# 360Roam: Real-Time Indoor Roaming Using Geometry-Aware $360°$ Radiance Fields

HUAJIAN HUANG, YINGSHU CHEN, TIANJIA ZHANG, and SAI-KIT YEUNG, Hong Kong University of Science and Technology, China

Fig. 1. *360Roam* is a system for immersive 6-DoF roaming in the indoor scenarios. It takes 360° images as input to learn geometry-aware omnidirectional radiance fields which consist of plentiful adaptive-assigned neural perceptrons (illustrated as different colors in the middle figure). It can render novel-view images at real-time speed on a single GPU GTX1080. The system also outputs the floorplan of the observed scene to enhance immersion for VR applications.

Neural radiance field (NeRF) has recently achieved impressive results in novel view synthesis. However, previous works on NeRF mainly focus on object-centric scenarios. In this work, we propose *360Roam*, a novel scene-level NeRF system that can synthesize images of large-scale indoor scenes in real time and support VR roaming. Our system first builds an omnidirectional neural radiance field (*360NeRF*) from multiple input 360° images. Using *360NeRF*, we then progressively estimate a 3D probabilistic occupancy map which represents the scene geometry in the form of spacial density. Skipping empty spaces and upsampling occupied voxels essentially allows us to accelerate volume rendering by using 360NeRF in a geometry-aware fashion. Furthermore, we use an adaptive divide-and-conquer strategy to slim and fine-tune the radiance fields for further improvement. The floorplan of the scene extracted from the occupancy map can provide guidance for ray sampling and facilitate a realistic roaming experience. To show the efficacy of our system, we collect a 360° image dataset in a large variety of scenes and conduct extensive experiments. Quantitative and qualitative comparisons among baselines illustrated our predominant performance in novel view synthesis for complex indoor scenes.

Authors' address: Huajian Huang, hhuangbg@connect.ust.hk; Yingshu Chen, yingshu2008@gmail.com; Tianjia Zhang, tzhangbl@connect.ust.hk; Sai-Kit Yeung, saikit@ust.hk, Hong Kong University of Science and Technology, Department of Computer Science and Engineering, Hong Kong, China.

CCS Concepts: • **Computing methodologies** → **Computational photography**; **Rendering**; **Virtual reality**.

Additional Key Words and Phrases: Novel view synthesis, Neural rendering, Photorealistic imagery

## 1 INTRODUCTION

Recently neural rendering has attracted great attention and demonstrated impressive rendering quality. Such learned view synthesis methods exploit neural networks to implicitly represent the structure and appearance of captured objects. The neural radiance field, pioneered by NeRF [Mildenhall et al. 2020] is currently the most promising path. Instead of pursuing precise geometry reconstruction, NeRF utilizes volume rendering techniques and multi-layer perceptrons to regress density and view-dependent color per ray. By densely sampling along each ray, it can generate photorealistic results even when encountered with light reflection, transparent objects, and thin structures. To reduce the long training and inference time, the following NeRF-based methods [Barron et al. 2021; Liu et al. 2020; Pumarola et al. 2021; Zhang et al. 2020] have made great efforts and modifications to increase rendering quality and speed. However, most of them focus on object-centric or small-scale scenes where placements of cameras are constrained during the

training and inference phases. GSN [DeVries et al. 2021] attempts to improve the limited performance of NeRF in open environments by introducing generative adversarial networks for scene-level radiance fields. Although it supports a four-degree-of-freedom (4-DoF) motion, the walk-though artifacts are noticeable, and the resolution of rendered images is too low.

In this work, we seek to explore the potential of neural radiance fields from object-centric scenarios into more challenging situations of complex indoor scene synthesis. Furthermore, we target real-time performance that will be essential for immersive applications such as VR roaming. Considering that the consumer-level 360° camera is becoming more and more accessible, we propose *360Roam*, an immersive roaming system that takes panorama sequences of complex scenes as input, generates high quality novel 360° images at real-time speed, and allows roaming with 6-DoF motion (Fig. 1).

*360Roam* consists of three stages following a coarse-to-fine process. First, a 360° neural radiance field (*360NeRF*) for the whole space is learned from all the input panoramas. Second, we use *360NeRF* to generate panoramic depth and uncertainty of the scene simultaneously, and then progressively update the occupied probability. This gives us a reliable occupancy map while reducing the reliance on unreliable hyperparameters (e.g., given the scene boundary). We further recover the floorplan (with objects) from the occupancy map. The floorplan enables collision detection and hence a more realistic indoor roaming, i.e., stops the camera from passing through the objects. More importantly, the geometric information effectively constrains the sampling range of the ray which allows us to slim a neural radiance field into numbers of sub-fields and fine-tune the radiance fields by skipping empty spaces and upsampling on occupied volume in the last stage. This geometry-aware radiance field essentially extends the render accelerating strategy from KiloNeRF [Reiser et al. 2021] to be scene-level and support adaptive radiance fields decomposing without manual intervention. See Fig. 2 for an overview of *360Roam*. To summarize, our contributions include:

- Proposing *360Roam*, an effective system for real-time 360° indoor-scene roaming that achieves 6-DoF flexibility.
- Proposing a practical pipeline that facilitates high-fidelity novel view synthesis of 360° panoramic images and reduces the reliance on hyperparameters.
- Collecting a dataset of real and synthetic scenes with 360° image sequences covering a diverse set of indoor environments, which facilitate future research of scene-level neural rendering. Relevant resources can be readily made available to the public.

## 2 RELATED WORK

### 2.1 Neural Radiance Field and Important Variants

NeRF [Mildenhall et al. 2020] is a promising and powerful neural scene representation for novel view synthesis. It uses the weight of a multi-layer perceptron to model the density and color of the scene as a function of continuous 5D coordinates and uses volume rendering to synthesize new views with no limitation on sampling rate. The input of NeRF is a combination of continuous coordinates without the need to discretize the scene and therefore, it

supports rendering images of arbitrary resolution. The introduction of the positional encoding mechanism could preserve high-frequency characteristics. Despite the outstanding visual performance, canonical NeRF suffers from many open issues which require more research such as dynamic scenes [Pumarola et al. 2021], varying lighting conditions [Martin-Brualla et al. 2021], fast training and rendering[Hedman et al. 2021; Piala and Clark 2021], general modeling[Wang et al. 2021] , few-shot learning [Yu et al. 2021b]. Until now, many researchers have proposed valuable variants to address these issues and have significantly expanded the application range and promoted robustness.

In our work, the most crucial issue is the efficiency of scene-level rendering to achieve real-time and immersive roaming. NeRF casts one camera ray per pixel and samples a set of points to query densities and radiance colors at those sampled points. For the pure implicit method, rendering the final color of each pixel requires hundreds of network queries, making it hard to complete the computation in a short time. Some methods[Lombardi et al. 2021; Xu et al. 2022; Yu et al. 2021a] exploit the sparse features of other representations such as mesh-based and primitive-based representations to improve efficiency. Other methods try to reduce the querying time and count. Neural Sparse Voxel Fields [Liu et al. 2020] speed up the rendering process using empty space skipping and early ray termination, which sample more points around entity surfaces and reduce queries. AutoInt [Lindell et al. 2021] replaced the numerical integration with a closed-form solution. DoNeRF [Neff et al. 2021] employs an extra network to predict sample locations directly and enforce those points closer to the intersection point of the surface and the camera ray. The Nvidia research team [Müller et al. 2022] reduced the training and rendering cost by embedding a multi-resolution hash table of learnable feature vectors into a small network which requires fewer floating-point operations. However, these methods are constrained to object-centric or relatively small scenes. In addition, KiloNeRF [Reiser et al. 2021] proves that using thousands of tiny independent MLPs to represent entire radiance fields can decrease the rendering time by three orders of magnitude. Their divide-and-conquer strategies have an underlying assumption that radiance fields are bounded by a box and its distribution is balanced, so they decompose the radiance fields into multiple volumes with same pre-defined resolution. However, different from object-centric scenes, the structures of complex scenes are irregular with the imbalanced spatial placement of objects. Balanced decomposing will waste a lot of neural perceptrons in an empty space. Therefore, our method introduces an adaptive partition scheme to make use of the perceptrons' capacity. Another related concurrent work, Block-NeRF [Tancik et al. 2022], also decomposes the whole scene into several blocks and uses one individual NeRF to represent each block. The main difference is that their rendered image is composited from rendering results of several block-NeRFs while we assign each block with a tinier NeRF network before the rendering process and synthesize the final image at once. Our method has less memory consumption and higher rendering efficiency in principle.
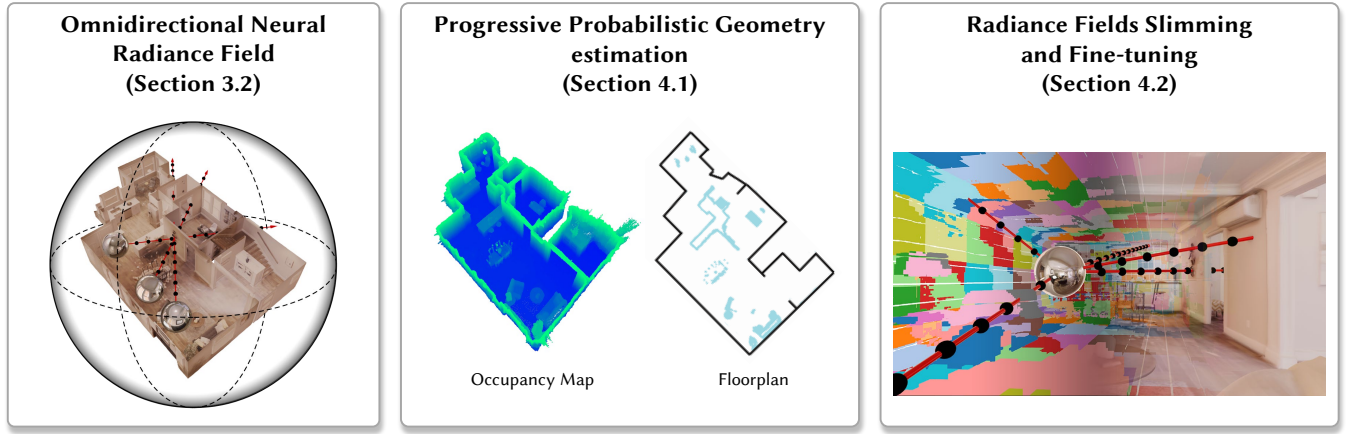
Fig. 2. Overview of the *360Roam* system. First, an omnidirectional neural radiance field (360NeRF) will be learned from a set of input 360° images. Then we sample 360NeRF at different camera positions to progressively estimate the probabilistic geometry and fuse them into an occupancy map from which a floorplan with objects is extracted. The *360NeRF* will then be slimmed and adaptively divided into multiple geometry-aware sub-fields that improves both the rendering quality and speed significantly.

