0.017	0.020	99.883
		90	0.024	0.025	0.024	99.500
KITTI 00-10 sequence	NeRF-LOAM	20	0.062	0.174	0.118	86.477
		67	0.060	0.182	0.121	86.253
		90	0.060	0.265	0.162	79.661
	PC-NeRF	20	0.037	0.090	0.063	96.141
		67	0.034	0.070	0.052	97.215
		90	0.032	0.069	0.051	96.868

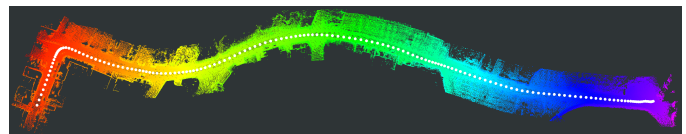
TABLE III
NOVEL LiDAR VIEW SYNTHESIS AND SINGLE-FRAME 3D RECONSTRUCTION ON LARGE-SCALE SCENES (FRAME SPARSITY: 20 %)

Method	Depth inference	Dep. Err. [m] ↓	Dep. Acc@ 0.2m [%] ↑	CD ↓ [m]	F@ 0.2m ↑
MapRayCasting	-	0.898	34.812	0.452	0.729
OriginalNeRF	one-step	2.441	8.290	1.438	0.271
PC-NeRF	two-step	0.658	43.556	0.273	0.834

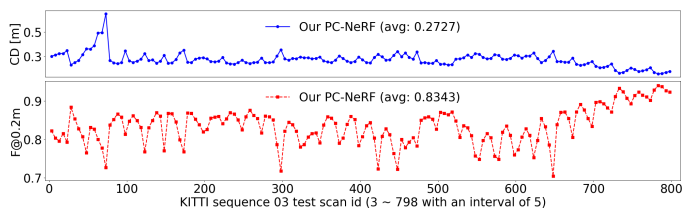
view synthesis and single-frame 3D reconstruction accuracy within a specific range of λ_{cf} variation. However, too high or too low λ_{cf} reduces the accuracy. Based on child NeRF free loss, child NeRF depth loss controlled by λ_{cd} and λ_{in} is used to further optimize the point-level and segment-level scene representations. λ_{cd} and λ_{in} are independent, so many pairings of λ_{cd} and λ_{in} need to be experimentally verified. Here, we simplify the pairing by taking $\lambda_{cd} = \lambda_{cf} \times \lambda_{in}$. As shown in rows 7 to 10 of Tab. VI, unfortunately, adding child NeRF depth loss does not improve the novel LiDAR



(a) Real LiDAR point clouds of all test frames.



(b) Reconstructed LiDAR point clouds by our proposed PC-NeRF.



(c) Single-frame 3D reconstruction accuracy of our proposed PC-NeRF.

Fig. 6. 3D reconstruction effect using our proposed PC-NeRF on KITTI 03 sequence. The white dots in Fig. 6(a) and Fig. 6(b) indicate the LiDAR position of each frame, corresponding to each data point in Fig. 6(c). As qualitative and quantitative results are shown in the three sub-figures, our proposed PC-NeRF has high 3D reconstruction accuracy in the KITTI 03 sequence.

view synthesis and 3D reconstruction accuracy. So, we further explore the smooth transition between child NeRF free loss and child NeRF depth loss in the following paragraph.

