# NeurAR: Neural Uncertainty for Autonomous 3D Reconstruction with Implicit Neural Representations

Yunlong Ran[1], Jing Zeng[1], Shibo He[2], Jiming Chen[2], Lincheng Li[3], Yingfeng Chen[3], Gimhee Lee[4], Qi Ye[2]

*Abstract*—**Implicit neural representations have shown compelling results in offline 3D reconstruction and also recently demonstrated the potential for online SLAM systems. However, applying them to autonomous 3D reconstruction, where a robot is required to explore a scene and plan a view path for the reconstruction, has not been studied. In this paper, we explore for the first time the possibility of using implicit neural representations for autonomous 3D scene reconstruction by addressing two key challenges: 1) seeking a criterion to measure the quality of the candidate viewpoints for the view planning based on the new representations, and 2) learning the criterion from data that can generalize to different scenes instead of a hand-crafting one. To solve the challenges, firstly, a proxy of Peak Signal-to-Noise Ratio (PSNR) is proposed to quantify a viewpoint quality; secondly, the proxy is optimized jointly with the parameters of an implicit neural network for the scene. With the proposed view quality criterion from neural networks (termed as Neural Uncertainty), we can then apply implicit representations to autonomous 3D reconstruction. Our method demonstrates significant improvements on various metrics for the rendered image quality and the geometry quality of the reconstructed 3D models when compared with variants using TSDF or reconstruction without view planning. Project webpage https://kingteeloki-ran.github.io/NeurAR/**

*Index Terms*—**Computer Vision for Automation; Motion and Path Planning; Planning under Uncertainty**

## I. INTRODUCTION

**A**UTONOMOUS 3D reconstruction has a wide range of applications, e.g. augmented/virtual reality, autonomous driving, filming, gaming, medicine, architecture. The problem requires a robot to make decisions about moving towards which viewpoint in each step to get the best reconstruction quality of an unknown scene with the lowest cost, i.e. view planning. In this work, we assume a robot can localize itself

[1]Yunlong Ran and Jing Zeng are with the College of Control Science and Engineering, Zhejiang University, Hangzhou, China. {`yunlong_ran`, `zengjing`}`@zju.edu.cn`

[2]Shibo He, Jiming Chen, and Qi Ye are with the College of Control Science and Engineering, the State Key Laboratory of Industrial Control Technology, Zhejiang University, and the Key Laboratory of Collaborative Sensing and Autonomous Unmanned Systems of Zhejiang Province. {`s18he`, `cjm qi.ye`}`@zju.edu.cn`

[3]Lincheng Li and Yingfeng Chen are with Fuxi AI Lab, NetEase Inc., Hangzhou, China. {`lilincheng`, `chenyingfeng1`}`@corp.netease.com`

[4]Gimhee Lee is with the Department of Computer Science, National University of Singapore, Singapore. `dcslgh@nus.edu.sg`

and at each viewpoint, the information of a scene is captured by an RGB image (with an optional depth image).

Implicit neural representations for 3D objects have shown their potential to be precise in geometry encoding, efficient in memory consumption (adaptive to scene size and complexity), predictive in filling unseen regions and flexible in the amount of training data. The reconstruction with offline images [1] or online images with cameras held by human [2] has achieved compelling results recently with implicit neural representations. However, leveraging these advancements to achieve high-quality autonomous 3D reconstruction has not been studied.

Previous 3D representations for autonomous 3D reconstruction include point cloud, volume, and surface. To plan the view without global information of a scene, previous work resorts to a greedy strategy: given the current position of a robot and the reconstruction status, they quantify the quality of the candidate viewpoints via information gain to plan the next best view (NBV). In these works, the information gain relies on hand-crafted criteria, each designed ad-hoc for a particular combination of a 3D representation and a reconstruction algorithm. For example, Mendez *et al.* [3] define view cost by triangulated uncertainty given by the algorithm inferring depth from stereo RGB images, Isler *et al.* [4] quantify the information gain for a view using entropy in voxels seen in this viewpoint, Wu *et al.* [5] identify the quality of the view by a Poisson field from point clouds and Song *et al.* [6] leverage mesh holes and boundaries to guide the view planning.

To use implicit neural representations for autonomous 3D reconstruction, a key capability is to quantify the quality of the candidate viewpoints. *For implicit neural representations, how to define the viewpoint quality? Is it possible for the neural network to learn a measurement of the quality from data instead of defining it heuristically?* In this paper, we make efforts to answer both questions.

The quality of reconstructed 3D models can be measured by the quality of images rendered from different viewpoints and the measurement is adopted in many offline 3D reconstruction works. PSNR is one popular measurement and it is defined according to the difference between the images rendered from the reconstructed model and the ground truth model.

The ground truth images, however, are unavailable for unvisited viewpoints to calculate PSNR during autonomous reconstruction. Is it possible to learn a proxy for PSNR? In [7], [8], the authors point out that if the target variable to regress is under a Gaussian noise model and the target distribution conditioned on the input is optimized by maximum likelihood, the optimum value of the noise variance is given
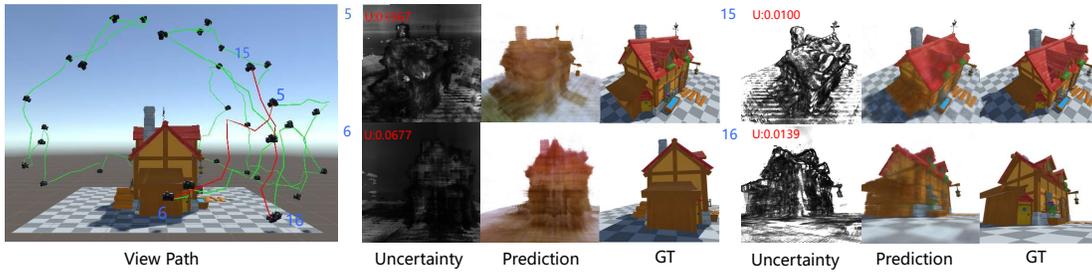
Fig. 1. View paths planned by NeurAR and the uncertainty maps used to guide the planner. Given current pose 5 and 15, paths are planned toward viewpoints having higher uncertainties, *i.e.* pose 6 and 16. Uncertainties of the viewpoints are shown in red text. Notice that darker regions in uncertainty maps relate to worse quality regions of rendered images.

by the residual variance of the target values and the regressed ones. Inspired by this, we assume the color to regress for a spatial point in a scene as a random variable modeled by a Gaussian distribution. The Gaussian distribution models the uncertainty of the reconstruction and the variance quantifies the uncertainty. When the regression network converges, the variance of the distribution is given by the squared error of the predicted color and the ground truth color; the integral of the uncertainty of points in the frustum of a viewpoint can be taken as a proxy of PSNR to measure the quality of candidate viewpoints.