## 2.2 Panorama Structure Estimation

The challenge of panorama view synthesis lies in the accurate estimation of the 3D structures of the scene. The entangled contents encoded in input images usually result in blurry artifacts on novel panoramas. We could exploit structure information as a precursor to regularize the model to output reasonable renderings. The depth and the room layout are two of the most common and easily obtained structure information. Methods for layout estimation are to either directly output the layout [Yang et al. 2019] or predict probability maps [Sun et al. 2019; Zou et al. 2018] to parse the scenes. Some works [Jin et al. 2020; Zeng et al. 2020; Zioulis et al. 2019] jointly learn the layout and panoramic depth and perform panoramic view synthesis as a proxy task. Inspired by [Xu et al. 2021], which uses room layout to guide the content generation process, we leverage the layout to clip the depth map rendered from NeRF, which could reduce the ambiguity and help render more photorealistic novel panoramas. As far as we know, we are the first to extend NeRF to panoramic view synthesis.

## 2.3 VR Photography

The omnidirectional image, which can cover 360° insight, is one of the most typical types of VR photography and thus can provide a more comprehensive view of the environment. It allows sparser views to perform novel view synthesis on panoramas. Some works have used RGBD panoramas [Serrano et al. 2019] including extra depth information obtained from other types of sensors for novel panorama synthesis. However, depth acquisition relies on expensive and fragile sensors and depth estimation suffers from noises and lighting conditions. We encourage synthesis methods that are supervised by only visual images. [Huang et al. 2017; Serrano et al. 2019] developed VR applications which support flexible viewing from 360° videos. Huang et al. [2017] utilized point clouds extracted from a 360° video to enable real-time video playback. Serrano et al. [2019] presented a layered representation for adding parallax

and real-time playback of 360° videos. Recently, inspired by the MPI representation[Zhou et al. 2018], Lin et al. [2020] and Attal et al. [2020] proposed multi-depth panorama (MDP) and multi-sphere image (MSI) representation, respectively, to conduct rendering from stereo 360° imagery. Previous works focused on processing existing captured 360° videos which cost lots of additional time to capture or post-process. Broxton et al. [2020] capture data using a customized low-cost hemispherical array made from 46 synchronized action sports cameras. The system can produce 6-DoF videos with a large baseline, fine resolution, wild field of view and high frame rates. The multi-sphere image (MSI) representation is used for light field video synthesis, yielding accurate results but involves a memory-heavy process. OmniPhotos [Bertel et al. 2020] reduced the capture time by attaching a 360° camera to a rotating selfie stick. It enables quick and casual capture of high-quality 360° panoramas with motion parallax. The visual rendering quality is improved by automatically and robustly reconstructing a scene-adaptive proxy geometry that reduces vertical distortions during image-based view synthesis. However, the user-centric VR methods often rely on dense capturing requirements and do not scale for room- or scene-level roaming.

## 3 OMNIDIRECTIONAL NEURAL RENDERING

### 3.1 Neural Radiance Field (NeRF)

NeRF [Mildenhall et al. 2020] represents the scene as a function of Cartesian coordinates, outputs densities and emitted radiance colors at those querying locations. The function could be modeled by simple MLP networks. Given a collection of $M$ images $\{I_m\}^M$ with corresponding camera poses $\{\Omega_{4\times4}^m\}^M$ in a world frame, the goal is to use the existing posed images to render new images from arbitrary known but unseen perspectives $\Omega^k \notin \{\Omega^m\}^M$. In order to get the representation, for each pixel in the input images, a camera ray is marched from its camera center through the scene within the near and far sampling range $[t_n, t_f]$. For any 3D point $\boldsymbol{p}_j(x, y, z)$ in the

ray, a global MLP model is applied to transform the concatenation of the point location $(x, y, z)$ and viewing direction (represented by a unit vector $\boldsymbol{d}$) into a 1-dimension density $(\sigma)$ and a 3-dimension RGB color $(c)$. Since the density is only dependent to point position, the model decomposes the mapping into predicting the density using only the location as well as predicting the color using both the location and viewing direction. Specifically, the mapping is represented as a two-stage MLP. The first 8 fully-connected layers of the network, denoted as $MLP_{1:8}$, take as input the 3D position $\boldsymbol{p}_j(x, y, z)$, and output the density $\sigma$ and a 256-dimension radiance feature $f_c$. The combination of the feature $f_c$ and viewing direction $\boldsymbol{d}$ is processed by the final fully connected layer $MLP_{9:}$ and decoded into the radiance $c(\boldsymbol{p}_j)$ in terms of RGB value. The MLP model can be represented as:

$$
\begin{aligned}
\sigma, f_c &= MLP_{1:8}(x, y, z) \\
c(\boldsymbol{p}_j) &= MLP_{9:}(f_c, \boldsymbol{d}) \, .
\end{aligned}
\tag{1}
$$

To improve rendering results, it utilizes positional encoding (which is borrowed from the Transformer[Vaswani et al. 2017]) and maps the input coordinates vectors into a high-dimension space. The mapping is useful in modeling high-frequency functions. For each coordinate component $x_i$ in both position and viewing direction, the encoding is formulated as:

$$
\mathcal{F}(x_i) = (\sin 2^0 \pi x_i, \cos 2^0 \pi x_i, ..., \sin 2^{D-1} \pi x_i, \cos 2^{D-1} \pi x_i) \, .
\tag{2}
$$

In addition, NeRF jointly optimizes two networks of same architecture as Eq. 1 with different parameters to increase rendering efficiency and avoid over-sampling in free spaces. A coarse network is evaluated on a set of $N_c$ locations using stratified sampling. Based on the output densities of the coarse network, the coefficients which could be represented as a function of density are normalized as weights to form a piecewise-constant PDF (probability density function). Additional $N_f$ points are sampled under this distribution. All $N_c + N_f$ points are used to query densities and colors through the fine network and render the final color of the ray.