Effect of smooth transition between child NeRF free loss and child NeRF depth loss: As shown in rows 11 to 16 of Tab. VI, enlarging the smooth transition interval γ between child NeRF free loss and child NeRF depth loss improves the novel LiDAR view synthesis and single-frame 3D reconstruction accuracy. At the same time, a sizeable smooth transition interval decreases the accuracy. The accuracy decreases because, when the smoothing transition interval is too large, the child NeRF depth loss needs to supervise

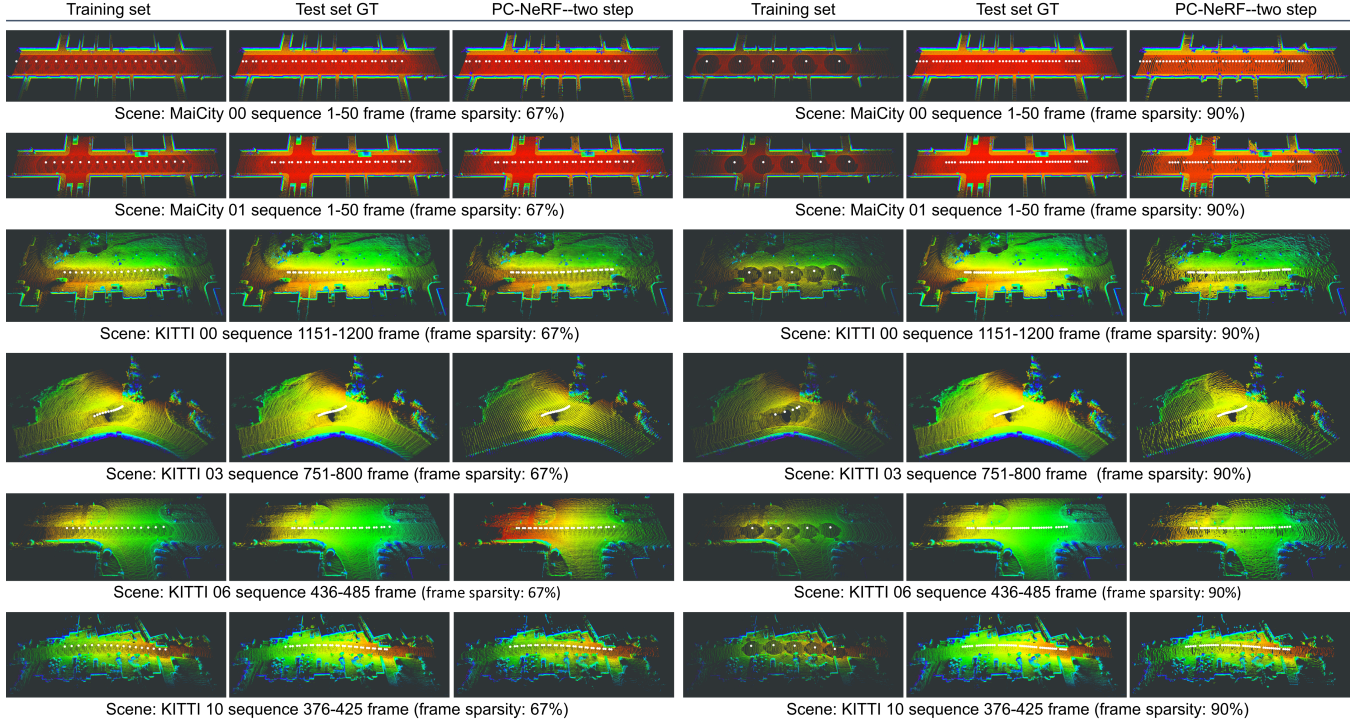


Fig. 7. 3D scene reconstruction using sparse LiDAR frames on the MaiCity and KITTI datasets. “two-step” denotes the two-step depth inference method. Each subfigure represents the result of concatenating multiple LiDAR frames using real poses. The white dots in each subfigure represent the LiDAR positions of each frame.

almost the sampling points on the entire LiDAR ray and cannot effectively supervise the sampling points around the child NeRF near and far bounds, i.e., around the real objects.

Effect of two-step depth inference: In Fig. 5 and Tab. I, two-step depth inference often outperforms one-step depth inference. Moreover, the extensive experiments in Sec. IV-B, Sec. IV-C, and Sec. IV-D demonstrate that our proposed PC-NeRF training and two-step depth inference methods are highly compatible and consistently robust.