With the key questions solved, we are able to build an autonomous 3D reconstruction system (NeurAR) using an implicit neural network. In summary, the contributions of the paper are:

- We propose the first autonomous 3D reconstruction system using an implicit neural representation.
- We propose a novel view quality criterion that learns online from continuously added input images per target scene instead of hand-engineering or learning from a large corpus of 3D scenes.
- Our proposed method significantly improves on various metrics from alternatives using voxel-based representations or using man-designed paths for the reconstruction.

## II. RELATED WORK

**View Planning** Most view planning methods focus on the NBV problem which uses feedback from current partial reconstruction to determine NBV. According to the representations for the 3D models, these methods can be divided into voxel-based method [3], [4] and surface-based method [5], [9].

The voxel-based methods are most commonly used due to their simplicity in representing space. Vasquez-Gomez *et al*. [10] analyze a set of boundary voxels and determine NBV for dense 3D modeling of small-scale objects. Stefan Isler *et al*. [4] provide several metrics to quantify the volume information contained in the voxels. Mendez *et al*. [3] define the information gain by the triangulation uncertainty of stereo images. Despite the simplicity, these methods suffer from memory consumption with growing scene complexity and higher spatial resolutions.

A complete volumetric map does not necessarily guarantee a perfect 3D surface. Therefore, researchers propose to analyze the shape and quality of the reconstructed surface for NBV [5], [9]. Wu *et al*. [5] estimate a confidence map representing the

completeness and smoothness of the constructed Poisson iso-surface and the confidence map is used to guide the calculation of NBV. Schmid *et al*. [9] propose information gain to evaluate the quality of observed surfaces and unknown voxels near the observed surfaces, then plan a path by RRT.

In contrast, some methods explore function learning to solve the NBV problem [11]–[13]. Supervised learning-based methods [13] learn information gain of a viewpoint given a partial occupancy map and [12] learns an informed distribution of high-utility viewpoints based on a partial occupancy map. For the reinforcement learning-based methods [11], no hand-crafted heuristics are required and an agent explores the viewpoints with high overall coverage. These learning-based methods require a large-scale dataset for training and may be hard to generalize to different scenes.

Different from the hand-crafted heuristics, the learned NBV policy, and the existing learning-based information gain, our proposed neural uncertainty is learned per target scene during the reconstruction, requiring no manual definition, no large training set, and being able to work in any new scene.

**Online Dense Reconstruction** For online 3D reconstruction from RGB images, most methods use Multi-View-Stereo (MVS) [14] to reconstruct dense models by first getting a sparse set of initial matches, iteratively expanding matches to nearby locations and performing surface reconstruction. With the release of commodity RGB-D sensors, the fusion of point clouds reprojected from depth images gains popularity. KinectFusion [15] achieves real-time 3D reconstruction with a moving depth camera by integrating points from depth images with Truncated Signed Distance Functions (TSDFs). OctoMap [16] builds a probabilistic occupancy volume based on octree. Recently, implicit representations have shown compelling results in 3D reconstruction, either for the radiance field approximation from RGB images [17] or the shape approximation from point clouds [18]. The novel representations are also studied for online dense reconstruction. iMAP [2] adopts MLPs as the scene representation and reconstructs the scene from RGBD images. In addition to using MLPs to represent a scene implicitly, NeRFusion [19] further combines a feature volume to fuse information from different views as a latent scene representation. Similarly, we use the implicit neural function to represent 3D models but we focus on how to leverage this representation for autonomous view planning.
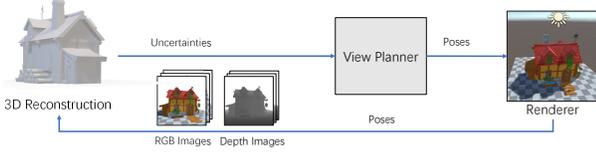
Fig. 2. The pipeline of our method

## III. METHOD

### A. Problem Description and System Overview

The problem considered is to generate a trajectory for a robot that yields high-quality 3D models of a bounded priori unknown scene and fulfills robot constraints. The trajectory is defined as a sequence of viewpoints $V = (v_1, ..., v_n)$ where $v_i \in \mathcal{R}^3 \times \mathcal{SO}(2)$ (usually roll is not considered). Let $S$ be the set of all the sequences. The problem can therefore be expressed as

$$V^* = \arg\max_{V \in S} I(V) \text{ s.t. } L(V) \leq L_{max}, \tag{1}$$

where $I(V)$ equals $\Sigma I(v_i)$, an objective function measuring the contributions from different viewpoints to the reconstruction quality and $L(V)$ is the length of the trajectory.

One difficulty of the problem is defining the objective function itself as there is no ground truth data to measure the quality of the reconstruction. Hence, one major line of work in this field is seeking a surrogate quantity that can be used for the view sequence optimization, which is also the focus of this paper.

Finding the best sequence for (1) is prohibitively expensive and most work resolves this problem by a greedy strategy. At each step of the reconstruction process, a number of candidate viewpoints are sampled according to some constraints, the contributions of these viewpoints to the reconstruction quality are evaluated based on the current reconstructed scene and the viewpoint with the most contribution is chosen as NBV. The process is repeated until a maximum length is reached. To move a robot to NBV, various path planning methods can be applied to navigate the robot.

In this work, we follow the greedy strategy and choose an existing RRT* based view planner. For each viewpoint, an RGB image (with an optional depth image) is captured. To achieve an autonomous 3D reconstruction system depending only on the input images, joint estimation of camera poses and the target scene is required. To focus on view planning with implicit scene representations, we assume the ground truth camera poses are given and leave the localization of the camera as future work.

Accordingly, our proposed autonomous 3D reconstruction system can be divided into three modules (shown in the top of Fig. 1): a 3D reconstruction module, a view planner module, and a simulation module. At each step, the 3D reconstruction module reconstructs a scene represented by an implicit neural network with new images from the simulation module and provides a quantity (neural uncertainty) measuring the reconstructed quality of each position (Section III-B). To adapt the implicit neural representation to online reconstruction, we further propose several strategies for online training and

acceleration in Section III-C. The view planner module (Section III-D) samples viewpoints from empty space, measures the contributions of these sampled viewpoints by composing the neural uncertainty, chooses NBV and plans a view path to NBV. The simulation module, built upon Unity Engine, renders RGBD images from new viewpoints and provides both the images and the camera poses to the 3D reconstruction module.