## 3.2 Omnidirectional Neural Radiance Field (*360NeRF*)

The involved coordinate systems used in NeRF are all Cartesian while a panorama uses the panoramic pixel grid coordinate system, in which each pixel $(u, v)$ in the panoramic image corresponds to a point $(\phi, \theta)$ in a sphere surface represented by the spherical polar coordinate system. The transformation between the two coordinate systems is described by:

$$
\begin{aligned}
u &= \phi * W/(2\pi) + W/2 \\
v &= -\theta * H/\pi + H/2 \, ,
\end{aligned}
\tag{3}
$$

where $u, v$ denote column and row of the panorama, $\phi, \theta$ denote longitude and latitude of the spherical surface, and $H, W$ denote the height and width of the panorama respectively. Additionally, the relation between spherical and 3D Cartesian coordinate systems is:

$$
\begin{aligned}
x &= \cos \theta \sin \phi \\
y &= -\sin \theta \\
z &= \cos \theta \cos \phi \, .
\end{aligned}
\tag{4}
$$

The camera ray cast from a panoramic pixel can be formulated as $r = \boldsymbol{o} + r\boldsymbol{d}(\phi/f, \theta/f, 1)$, where $\boldsymbol{o}$ is the camera center, $f$ is the focal length, $\boldsymbol{d}$ is the ray direction and $r$ is the radial distance of a point in the ray. However, instead of transforming the coordinate to a 3D Cartesian coordinate, we use an alternative parameterization to model the omnidirectional neural radiance field, referred as *360NeRF*. 360NeRF conducts sampling in the spherical coordinate system and represents the points in inverse distance behavior, $\boldsymbol{p}_j(\phi, \theta, 1/r)$. And then 360NeRF queries a MLP for corresponding color and density. Therefore, combining the initial MLP model (Eq. 1) and positional encoding (Eq. 2), our model could be rewritten as:

$$
\begin{aligned}
\sigma, f_c &= MLP_{1:8}(\mathcal{F}(\boldsymbol{p}_j)) \\
c(\boldsymbol{p}_j) &= MLP_{9:}(f_c, \mathcal{F}(\boldsymbol{d})) \, ,
\end{aligned}
\tag{5}
$$

where the positional embedding function $\mathcal{F}$ is applied individually on each component of a vector.

Finally, we accumulate the densities and radiance of $N$ samplings along the ray $\mathcal{R}_i$ to get the final color $\hat{C}(\mathcal{R}_i)$ and depth $\hat{D}(\mathcal{R}_i)$ estimation of the pixel, similar to the volume rendering formulated in NeRF. Supposing the sampling position along the ray is $t_i$, the corresponding color and depth can be calculated by :

$$
\begin{aligned}
\hat{C}(\mathcal{R}_i) &= \sum_{i=1}^{N} T_i \boldsymbol{c}_i \\
\hat{D}(\mathcal{R}_i) &= \sum_{i=1}^{N} T_i t_i \, ,
\end{aligned}
\tag{6}
$$

where the weight
$$
T_i = \exp\left(-\sum_{k=1}^{i-1} \sigma_k(t_{k+1} - t_k)\right)(1 - \exp(-\sigma_i(t_{i+1} - t_i))).
$$

## 4 GEOMETRY-AWARE *360NERFS*

In structured scenes, there are a lot of open spaces and the distributions of objects are spatially sparse. It means that the density of radiance fields is highly imbalanced. The perspectives from the camera are always surrounded by the fields. If we uniformly sample along the rays, it is highly inefficient and easy to get trapped in the local minimum. Conversely, by skipping empty spaces and upsampling, occupied volumes are able to accelerate volume rendering and effectively improve the rendering quality. In our systems, we introduce a probabilistic progressive layout estimation method to recover geometric information from *360NeRF*. Furthermore, based on the density of radiance fields, *360NeRF* will be adaptively slimmed and further fine-tuned to become multiple tiny geometry-aware radiance fields.

## 4.1 Progressive Probabilistic Geometry Estimation

*4.1.1 Occupancy Map.* To extract geometric information from NeRF, the existing approach relying on global density sampling has two limitations. First, it assumes the radiance field is covered by a cuboid and based on the known boundary to estimate the geometry which is generally not practical especially for outward-facing scenes. Although we can utilize a structure-from-motion algorithm to roughly estimate a bounding box of radiance field, global density sampling is sensitive to outliers and difficult to handle polyhedral scenes due to overfitting, as shown in Fig.3a. It will lead to a waste of neural
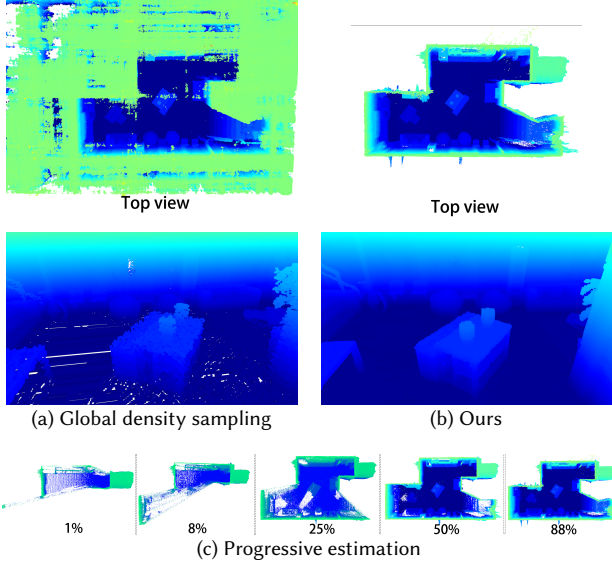
Top view

Top view

(a) Global density sampling

(b) Ours

1%        8%        25%        50%        88%

(c) Progressive estimation

Fig. 3. The comparison of occupancy maps reconstructed by different methods. Exiting method (a) relies on a global sampling of radiance fields which results in overfitting. We propose a progressive method to recover the occupancy map and achieve better performance (b), while the iterative results are demonstrated on (c).



Upper slicing region: height of 60-65% above floor

Gravity Direction

Noisy voxels outside the wall

Lower slicing region: height of 20-30% above floor

Upper-slice Floorplan

Lower-slice Floorplan

(a) Upper & lower slicing.

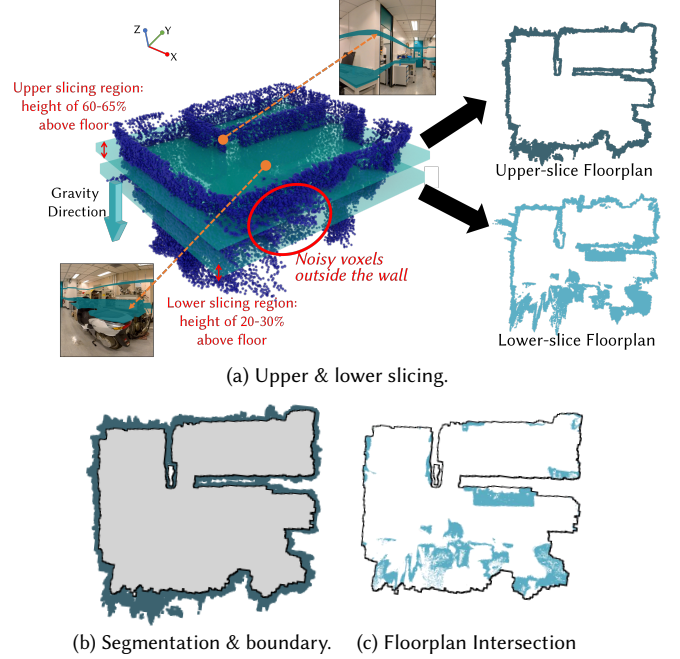(b) Segmentation & boundary.        (c) Floorplan Intersection

Fig. 4. Floorplan Estimation on scene Lab. (a) We collect occupied voxels at the height of upper and lower ranges and project them to floor-parallel plane as the initial floorplan. (b) We apply simple morphological operations on the upper-slice floorplan to extract the interior segment (gray area) and hence the floorplan inner boundary (black line). (c) Intersecting the inner boundary with the lower-slice floorplan gives the final floorplan.

perceptrons since we rely on recovered geometry to decompose radiance fields and assign networks as elaborated in the Section 4.2. The other limitation is being highly reliant on the prior information of density value. It is vulnerable in terms of dealings with transparent objects and mirror refection.

To bridge this gap, we formulated occupancy recovery as a progressive probability estimation problem. Intuitively, the present state of each voxel depends on the current measurement and the previous estimation. The probabilistic model [Yguel et al. 2008] could be formulated in logit notation:

$$logit(p|z_{1:t}) = logit(p|z_{1:t-1}) + logit(p|z_t),  \quad (7)$$

where $logit(p) = log\frac{p}{1-p}$ and $p$ is the probability of occupancy; $z_{1:t}$ represents the accumulated observation from the start to time $t$ and $z_t$ is the current observations at time $t$; With the incoming observations, it can progressively update the occupied probability of each voxel. The range of $logit(p)$ is from -2 and 3.5 which indicates the occupied probability of 0.12 and 0.97. When the probability is larger than the threshold, the voxel will be considered occupied. To facilitate these processes, we adapt a typical reconstruction framework Octomap [Hornung et al. 2013] which utilizes an octree to store and maintain the occupied probability of voxel. Based on predefined measurement uncertainty, Octomap only takes depth maps as input to generate occupancy voxel. However, we can take advantage of *360NeRF* to obtain a panorama distance map at each sampling location as well as the uncertainty estimation by normalizing the weight of Eq. 6. Therefore, we shall make use of the information to enhance occupancy map reconstruction. The measurement probability will be changed according to the current 360NeRF estimation. To reduce the use of hyperparameters, the resolution of occupancy

voxel depends on the ray sample rate of *360NeRF*. In addition, the sampling locations can base on the pose of the training images or we can densely sample more points around these prior locations to further increase the accuracy. In our system, it simply goes through locations in the training set which is generally sufficient to reconstruct an appropriate occupancy map of the scene. This process is depicted in the first part of Algorithm 1. As Fig. 3 shows, this progressive probabilistic geometry estimation method is more effective compared to the global density sampling.