V. CONCLUSION

This paper proposes a parent-child neural radiance fields (PC-NeRF) framework for large-scale 3D scene reconstruction and novel LiDAR view synthesis optimized for efficiently utilizing temporally sparse LiDAR frames in outdoor autonomous driving. PC-NeRF proposes a hierarchical spatial partitioning approach to divide the autonomous vehicle driving environment into large blocks, referred to as parent NeRFs, and subsequently subdivide each block into geometric segments, represented by child NeRFs. A parent NeRF shares a network with the child NeRFs within it. Leveraging the hierarchical spatial partitioning approach, PC-NeRF introduces a comprehensive multi-level scene representation. This representation is crafted to collectively optimize scene-level, segment-level, and point-level features, enabling efficient utilization of sparse LiDAR frames. Besides, we propose a two-step depth inference method to realize segment-to-point inference. Our proposed PC-NeRF is validated with extensive experiments to achieve high-precision novel LiDAR view synthesis and 3D reconstruction in large-scale scenes.

Furthermore, PC-NeRF demonstrates significant deployment potential, achieving notable accuracy in novel LiDAR view synthesis and 3D reconstruction with just one epoch of training on most test scenes from the KITTI and MaiCity datasets. Most importantly, PC-NeRF can tackle practical situations with temporally sparse LiDAR frames. Our future work will further explore the potential of PC-NeRF in object detection and localization for autonomous driving.

REFERENCES

- [1] L. Chen, Y. Li, C. Huang, B. Li, Y. Xing, D. Tian, L. Li, Z. Hu, X. Na, Z. Li, C. Lv, J. Wang, D. Cao, N. Zheng, and F.-Y. Wang, “Milestones in autonomous driving and intelligent vehicles: Survey of surveys,” *IEEE Transactions on Intelligent Vehicles*, vol. 8, no. 2, pp. 1046–1056, 2023, doi:[10.1109/TIV.2022.3223131](https://doi.org/10.1109/TIV.2022.3223131).
- [2] S. Teng, X. Hu, P. Deng, B. Li, Y. Li, Y. Ai, D. Yang, L. Li, Z. Xuanyuan, F. Zhu, and L. Chen, “Motion planning for autonomous driving: The state of the art and future perspectives,” *IEEE Transactions on Intelligent Vehicles*, vol. 8, no. 6, pp. 3692–3711, 2023, doi:[10.1109/TIV.2023.3274536](https://doi.org/10.1109/TIV.2023.3274536).
- [3] Z. Li and J. Zhu, “Point-based neural scene rendering for street views,” *IEEE Transactions on Intelligent Vehicles*, 2023, doi:[10.1109/TIV.2023.3304347](https://doi.org/10.1109/TIV.2023.3304347).
- [4] 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, doi:[10.1109/ICRA48891.2023.10160907](https://doi.org/10.1109/ICRA48891.2023.10160907).
- [5] Y. Ran, J. Zeng, S. He, J. Chen, L. Li, Y. Chen, G. Lee, and Q. Ye, “Neurar: Neural uncertainty for autonomous 3d reconstruction with implicit neural representations,” *IEEE Robotics and Automation Letters*, vol. 8, no. 2, pp. 1125–1132, 2023, doi:[10.1109/LRA.2023.3235686](https://doi.org/10.1109/LRA.2023.3235686).
- [6] 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 (ICCV)*, October 2023, pp. 8218–8227, doi:[10.48550/arXiv.2303.10709](https://doi.org/10.48550/arXiv.2303.10709).