### B. 3D Reconstruction with Neural Uncertainty

In the section, we first formulate Neural Uncertainty when a target scene is approximated by an implicit neural network. Then we show how to optimize the uncertainty together with the network parameters. With the neural uncertainty formulation, a strong linear relationship is established between the uncertainty and the image quality metric, i.e. PSNR and therefore a proxy for PSNR is proposed.

*1) Problem formulation:* NeRF [1] represents a continuous scene with an implicit neural function by taking the location of a point $\mathbf{x}$ on a ray of direction $\mathbf{d}$ as inputs and its color value $c$, density $\rho$ as outputs, *i.e.* $\mathbf{F}_\theta(\mathbf{x}, \mathbf{d}) = (c, \rho)$. The function $\mathbf{F}_\theta$ is represented by an MLP.

To learn the neural uncertainty for view planning, we treat the color as a random variable under a Gaussian distribution instead of taking it as a deterministic one. The color distribution can be represented as $p(C) = \mathcal{N}(\mu, \sigma^2)$, where the RGB channels of $C$ share the same $\sigma$. To estimate the distribution, we map the inputs $(\mathbf{x}, \mathbf{d})$ to its parameters $(\mu, \sigma)$ of the distribution for the color variable $C$; in other words, the outputs of the MLP is $\mu$, $\sigma$ rather than $c$, so we have $\mathbf{F}_\theta(\mathbf{x}, \mathbf{d}) = (\mu, \sigma, \rho)$. Having defined the color distribution for a point $\mathbf{x}$ on a ray $\mathbf{d}$, we now deduce the rendered color distribution for a camera ray. Following NeRF, we define the rendered color as

$$C_r = \sum_{i=1}^{N} \omega_i C_i, \tag{2}$$

where $\omega_i$ is $o_i \prod_{j=1}^{i-1}(1 - o_j)$, $o_i = (1 - \exp(-\rho_i \delta_i))$ and $\delta_i = d_{i+1} - d_i$ which represents the inter-sample distance. $p(C_i) = \mathcal{N}(\mu_i, \sigma_i^2)$ is the color distribution for a point on the ray. Assuming $\{C_i\}$ for points on the ray are independent from each other, the distribution for $C_r$ is a Gaussian distribution. The probability of the rendered color value being $c_r$ is $p(C_r = c_r) = \mathcal{N}(c_r|\mu_r, \sigma_r)$, where the mean and variance are $\mu_r = \sum_{i=1}^{N} \omega_i \mu_i$ and $\sigma_r^2 = \sum_{i=1}^{N} \omega_i \sigma_i^2$.

Similarly, assuming $\{C_r\}$ for different rays are independent, the random variable $C_I$ for the mean rendered color of rays sampled from an image pool is also under a Gaussian distribution, whose mean and variance are

$$\mu_I = \frac{1}{R} \sum_{r=1}^{R} \mu_r = \frac{1}{R} \sum_{r=1}^{R} \sum_{i=1}^{N} \omega_{ri} \mu_{ri}, \tag{3}$$

$$\sigma_I^2 = \frac{1}{R} \sum_{r=1}^{R} \sigma_r^2 = \frac{1}{R} \sum_{r=1}^{R} \sum_{i=1}^{N} \omega_{ri} \sigma_{ri}^2, \tag{4}$$

where $r$ represents a camera ray tracing through a pixel from an image pool and $R$ the number of rays sampled for the image pool.

*2) Optimization:* We now consider how to optimize the network parameters $\theta$ and determine $\mu, \sigma, \rho$ for each point on a ray (we ignore the subscript $ri$ for brevity). Consider a pool of images, the likelihood for the mean color value $c_I$ of a set of rays shooting from the pixels sampled from the image pool are obtained as $p(C_I = c_I) = \mathcal{N}(c_I | \mu_I, \sigma_I)$. Taking the negative logarithm of the likelihood and ignoring the constant, we have

$$L_{color} = \log \sigma_I + \frac{L_{mean}}{2\sigma_I^2} \le \log \sigma_I + \frac{L_I}{2\sigma_I^2}, \qquad (5)$$

$$L_{mean} = \|c_I - \mu_I\|_2^2 \le L_I = \frac{1}{R} \sum_{r=1}^{R} \|c_r - \mu_r\|_2^2. \qquad (6)$$

For $L_{mean}$ of (6), the constraint for each pixel is too weak and the network is not able to converge a meaningful result (PSNR about 10). As we have supervision for the color of each pixel and $L_{mean}$ is smaller than or equal to $L_I$, we choose to minimize $L_I$ instead. The loss function above is differentiable w.r.t $\theta$ and $\mu, \sigma, \rho$ are the outputs of MLPs ($\mathbf{F}_\theta$); we can use gradient descent to determine the network parameters $\theta$ and $\mu, \sigma, \rho$. To make $\sigma$ positive, we let the network estimate $\sigma^2$ and the output of MLPs $s$ is activated by $e^s$ to get $\sigma^2$.

Consider the minimization with respect to $\theta$. Given $\sigma_I$, we can see the minimization of the maximum likelihood under a conditional Gaussian distribution for each point is equivalent to minimizing a mean-of-squares error function given by $L_I$ in (6). Apply $\mathbf{F}_\theta$ on a point on a ray and $\mu, \rho, \sigma$ can be obtained.

*3) Neural Uncertainty and PSNR:* For the variance $\sigma_I^2$, or Neural Uncertainty, the optimum value can be achieved by setting the derivative of $L_{color}$ with respect to $\sigma_I$ to zero, giving

$$\sigma_I^2 = L_I. \qquad (7)$$

The equation above indicates that the optimal solution for $\sigma_I^2$ is the squared errors between the predicted image and the ground truth image.

On the other hand, PSNR is defined as $10 \log_{10} \frac{MAX_I^2}{MSE}$, where $MAX_I$ is the maximum possible pixel value of the image. When a pixel is represented using 8 bits, it is 255 and MSE is the mean squared error between two images, the same as $L_I$. Then we establish a linear relationship between the logarithm of Neural Uncertainty and PSNR, i.e.

$$PSNR = A \log \sigma_I^2, \qquad (8)$$

where $A$ is a constant coefficient.

