

---

# Learning Depth-regularized Radiance Fields from Asynchronous RGB-D Sequences

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Recently it is shown that learning radiance fields with depth rendering and depth  
2 supervision can effectively promote the view synthesis quality and convergence.  
3 But this paradigm requires input RGB-D sequences to be synchronized, hindering  
4 its usage in the UAV city modeling scenario. To this end, we propose to jointly  
5 learn large-scale depth-regularized radiance fields and calibrate the mismatch  
6 between RGB-D frames. Although this joint learning problem can be simply  
7 addressed by adding new variables, we exploit the prior that RGB-D frames are  
8 actually sampled from the same physical trajectory. As such, we propose a novel  
9 **time-pose function**, which is an implicit network that maps timestamps to SE(3)  
10 elements. Our algorithm is designed in an alternative way consisting of three  
11 steps: (1) time-pose function fitting; (2) radiance field bootstrapping; (3) joint  
12 pose error compensation and radiance field refinement. In order to systematically  
13 evaluate under this new problem setting, we propose a large synthetic dataset with  
14 diverse controlled mismatch and ground truth. Through extensive experiments,  
15 we demonstrate that our method outperforms strong baselines. We also show  
16 qualitatively improved results on a real-world asynchronous RGB-D sequence  
17 captured by drones. Codes, data, and models will be made publicly available.

## 18 1 Introduction

19 Incorporating depth rendering and depth supervision into radiance fields has been demonstrated as  
20 a helpful regularization technique in several recent studies [2, 21, 32, 20]. However, this technique  
21 has not yet been successfully introduced into radiance field learning from UAV (Unmanned Aerial  
22 Vehicle) images, despite it's a typical choice in city modeling. A closer look at the aforementioned  
23 works reveals that they assume synchronized RGB and depth signals, which is hard to guarantee in  
24 UAV vision due to the lack of suitable synchronized sensors for long sensing ranges. So we study the  
25 *problem* of learning depth-regularized radiance fields from asynchronous RGB-D sequences.

26 As a recap, the canonical radiance field [15] learns a neural network parameterized by  $\theta$  that represents  
27 a 3D scene, from input images  $I$  and their intrinsic/extrinsic parameters  $\mathcal{T}_I$ . To alleviate the reliance  
28 on  $\mathcal{T}_I$ , some works [30, 11, 8] aim to resolve a different problem that self-calibrates  $\mathcal{T}_I$ . In other  
29 words, they jointly learn  $\theta$  and  $\mathcal{T}_I$  from input images  $I$ . Similarly, the *formulation* considered here is  
30 to learn scene representation  $\theta$ , camera parameters  $\mathcal{T}_I$  and  $\mathcal{T}_D$  from inputs images  $I$  and depths  $D$ .

31 It is natural to develop the joint learning *formulation* to resolve the *problem*, as it amounts to adding  
32 new parameters to existing methods [30, 11, 8]. However, an important prior is ignored that RGB-D  
33 frames are actually sampled from the same physical trajectory. As conceptually shown in Fig. 1-a/b,  
34  $\mathcal{T}_I$  and  $\mathcal{T}_D$  can be considered as samples from a function that maps timestamps to SE(3) elements.  
35 We name this function as **time-pose function** and model it with a neural network parameterized by  
36  $\phi$ . As such, we address the *problem* with a *new formulation* that learns scene representation  $\theta$  and

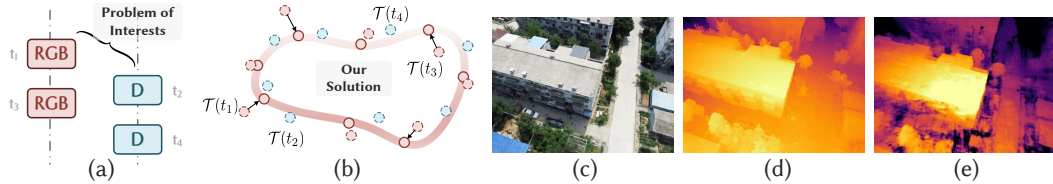


Figure 1: The problem of interest is to learn a depth-regularized radiance field using asynchronous RGB-D sequences (a). We proposed a time-pose function as conceptually shown in (b) to leverage the prior that RGB-D sequences are actually sampled from the same physical underlying trajectory. For a novel view (c), our method can render a better depth map (d) than Mega-NeRF (e).

37 time-pose function  $\phi$  from inputs RGB images  $I$  and depths  $D$ . An interesting fact is that both  $\theta$  and  
 38  $\phi$  are implicit neural representation networks (or say coordinate-based networks) that allow fully  
 39 differentiable training. To our knowledge, this *new formulation* has not been proposed before.

40 We also propose an effective learning scheme designed in an alternative manner. In the first stage, we  
 41 fit the time-pose function  $\phi$  using one modality (e.g., RGB images) and infer the poses of the other,  
 42 using a balanced pose regression loss and a speed regularization term. Secondly, we bootstrap a large-  
 43 scale radiance field  $\theta$  based upon Mega-NeRF [29] using the outputs of the trained time-pose function.  
 44 Thanks to the first step, depth regularization can be imposed here in spite of RGB-D misalignment.  
 45 Finally, thanks to the cascade of two fully differentiable implicit representation networks, we jointly  
 46 optimize the 3D scene representation  $\theta$  and compensate pose errors by updating  $\phi$ .

47 Since the *problem* considered is new, we contribute a synthetic dataset (named AUS) for systematic  
 48 evaluation. Using six large-scale 3D scenes, realistic drone trajectories of different difficulty levels  
 49 are generated. Specifically speaking, simple trajectories are heuristically designed with a zig-zag  
 50 pattern while complicated ones are generated by manual control signals in simulation. We also  
 51 control the mismatch between RGB-D sequences using different protocols, to cover as many as  
 52 possible scenarios that the algorithm may encounter in reality. Through a set of comprehensive  
 53 experiments, we show the proposed method outperforms several state-of-the-art counterparts and  
 54 our design choices contribute positively to performance. Last but not least, we present a real-world  
 55 evaluation using asynchronous sensors on drones. Our depth rendering results (on unseen viewpoint)  
 56 is shown in Fig. 1-d, which is much better than the result of Mega-NeRF shown in Fig. 1-e. This  
 57 success is credited to the usage of depth regularization as made possible by our novel algorithm.