TABLE IV
NOVEL LiDAR VIEW SYNTHESIS AND SINGLE-FRAME 3D RECONSTRUCTION OF PC-NeRF USING SPARSE LiDAR FRAMES ON THE MAICITY AND KITTI DATASETS (TYPE OF DEPTH INFERENCE FOR PC-NeRF: TWO-STEP, BOLDED: BEST RESULTS WHEN TRAINING ONE EPOCH)

Scene/Dataset		Training set proportion	Frame sparsity [%]	Train epoch	Train time/ epoch [min]	Memory consumption	Dep. Err. [m] ↓	Dep. Acc@ 0.2m [%] ↑	CD [m] ↓	F@ 0.2m ↑	
Scene	MaiCity 00 sequence 1-50 frame	4/5	20	1	90.0	12.1 MB + 275.1 KB	0.303	88.956	0.172	0.955	
		3/4	25	1	75.7	12.1 MB + 274.1 KB	0.287	89.159	0.166	0.957	
		2/3	33	1	66.3	12.1MB + 273.0 KB	0.328	88.343	0.192	0.946	
		1/2	50	1	56.1	12.1 MB + 268.9 KB	0.245	88.810	0.143	0.953	
		1/3	67	1	38.0	12.1 MB + 258.7 KB	0.189	90.168	0.109	0.961	
		1/4	75	1	31.0	12.1 MB + 248.0 KB	0.160	91.219	0.097	0.969	
	KITTI 00 sequence 1151-1200 frame	1/5	80	5	22.4	12.1 MB + 231.2 KB	0.177	90.283	0.106	0.962	
		1/10	90	5	11.3	12.1 MB + 222.3 KB	0.163	88.747	0.108	0.946	
		4/5	20	1	76.0	12.1 MB + 736.1 KB	0.488	66.654	0.224	0.891	
		3/4	25	1	55.0	12.1 MB + 734.8 KB	0.443	69.213	0.206	0.900	
		2/3	33	1	49.0	12.1 MB + 735.0 KB	0.447	69.547	0.206	0.900	
		1/2	50	1	35.6	12.1 MB + 702.8 KB	0.461	68.293	0.213	0.898	
	KITTI 06 sequence 436-485 frame	1/3	67	1	24.5	12.1 MB + 676.3 KB	0.439	72.185	0.197	0.909	
		1/4	75	2	24.0	12.1 MB + 662.2 KB	0.433	71.963	0.202	0.906	
		1/5	80	5	14.9	12.1 MB + 621.3 KB	0.423	72.895	0.197	0.908	
		1/10	90	10	7.3	12.1 MB + 550.2 KB	0.412	73.620	0.197	0.904	
		4/5	20	1	68.0	12.1 MB + 693.7 KB	0.350	72.420	0.209	0.894	
		3/4	25	1	64.4	12.1 MB + 686.6 KB	0.337	73.829	0.206	0.898	
	Dataset	MaiCity 00-01 sequence	2/3	33	1	44.6	12.1 MB + 687.6 KB	0.330	74.303	0.204	0.902
			1/2	50	1	32.8	12.1 MB + 653.2 KB	0.332	74.564	0.203	0.901
1/3			67	1	22.4	12.1 MB + 639.1 KB	0.393	74.556	0.216	0.898	
KITTI 00-10 sequence		1/4	75	2	22.0	12.1 MB + 615.1 KB	0.312	75.942	0.196	0.905	
		1/5	80	10	12.3	12.1 MB + 589.5 KB	0.323	74.808	0.200	0.905	
		1/10	90	10	7.6	12.1 MB + 505.5 KB	0.311	75.424	0.196	0.902	