*4.1.2 Floorplan.* The estimated occupancy map provides useful geometric information to refine the radiance fields and we further exploit it to estimate the floorplan of the scene. The floorplan can be used to guide ray sampling and restrict the walking path at inference time. The floorplan we are estimating here is essentially the inner surface of the wall and any objects in the scene that obstruct the navigation path. To better filter out noisy voxels from the occupancy map, we project a range of occupied voxels in the upper slicing regions (60% to 65% of the height to avoid head jamb for multi-room contours) to a floor-parallel plane to form an initial layout (Fig.4a). By using some simple morphological operations we segment out the inner boundary as show in Fig.4b. The most important part here is to determine the slicing direction which should be floor-parallel. Namely, we need to estimate the direction of the gravity of the reconstructed scene. Benefit from 360° images, we can get the gravity direction of the structured scene via calculating three

**ALGORITHM 1:**

Density distribution estimation and networks assignment

**Input:** 360NeRF and the amount of networks $n$

/* Stage I */

Sample positions: $Q \leftarrow$ training set;

**for** $q \in Q$ **do**

> distance map, uncertainty map $\leftarrow$ 360NeRF ($q$);
>
> Updating occupied probability map using Eq.7 ;

**end**

Obtained occupancy map $O$ and the amount of occupied voxel $S$;

/* Stage II */

Size of each sub-field $s \leftarrow S/n$ ;

**for** $k \leftarrow 1$ **to** $n$ **do**

> Pop the unlabeled voxel with minimum coordinate from $O$ ;
>
> Perform Nearest Neighbor Search algorithm and cluster $s$
> occupied voxel into $k^{th}$ network;

**end**



(a) NeRF          (b) KiloNeRF
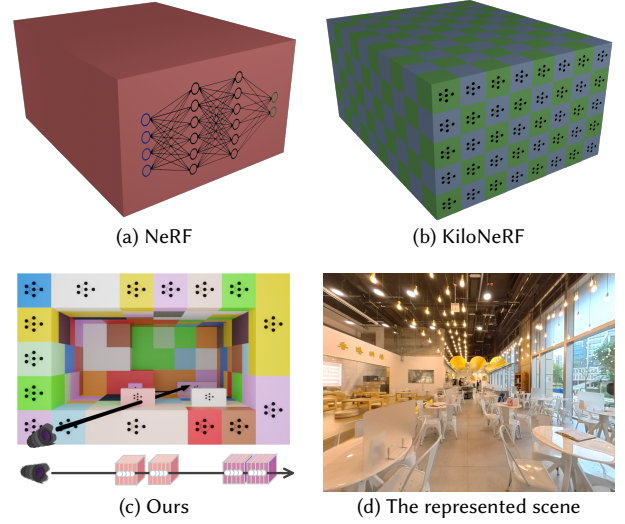
(c) Ours          (d) The represented scene

Fig. 5. Different representations of radiance fields. (a) NeRF uses a single deep MLP to represent the entire field while (b) KiloNeRF uniformly decomposes the field into thousands of tiny cubes (MLP) with fixed resolution and handles rays parallelly. The inference speed of KiloNeRF is faster than NeRF, but their renderings quality are smiliar. (c) By contrast, we use the occupancy map to adaptively decompose the *360NeRF* into $n$ sub-fields and assign the tiny MLP with the same network structure to handle different regions. Each sub-field has similar amount of occupied voxels. The actual number depends on the scene geometry in practice. For illustration purpose here each sub-field corresponds to two occupied voxels. The resultant geometry-aware *360NeRFs* allow effective empty space skipping and dense sampling in occupied voxels to improve rendering quality.

vanishing points [Sun et al. 2021]. After aligning the 360° images with gravity, we are able to estimate the height and extract the floor-parallel plane. Finally, to complete the floorplan with interior items such as chairs, tables, etc. which obstruct the indoor roaming path, we project occupied voxels in the lower slicing region (20% to 30% of the height that covers most of the obstructing objects) to the floor-parallel plane (Fig.4a) and intersect the projected voxels with the inner boundary (Fig.4c). In *360Roam*, the intersected clean floorplan can enhance the user experience with visualized scene boundary and item obstacles.

### 4.2 Adaptive Radiance Fields Slimming and Fine-tuning

Rendering an image from the NeRF model for one pixel costs hundreds of network queries. In order to accelerate inference speed and meet the need for real-time rendering, we can divide and decouple the entire neural radiance field and parallel processing for each bundle of rays. However, in a large-scale scene, the distribution of objects and their corresponding radiance fields are highly imbalanced. As the empty space does not contribute to volume rendering, we should avoid sampling and wasting the capacity of neural perceptron in empty spaces. Therefore, in contrary to the uniform divide-and-conquer strategy used in KiloNeRF [Reiser et al. 2021], we no longer decompose the radiance field simply based on identical volume. Instead, our decomposition strategy is geometry-aware and takes scene geometry information into account. Specifically, based on the extracted occupancy map, we can obtain the density distribution of the radiance fields which can be represented by the number of the occupied voxels. We then decompose the field into $n$ sub-fields with a similar amount of occupied voxels as described in the second stage of Algorithm 1 . Since the occupied voxels in *360NeRF* are spatially imbalanced, each sub-field has a distinct volume. This step only requires one hyperparameter $n$, the number of tiny MLPs, which can be adjusted according to computational resources and demand in terms of rendering performance. After this, each of the $n$ tiny MLPs is used to cover each sub-field which is illustrated in the Fig.5. Eventually, our system can make use of the

capability of neural networks and reduce the usage of MLPs while maintaining high-fidelity rendering.

To decouple each sub-field while persisting global consistency, a fine-tuning process is conducted. The initial weights of tiny MLPs are distilled from the *360NeRF* to increase the fine-tuning speed. Similar to distillation procedures [Reiser et al. 2021], each MLP $\mathcal{M}_k$ will sample $I$ corresponding points in the *360NeRF* to obtain the global domain radiance values $c_i$ and densities $\sigma_i$. $\mathcal{M}_k$'s parameters are optimized via minimizing the mean squared errors between global and local values $(c_i^k, \sigma_i^k)$, which is formulated as:

$$\mathcal{L}_k = \frac{1}{I} \sum_{i \in I} \|\sigma_i^k - \sigma_i\|_2^2 + \|c_i^k - c_i\|_2^2 . \tag{8}$$

After initialization, the parameters of all tiny neural networks will be fine-tuned by using the original training images. During the fine-tuning phase, it will densely sample the occupied voxels instead of uniform sampling along the rays. As radiance fields are slimmed and a vast amount of empty sampling is being avoided, the fine-tuning process can be converged quickly.

### 4.3 Real-time *360Roam*

To achieve real-time rendering performance for immersive 360° roaming, we have taken advantage of ray casting techniques to quickly localize occupied voxels. After partitioning the radiance
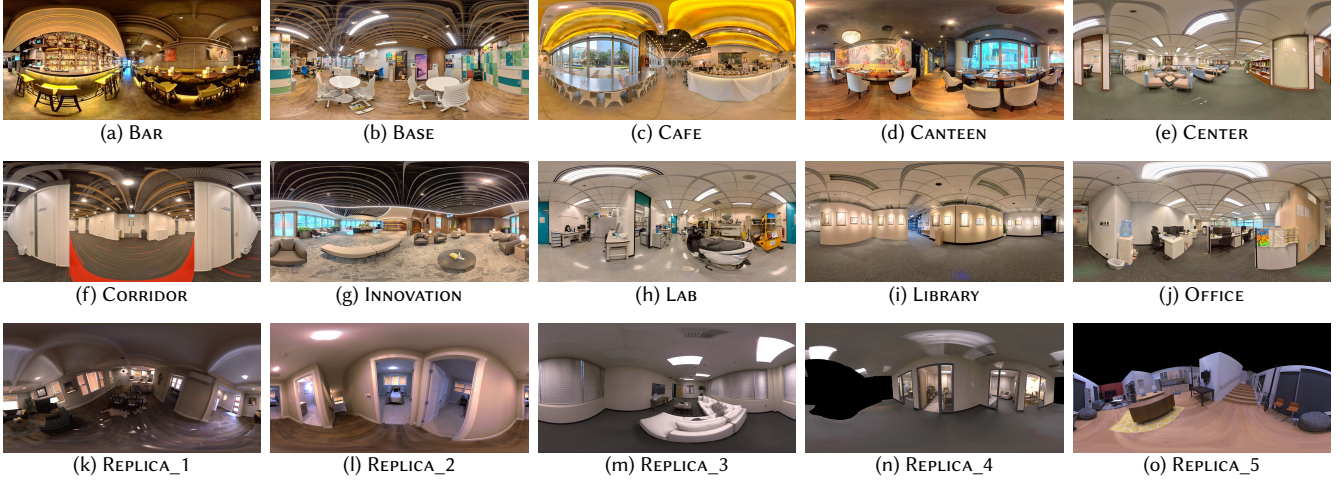
Fig. 6. The dataset contains 15 indoor scenes in which 10 sequences are captured in real world using a 360° camera, and another 5 sequences are rendered from synthetic scenes in Replica [Straub et al. 2019].

field as described in Section 4.2, each occupied voxel in the occupancy map has been labeled to the corresponding tiny MLP $\mathcal{M}_k$'s. In addition, we converted the occupancy map into a voxel hashing data structure and stored it in a GPU such that each ray can be concurrently processed. Due to the support of data structure, efficient empty space skipping, and parallel MLPs queries, our proposed system can run at real-time speed on a GPU GTX1080.