To verify the linear relationship, we scatter data pairs of $(PSNR, \log \sigma_I^2)$ for images in the testing set evaluated at different iterations when optimizing $\mathbf{F}_\theta$ for a cabin scene (the scene is shown in Section IV). Two different training strategies are conducted: online training with images captured along a planned trajectory added sequentially and offline batch training using all the images precaptured from the trajectory. As can be seen Fig. 3(c-d), a strong correlation exists between $PSNR$ and $\log \sigma^2$, whose Pearson Correlation Coefficient (PCC) is -0.96 for online and -0.92 for offline. The two variables are almost perfectly negatively linearly related. During the training, $\sigma_I$ and $L_I$ are jointly optimized. The loss curves of
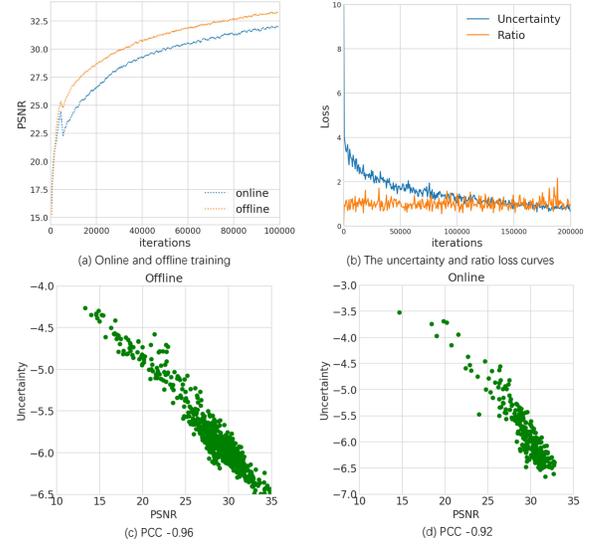


Fig. 3. Loss curves and linear relationship between $\log \sigma_I^2$ and PSNR. PCC value of -1 signifies strong negative correlation. (a) Training loss curves for online training and offline training. (b) The loss curves for the uncertainty part $\log \sigma_I$ and the ratio part $\frac{L_I}{\sigma_I^2}$ in (5). (c) Linearity when the scene is optimized using offline images. (d) Linearity when the scene is optimized using online images.

the uncertainty part $\log \sigma_I$ and the ratio part $\frac{L_I}{\sigma_I^2}$ in (5) for online are shown in Fig. 3(b). Notice that the loss curve for the ratio part stays almost constant during training, which also verifies the effectiveness using uncertainty as a proxy of PSNR.

For the verification and the usage of the linear relationship for autonomous reconstruction, we adopt NeRF [1] as our implicit representation for a scene. From the formulation and the derivation of the relation between neural uncertainty and PSNR, our neural uncertainty is agnostic to the underlying function $\mathbf{F}_\theta$, which can be MLPs like NeRF or networks based on trainable feature vectors with MLP Decoder [19]

### C. Online Training and Acceleration

Though online reconstruction with Neural Uncertainty supervised by images can achieve similar accuracy with NeRF, the convergence is too slow for view planning. We accelerate training by introducing particle filter, depth supervision and a keyframe strategy.

**Particle filter** keeps particles (rays) active in high loss region, which helps the network optimize details faster. At each step, when a new image is added to an image pool for the training, a set of particles are randomly sampled from the image. After an iteration of training, a quarter of particles are resampled according to the weight of particles, which is defined according to the loss of a ray, and the other particles are uniformly sampled from the image.

Particle filter is applied after coarse learning of the whole scene is done. At the early stage of training, the model knows little information about the scene and tends to have a higher loss for rays shooting at objects than that for empty space. This results in particles always staying at the surfaces of objects and therefore the network learns the whole space slower. Considering this, at the beginning iterations of the training, we use random sampling.

**Depth supervision** can greatly speed up training. Depth images for NeRF are rendered similar to color images in [1]. We define depth loss and our final loss as

$$L_{depth} = \frac{1}{R}\sum_{r=1}^{R}\|\hat{z}_r - z_r\|_2^2 , \quad L = L_{color} + \lambda_d L_{depth}, \quad (9)$$

where $\hat{z}_r$ and $z_r$ represent the rendered depth from reconstructed 3D model $\mathbf{F}_\theta$ and the ground truth depth for a pixel. As depth captured from real sensors typically has noise, we find using depth supervision may make the model not able to converge well due to the conflict between noisy depth and the depth inferred from multiview RGB images. A balance needs to be made between the two cases. We strengthen depth supervision at the early stage of training to accelerate training and decrease it after getting a coarse 3D structure. The weight of depth loss is decreased from 1 to 1/10 after $N_d$ iterations to emphasize structure from multiview images.

**Keyframe pool** We follow iMAP [2] to maintain a keyframe pool containing 4 images for continual training. The pool is initialized with the first four views and during training, the image with minimum image loss in the pool is replaced with a new image or an image from a seen images set.

### D. View Planning with Neural Uncertainty

As explained in Section III-A, we follow the greedy strategy. For the view path planning, we exploit the scenic path planner based on RRTs* [3] to evaluate the efficacy of the proposed Neural Uncertainty in guiding the view planning. In the scenic path planner, the RRTs* algorithm samples from a prior distribution for the view quality cost-space approximated by Sequential Monte-Carlo (SMC) instead from $SE(3)$, which baises the growth of the tree towards areas with good NBV cost while also aims at a shortest Euclidean path. Our view planner replaces their view cost with

$$C_{view} = \sigma_I^2, \quad (10)$$

where $\sigma_I^2$ is defined in (4) and here the image pool contains only one image.

For the view planning, we introduce a prior of objects in the center of the space for the candidate viewpoint sampling: the camera moves in a band with the nearest distance to the object $D_{near}$ meters and the farthest distance $D_{far}$ meters; the camera looks at directions pointing roughly to the center of the object; candidate viewpoints are sampled in a sphere with the center at the current camera position and the radius $R_{sample}$ meters.

## IV. RESULTS

**Data** The experiments are conducted on five 3D models, $drums$ from NeRF synthetic dataset, $Alexander^1$, $cabin$, $tank$ and $monsters$ we collect online. The scenes used in the paper are mostly of the size $5m\times5m\times5m$, and the larger scene $Alexander$ is about $80m\times80m\times80m$. $D_{near}$, $D_{far}$, $R_{sample}$ are 80m, 90m, 80m for $Alexander$ and 3m, 4m, 3m for

[1]https://3dwarehouse.sketchup.com/

other scenes. RGBD images are rendered by Unity Engine and we assume their corresponding camera poses are known. To simulate depth noise, all rendered depth images are added with noise scaling approximately quadratically with depth $z$ [20]. The depth noise model is $\epsilon = N(\mu(z), \sigma(z))$ where $\mu(z) = 0.0001125z^2 + 0.0048875$, $\sigma(z) = 0.002925z^2 + 0.003325$ and the constant parameters are acquired by fitting the model to the noise reported for Intel Realsense L515 [2]. The depth noise model for $Alexander$ of large size, $\mu(z) = 0.00001235z^2 + 0.00004651$, $\sigma(z) = 0.00001228z^2 + 0.00001571$, the constant parameters acquired by fitting the model to the noise reported for Lidar VLP16[3].