58 To summarize, we have the following contributions in this paper: (1) We formalize the new *problem*  
 59 of learning depth-regularized radiance fields from asynchronous RGB-D sequences, which is rooted  
 60 in UAV city modeling. (2) We identify an important domain-specific prior in this problem: RGB-D  
 61 frames are sampled from the same underlying trajectory. We instantiate this prior into a novel time-  
 62 pose function and develop a cascaded fully differentiable implicit representation network. (3) In order  
 63 to systematically study the problem, we contribute a photo-realistically rendered synthetic dataset  
 64 that simulates different types of mismatch. (4) Through a comprehensive benchmarking on this new  
 65 dataset and real-world asynchronous RGB-D sequences, we demonstrate that our method can promote  
 66 performance over strong prior arts. Anonymous code: <https://anonymous.4open.science/r/async-nerf>

## 67 2 Related Works

68 **Large-scale Radiance Fields.** Neural Radiance Field (NeRF) [15] has shown impressive results in  
 69 neural reconstruction and rendering. However, its capacity to model large-scale unbounded 3D scenes  
 70 is limited. Several strategies [29, 26, 32, 14, 35] have been proposed to address this limitation, with a  
 71 common principle of dividing large scenes into blocks or decomposing the scene into multiple levels.  
 72 Block-NeRF [26] clusters images by dividing the whole scene according to street blocks. Mega-NeRF  
 73 [29] utilizes a clustering algorithm that partitions sampled 3D points into different NeRF submodules.  
 74 BungeeNeRF [32] trains NeRFs using a growing model of residual blocks with predefined multiple  
 75 scales of data. Switch-NeRF [14] designs a gating network to jointly learn the scene decomposition  
 76 and NeRFs without any priors of 3D scene shape or geometric distribution. However, these prior  
 77 works fail to leverage the rich geometric information in depth images for effective regularization.

78 **Depth-regularized Radiance Fields.** Volumetric rendering requires extensive samples and sufficient  
 79 views to effectively differentiate between empty space and opaque surfaces. Depth maps can  
 80 serve as geometric cues, providing regularization constraints and sampling prior, which accelerates  
 81 NeRF’s convergence towards the correct geometry. DS-NeRF [3] enhances this process using depth  
 82 supervision from 3D point clouds, estimated by structure-from-motion, and a specific loss for rendered  
 83 ray termination distribution. Mono-SDF [38] and Dense-Depth Prior [21] further supplement this  
 84 with a pretrained dense monocular depth estimator for less-observed and textureless areas. To adapt  
 85 NeRF for outdoor scenarios, URF [20] rasterizes a pre-built LiDAR point cloud map to generate  
 86 dense depth images and alleviates floating elements by penalizing floaters in the free space. Moreover,  
 87 S-NeRF [34] completes depth on sparse LiDAR point clouds using a confidence map, effectively  
 88 handling street-view scenes with limited perspectives. However, those methods are not readily  
 89 applicable to UAV captured images due to the lack of suitable synchronized sensors for long ranges.

90 **Broader UAV Vision and Synchronization.** Like autonomous driving, UAV vision is drawing  
 91 increasing attention due to its unique characteristics. Broader UAV vision covers many topics like  
 92 counting [31][7], trajectory forecasting [18], intention prediction [33], object tracking [16], physics  
 93 understanding [39], next-best-view prediction [6], 3D reconstruction [41], and calibration [19].  
 94 Sensor synchronization is challenging for UAV vision (and other settings) and several works address  
 95 the problem from an algorithmic perspective. One possibility is to adopt tailored hardware designs or  
 96 software protocols [1] to synchronize all the devices. Another branch of sensor-agnostic methods  
 97 utilizes temporal priors by using Sum-of-Gaussians [4] or parametric interpolation functions [36].

98

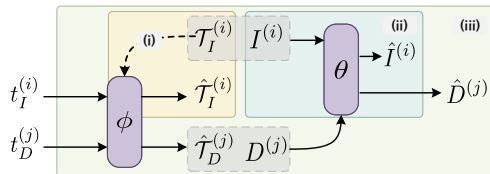


Figure 2: **Three-step Optimization.** (i) A time-pose function parameterized by  $\phi$  is trained to predict camera poses from timestamps; (ii) The neural radiance field parameterized by  $\theta$  is bootstrapped with pure RGB losses; (iii) Both of the parameters  $\theta, \phi$  are jointly optimized with RGB-D supervision.

### 99 3 Problem Formulation and Optimization Pipeline

100 **Problem & Challenge.** Our goal is to learn a neural radiance field parameterized by  $\theta$  for large-scale  
 101 scene representation from UAV images as done in prior works [29, 32]. However, these prior works  
 102 fail to leverage depth supervision, which is known [3, 21] as useful for training floater-less NeRFs.  
 103 To our knowledge, there are no easily accessible synchronized RGB-D sensor suites for **large-scale**  
 104 outdoor scenes, and trivially synchronizing them according to timestamp<sup>1</sup> cannot fully address the  
 105 misalignment issue. Instead of using expensive hardware, we take an algorithmic perspective.

106 **Input & output.** There are some prior works on large-scale scene modeling using aerial images  
 107 [32, 29, 5, 6]. In this study, we assume an input RGB-D stream captured by drones: a set of  
 108 RGB camera images  $\{I^{(i)}\}_{i=1}^{N_I}$  and a set of depth maps  $\{D^{(j)}\}_{j=1}^{N_D}$  (shown in Fig. 1-a) and we aim  
 109 to recover the spatiotemporal transformations between them. Given that our focus is on relative  
 110 transformation, it is viable to consider either the RGB or the depth stream as the reference without  
 111 compromising generality. For convenience, we assume a set of camera poses  $\{\mathcal{T}_I^{(i)}\}_{i=1}^{N_I}$  for color  
 112 images are obtained by an SfM algorithm. The neural scene representation parameterized by  $\theta$   
 113 outputs an image  $\hat{I}$  as well as a depth map  $\hat{D}$  at a given perspective camera pose  $\mathcal{T}$ .