## 5 EXPERIMENT

In this section, we detail how to collect 360° panoramic data and perform the novel view synthesis from scratch. We describe the applied network architectures and important hyper-parameters. Moreover, we conducted thorough qualitative and quantitative evaluations to show our results outperform previous baseline methods and achieve SOTA performance. Finally, the ablation studies prove the validity of our proposed method.

### 5.1 Dataset

To evaluate our proposed methods, we employed both synthetic and real-world data which are a collection of 360° image sequences captured at various indoor scenes as Fig.6 shows.

In terms of real-world data, we used an 'Insta360 ONE X2' 360° camera to capture real-world panoramas. We carefully selected 10 scenes named BAR, BASE, CAFE, CANTEEN, CENTER, CORRIDOR, INNOVATION, LAB, LIBRARY and OFFICE respectively and picked 140 panoramas on average for each scene. They cover diverse indoor environments which have different styles, furniture arrangements, and layouts with multiple relevant separated spaces. To facilitate data collection, we fixed the camera rig on a tracked mobile robot with a height of 1m and controlled the robot remotely. In addition, we used HDR capture model to preserve image quality. The resolution of the image is $6080 \times 3040$. Finally, we used a 360 SFM pipeline [Huang and Yeung 2022] to process 360° images, while we only reserved keyframes extracted by SFM pipeline. In addition, we

separated the keyframes into two sets, one set serves as training data and another serves as testing data with ratio 4:1.

The synthetic data is better for quantitative comparison since we can get error-free testing ground-truth data. Therefore, we applied the off-the-shelf Replica [Straub et al. 2019] dataset to render images from reconstructed 3D models of indoor spaces. The related SDK provides a user interface to generate perspective images. We rendered a cube map containing 6 perspective images at each capturing location and transformed the cube maps into equirectangular images with the resolution of $3840 \times 1920$. We have collected 5 sets of synthetic data, named as REPLICA_1, REPLICA_2, REPLICA_3, REPLICA_4, REPLICA_5. We collected around 80 panoramas per scene with on average 66 for training and 16 for testing.

### 5.2 Implementations

*5.2.1 Baselines.* To evaluate our model, we compared it to the canonical NeRF [Mildenhall et al. 2020], Mip-NeRF [Barron et al. 2021], KiloNeRF [Reiser et al. 2021] and NSVF [Liu et al. 2020]. NSVF aforehand defines a set of voxel-bounded implicit fields organized in a voxel octree to model local properties in each cell and progressively learn the underlying voxel structures during training which achieves good performance in object-centric scenarios. KiloNeRF dramatically improves inference speed by introducing uniform divide-and-conquer strategy while maintaining similar rendering quality to vanilla NeRF. Mip-NeRF found that the uneven placement of objects results in blurred renderings of NeRF in near views and aliasing artifacts in far views. Mip-NeRF resolved this problem by casting a cone for each pixel as well as decomposing the cone into several conical frustums. Instead of applying positional encoding features on those frustums, Mip-NeRF proposed an integrated positional encoding feature based on a multivariate Gaussian. These strategies already take scale into consideration and thereby a unified MLP network is able to represent the scene. The hierarchical

| layer | input | channels | activation |
|---|---|---|---|
| pe1 | position vector | 3/60 | - |
| layer1 | pe1 | 60/256 | ReLU |
| layer2 | layer1 | 256/256 | ReLU |
| layer3 | layer2 | 256/256 | ReLU |
| layer4 | layer3 | 256/256 | ReLU |
| layer5 | layer4 + pe1 | 256/256 | ReLU |
| layer6 | layer5 | 256/256 | ReLU |
| layer7 | layer6 | 256/256 | ReLU |
| layer8 | layer7 | 256/256 | - |
| layer9-density | layer8 | 256/1 | ReLU |
| layer9-feature | layer8 | 256/256 | ReLU |
| pe2 | viewing direction | 3/24 | - |
| layer10 | layer9-feature+pe2 | (256+24)/128 | ReLU |
| layer11 | layer10 | 128/3 | sigmoid |

Table 1. The MLP network architecture of *360NeRF*

| layer | input | channels | activation |
|---|---|---|---|
| pe1 | position vector | 3/60 | - |
| layer1 | pe1 | 60/32 | ReLU |
| layer2-density | layer1 | 32/1 | ReLU |
| layer2-feature | layer1 | 32/32 | - |
| pe2 | viewing direction | 3/24 | - |
| layer3 | layer2-feature+pe2 | (32+24)/32 | ReLU |
| layer4 | layer3 | 32/3 | sigmoid |

Table 2. The MLP network architecture of slimming radiance fields

sampling strategy is reserved while the query of coarse points and fine points is performed through the same network.

*5.2.2 Network Architectures Used in Our Pipeline.* For 360NeRF, we applied a single normal MLP network detailed in Table 1 to model the scene. All layers are fully-connected layers with ReLU activation unless specified otherwise. The input position vector is transformed into positional encoding vector and then passed through 8 layers. A skip connection is used in the fifth layer to preserve information from former layers by concatenating the positional encoding and the output of the fourth layer. The ninth layer outputs a density rectified by a ReLU to guarantee non-negativity and a feature vector. Viewing direction with positional encoding is concatenated with the feature vector to contribute to producing radiance. The final layer is activated by sigmoid to restrict the values within range $[0, 1]$. The count of network parameters is about 0.596M. We do not use hierarchical sampling strategy, two-pass MLP query, in order to reduce training and inference time.

In terms of tiny MLP networks in the final system, the network architecture is detailed in Table 2. The skip connection is excluded from the network. Compared to the normal MLP, the width and the depth of the tiny network are significantly compressed while the amount of network parameters is 6212.

*5.2.3 Training Details and Parameters.* As the training of NeRF or other NeRF-based methods is commonly time-comsuming, we resized real-scene images into $1520 \times 760$ and synthetic images into

$1920 \times 960$ for training and evaluation. For NeRF, it sampled 64 points per ray for the coarse query. And then it sampled another 128 points and used all 192 points for the fine query. Therefore, it took 254 network queries in total to render a pixel. To ensure a fair comparison, we kept the number of queries and training epochs of all approaches the same. However, NSVF is extremely GPU-memory hungry, and the demand on memory would continuously increase as it iteratively subdivides and optimizes voxel during training. Following the default parameters configuration, all the models of NSVF were trained in a multi-process distributed manner using 4 NVIDIA RTX 3090 GPUs with 24 GB memory. NSVF would stop training once out of memory error occurs. Moreover, both NSVF and KiloNeRF requires bounding boxes of scenes and fixed voxel resolutions for training. These prerequisite parameters were estimated from the occupancy maps reconstructed by our methods.

Our proposed method is implemented in Pytorch and CUDA. We applied the Adam optimizer with a learning rate that begins at 1e-3 and decays exponentially to 5e-5 during optimization. The hyperparameters of the optimizer are set to the default values with $\beta_1 = 0.9, \beta_2 = 0.99$. For every backward stage, we used a batch of 4096 rays per process. The coarse training of a 360NeRF takes about 15 hours on 4 RTX 2080 Ti GPUs while the fine-tuning of thousands of tiny networks generally takes 4 hours after slimming.

### 5.3 Floorplans

From the geometric estimation (Sec. 4.1), we estimated the floorplans for all real and synthetic scenes, parts of the results are illustrated in Fig. 9 and 10. Our system can support diverse shapes of indoor scenes which are not only limited to single room or cuboid combinations. The estimated floorplan with objects can be used to provide users with a visual guidance for indoor roaming and avoiding occlusions. Please refer to supplementary for more visual floorplan samples.

### 5.4 Novel View Synthesis

We evaluated the novel view synthesis performance for both real and synthetic scenes quantitatively and qualitatively. Quantitative measurements in terms of appearance similarity including Peak Signal-to-noise Ratio (PSNR), Structural Similarity Index Measure (SSIM) and Learned Perceptual Image Patch Similarity (LPIPS) [Zhang et al. 2018] are adopted to analyze performance. We compared performance between baselines and our methods, and conducted ablation study to analyze the influence of the number of tiny networks used to conquer the radiance fields and validate the effectiveness of making use of probabilistic geometry for fine-tuning and adaptive slimming.