**Implementation details** The networks and the hyperparameters of all the experiments for different scenes below are set to the same value. Most hyper parameters use the default setting in NeRF, including parameters of Adam optimizer (with hyperparameters $\beta_1 = 0.9$, $\beta_2 = 0.999$, $\epsilon = 10^{-7}$), 64/128 sampling points on a ray for coarse/fine sampling, a batch size of 1024 *etc*.
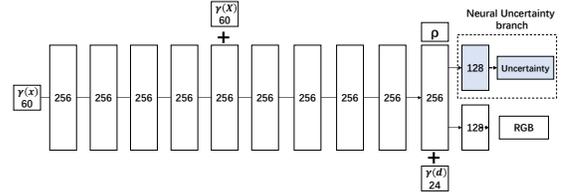


Fig. 4. The network architecture of the 3D reconstruction module of NeurAR. Similar to NeRF. The only difference is an additional branch of a 128 fully connected layer for the estimation of uncertainty
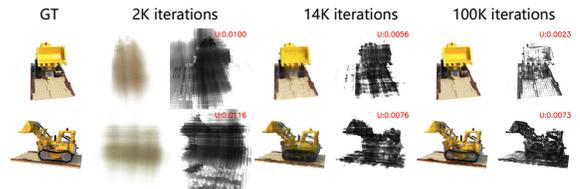


Fig. 5. Uncertainty maps from different viewpoints and training iterations. Top: Uncertainty maps for seen viewpoints; Bottom: uncertainty maps for an unseen viewpoint. The overall uncertainty decreases with training while the uncertainty for the unseen viewpoint stays high even when the network converges.

NeurAR reconstructs scenes and plans the view path simultaneously, running parallel on two RTX3090 GPUs: the 3D reconstruction module, *i.e.* optimization of NeRF on a GPU; the view planner module and the renderer module on the other one. After a step of the view planner (also 700 iterations for the NeRF optimization), the optimization receives one or two images collected on the view path from cameras. We set the maximum views for planning to be 28 views, *i.e.* 11k training iterations and about 15 steps for the optimization during the view planning. In each step, 1 or 2 views are selected along the planned path. The network architecture of our NeurAR is shown in Fig. 4.

**Algorithmic Variants** To evaluate the efficacy of the proposed method, we designed baselines and variants of our method.

[2]https://www.intelrealsense.com/lidar-camera-l515/
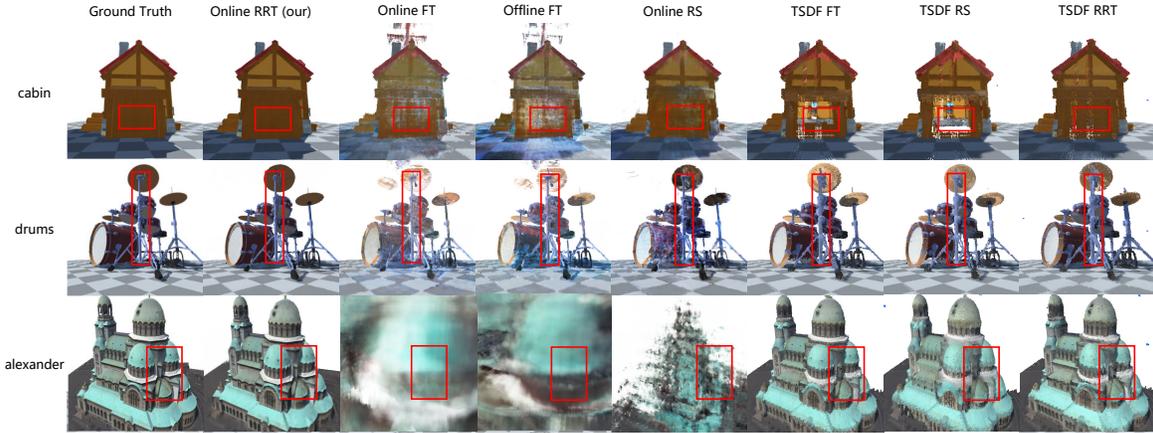[3]https://usermanual.wiki/Pdf/VLP16Manual.1719942037/view/

Fig. 6. Comparison of the reconstruction models with different methods. Refer to the supplementary video for higher resolution, more comparison and more viewpoints.

TABLE I
EVALUATIONS ON THE RECONSTRUCTED 3D MODELS USING DIFFERENT METHODS

| Method | cabin | | | drums | | | alexander | | | tank | | | monsters | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | PSNR↑ | SSIM↑ | LPIPS↓ | PSNR↑ | SSIM↑ | LPIPS↓ | PSNR↑ | SSIM↑ | LPIPS↓ | PSNR↑ | SSIM↑ | LPIPS↓ | PSNR↑ | SSIM↑ | LPIPS↓ |
| TSDF FT | 21.17 | 0.768 | 0.140 | 18.67 | 0.746 | 0.150 | 19.30 | 0.662 | 0.190 | 20.39 | 0.782 | 0.151 | 20.42 | 0.786 | 0.126 |
| TSDF RS | 20.87 | 0.75 | 0.151 | 18.68 | 0.738 | 0.160 | 19.15 | 0.649 | 0.195 | 18.64 | 0.749 | 0.171 | 20.19 | 0.777 | 0.131 |
| TSDF RRT [9] | 20.34 | 0.739 | 0.171 | 17.95 | 0.716 | 0.181 | 18.86 | 0.644 | 0.199 | 21.36 | 0.772 | 0.172 | 20.14 | 0.773 | 0.146 |
| Offline FT | 24.31 | 0.842 | 0.108 | 21.17 | 0.831 | 0.117 | 8.40 | 0.431 | 0.606 | 21.71 | 0.819 | 0.142 | 25.01 | 0.880 | 0.076 |
| Online FT | 23.42 | 0.833 | 0.114 | 21.49 | 0.831 | 0.117 | 8.71 | 0.419 | 0.597 | 23.03 | 0.839 | 0.124 | 24.42 | 0.878 | 0.077 |
| Online RS | 26.74 | 0.864 | 0.082 | 23.66 | 0.859 | 0.0838 | 13.08 | 0.533 | 0.373 | 21.97 | 0.789 | 0.162 | 24.68 | 0.871 | 0.063 |
| Online even cover. | 26.56 | 0.868 | 0.106 | 24.82 | 0.913 | 0.049 | **24.30** | **0.767** | **0.182** | 24.81 | 0.875 | 0.108 | 25.90 | 0.909 | 0.065 |
| Ours Offline | **28.86** | **0.917** | **0.048** | **26.45** | **0.916** | **0.051** | 23.98 | 0.758 | 0.188 | **27.54** | **0.909** | **0.064** | **27.57** | **0.927** | **0.039** |
| Ours | 28.35 | 0.902 | 0.062 | 25.73 | 0.905 | 0.058 | 24.07 | 0.757 | 0.199 | 25.83 | 0.874 | 0.097 | 26.57 | 0.908 | 0.054 |