114 **Observation.** Note that all the sensor data are captured with a drone on the **same trajectory**<sup>2</sup>, we  
 115 can model the relationship between capture time  $t$  and sensor poses  $\mathcal{T}$  with an implicit **time-pose**  
 116 **function** as  $\phi : t \rightarrow \hat{\mathcal{T}} = [\hat{\mathbf{x}}, \hat{\mathbf{q}}]$ , where  $t$  is the timestamp of capture, and the estimated pose  $\hat{\mathcal{T}}$  is  
 117 represented by a translation vector  $\hat{\mathbf{x}} \in \mathbb{R}^3$  and a quaternion  $\hat{\mathbf{q}} \in \mathbb{R}^4$ .

118 **Pipeline overview.** We formulate our method as a 3-step optimization problem (as shown in Fig. 2).  
 119 First, since the time-pose relationship for the RGB captures are given, we can train a time-pose

<sup>1</sup>The so-called soft synchronization.

<sup>2</sup>but they are not necessarily synchronized in terms of acquisition time.

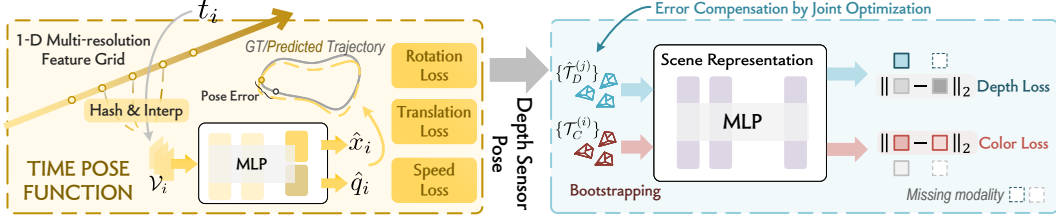


Figure 3: **Method Pipeline.** The time-pose function is modeled using a 1-D multi-resolution hash grid with direct and speed losses. After bootstrapping the scene representation networks with pure RGB signals, the predicted depth sensor poses are used for jointly optimizing the NeRFs’ parameters  $\theta$ . At each timestamp ( $t_i$  from RGB sequence or  $t_j$  from depth sequence), only one modality of sensor signals is provided, thus only one loss term is activated (shown on the right).

120 function on the RGB sequence (Fig. 2-(i)). Then, to train the neural radiance field, we first bootstrap  
 121 the network with pure RGB supervision (Fig. 2-(ii)). To further enable training with RGB-D  
 122 supervision, we can use the previously trained time-pose function to estimate the corresponding depth  
 123 camera poses  $\{\mathcal{T}_D^{(j)}\}$  of the depth timestamps  $\{t_D^{(j)}\}$ . Since both of the networks are differentiable,  
 124 we jointly optimize the networks in an end-to-end manner in the third stage (Fig. 2-(iii)).

## 125 4 Method

126 We introduce in Section 4.1 the details of learning an implicit time-pose function. In Section 4.2, we  
 127 describe our neural scene representation networks and the bootstrapping strategy. In Section 4.3, we  
 128 adopt depth supervision and jointly train the time-pose function with RGB-D pairs.

### 129 4.1 Time-Pose Function

130 We represent the camera trajectory as an implicit time-pose function  $\phi$  whose input is a timestamp  $t$ ,  
 131 and whose output is a 6-DoF pose  $\mathcal{T}$  that consists of a 3-D translation  $x_i$  and a 4-D quaternion  $q_i$ .

132 **Network Overview** The time-pose function (shown in the left part of Fig. 3) is approximated with  
 133 a compact 1-D multi-resolution hash grid  $\{\mathcal{G}^{(l)}\}_{l=1}^L$ , followed by an MLP decoder. The hash grid  
 134 consists of  $L$  levels of separate feature grids with trainable hash encodings [17]. The reason why we  
 135 choose this architecture is as follows: The time-pose function is a coordinate-based function that may  
 136 contain coarse and fine-level feature components <sup>3</sup>[27], and this architecture allows us to sample the  
 137 hash encodings from each grid layer with different resolutions and perform quadratic interpolation on  
 138 the extracted encodings to obtain a feature vector  $\mathcal{V}_i$  when querying a specific timestamp  $t$  that is in  
 139 the range of all timestamps. This design choice is also empirically validated in Table. 3.

140 After obtaining the interpolated feature vector, an MLP with two separated decoder heads is used  
 141 to predict the output translation  $\hat{x}_i$  and rotation  $\hat{q}_i$  vectors respectively. The forward pass can be  
 142 expressed in the following equations:

$$\mathcal{V}_i = \mathcal{F}_{\text{MLP}} \left( \text{concat}\{\text{interp}(\text{h}(t; \pi_l), \mathcal{G}_\theta^l)\}_{l=1}^L; \Phi_{\text{MLP}} \right), \quad (1)$$

$$\hat{\mathcal{T}}_i = [\hat{x}_i, \hat{q}_i] = l_{\text{trans}}(\mathcal{V}_i; \Phi_{\text{trans}}), l_{\text{rot}}(\mathcal{V}_i; \Phi_{\text{rot}}), \quad (2)$$

143 where  $\text{interp}$  denotes interpolation,  $\text{h}$  is the hash function parameterized by  $\pi_l$ ,  $\mathcal{F}_{\text{MLP}}$ ,  $l_{\text{trans}}$ ,  $l_{\text{rot}}$  are the  
 144 MLP networks and the decoder heads, with  $\Phi_{\text{MLP}}$ ,  $\Phi_{\text{trans}}$ ,  $\Phi_{\text{rot}}$  representing their parameters.

145 **Depth-pose Prediction** Since both the depth maps and the RGB images are collected by the same  
 146 drone on the same flight, they cover the same spatial-temporal footprints except for the difference  
 147 in the placement of the two sensors on the aircraft. For every depth frame, we first predict the  
 148 RGB camera pose using the capture timestamps of the depth sensor with the time-pose function  
 149 then transform the predicted RGB camera pose to the depth sensor pose with a pre-calibrated pose  
 150 transformation  $\mathcal{T}_{\text{RGB} \rightarrow \text{D}}$  between sensors.