- [7] H. Turki, D. Ramanan, and M. Satyanarayanan, “Mega-nerf: Scalable construction of large-scale nerfs for virtual fly-throughs,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 12 922–12 931, doi:[10.1109/CVPR52688.2022.01258](https://doi.org/10.1109/CVPR52688.2022.01258).
- [8] Y. Chang, K. Ebadi, C. E. Denniston, M. F. Ginting, A. Rosinol, A. Reinke, M. Palieri, J. Shi, A. Chatterjee, B. Morrell *et al.*, “Lamp 2.0: A robust multi-robot slam system for operation in challenging large-scale underground environments,” *IEEE Robotics and Automation Letters*, vol. 7, no. 4, pp. 9175–9182, 2022, doi:[10.1109/LRA.2022.3191204](https://doi.org/10.1109/LRA.2022.3191204).
- [9] H. Parr, C. Harvey, and G. Burnett, “Investigating levels of remote operation in high-level on-road autonomous vehicles using operator sequence diagrams,” 2023, doi:[10.21203/rs.3.rs-2510863/v1](https://doi.org/10.21203/rs.3.rs-2510863/v1).
- [10] J. Zhang, J. Pu, J. Xue, M. Yang, X. Xu, X. Wang, and F.-Y. Wang, “Hivegpt: human-machine-augmented intelligent vehicles with generative pre-trained transformer,” *IEEE Transactions on Intelligent Vehicles*, 2023, doi:[10.1109/TIV.2023.3256982](https://doi.org/10.1109/TIV.2023.3256982).
- [11] J. Wang, Z. Wang, B. Yu, J. Tang, S. L. Song, C. Liu, and Y. Hu, “Data fusion in infrastructure-augmented autonomous driving system: Why? where? and how?” *IEEE Internet of Things Journal*, 2023, doi:[10.1109/IJOT.2023.3266247](https://doi.org/10.1109/IJOT.2023.3266247).
- [12] Z. Song, Z. He, X. Li, Q. Ma, R. Ming, Z. Mao, H. Pei, L. Peng, J. Hu, D. Yao *et al.*, “Synthetic datasets for autonomous driving: A survey,” *arXiv preprint arXiv:2304.12205*, 2023, doi:[10.48550/arXiv.2304.12205](https://doi.org/10.48550/arXiv.2304.12205).
- [13] Y. Wang, L. Xu, F. Zhang, H. Dong, Y. Liu, and G. Yin, “An adaptive fault-tolerant ekf for vehicle state estimation with partial missing measurements,” *IEEE/ASME Transactions on Mechatronics*, vol. 26, no. 3, pp. 1318–1327, 2021, doi:[10.1109/TMECH.2021.3065210](https://doi.org/10.1109/TMECH.2021.3065210).
- [14] I. Raouf, A. Khan, S. Khalid, M. Sohail, M. M. Azad, and H. S. Kim, “Sensor-based prognostic health management of advanced driver assistance system for autonomous vehicles: A recent survey,” *Mathematics*, vol. 10, no. 18, p. 3233, 2022, doi:[10.3390/math10183233](https://doi.org/10.3390/math10183233).
- [15] M. Waqas and P. Ioannou, “Automatic vehicle following under safety, comfort, and road geometry constraints,” *IEEE Transactions on Intelligent Vehicles*, vol. 8, no. 1, pp. 531–546, 2022, doi:[10.1109/TIV.2022.3177176](https://doi.org/10.1109/TIV.2022.3177176).
- [16] A. Hornung, K. M. Wurm, M. Bennewitz, C. Stachniss, and W. Burgard, “Octomap: An efficient probabilistic 3d mapping framework based on octrees,” *Autonomous robots*, vol. 34, pp. 189–206, 2013, doi:[10.1007/s10514-012-9321-0](https://doi.org/10.1007/s10514-012-9321-0).
- [17] X. Hu, G. Xiong, J. Ma, G. Cui, Q. Yu, S. Li, and Z. Zhou, “A non-uniform quadtree map building method including dead-end semantics extraction,” *Green Energy and Intelligent Transportation*, vol. 2, no. 2, p. 100071, 2023, doi:[10.1016/j.geits.2023.100071](https://doi.org/10.1016/j.geits.2023.100071).
- [18] 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, doi:[10.1109/ICRA48506.2021.9562069](https://doi.org/10.1109/ICRA48506.2021.9562069).
- [19] X. Yang, G. Lin, Z. Chen, and L. Zhou, “Neural vector fields: Implicit representation by explicit learning,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 16 727–16 738, doi:[10.1109/CVPR52729.2023.01605](https://doi.org/10.1109/CVPR52729.2023.01605).
- [20] D. Yu, M. Lau, L. Gao, and H. Fu, “Sketch beautification: Learning part beautification and structure refinement for sketches of man-made objects,” *arXiv preprint arXiv:2306.05832*, 2023, doi:[10.48550/arXiv.2306.05832](https://doi.org/10.48550/arXiv.2306.05832).
- [21] P. Wang, L. Liu, Y. Liu, C. Theobalt, T. Komura, and W. Wang,