| Method | Acc↓ | Comp↓ | C.R.↑ | Acc↓ | Comp↓ | C.R.↑ | Acc↓ | Comp↓ | C.R.↑ | Acc↓ | Comp↓ | C.R.↑ | Acc.↓ | Comp↓ | C.R.↑ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| TSDF FT | 2.47 | 2.68 | 0.39 | 2.51 | 1.64 | 0.21 | 87.16 | 164.30 | 0.44 | 2.46 | 3.14 | 0.40 | 2.44 | 1.21 | 0.43 |
| TSDF RS | 2.83 | 2.81 | 0.42 | 3.09 | 1.56 | 0.22 | 107.85 | 159.78 | 0.44 | 3.04 | 3.53 | 0.40 | 2.93 | 1.19 | 0.48 |
| TSDF RRT [9] | 2.05 | 2.53 | 0.44 | 2.90 | 1.33 | 0.25 | 100.85 | 169.56 | 0.43 | 2.01 | 1.62 | 0.48 | 2.97 | 1.14 | 0.48 |
| Offline FT | 1.09 | 1.02 | 0.70 | 1.77 | 1.21 | 0.66 | 1582.74 | 193.77 | 0.04 | 1.38 | 1.28 | 0.63 | **1.10** | 1.10 | 0.71 |
| Online FT | 1.19 | 1.06 | 0.67 | 1.79 | 1.22 | 0.65 | 1506.59 | 202.52 | 0.05 | 1.37 | 1.30 | 0.62 | 1.19 | 1.09 | 0.71 |
| Online RS | 1.88 | 1.20 | 0.57 | 1.53 | 1.10 | 0.70 | 652.66 | 211.00 | 0.27 | 2.62 | 2.25 | 0.37 | 1.89 | 1.11 | 0.66 |
| Online even cover. | 1.21 | 1.01 | 0.74 | **1.15** | 1.08 | **0.77** | 54.24 | 35.69 | 0.48 | **1.50** | 1.29 | 0.64 | 1.17 | 1.06 | 0.76 |
| Ours Offline | **0.96** | **0.93** | **0.76** | 1.31 | 1.07 | 0.74 | **31.66** | **21.90** | **0.73** | 1.51 | **1.19** | **0.68** | 1.27 | **1.01** | **0.76** |
| Ours | 1.07 | 0.95 | 0.74 | 1.29 | **1.04** | 0.75 | 48.48 | 34.16 | 0.60 | 1.37 | 1.20 | 0.65 | 1.43 | 1.02 | 0.74 |

For the implicit scene representation, we adopt TSDF as the baseline as it is oen of the most used representation for SLAM and autonomous reconstruction [4], [9], [21]. For view path planning based on the proposed Neural Uncertainty, we construct two variants: one using a pre-defined circular trajectory with which existing work [21], [22] usually compares and the other one randomly sampling a viewpoint instead of using the view cost to choose NBV at each step.

The variants are 1) **TSDF FT**: RGBD images and corresponding poses are collected from a **F**ixed circular **T**rajectory and are fed into the system sequentially. The voxel resolution of TSDF is 1cm. 2) **TSDF RS**: replace the fixed trajectory in **TSDF FT** with **R**andomly **S**ampled NBVs. 3)**TSDF RRT**: As for autonomous scene reconstruction most work is not open source, we re-implement the view cost defined in [9] which adopts TSDF representation. The online trajectory is planned with **RRT** and the view cost is defined according to the quality of a reconstructed voxel, measured by the number of rays traversing through it. Images and corresponding poses for the fusion are collected from the planned views and are fed into the system sequentially. 4) **Offline FT**: The variant is NeurAR with views from the fixed trajectory and trained offline. 5) **Online FT**: The variant is NeurAR with views from the fixed trajectory. 6) **Online RS**: NeurAR with randomly

sampled NBVs. 7) **Online even cover.**: NeurAR with views evenly distributed in a dome to cover the whole space. 8) **Ours offline**: NeurAR trained offline, *i.e.* trained from scratch with images precaptured from all the planned views from our proposed method.

**Metrics** The quality of the reconstructed models can be measured in two aspects: the quality of the rendered images (measure both the geometry and the texture) and the quality of the geometry of the constructed surface. For the former, we evenly distribute about 200 testing views 80m for $Alexander$, 3m and 3.4m from the center for other scenes, render images for the reconstructed models and evaluate PSNR, SSIM and LPIPS for these images. For the latter, we adopt metrics from iMAP [2]: $Accuracy$ (cm), $Completion$ (cm), $Completion Ratio$ (the percentage of points in the reconstructed mesh with $Completion$ under a threshold, 30 cm for $Alexander$ and 1cm for others). For the geometry metrics, about 300k points are sampled from the surfaces.

### A. Neural Uncertainty

Fig. 3 has quantitatively verified the correlation between Neural Uncertainty and the image quality. Here, we further demonstrate the correlation with examples. In this experiment, the reconstruction network is optimized by the loss

defined in (5) with 73 images collected from cameras placed roughly at one hemisphere. Fig. 5 shows the images rendered from different viewpoints for the reconstructed 3D models and their corresponding uncertainty map during training. The uncertainty map is rendered from an uncertainty field where the value of each point in the field is $\frac{1}{log\sigma}$. Viewpoints of Row 1 are in the hemisphere seen in the training and the viewpoint of the last row is in the other hemisphere. At the beginning, for all viewpoints, the object region and the empty space both have high uncertainty. For the seen viewpoints, with training continuing, the uncertainty of the whole space decreases and the quality of the rendered images improve. When the network converges, the uncertainty only remains high in the local areas having complicated geometries. For the unseen viewpoint, though the uncertainty for the empty space decreases dramatically with training, the uncertainty is still very high on the whole surface of the object even when the network converges.