151 **To optimize the Time-Pose Function**, we propose the following objective function:

$$\mathcal{L} = \lambda_{\text{trans}} \mathcal{L}_{\text{trans}} + \lambda_{\text{rot}} \mathcal{L}_{\text{rot}} + \lambda_{\text{speed}} \mathcal{L}_{\text{speed}}, \quad (3)$$

<sup>3</sup>This is shown by the ground truth trajectory in Fig. 4 and the predicted trajectory in the supplementary.

152 where  $\mathcal{L}_{\text{trans}}, \mathcal{L}_{\text{rot}}, \mathcal{L}_{\text{speed}}$  are translation, rotation and speed losses respectively as shown in the left  
 153 panel of Fig. 3. and  $\lambda_{\text{trans}}, \lambda_{\text{rot}}, \lambda_{\text{speed}}$  are the weighting parameters. Note that  $\lambda_{\text{trans}}$  and  $\lambda_{\text{rot}}$  are  
 154 automatically adjusted as explained in a later paragraph.

155 **Pose Representation.** There are some common choices to represent rotation for optimizing camera  
 156 poses like rotation matrices [37] or Euler-angles [25, 28]. However, they are not continuous for  
 157 representing rotation [42] due to their non-homeomorphic representation space to  $SO(3)$ . We choose  
 158 to use unit quaternion as our raw rotation representation because arbitrary 4-D vectors can be easily  
 159 mapped to legitimate rotations by normalizing them to the unit length [10].

160 **Optimization of Translation and Rotation.** We optimize the translation and the rotation vectors by  
 161 minimizing the mean square error (MSE) between the estimated and ground-truth camera poses:

$$\mathcal{L}_{\text{trans}} = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{x}_i)^2, \mathcal{L}_{\text{rot}} = \frac{1}{n} \sum_{i=1}^n (q_i - \hat{q}_i)^2. \quad (4)$$

162 Since  $x$  and  $q$  are in different units, the scaling factor  $\lambda_{\text{trans}}$  and  $\lambda_{\text{rot}}$  play an important role in balancing  
 163 the losses. To prevent translation and rotation from negatively influencing each other in training and to  
 164 tap into possible mutual facilitation, we make the weighting factors learnable by using homoscedastic  
 165 uncertainty [9] as  $\mathcal{L}_{\sigma} = \mathcal{L}_{\text{trans}} \exp(-\hat{s}_{\text{trans}}) + \hat{s}_{\text{trans}} + \mathcal{L}_{\text{rot}} \exp(-\hat{s}_{\text{rot}}) + \hat{s}_{\text{rot}}$ , where  $\hat{s}$  are learnable  
 166 parameters, thus the loss terms are balanced during training course<sup>4</sup>.

167 **Optimization of Motion Speed.** Observing that the time-pose function is essentially a function of  
 168 translational displacement and angular displacement with respect to time, we can use the average  
 169 linear velocity<sup>5</sup> to supervise the gradient of the network output, with regard to the input vectors.  
 170 Since the linear velocity variation is small and the angular velocity variation is relatively larger in the  
 171 scenes captured by the drone, only the average linear velocity is used to supervise the neural network  
 172 and the latter is not supervised in our method:

$$\mathcal{L}_{\text{speed}} = \text{MSE}(v(t_i), \hat{v}(t_i)) = \frac{1}{n} \sum_{i=1}^n (v(t_i) - \frac{\partial \hat{x}}{\partial t}(t_i))^2, v(t_i) = \frac{\partial x}{\partial t} \Big|_{t=t_i} \approx \frac{x_i - x_{i-1}}{t_i - t_{i-1}} \quad (5)$$

## 173 4.2 Bootstrapping Large-scale Neural Radiance Fields

174 In this part, we introduce our proposed scene representation (right half of Fig. 3) that is bootstrapped  
 175 in the second phase of the optimization process (Fig. 2-(ii)). Due to the limited capacity of MLPs,  
 176 we follow Mega-NeRF[29] and partition the scene map into a series of equal-sized blocks in terms  
 177 of spatial scope, and each block learns its individual scene representation with an implicit field. In  
 178 this stage, we optimize the scene representation with pure RGB data. Specifically, the radiance field  
 179 is denoted as  $\{f_{\text{NeRF}}^{(i)}\}_{i=1}^{N_x \times N_y}$ , where  $N_x, N_y$  denotes the spatial grid size. Each implicit function  
 180 represents a geographic region with  $\mathbf{x}_i^{\text{centroid}}$  as its centroid. The  $k$ th scene model can be written as:

$$f_{\text{NeRF}}^{(k)}(\gamma(\mathbf{x}_{\text{pts}}), \gamma(\mathbf{d})) \rightarrow (\hat{c}, \sigma), \quad (6)$$

181 where  $k = \arg \min_j \|\mathbf{x}_{\text{pts}} - \mathbf{x}_j^{\text{centroid}}\|_2$  and  $\gamma$  is the positional encoding function.

182 For view synthesis, we adopt volume rendering techniques to synthesize color image  $\hat{I}$  and depth  
 183 map  $\hat{D}$ . To be specific, we sample a set of points for each emitted camera ray in a coarse-to-fine  
 184 manner [15] and accumulate the radiance and the distance along the corresponding ray to calculate the  
 185 rendered color  $\hat{I}$  and depth  $\hat{D}$ . To obtain the radiance of a spatial point  $\mathbf{x}_{\text{pts}}$ , we use the nearest scene  
 186 model for prediction. A set of per-image appearance embedding [13] is also optimized simultaneously  
 187 in the training.

$$\hat{I}(\mathbf{o}, \mathbf{d}) = \int_{\text{near}}^{\text{far}} T(t) \sigma^{(k)}(\mathbf{x}(t)) \cdot c^{(k)}(\mathbf{x}(t), \mathbf{d}) dt, \hat{D}(\mathbf{o}, \mathbf{d}) = \int_{\text{near}}^{\text{far}} T(t) \sigma^{(k)}(\mathbf{x}(t)) \cdot t dt, \quad (7)$$

188 where  $\mathbf{o}$  and  $\mathbf{d}$  denote the position and orientation of the sampled ray,  $\mathbf{x}(t) = \mathbf{o} + t\mathbf{d}$  represents  
 189 the sampled point coordinates in the world space, and  $T(t) = \exp\left(-\int_{\text{near}}^t \sigma^{(k)}(\mathbf{x}(s)) ds\right)$  is the

<sup>4</sup>Manual selection of weights requires laborious tuning, but comparable performance can be achieved.