*5.4.1 Influence of number of tiny networks.* 360Roam relies on geometry-aware $360°$ radiance fields which are essentially generated via slimming and decomposing the initial 360NeRF into *n* slimming fields. To analyze the effect and identify the necessary amount to achieve a balance between the rendering quality and computational cost, we measured the rendering quality of 5 synthetic scenes using different *n*, including 100, 200, 300, 400, 512, 1024, 2048, 3072 and 4096. The evaluated results are illustrated as charts in Fig. 7. By thousands of tiny networks, the more networks used, the better performance is. When amount exceeds a thousand, there is no obvious improvement and may even cause degradation in the
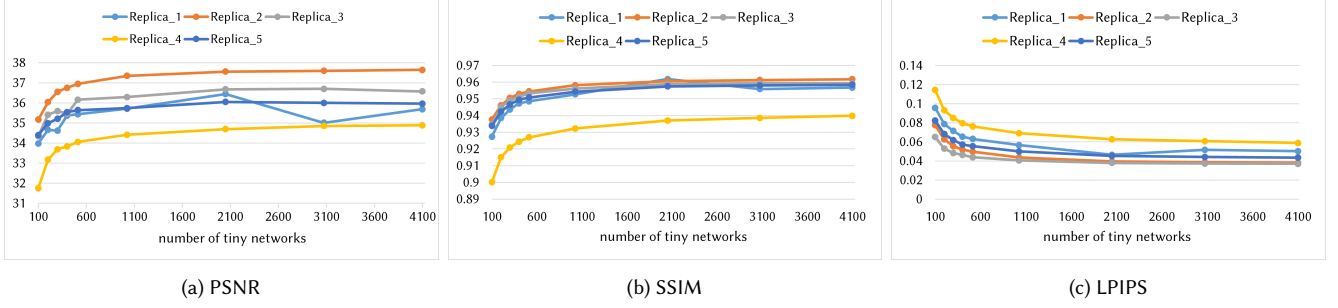
(a) PSNR

(b) SSIM

(c) LPIPS

Fig. 7. The influence of number of the tiny networks on the synthetic scenes.

| Scene | | NSVF | NeRF | Mip-NeRF | KiloNeRF | 360NeRF | 360Roam-100 | 360Roam-200 | 360Roam-u | 360Roam |
|---|---|---|---|---|---|---|---|---|---|---|
| Real Scenes | PSNR↑ | 17.050 | 22.443 | 22.530 | 21.426 | 23.048 | 21.961 | 22.157 | 22.415 | **24.679** |
| | SSIM↑ | 0.516 | 0.672 | 0.670 | 0.690 | 0.720 | 0.659 | 0.677 | 0.704 | **0.755** |
| | LPIPS↓ | 0.758 | 0.339 | 0.336 | 0.283 | 0.297 | 0.344 | 0.308 | 0.268 | **0.209** |
| Synthesis | PSNR↑ | 16.237 | 34.018 | 33.568 | 36.229 | 34.839 | 33.927 | 34.852 | 36.270 | **36.284** |
| | SSIM↑ | 0.745 | 0.931 | 0.925 | 0.951 | 0.938 | 0.927 | 0.937 | 0.951 | **0.955** |
| | LPIPS↓ | 0.570 | 0.096 | 0.103 | 0.058 | 0.085 | 0.087 | 0.071 | 0.057 | **0.046** |
| Inference Speed (s) | | 3 | 60 | 60 | 0.04 | 31 | 0.7 | 0.4 | 0.03 | 0.03 |

Table 3. Comparison of Novel View Synthesis. 360Roam contains 2048 tiny MLPs and achieves the best performance. 360Roam-u denotes the results of uniformly dividing the radiance field without consideration of density distribution. 360Roam-# denotes the system is composed of # MLPs. The counts of network parameters of 360Roam-200 and NeRF are similar, while storage memory of 360Roam-100 and NeRF is similar. 360NeRF is the initial model which is not slimmed and fine-tuned via making use of probabilistic geometry. Approximate time to render one image is also reported.

final rendering quality. In general, a larger amount of networks can have a faster rendering speed until reaching GPU capability of parallel computing. The memory and the number of tiny networks are linearly dependent, while 100 and 4096 tiny network models requires $7.13MB$ and $291.21MB$ memory space respectively. It is worth noting that the memory cost of 360Roam-100 using 100 tiny networks is similar to NeRF while the capacities of neural perceptron are comparable in 360Roam-200 and NeRF, having about 1.2M parameters. Overall, we used 2048 as a typical parameter to conduct the below evaluation.

*5.4.2 Quantitative Evaluation.* The average metrics on real and synthetic sets are reported in Table. 3 and you could refer to the supplementary for the detailed results of individual scenes. From the results reported in Table. 3, our complete system, 360Roam, achieves the best quality on real and synthetic data while the render speed is two thousands times faster than NeRF [Mildenhall et al. 2020] and Mip-NeRF [Barron et al. 2021]. If we decompose the global radiance field into uniform voxel and use similar amount of networks, 360Roam-u can achieve comparable performance on synthetic scenes. But the performance of uniform decomposition is significantly degraded on real scenes. Another method relying on uniform decomposition, KiloNeRF [Reiser et al. 2021], has the same phenomenon. This is because the layout of real scenes are more complicated and the density distribution of the radiance field is highly imbalanced. Uniform decomposition leads to capacity of neural perceptrons wasted in insignificant regions. Obviously, the

proposed adaptive slimming method made the model more stable. When there is no exploit of geometry to fine-tune the radiance filed, 360NeRF sightly suppresses the original NeRF and Mip-NeRF. Furthermore, with the increasing use of tiny MLP, the rendering quality of the system is enhanced. In generally, it requires more tiny MLPs and memory to reach similar capacity of a deep MLP. Although 360Roam-200 achieves competitive results against NeRF, its performance still degrades compared to 360NeRF. As a method focusing on object-centric or small-scale scenes, NSVF [Liu et al. 2020] cannot converge on these dataset. But we would like to note that results of NSVF could possibly be improved by providing the depth to explicitly shape the neural voxels or using more powerful GPU during training.

It is noticeable that results of real scenes generally have a worse performance compared to synthetic scenes. Since real scenes suffer from nonuniform lighting exposure cross views and camera pose approximation, while synthetic scenes provide ground-truth camera parameters and are rendered under the same illumination settings.

*5.4.3 Qualitative Evaluation.* Visual comparisons between baselines and our method are displayed in Fig. 8. Our novel view images for either the real scene (top scene) or the synthetic scene (bottom scene) always have higher fidelity of the texture at any distance. For example, the complex wall texture in the red bounded area in the real scene (first row), and the paintings with textual elements in the blue bounded area in the synthetic scene (fourth row) are generated with clear and recognizable details by our approach. The

|Ground Truth|NSVF [Liu et al. 2020]|NeRF [Mildenhall et al. 2020]|Mip-NeRF [Barron et al. 2021]|Our approach|



Fig. 8. Comparison of novel view synthesis result on scenes CANTEEN and REPLICA_2. Our method outperforms other methods and has higher fidelity in terms of geometry and texture.Note: zoom in for better view.



Fig. 9. Roaming in the scene LAB using perspective view.

baselines can only synthesize blurry and non-photorealistic texture. In addition, ours supports rendering fine details at near and far distances, while the baselines produce blurry objects in the near area and a hazy effect for distant area. In the far blue bounded area of the real scene (second row), our images contain subtle details for the ceiling lamp and tableware at a closer distance, and clear parasols and buildings at a farther distance. At such ranges, the baselines fail to synthesize any clear objects. At closer ranges in the red bounded area of the synthetic scene (third row), our approach produces delicate textures from the near bedside to pillows while baselines only render indistinct texture.

Fig. 10. Our Results of Novel View Synthesis on four real scenes. In the first two columns are synthesized novel views. The corresponding floorplans display on the rightmost column with circled numbers and arrows. The circled numbers indicate relative locations and in which column the rendered image locates, and arrows represent the viewing directions.
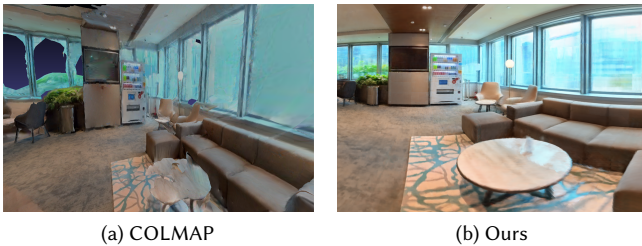


(a) COLMAP

(b) Ours

Fig. 11. Comparison to 3D reconstruction and texturing. We show a perspective result of typical 3D reconstruction (COLMAP [Schönberger and Frahm 2016]) and texture mapping (MVS-Texturing [Waechter et al. 2014]) on the left, and corresponding rendered view from our system on the right. Reconstructed textured model has obvious geometry and texture artifacts and leads to non-photorealism, while ours is much more photorealistic and suitable for VR roaming.