TABLE V
ABLATION STUDY OF PARENT NERF DEPTH LOSS

Scene	λ_{pd}	Dep. Err. [m] ↓	Dep. Acc@ 0.2m [%] ↑	CD [m] ↓	F@ 0.2m ↑
KITTI	0	0.4892	66.6861	0.2243	0.8910
00 sequence	1 (our)	0.4884	66.6541	0.2239	0.8908
1151-1200	10	0.4895	66.6587	0.2241	0.8908
frame	100	0.4675	69.5239	0.2077	0.8993
(frame sparsity = 20 %)	1.0e3	0.4683	69.7551	0.2089	0.8989
	1.0e4	0.4729	69.9102	0.2079	0.8992
	1.0e5	0.5307	67.6254	0.2240	0.8931
KITTI	0	0.4384	72.1457	0.1966	0.9086
00 sequence	1 (our)	0.4390	72.1848	0.1967	0.9086
1151-1200	10	0.4400	72.1695	0.1969	0.9085
frame	100	0.4454	72.1620	0.1988	0.9078
(frame sparsity = 67 %)	1.0e3	0.4540	72.4099	0.2005	0.9062
	1.0e4	0.4443	72.5878	0.1969	0.9091
	1.0e5	0.4977	70.7471	0.2112	0.9035
KITTI	0	1.3424	39.6964	0.3521	0.7422
01 sequence	1 (our)	1.3358	39.6509	0.3501	0.7439
1001-1050	10	1.3435	39.6998	0.3526	0.7417
frame	100	1.3490	39.7598	0.3557	0.7384
(frame sparsity = 20 %)	1.0e3	1.3259	39.9901	0.3567	0.7361
	1.0e4	1.2874	38.8882	0.3596	0.7266
	1.0e5	1.2785	38.6314	0.3580	0.7268

Note: $\lambda_{cf} = 1.0e6$, $\lambda_{cd} = 1.0e5$, $\lambda_{in} = 0.1$, $\gamma = 2$ m

TABLE VI
ABLATION STUDY OF CHILD NERF FREE LOSS AND CHILD NERF DEPTH LOSS ON THE KITTI 00 SEQUENCE 1151-1200 FRAME SCENE ($\lambda_{pd} = 1$, FRAME SPARSITY: 20 %.)

λ_{cf}	λ_{cd}	λ_{in}	γ [m]	Dep. Err. [m] ↓	Dep. Acc@ 0.2m [%] ↑	CD [m] ↓	F@ 0.2m ↑
0				0.511	65.232	0.220	0.890
1				0.507	66.129	0.217	0.894
1.0e3	0	0	0	0.472	69.548	0.209	0.899
1.0e6				0.465	67.109	0.219	0.894
1.0e9				0.466	67.093	0.219	0.893
	2.5e4	0.025		0.480	67.001	0.222	0.891
	5.0e4	0.05		0.494	66.819	0.224	0.890
	1.0e5	0.1		0.509	66.548	0.228	0.887
	2.0e5	0.2		0.538	65.988	0.234	0.884
			0.5	0.511	66.430	0.228	0.887
			1	0.513	65.970	0.229	0.887
			2 (our)	0.488	66.654	0.224	0.891
			3	0.447	68.697	0.207	0.899
			5	0.452	69.344	0.206	0.899
			10	0.468	70.269	0.209	0.899

“Neus: Learning neural implicit surfaces by volume rendering for multi-view reconstruction,” *arXiv preprint arXiv:2106.10689*, 2021, doi:10.48550/arXiv.2106.10689.