<sup>5</sup>Note that the average velocity refers to the mean value calculated from the ground-truth camera pose of the current frame and the two adjacent frames, rather than the average value in the whole sequence.

190 accumulated transmittance. We optimize the scene representation model with only the photometric  
 191 error as  $\mathcal{L}_{\text{bootstrap}} = \text{MSE}(I, \hat{I})$ . We empirically observe that this bootstrapping is critical to the  
 192 challenging third stage which jointly learns  $\theta$  and  $\phi$  using asynchronous RGB-D data.

### 193 4.3 Joint Optimization

194 While the time-pose function learns a good initialization from the RGB sequence, there are still errors  
 195 to be compensated. In this section, we describe how we perform simultaneous mapping and pose  
 196 optimization, which compensates for the initial error of the time-pose function.

197 We jointly optimize the inaccurate camera poses and the implicit maps: when fitting parameters  
 198  $\Theta_{\text{NeRF}}^{(k)}$  of the scene representation, the estimated depth camera poses  $\hat{\mathcal{T}}_D^{(j)} \in \text{SE}(3)$  (where  $t \in \mathbb{R}^3$   
 199 and  $q \in \text{SO}(3)$ ) will be simultaneously optimized on the manifold:

$$\theta, \{\hat{\mathcal{T}}_D\} = \underset{\theta, \mathcal{T} \in \text{SE}(3)}{\text{argmin}} \mathcal{L}(\{I^{(i)}\}, \{D^{(j)}\} \mid \theta, \{\hat{\mathcal{T}}_D\}), \quad (8)$$

200 where  $\mathcal{L}$  is the objective function we demonstrate in the next paragraph.

201 To train the implicit representation to obtain photo-realistic RGB rendering maps and accurate depth  
 202 map estimation, we update the mapping losses as:

$$\mathcal{L} = \lambda_{\text{color}} \sum_i \text{MSE}(I^{(i)}, \hat{I}^{(i)}) + \lambda_{\text{depth}}(\alpha) \sum_j \text{MSE}(D^{(j)}, \hat{D}^{(j)}), \quad (9)$$

203 where  $\lambda_{\text{color}}$  and  $\lambda_{\text{depth}}(\alpha)$  are weighting hyper-parameters for color and depth loss, in which the  
 204 depth loss weight starts to grow from zero gradually with the training process  $\alpha$ .

205 To compensate for the error from the time-pose function extracted poses, we jointly optimize two  
 206 implicit representation networks thanks to the end-to-end differentiable nature.

## 207 5 Asynchronous Urban Scene (AUS) Dataset

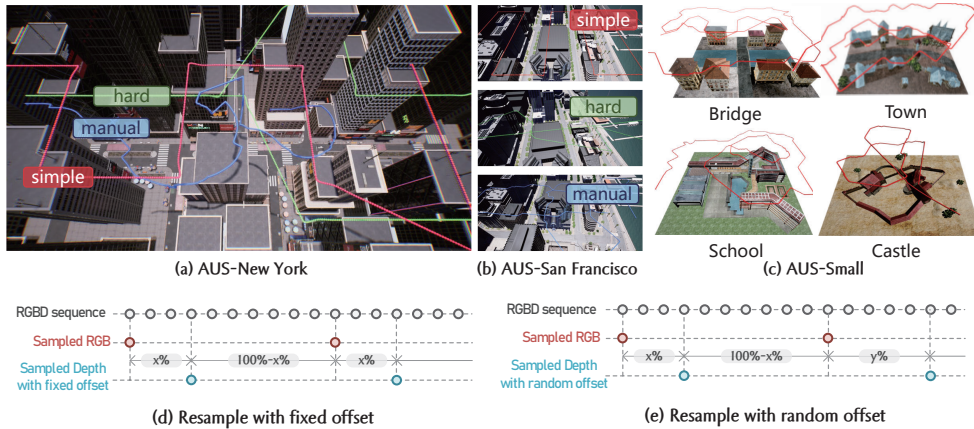


Figure 4: We propose a photo-realistically rendered dataset named Asynchronous Urban Scene (AUS) for evaluation. (a/b) are large-scale city scenes designed according to New York and San Francisco while (c) is (relatively) small-scale scenes provided by UrbanScene3D. Drone trajectories of different difficulty levels are visualized in (a-c). On these trajectories, we first capture an RGB-D sequence with an enough high framerate. Then we exploit two resampling strategies: fixed offset (d) and random offset (e).  $x$  equals 30 in (d) for every RGB-D pair.  $x$  equals 30 while  $y$  equals 50 in (e).

208 **Dataset Collection.** Our Asynchronous Urban Scene (AUS) dataset as illustrated in Fig.4 is generated  
 209 using Airsim [23], a simulator plug-in for Unreal Engine. With 3D city models loaded in Unreal  
 210 Engine, the simulator can output photorealistic and high-resolution RGB images with synchronized  
 211 depth images (resampled later) according to the a drone trajectory and a capture framerate. We choose  
 212 Airsim as it strikes a good balance between rendering quality and dynamics modeling flexibility.