Fig. 9 shows the a roaming example with long trajectories in LAB. Fig. 10 is a gallery of some novel view 360° images of real scenes with their estimated floorplans on the right using our system. We show two novel views for each scene, and the relative locations and viewing directions are marked in the floorplan with circled numbers and arrows. Additional roaming results and qualitative comparisons are shown in the supplementary video.

## 6  DISCUSSION, LIMITATION, AND FUTURE WORK

We demonstrated how *360Roam* effectively extends the conventional NeRF into the geometry-aware *360NeRFs* to generate high quality 360° images for real-time indoor roaming. We also conduct an additional comparison with a typical way which uses 3D reconstruction [Schönberger and Frahm 2016; Schönberger et al. 2016] and texture mapping [Waechter et al. 2014] techniques to recover
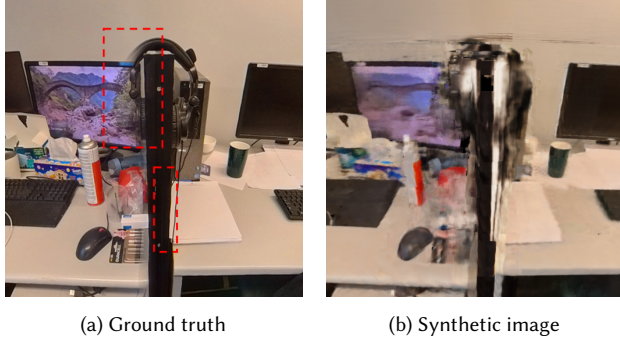
(a) Ground truth      (b) Synthetic image

Fig. 12. Stitching artifacts of 360° images affect rendering results.



(a) Synthetic image      (b) Corresponding geometry

Fig. 13. Incorrect geometry recovery would impair performance of 360Roam.

meshes for real scene roaming. As the COLMAP does not support 360° images, we crop the training panoramas into cube maps and then use the same camera poses to run the COLMAP. An illustrative example is shown in Fig.11. Our approach produces real-time photo-realisitc renderings and has a dramatic advantage over conventional 3D reconstruction which shows obvious geometry and texture artifacts. 360Roam has large potential for VR roaming for large-scale indoor real scenes with less resource requirements and manpower. However, by doing so we also inherit the limitations of NeRF and panorama scene understanding. We list the most important issues and potential directions to further improve the system.

**Camera model.** We use an ideal spherical camera model to describe the projection of 360° camera. However, a consumer-grade 360° camera outputs 360° images by optimizing the stitching of two fish-eye images. Due to manufacturing deficiency and the imperfection of factory camera calibration, original images sometimes have obvious stitching artifacts which will affect the rendering quality, as Fig. 12 shows. This is one key reason why the performance on real scenes are noticeably worse than that of synthetic scenes. Although we can exploit a professional 360° camera for capturing, it is necessary to take distortion into account and optimize camera intrinsic parameters during training.

**Influence of the estimated geometry**. Different from vanilla NeRF or 360NeRF which models geometry implicitly, the complete 360Roam is a hybrid system where the ray sampling is based on explicit geometry for increasing synthetic quality and speed. Our progressive probabilistic geometry estimation method allows us to obtain a proper occupancy map fitting the scene rather than to get an exquisite 3D reconstruction. It is capable of modeling transparent objects, as shown in the first row of Fig. 13. However, once the recovered geometry is deficient, the quality of the synthetic image would suffer. The bottom row of Fig. 13 provides an example.

**Training time and generalization**. NeRF-based methods require a long training time to train MLP networks overfitting the scene. Consequently, it cannot generalize to predict novel images of arbitrary scenes. Each time the data of a new scene is collected, it usually takes tens of hours to converge which brings inconvenience. If we want to use different training resolutions, add newly captured images or fix flaws in a few panoramas, it requires the parameters of the model to be retrained from scratch. It is important to either reduce the training time or expand the model application range.
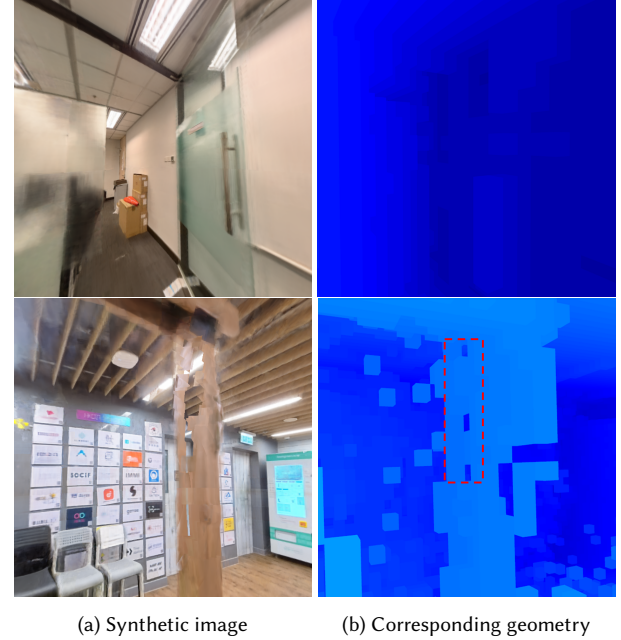
**Dynamic and deformable objects.** Currently, our method assumes the scene is static and no motion blur will occur in the images. When we collected real-world scenes, we avoided human interruption or furniture arrangement changed. It will be useful if we can detect moving and non-rigid regions from the captured panoramas and apply additional models [Li et al. 2021] or interpolation methods [Jiang et al. 2018] to deal with those areas.

**Few-shot sampling.** Most NeRF-based methods require dense scene sampling to collect as much environmental information as possible. And most existing datasets for novel view synthesis task usually contains tens of or even hundreds of images representing an object or small scenes. However, our dataset only uses about hundreds of images to represent a large-scale scenes. The distributions of training panoramas of all scenes are displayed in the supplementary. Although panoramas having a wider field of view are beneficial for representing large scenes, the rendering quality can still be poor in under-sampled areas. How to provide a more realistic view-dependent effect with few-shot training will be an interesting future direction.

## 7 CONCLUSION

In this paper, we seek to extend neural radiance fields to handle large indoor scenes while having real-time performances for novel view synthesis. Therefore, we propose *360Roam* which is a novel neural rendering pipeline taking 360° images as input to learn omnidirectional radiance fields *360NeRF* first. The *360NeRF* is further slimmed and fine-tuned according to the estimated probabilistic geometry. The geometry-aware 360° radiance fields not only accelerate the rendering speed but also improve the rendering quality. It can maintain fidelity even if the perspective of view undergoes a large translational motion. With the floorplan guidance, our system

is capable of providing an appealing and immersive indoor roaming experience which is promising for various VR applications.

## REFERENCES

Benjamin Attal, Selena Ling, Aaron Gokaslan, Christian Richardt, and James Tompkin. 2020. MatryODShka: Real-time 6DoF Video View Synthesis using Multi-Sphere Images. In *European Conference on Computer Vision (ECCV)*. https://visual.cs.brown.edu/matryodshka

Jonathan T Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and Pratul P Srinivasan. 2021. Mip-NeRF: A Multiscale Representation for Anti-Aliasing Neural Radiance Fields. *arXiv preprint arXiv:2103.13415* (2021).

Tobias Bertel, Mingze Yuan, Reuben Lindroos, and Christian Richardt. 2020. OmniPhotos: casual 360° VR photography. *ACM TOG* 39, 6 (2020), 1–12.

Michael Broxton, John Flynn, Ryan Overbeck, Daniel Erickson, Peter Hedman, Matthew Duvall, Jason Dourgarian, Jay Busch, Matt Whalen, and Paul Debevec. 2020. Immersive light field video with a layered mesh representation. *ACM Transactions on Graphics (TOG)* 39, 4 (2020), 86–1.

Terrance DeVries, Miguel Angel Bautista, Nitish Srivastava, Graham W. Taylor, and Joshua M. Susskind. 2021. Unconstrained Scene Generation with Locally Conditioned Radiance Fields. *arXiv* (2021).

Peter Hedman, Pratul P Srinivasan, Ben Mildenhall, Jonathan T Barron, and Paul Debevec. 2021. Baking neural radiance fields for real-time view synthesis. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 5875–5884.

Armin Hornung, Kai M Wurm, Maren Bennewitz, Cyrill Stachniss, and Wolfram Burgard. 2013. OctoMap: An efficient probabilistic 3D mapping framework based on octrees. *Autonomous robots* 34, 3 (2013), 189–206.