- [22] 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, doi:10.1145/3503250.
- [23] S. Liu and J. Zhu, “Efficient map fusion for multiple implicit slam agents,” *IEEE Transactions on Intelligent Vehicles*, 2023, doi:10.1109/ITV.2023.3297194.

- [24] X. Chen, Z. Song, J. Zhou, D. Xie, and J. Lu, “Camera and lidar fusion for urban scene reconstruction and novel view synthesis via voxel-based neural radiance fields,” *Remote Sensing*, vol. 15, no. 18, p. 4628, 2023, doi:10.3390/rs15184628.
- [25] 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, doi:10.1109/ICCV48922.2021.00617.
- [26] A. Moreau, N. Piasco, D. Tsishkou, B. Stanculescu, and A. de La Fortelle, “Lens: Localization enhanced by nerf synthesis,” in *Conference on Robot Learning*. PMLR, 2022, pp. 1347–1356, doi:10.48550/arXiv.2110.06558.
- [27] 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. 12786–12796, doi:10.1109/CVPR52688.2022.01245.
- [28] X. Yu, Y. Liu, S. Mao, S. Zhou, R. Xiong, Y. Liao, and Y. Wang, “Nf-atlas: Multi-volume neural feature fields for large scale lidar mapping,” *arXiv preprint arXiv:2304.04624*, 2023, doi:10.48550/arXiv.2304.04624.
- [29] T. Tao, L. Gao, G. Wang, P. Chen, D. Hao, X. Liang, M. Salzmann, and K. Yu, “Lidar-nerf: Novel lidar view synthesis via neural radiance fields,” 2023, doi:10.48550/arXiv.2304.10406.
- [30] J. Zhang, F. Zhang, S. Kuang, and L. Zhang, “Nerf-lidar: Generating realistic lidar point clouds with neural radiance fields,” *arXiv preprint arXiv:2304.14811*, 2023, doi:10.48550/arXiv.2304.14811.
- [31] S. Huang, Z. Gojcic, Z. Wang, F. Williams, Y. Kasten, S. Fidler, K. Schindler, and O. Litany, “Neural lidar fields for novel view synthesis,” *arXiv preprint arXiv:2305.01643*, 2023, doi:10.48550/arXiv.2305.01643.
- [32] H. Kuang, X. Chen, T. Guadagnino, N. Zimmerman, J. Behley, and C. Stachniss, “Ir-mcl: Implicit representation-based online global localization,” *IEEE Robotics and Automation Letters*, vol. 8, no. 3, pp. 1627–1634, 2023, doi:10.1109/LRA.2023.3239318.
- [33] L. Wiesmann, T. Guadagnino, I. Vizzo, N. Zimmerman, Y. Pan, H. Kuang, J. Behley, and C. Stachniss, “Locndf: Neural distance field mapping for robot localization,” *IEEE Robotics and Automation Letters*, 2023, doi:10.1109/LRA.2023.3291274.
- [34] K. Rematas, A. Liu, P. P. Srinivasan, J. T. Barron, A. Tagliasacchi, T. Funkhouser, and V. Ferrari, “Urban radiance fields,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 12932–12942, doi:10.1109/CVPR52688.2022.01259.
- [35] M. Tancik, V. Casser, X. Yan, S. Pradhan, B. Mildenhall, P. P. Srinivasan, J. T. Barron, and H. Kretschmar, “Block-nerf: Scalable large scene neural view synthesis,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 8248–8258, doi:10.1109/CVPR52688.2022.00807.
- [36] M. Zhenxing and D. Xu, “Switch-nerf: Learning scene decomposition with mixture of experts for large-scale neural radiance fields,” in *The Eleventh International Conference on Learning Representations*, 2022. [Online]. Available: <https://openreview.net/forum?id=PQ2zoIZqvm>
- [37] D. Rebain, W. Jiang, S. Yazdani, K. Li, K. M. Yi, and A. Tagliasacchi, “Derf: Decomposed radiance fields,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 14153–14161, doi:10.1109/CVPR46437.2021.01393.
- [38] C. Reiser, S. Peng, Y. Liao, and A. Geiger, “Kilonerf: Speeding up neural radiance fields with thousands of tiny mlps,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 14335–14345, doi:10.1109/ICCV48922.2021.01407.
- [39] L. Liu, J. Gu, K. Zaw Lin, T.-S. Chua, and C. Theobalt, “Neural sparse voxel fields,” *Advances in Neural Information Processing Systems*, vol. 33, pp. 15651–15663, 2020, doi:10.5555/3495724.3497037.
- [40] Y. Xiangli, L. Xu, X. Pan, N. Zhao, A. Rao, C. Theobalt, B. Dai, and D. Lin, “Bungeenerf: Progressive neural radiance field for extreme multi-scale scene rendering,” in *European conference on computer vision*. Springer, 2022, pp. 106–122, doi:10.1007/978-3-031-19824-3_7.
- [41] R. Martin-Brualla, N. Radwan, M. S. Sajjadi, J. T. Barron, A. Dosovitskiy, and D. Duckworth, “Nerf in the wild: Neural radiance fields for unconstrained photo collections,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 7210–7219.
- [42] J. Ma, J. Zhang, J. Xu, R. Ai, W. Gu, and X. Chen, “Overlaptransformer: An efficient and yaw-angle-invariant transformer network for lidar-based place recognition,” *IEEE Robotics and Automation Letters*, vol. 7, no. 3, pp. 6958–6965, 2022, doi:10.1109/LRA.2022.3178797.
- [43] K. Wang, T. Zhou, X. Li, and F. Ren, “Performance and challenges of 3d object detection methods in complex scenes for autonomous driving,”