Scene	Method	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	RMSE $\downarrow$	RMSE log $\downarrow$	$\delta_1(\%) \uparrow$	$\delta_2(\%) \uparrow$	$\delta_3(\%) \uparrow$
NY	NeRF-W	23.32	0.8105	0.2249	17.40	0.2630	80.17	90.11	94.72
	Mega-NeRF	23.53	<b>0.8375</b>	0.1920	23.99	0.2943	80.11	88.77	92.87
	Ours	<b>24.33</b>	0.8346	<b>0.1833</b>	<b>6.15</b>	<b>0.0816</b>	<b>94.85</b>	<b>98.23</b>	<b>99.22</b>
SF	NeRF-W	19.21	0.6610	0.3632	24.93	0.1877	81.81	91.54	96.93
	Mega-NeRF	20.53	0.7334	<b>0.2619</b>	23.56	0.1713	88.58	94.74	96.83
	Ours	<b>22.14</b>	<b>0.7930</b>	0.2620	<b>7.64</b>	<b>0.0789</b>	<b>96.34</b>	<b>98.80</b>	<b>99.70</b>
Bridge	NeRF-W	26.79	0.8053	0.2438	131.88	1.2277	53.76	61.57	58.98
	Mega-NeRF	27.98	0.8674	<b>0.1548</b>	120.41	1.3246	69.10	72.54	73.17
	Ours	<b>29.06</b>	<b>0.8751</b>	0.1952	<b>26.56</b>	<b>0.3248</b>	<b>93.24</b>	<b>96.32</b>	<b>98.26</b>
Town	NeRF-W	21.32	0.6208	0.4088	132.70	1.4640	44.89	55.68	57.90
	Mega-NeRF	24.69	0.7305	0.3103	129.50	1.4240	54.54	59.18	57.90
	Ours	<b>25.32</b>	<b>0.7675</b>	<b>0.2631</b>	<b>15.61</b>	<b>0.4632</b>	<b>91.92</b>	<b>96.89</b>	<b>98.49</b>
School	NeRF-W	19.69	0.5715	0.4453	88.83	0.9365	61.73	72.58	75.80
	Mega-NeRF	25.57	0.7739	0.3191	63.10	0.7651	77.18	85.02	86.69
	Ours	<b>26.51</b>	<b>0.7971</b>	<b>0.3175</b>	<b>21.19</b>	<b>0.2083</b>	<b>92.87</b>	<b>95.78</b>	<b>97.51</b>
Castle	NeRF-W	22.63	0.7443	0.2557	78.18	0.8651	75.72	79.26	81.11
	Mega-NeRF	28.06	<b>0.9053</b>	0.1159	54.99	0.6167	79.69	83.59	87.43
	Ours	<b>28.21</b>	0.8976	<b>0.1113</b>	<b>16.66</b>	<b>0.3565</b>	<b>93.12</b>	<b>97.23</b>	<b>98.45</b>
Mean	NeRF-W	21.80	0.7156	0.3118	55.86	0.5846	72.20	81.40	84.88
	Mega-NeRF	23.85	0.7990	0.2262	51.06	0.5527	78.66	85.09	87.43
	Ours	<b>24.85</b>	<b>0.8220</b>	<b>0.2223</b>	<b>12.14</b>	<b>0.1834</b>	<b>94.47</b>	<b>97.73</b>	<b>98.95</b>
		(+1.00)	(+0.0230)	(-0.0039)	(-38.92)	(-0.3693)	(+15.11)	(+12.64)	(+11.52)

Table 1: We quantitatively evaluate our method on the AUS dataset. Our method can synthesize more realistic images and more accurate depth maps than the baseline methods. For the **NY** and **SF**, we only report the mean performances on all sequences (Simple / Hard / Manual) due to limited space and more detailed results are in the supplementary materials.

213 **3D City Scene Models.** To generate the AUS dataset, we exploit a total of six scene models, covering  
214 two large-scale ones shown in Fig. 4-a/b and four (relatively) small-scale ones shown in Fig. 4-c. The  
215 former uses the New York and San Francisco city scenes provided by Kirill Sibiriyakov [24], in which  
216 AUS-NewYork covers a  $250 \times 150m^2$  area with many detailed buildings and AUS-SanFrancisco  
217 consists of a  $500 \times 250m^2$  area near the Golden Gate Bridge. The latter uses four model files provided  
218 in the UrbanScene3D dataset [12]. As such, at the scene level, AUS features a good coverage of both  
219 large-scale modern cities and smaller cultural heritage sites.

220 **Trajectory Generation.** Trajectory complexity matters for our problem. On one hand, in many  
221 real-world applications, photographers may manually control drones to capture a city. On the other  
222 hand, simple trajectories can be modeled by simple functions, rendering the neural time-pose function  
223 unnecessary. To build a meaningful and comprehensive benchmark, we use three types of trajectories:  
224 a trivial Zig-Zag trajectory (simple in Fig. 4), a more complex randomly generated trajectory (hard in  
225 Fig. 4), and a very complex manually controlled trajectory (manual in Fig. 4). In AUS-Small, we  
226 only provide manually controlled trajectories, since the scene sizes are relatively small and using the  
227 former two trajectory strategies leads to an unrealistically large overlap between frames.

228 **Mismatch Resampling.** We first sample synchronized RGB-D sequences in the simulator at a high  
229 frequency (50fps) then re-sample RGB and depth images with various offsets to create asynchronous  
230 RGB-D sequences. As shown in Fig. 4-d/e, we exploit two settings for the AUS dataset. In Fig. 4-d,  
231 every RGB-D pair is resampled according to a fixed offset denoted by  $x\%$ . For example, we sample  
232 the RGB image at 5fps or say every 10 frames and  $x = 30$  means every depth image is 3 frames  
233 later than the RGB counterpart. In Fig. 4-e, the offset between an RGB-D pair is randomly selected,  
234 simulating a challenging real-world asynchronous sequence. Offset ablation will be shown later.

## 235 6 Experiments

236 In this section, we show the effectiveness of our 3-step optimization pipeline by qualitatively and  
237 quantitatively evaluating our proposed methods and comparing with baseline methods.

### 238 6.1 Results

239 We evaluate our proposed method against NeRF-W [13] and city-scale Mega-NeRF [29] and present  
240 the quantitative results in Table 1 and Table 2. NeRF-W is the baseline from which we borrow  
241 the aforementioned idea of per-image appearance embedding and Mega-NeRF is a state-of-the-art  
242 (SOTA) large-scale scene modeling framework which our network is built upon.

Scene	Time-Pose Function		Joint Optimization	
	Rot. ( $^{\circ}$ )	Trans. ( $m$ )	Rot. ( $^{\circ}$ )	Trans. ( $m$ )
<b>NY Full</b>	0.66 / 0.59 / 3.70	1.84 / 1.12 / 0.46	0.13 / 0.09 / 1.47	0.34 / 0.56 / 0.20
<b>SF Full</b>	0.17 / 0.67 / 0.65	1.34 / 1.45 / 0.94	0.05 / 0.41 / 0.02	0.32 / 1.09 / 0.66
<b>Small</b>	1.51 / 0.68 / 0.70 / 1.05	0.95 / 1.35 / 0.89 / 0.38	0.49 / 0.36 / 0.68 / 0.38	0.57 / 0.85 / 0.56 / 0.12
<b>Mean</b>	<b>1.04</b>	<b>1.07</b>	<b>0.41 (-0.63)</b>	<b>0.53 (-0.54)</b>

Table 2: We show that the time-pose function learns an accurate implicit trajectory from the RGB sequence that can estimate accurate poses for depth frames. By further tuning the time-pose function jointly with the scene representation network, the accuracy of the predicted depth sensor poses can be improved. The results of Simple / Hard / Manual on **NY** and **SF** are shown in the first two lines. The results of the 4 small scenes of (Bridge / Town / School / Castle) are shown at the bottom.