Huajian Huang and Sai-Kit Yeung. 2022. 360VO: Visual Odometry Using A Single 360 Camera. In *International Conference on Robotics and Automation (ICRA)*. IEEE.

Jingwei Huang, Zhili Chen, Duygu Ceylan, and Hailin Jin. 2017. 6-DOF VR videos with a single 360-camera. In *2017 IEEE Virtual Reality (VR)*. IEEE, 37–44.

Huaizu Jiang, Deqing Sun, Varun Jampani, Ming-Hsuan Yang, Erik Learned-Miller, and Jan Kautz. 2018. Super slomo: High quality estimation of multiple intermediate frames for video interpolation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 9000–9008.

Lei Jin, Yanyu Xu, Jia Zheng, Junfei Zhang, Rui Tang, Shugong Xu, Jingyi Yu, and Shenghua Gao. 2020. Geometric structure based and regularized depth estimation from 360 indoor imagery. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 889–898.

Zhengqi Li, Simon Niklaus, Noah Snavely, and Oliver Wang. 2021. Neural scene flow fields for space-time view synthesis of dynamic scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 6498–6508.

Kai-En Lin, Zexiang Xu, Ben Mildenhall, Pratul P Srinivasan, Yannick Hold-Geoffroy, Stephen DiVerdi, Qi Sun, Kalyan Sunkavalli, and Ravi Ramamoorthi. 2020. Deep multi depth panoramas for view synthesis. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XIII 16*. Springer, 328–344.

David B Lindell, Julien NP Martel, and Gordon Wetzstein. 2021. Autoint: Automatic integration for fast neural volume rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 14556–14565.

Lingjie Liu, Jiatao Gu, Kyaw Zaw Lin, Tat-Seng Chua, and Christian Theobalt. 2020. Neural Sparse Voxel Fields. *NeurIPS* (2020).

Stephen Lombardi, Tomas Simon, Gabriel Schwartz, Michael Zollhoefer, Yaser Sheikh, and Jason Saragih. 2021. Mixture of volumetric primitives for efficient neural rendering. *ACM Transactions on Graphics (TOG)* 40, 4 (2021), 1–13.

Ricardo Martin-Brualla, Noha Radwan, Mehdi S. M. Sajjadi, Jonathan T. Barron, Alexey Dosovitskiy, and Daniel Duckworth. 2021. NeRF in the Wild: Neural Radiance Fields for Unconstrained Photo Collections. In *CVPR*.

Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. 2020. Nerf: Representing scenes as neural radiance fields for view synthesis. In *European conference on computer vision*. Springer, 405–421.

Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. 2022. Instant Neural Graphics Primitives with a Multiresolution Hash Encoding. *arXiv preprint arXiv:2201.05989* (2022).

Thomas Neff, Pascal Stadlbauer, Mathias Parger, Andreas Kurz, Joerg H. Mueller, Chakravarty R. Alla Chaitanya, Anton S. Kaplanyan, and Markus Steinberger. 2021. DONeRF: Towards Real-Time Rendering of Compact Neural Radiance Fields using Depth Oracle Networks. *Computer Graphics Forum* 40, 4 (2021). https://doi.org/10.1111/cgf.14340

Martin Piala and Ronald Clark. 2021. Terminerf: Ray termination prediction for efficient neural rendering. In *2021 International Conference on 3D Vision (3DV)*. IEEE, 1106–1114.

Albert Pumarola, Enric Corona, Gerard Pons-Moll, and Francesc Moreno-Noguer. 2021. D-nerf: Neural radiance fields for dynamic scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 10318–10327.

Christian Reiser, Songyou Peng, Yiyi Liao, and Andreas Geiger. 2021. KiloNeRF: Speeding up Neural Radiance Fields with Thousands of Tiny MLPs. *arXiv preprint arXiv:2103.13744* (2021).

Johannes Lutz Schönberger and Jan-Michael Frahm. 2016. Structure-from-Motion Revisited. In *Conference on Computer Vision and Pattern Recognition (CVPR)*.

Johannes Lutz Schönberger, Enliang Zheng, Marc Pollefeys, and Jan-Michael Frahm. 2016. Pixelwise View Selection for Unstructured Multi-View Stereo. In *European Conference on Computer Vision (ECCV)*.

Ana Serrano, Incheol Kim, Zhili Chen, Stephen DiVerdi, Diego Gutierrez, Aaron Hertzmann, and Belen Masia. 2019. Motion parallax for 360° RGBD video. *IEEE Transactions on Visualization and Computer Graphics* 25, 5 (2019), 1817–1827. https://doi.org/10.1109/TVCG.2019.2898757

Julian Straub, Thomas Whelan, Lingni Ma, Yufan Chen, Erik Wijmans, Simon Green, Jakob J. Engel, Raul Mur-Artal, Carl Ren, Shobhit Verma, Anton Clarkson, Mingfei Yan, Brian Budge, Yajie Yan, Xiaqing Pan, June Yon, Yuyang Zou, Kimberly Leon, Nigel Carter, Jesus Briales, Tyler Gillingham, Elias Mueggler, Luis Pesqueira, Manolis Savva, Dhruv Batra, Hauke M. Strasdat, Renzo De Nardi, Michael Goesele, Steven Lovegrove, and Richard Newcombe. 2019. The Replica Dataset: A Digital Replica of Indoor Spaces. *arXiv preprint arXiv:1906.05797* (2019).

Cheng Sun, Chi-Wei Hsiao, Min Sun, and Hwann-Tzong Chen. 2019. Horizonnet: Learning room layout with 1d representation and pano stretch data augmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 1047–1056.

Cheng Sun, Min Sun, and Hwann-Tzong Chen. 2021. Hohonet: 360 indoor holistic understanding with latent horizontal features. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2573–2582.

Matthew Tancik, Vincent Casser, Xinchen Yan, Sabeek Pradhan, Ben Mildenhall, Pratul P Srinivasan, Jonathan T Barron, and Henrik Kretzschmar. 2022. Block-NeRF: Scalable Large Scene Neural View Synthesis. *arXiv preprint arXiv:2202.05263* (2022).

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in neural information processing systems*. 5998–6008.

Michael Waechter, Nils Moehrle, and Michael Goesele. 2014. Let there be color! Large-scale texturing of 3D reconstructions. In *ECCV*. Springer, 836–850.

Qianqian Wang, Zhicheng Wang, Kyle Genova, Pratul P Srinivasan, Howard Zhou, Jonathan T Barron, Ricardo Martin-Brualla, Noah Snavely, and Thomas Funkhouser. 2021. Ibrnet: Learning multi-view image-based rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 4690–4699.

Jiale Xu, Jia Zheng, Yanyu Xu, Rui Tang, and Shenghua Gao. 2021. Layout-Guided Novel View Synthesis from a Single Indoor Panorama. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 16438–16447.

Qiangeng Xu, Zexiang Xu, Julien Philip, Sai Bi, Zhixin Shu, Kalyan Sunkavalli, and Ulrich Neumann. 2022. Point-NeRF: Point-based Neural Radiance Fields. *arXiv preprint arXiv:2201.08845* (2022).

Shang-Ta Yang, Fu-En Wang, Chi-Han Peng, Peter Wonka, Min Sun, and Hung-Kuo Chu. 2019. Dula-net: A dual-projection network for estimating room layouts from a single rgb panorama. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 3363–3372.

Manuel Yguel, Olivier Aycard, and Christian Laugier. 2008. Update policy of dense maps: Efficient algorithms and sparse representation. In *Field and Service Robotics*. Springer, 23–33.

Alex Yu, Ruilong Li, Matthew Tancik, Hao Li, Ren Ng, and Angjoo Kanazawa. 2021a. Plenoctrees for real-time rendering of neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 5752–5761.

Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. 2021b. pixelNeRF: Neural Radiance Fields from One or Few Images. In *CVPR*.

Wei Zeng, Sezer Karaoglu, and Theo Gevers. 2020. Joint 3d layout and depth prediction from a single indoor panorama image. In *European Conference on Computer Vision*. Springer, 666–682.

Kai Zhang, Gernot Riegler, Noah Snavely, and Vladlen Koltun. 2020. NeRF++: Analyzing and Improving Neural Radiance Fields. arXiv:2010.07492 [cs.CV]

Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. 2018. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 586–595.

Tinghui Zhou, Richard Tucker, John Flynn, Graham Fyffe, and Noah Snavely. 2018. Stereo Magnification: Learning View Synthesis using Multiplane Images. In *SIGGRAPH*.

Nikolaos Zioulis, Antonis Karakottas, Dimitrios Zarpalas, Federico Alvarez, and Petros Daras. 2019. Spherical view synthesis for self-supervised 360 depth estimation. In *2019 International Conference on 3D Vision (3DV)*. IEEE, 690–699.

Chuhang Zou, Alex Colburn, Qi Shan, and Derek Hoiem. 2018. Layoutnet: Reconstructing the 3d room layout from a single rgb image. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2051–2059.