Fig. 1 shows the Neural Uncertainty can successfully guide the planner to planning view paths for cameras to look at regions that are not well reconstructed. For example, given current pose 15 and $\mathbf{F}_\theta$, a path is planned toward viewpoint 16 having a higher uncertainty map. The rendered image from viewpoint 16 from $\mathbf{F}_\theta$ exhibits poor quality (zoom in the image and the uncertainty map for details under the roof).

### B. View planning with Neural Uncertainty

To show the efficacy of our proposed method, we compare metrics in Table I and the rendered images of reconstructed 3D models using different methods in Fig. 6. Table I demonstrates except for the even coverage, our method outperforms all variants on all metrics significantly. Our proposed NeurAR achieves better reconstruction results with shorter view paths compared with other variants and existing work.

Compared with methods using the implicit representation without path planning (**Offline FT**, **Online FT**, **Online RS**), our method demonstrates significant improvements in the image quality and geometry accuracy. Methods without the planned views 1) are even unable to converge to a reasonable result in the large scene (NeRF will fail due to overfitting without carefully planned input viewpoint coverage), 2) have many holes in the objects as these regions are unseen in the images and 3) inferior image details on the surface of objects as these regions have a little overlap of different views, making inferring the 3D geometry hard (check red box in Fig. 6 for visual comparison). In addition, these variants tend to have ghost effects in the empty space. This is largely because the viewpoints are designed to make the camera look at the objects and the empty space is not considered. Notice that placing viewpoints covering the whole space without the aid of visualization tools as feedback is non-trivial for humans.

Assuming we have the 3D shape of a scene with the target object in the center (which is our goal), to cover the whole space, we distribute viewpoints evenly at a dome centered at the scene center with a radius of 80 meters for *Alexander* and 3 meters for others. All viewpoints look at the center. An even coverage path can provide a good reference as even coverage

for many scenes of simple shape and texture is an optimum solution. Metrics for **Online even cover.** and **Ours** in Table I demonstrate that our planned views can achieve similar or even better reconstruction results and our method has no prior or only a minor knowledge about the scene.

**Ours offline** in most scenes gives better results than NeurAR trained online as using all views enforces multiview constraints at the start of the training.

In addition to the lower mean PSNR shown in the Table I, our method achieves much smaller variance regarding the PSNR of the rendered images from different viewpoints. For example, the PSNR variance of the rendered images for reconstructed cabins using **Online FT**, **Online RS** and our method is 32.55dB, 15.44dB to 4.78dB, indicating our method more even. Further, for the average path length in the smaller scenes, our NeurAR traverses 43.27m while **Online RS** about 70.24m; in the larger scene, our NeurAR traverses 907.60m while **Online RS** about 1329.75m.

For models using the reconstruction with TSDF (**TSDF FT**, **TSDF RS**, **TSDF RRT**), we render images via volume rendering. From the images in the last three columns in Fig. 6, the surfaces exhibit many holes; in addition, the surfaces of the reconstructed models are rugged due to the noise in the depth images. Though NeurAR uses depth too, it depends on the depth images to accelerate convergence only at the early stage of training and decreases its effect after coarse structures have been learned. The finer geometry is acquired by multiview image supervision. NeurAR fills holes in **TSDF RRT** and provide finer details. It outperforms **TSDF RRT** in the image quality, geometry quality and path length (43.27m vs 57.39m). Though post-processing can be applied to extra finer meshes for the reconstructions using TSDF and get smoother images than our rendering images from TSDF directly, denoising and filling the holes are non-trivial tasks, particularly in scenes having complex structures.

### C. Ablation study

**Training Iterations between Steps** The number of iterations for NeRF optimization $iter$ allowed for planning a view step affects final results. We choose different iterations to run our NeurAR system in the scene $cabin$ and compare PSNR of the rendered images of the reconstructed models. PSNR for the models optimized using 300, 700, 1400 iterations is 26.16, 28.58, 26.91. This is because too few iterations may lead to uncertainty not being optimized well while too many steps may lead to overfitting to images added.

**Depth noise** We construct variants of NeurAR using different noise magnitude from no noise, the noise magnitude of L515, to the double ($2 \times \mu(z)$, $2 \times \sigma(z)$) and the triple ($3 \times \mu(z)$, $3 \times \sigma(z)$) noise magnitude of L515. Table II shows the metrics measuring the cabin models reconstructed by these variants: NeurAR can be accelerated with very noisy depth at a cost of only a minor drop in the reconstruction quality. The robustness verifies the effect of dampening the depth supervision during the training in Section III-C .

**Number of views** To further study the influence, we set the maximum allowed views from 18 to 58 for the $cabin$ scene and

TABLE II
INFLUENCE OF DIFFERENT NOISE MAGNITUDE

| Method | PSNR↑ | SSIM↑ | LPIPS↓ | Acc↓ | Comp↓ | C.R.↑ |
|---|---|---|---|---|---|---|
| Ours (no noise) | 28.71 | 0.926 | 0.043 | 0.92 | 0.90 | 0.75 |
| Ours (with noise) | 28.35 | 0.902 | 0.062 | 1.07 | 0.95 | 0.74 |
| Ours (noise x2) | 28.19 | 0.903 | 0.062 | 0.99 | 0.97 | 0.74 |
| Ours (noise x3) | 27.65 | 0.900 | 0.063 | 1.00 | 0.94 | 0.77 |

provide the PSNR for the reconstructed models under different view settings.

TABLE III
EVALUATIONS ON THE RECONSTRUCTED 3D MODELS USING DIFFERENT NUMBER OF VIEWS

| Views | PSNR↑ | SSIM↑ | LPIPS↓ | Acc↓ | Comp↓ | C.R.↑ |
|---|---|---|---|---|---|---|
| 18 | 22.47 | 0.802 | 0.155 | 2.18 | 1.31 | 0.59 |
| 28 | 28.35 | 0.902 | 0.062 | 1.07 | 0.95 | 0.74 |
| 38 | 30.03 | 0.919 | 0.053 | 0.90 | 0.86 | 0.77 |
| 58 | 29.39 | 0.926 | 0.050 | 0.82 | 0.90 | 0.79 |

## V. CONCLUSION

In this paper, we have presented an autonomous 3D reconstruction method based on an implicit neural representation, with the view path planned according to a proxy of PSNR. The proxy is acquired by learning uncertainty for the estimation of the color of a spatial point on a ray during the reconstruction. A strong correlation is discovered between the uncertainty and PSNR, which is verified by extensive experiments. Compared with variants with a pre-defined trajectory and variants using TSDF, our method demonstrates significant improvements. Our work shows a promising potential of implicit neural representations for high-quality autonomous 3D reconstruction.