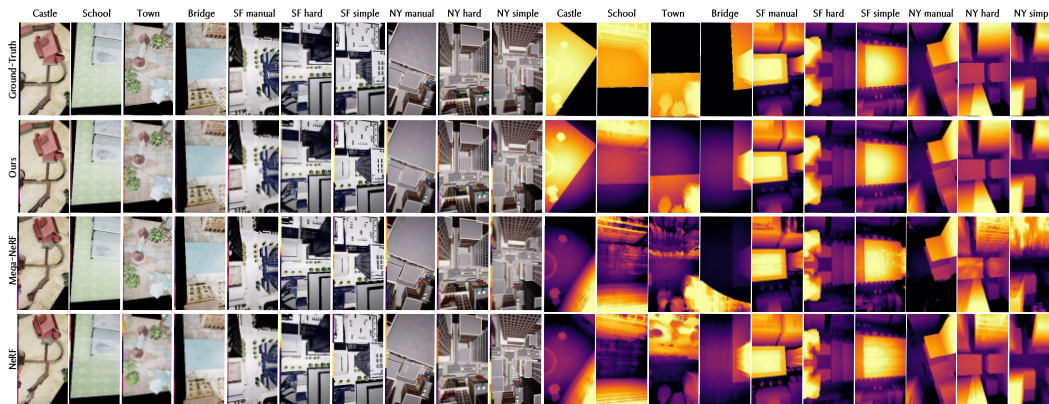


Figure 5: **Qualitative Results.** Our method can render photo-realistic novel views and the best depth estimation results. Please zoom in to see details using an electronic version.

243 **RGB-D View Synthesis** The standard metrics for novel view synthesis and depth estimation are  
 244 used for evaluation. For RGB view synthesis, metrics including PSNR, SSIM, and the VGG  
 245 implementation of LPIPS [40] are used. For depth estimation, RMSE, RMSE log,  $\delta_{1,2,3}$  are used. In  
 246 Table 1, we quantitatively show on the RGB view synthesis task that, together with depth supervision,  
 247 our method can generate more photo-realistic images than two SOTA baselines, and show on the  
 248 depth estimation tasks that our method significantly improves the learned geometry of the scene  
 249 representation network. We also present the RGB-D view synthesis results qualitatively in Fig.5 in  
 250 which our method synthesizes photo-realistic images and accurate depth maps, while baseline methods  
 251 fail at predicting reliable depth maps (e.g. in the **School** scene, they mispredict the void space as a  
 252 dense surface; in **NY hard**, the depth values around glasses are obviously inaccurate).

253 **Depth Pose Estimation** We evaluate the performance of our time-pose function in Table 2, or say  
 254 specifically the accuracy of our method to localize depth sensor poses. As shown quantitatively, our  
 255 method can achieve an average pose error of  $1.04m$  and  $1.07^{\circ}$  in the first stage, where only time-pose  
 256 pairs from the RGB sequence are used to optimize the network. After joint optimization in the third  
 257 stage, our method cuts half the errors to  $0.53m$  and  $0.41^{\circ}$ .

258 **Real-world Evaluation.** In the real-world experiments, we use the DJI M300 UAV (equipped with a  
 259 high-definition RGB camera and LiDAR to collect real data, where the RGB camera collects images  
 260 at the frame rate of 30fps and the LiDAR collects depth information at 240Hz. The poses of the RGB  
 261 images are provided by COLMAP [22]. The fixed transformations between sensors are provided by  
 262 the producer or can be calibrated manually. A qualitative comparison is provided in Fig. 1 and more  
 263 results are in the supplementary.

## 264 6.2 Ablation Studies

265 **Time-Pose Function Network Structure.** We compared our proposed time-pose function implemen-  
 266 tation with several commonly used implicit representation architectures on the localization accuracy



Method	Rotation ( $^{\circ}$ )		Translation ( $m$ )	
	Mean	Median	Mean	Median
MLP	26.62	17.53	8.23	7.50
Feature Grid	15.86	14.56	8.53	7.48
Ours (L=1)	24.24	12.99	9.20	8.01
Ours w/o speed	11.96	11.21	19.95	12.28
Ours	<b>11.36</b>	<b>11.17</b>	<b>6.29</b>	<b>4.03</b>

Table 3: Ablation on different network structures and the use of speed optimization.

Offset	Rotation ( $^{\circ}$ )		Translation ( $m$ )	
	Mean	Median	Mean	Median
10%	<b>0.66</b>	<b>0.26</b>	<b>3.42</b>	<b>1.88</b>
20%	0.94	0.55	4.96	3.90
30%	1.24	0.79	6.41	5.78
40%	1.41	0.75	7.26	6.61
50%	1.50	0.84	7.53	6.51
Random	1.12	0.52	5.17	4.05

Table 4: Results on the ablation of different sampling offset strategies.

(Tab. 3). In order to highlight the gap between different methods, we downsample the dataset to increase the difficulty. (a) **Pure MLP architecture** : processing the positional-encoded [15] input timestamps with an MLP; (b) **1-D Feature Grid**: storing a feature vector for each second in the timestamp span and performing linear feature interpolation in the query’s neighborhood. (c) **Ours**: our proposed 1-D multi-resolution hash grid with different resolution layers. The results (Table 3) show that our proposed multi-resolution architecture outperforms other network architectures in accuracy. The network structures are further detailed in the supplementary materials.

**Speed Loss.** We compared our method’s localization accuracy with and without the optimization of motion speed (see the comparison of ‘Ours’ and ‘Ours w/o speed’ in Tab. 3). The results show that minimizing the gradient error (i.e., speed loss) help a lot in improving the accuracy of the translation (from  $20m$  to  $6.3m$ ).