- IEEE Transactions on Intelligent Vehicles*, vol. 8, no. 2, pp. 1699–1716, 2022, doi:[10.1109/TIV.2022.3213796](https://doi.org/10.1109/TIV.2022.3213796).
- [44] A. Milioto, I. Vizzo, J. Behley, and C. Stachniss, “Rangenet++: Fast and accurate lidar semantic segmentation,” in *2019 IEEE/RSJ international conference on intelligent robots and systems (IROS)*. IEEE, 2019, pp. 4213–4220, doi:[10.1109/IROS40897.2019.8967762](https://doi.org/10.1109/IROS40897.2019.8967762).
 - [45] 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.
 - [46] E. Haines, “Essential ray tracing,” *Glas89*, pp. 33–77, 1989.
 - [47] 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, doi:[10.1109/CVPR.2012.6248074](https://doi.org/10.1109/CVPR.2012.6248074).
 - [48] J. Behley, M. Garbade, A. Milioto, J. Quenzel, S. Behnke, C. Stachniss, and J. Gall, “Semantickitti: A dataset for semantic scene understanding of lidar sequences,” in *Proceedings of the IEEE/CVF international conference on computer vision*, 2019, pp. 9297–9307, doi:[10.1109/ICCV.2019.00939](https://doi.org/10.1109/ICCV.2019.00939).
 - [49] 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.
 - [50] S. Thrun, “Probabilistic robotics,” *Communications of the ACM*, vol. 45, no. 3, pp. 52–57, 2002, doi:[10.1145/504729.504754](https://doi.org/10.1145/504729.504754).
 - [51] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” *arXiv preprint arXiv:1412.6980*, 2014, doi:[10.48550/arXiv.1412.6980](https://doi.org/10.48550/arXiv.1412.6980).