One limitation of NeurAR is the optimization of implicit neural representation is slow and the computation consumption is high. The optimization of the model between view steps takes about 50-120 seconds depending on the number of iterations, which is much slower than existing unknown area explorations. However, our method is agnostic to the underlying NeRF realization and therefore, we can easily swap the module for more advanced NeRF variants emerging lately [23]. Future work includes achieving acceleration of the training and rendering for NeurAR.

The other limitation is that we assume known camera poses, which require extra location systems to make NeurAR work. Future work includes joint optimization of the camera poses, reconstruction, and view planning. Also, extending NeurAR to multiple agents for large scene reconstructions with distributed learning [24] and adverse environments with the mmWave radar [25] are interesting directions.

## REFERENCES

[1] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng, "Nerf: Representing scenes as neural radiance fields for view synthesis," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2020, pp. 405–421.

[2] E. Sucar, S. Liu, J. Ortiz, and A. J. Davison, "imap: Implicit mapping and positioning in real-time," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021, pp. 6229–6238.

[3] O. Mendez, S. Hadfield, N. Pugeault, and R. Bowden, "Taking the scenic route to 3d: Optimising reconstruction from moving cameras," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2017, pp. 4677–4685.

[4] S. Isler, R. Sabzevari, J. Delmerico, and D. Scaramuzza, "An information gain formulation for active volumetric 3d reconstruction," in *2016 IEEE International Conference on Robotics and Automation (ICRA)*, 2016, pp. 3477–3484.

[5] S. Wu, W. Sun, P. Long, H. Huang, D. Cohen-Or, M. Gong, O. Deussen, and B. Chen, "Quality-driven poisson-guided autoscanning," *ACM Transactions on Graphics (TOG)*, vol. 33, no. 6, pp. 1–12, 2014.

[6] S. Song and S. Jo, "Surface-based exploration for autonomous 3d modeling," in *2018 IEEE International Conference on Robotics and Automation (ICRA)*, 2018, pp. 4319–4326.

[7] C. M. Bishop, *Pattern Recognition and Machine Learning*. Springer-Verlag New York, 2006.

[8] Q. Ye and T.-K. Kim, "Occlusion-aware hand pose estimation using hierarchical mixture density network," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 801–817.

[9] L. Schmid, M. Pantic, R. Khanna, L. Ott, R. Siegwart, and J. Nieto, "An efficient sampling-based method for online informative path planning in unknown environments," *IEEE Robotics and Automation Letters (RAL)*, vol. 5, no. 2, pp. 1500–1507, 2020.

[10] J. I. Vasquez-Gomez, L. E. Sucar, and R. Murrieta-Cid, "View/state planning for three-dimensional object reconstruction under uncertainty," *Autonomous Robots*, vol. 41, no. 1, pp. 89–109, 2017.

[11] M. Devrim Kaba, M. Gokhan Uzunbas, and S. Nam Lim, "A reinforcement learning approach to the view planning problem," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 6933–6941.

[12] L. Schmid, C. Ni, Y. Zhong, R. Siegwart, and O. Andersson, "Fast and compute-efficient sampling-based local exploration planning via distribution learning," *IEEE Robotics and Automation Letters (RAL)*, vol. 7, no. 3, pp. 7810–7817, 2022.

[13] B. Hepp, D. Dey, S. N. Sinha, A. Kapoor, N. Joshi, and O. Hilliges, "Learn-to-score: Efficient 3d scene exploration by predicting view utility," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 437–452.

[14] J. L. Schönberger, E. Zheng, J.-M. Frahm, and M. Pollefeys, "Pixelwise view selection for unstructured multi-view stereo," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2016, pp. 501–518.

[15] S. Izadi, D. Kim, O. Hilliges, D. Molyneaux, R. Newcombe, P. Kohli, J. Shotton, S. Hodges, D. Freeman, A. Davison *et al.*, "Kinectfusion: real-time 3d reconstruction and interaction using a moving depth camera," in *Proceedings of the 24th annual ACM symposium on User interface software and technology (UIST)*, 2011, pp. 559–568.

[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, no. 3, pp. 189–206, 2013.

[17] "Nerf project page," https://www.matthewtancik.com/nerf, 2020.

[18] A. Gropp, L. Yariv, N. Haim, M. Atzmon, and Y. Lipman, "Implicit geometric regularization for learning shapes," in *Proceedings of the International Conference on Machine Learning (ICML)*, 2020, pp. 3789–3799.

[19] X. Zhang, S. Bi, K. Sunkavalli, H. Su, and Z. Xu, "Nerfusion: Fusing radiance fields for large-scale scene reconstruction," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022, pp. 5449–5458.

[20] C. V. Nguyen, S. Izadi, and D. Lovell, "Modeling kinect sensor noise for improved 3d reconstruction and tracking," in *2012 Second International Conference on 3D Imaging, Modeling, Processing, Visualization & Transmission*, 2012, pp. 524–530.

[21] S. Song, D. Kim, and S. Choi, "View path planning via online multiview stereo for 3-d modeling of large-scale structures," *IEEE Transactions on Robotics*, 2021.

[22] D. Peralta, J. Casimiro, A. M. Nilles, J. A. Aguilar, R. Atienza, and R. Cajote, "Next-best view policy for 3d reconstruction," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2020, pp. 558–573.

[23] T. Müller, A. Evans, C. Schied, and A. Keller, "Instant neural graphics primitives with a multiresolution hash encoding," *ACM Transactions on Graphics (TOG)*, vol. 41, no. 4, pp. 1–15, 2022.

[24] Y. Huang, Y. Sun, Z. Zhu, C. Yan, and J. Xu, "Tackling data heterogeneity: A new unified framework for decentralized sgd with sample-induced topology," in *Proceedings of the International Conference on Machine Learning (ICML)*, 2022, pp. 9310–9345.

[25] A. Chen, X. Wang, S. Zhu, Y. Li, J. Chen, and Q. Ye, "mmbody benchmark: 3d body reconstruction dataset and analysis for millimeter wave radar," in *Proceedings of the 30th ACM International Conference on Multimedia*, 2022, pp. 3501–3510.