**Different Sampling Offsets.** To show the robustness of our proposed time-pose function, we compare the localization accuracy of implicit trajectory representations under different sampling offsets, whose definition is described in Fig. 4-d/e. The quantitative results (Table. 4) indicate that the network’s output exhibits a controllable margin of error as the data offset increases.

Scene	Ours		Mega-NeRF		Mega-NeRF-Depth	
	PSNR $\uparrow$	RMSE $\downarrow$	PSNR $\uparrow$	RMSE $\downarrow$	PSNR $\uparrow$	RMSE $\downarrow$
NY Mean	<b>24.24</b>	<b>5.93</b>	24.03	42.15	19.70	15.94
SF Mean	<b>22.70</b>	<b>7.26</b>	20.00	32.17	19.07	11.39
Bridge	<b>29.06</b>	<b>26.55</b>	27.98	120.41	22.35	96.16
Town	<b>25.32</b>	<b>15.61</b>	24.69	129.50	20.14	81.99
School	<b>26.51</b>	<b>21.19</b>	25.57	63.10	21.91	42.74
Castle	<b>28.22</b>	<b>16.66</b>	28.06	54.99	23.23	38.90
Mean	<b>26.01</b>	<b>15.53</b>	25.01	73.72	21.07	47.85

Table 5: Ablation on the joint optimization stage. We show that jointly optimizing the time-pose function and the scene representation significantly helps reduce geometric error.

**Joint Optimization for Pose Error Compensation.** To demonstrate the importance of rectifying erroneous poses of depth images in asynchronous RGB-D sequences using the time-pose function, we train a Mega-NeRF[29] with depth supervision but disabled the joint optimization stage. From the evaluation results (Table 5), we observe its substantial impact on the rendering quality (PSNR for RGB and RMSE for depth). Due to limited space, qualitative results are in the supplementary.

## 7 Conclusion

In this paper, we present a method to learn depth-supervised neural radiance fields from asynchronous RGB-D sequences. We leverage an important prior that the sensors cover the same spatial-temporal footprints and propose to utilize this prior with an implicit time-pose function. With a 3-staged optimization pipeline, our method calibrates the RGB-D poses and trains a large-scale implicit scene representation. Our experiments on a newly proposed large-scale dataset show that our method can effectively register depth camera poses and learns the 3D scene representation for photo-realistic novel view synthesis and accurate depth estimations. **Broader impact and limitations:** Large-scale scene modelling can be used for potential military use, which the method is not intended for.

## References

- 296
- 297 [1] Sameer Ansari, Neal Wadhwa, Rahul Garg, and Jiawen Chen. Wireless Software Synchronization of  
298 Multiple Distributed Cameras. In *2019 IEEE International Conference on Computational Photography*  
299 (*ICCP*), pages 1–9, May 2019. ISSN: 2472-7636.
- 300 [2] Haowen Deng, Mai Bui, Nassir Navab, Leonidas Guibas, Slobodan Ilic, and Tolga Birdal. Deep bingham  
301 networks: Dealing with uncertainty and ambiguity in pose estimation. *International Journal of Computer*  
302 *Vision*, pages 1–28, 2022.
- 303 [3] Kangle Deng, Andrew Liu, Jun-Yan Zhu, and Deva Ramanan. Depth-supervised nerf: Fewer views and  
304 faster training for free. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*  
305 *Recognition*, pages 12882–12891, 2022.
- 306 [4] A. Elhayek, C. Stoll, N. Hasler, K. I. Kim, H.-P. Seidel, and C. Theobalt. Spatio-temporal motion tracking  
307 with unsynchronized cameras. In *2012 IEEE Conference on Computer Vision and Pattern Recognition*,  
308 pages 1870–1877, June 2012. ISSN: 1063-6919.
- 309 [5] Antoine Guedon, Pascal Monasse, and Vincent Lepetit. Scone: Surface coverage optimization in unknown  
310 environments by volumetric integration. In *Advances in Neural Information Processing Systems*, 2022.
- 311 [6] Antoine Guédon, Tom Monnier, Pascal Monasse, and Vincent Lepetit. Macarons: Mapping and coverage  
312 anticipation with rgb online self-supervision. *arXiv preprint arXiv:2303.03315*, 2023.
- 313 [7] Meng-Ru Hsieh, Yen-Liang Lin, and Winston H Hsu. Drone-based object counting by spatially regularized  
314 regional proposal network. In *Proceedings of the IEEE international conference on computer vision*, pages  
315 4145–4153, 2017.
- 316 [8] Yoonwoo Jeong, Seokjun Ahn, Christopher Choy, Anima Anandkumar, Minsu Cho, and Jaesik Park.  
317 Self-calibrating neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on*  
318 *Computer Vision*, pages 5846–5854, 2021.
- 319 [9] Alex Kendall and Roberto Cipolla. Geometric loss functions for camera pose regression with deep learning.  
320 In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5974–5983,  
321 2017.
- 322 [10] Alex Kendall, Matthew Grimes, and Roberto Cipolla. Posenet: A convolutional network for real-time  
323 6-dof camera relocalization. In *Proceedings of the IEEE international conference on computer vision*,  
324 pages 2938–2946, 2015.
- 325 [11] Chen-Hsuan Lin, Wei-Chiu Ma, Antonio Torralba, and Simon Lucey. Barf: Bundle-adjusting neural  
326 radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages  
327 5741–5751, 2021.
- 328 [12] Liqiang Lin, Yilin Liu, Yue Hu, Xingguang Yan, Ke Xie, and Hui Huang. Capturing, reconstructing, and  
329 simulating: the urbanscene3d dataset. In *ECCV*, 2022.
- 330 [13] Ricardo Martin-Brualla, Noha Radwan, Mehdi S. M. Sajjadi, Jonathan T. Barron, Alexey Dosovitskiy, and  
331 Daniel Duckworth. NeRF in the Wild: Neural Radiance Fields for Unconstrained Photo Collections. In  
332 *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7206–7215, June  
333 2021. ISSN: 2575-7075.
- 334 [14] Zhenxing Mi and Dan Xu. Switch-nerf: Learning scene decomposition with mixture of experts for  
335 large-scale neural radiance fields. In *International Conference on Learning Representations (ICLR)*, 2023.
- 336 [15] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and Ren  
337 Ng. NeRF: Representing Scenes as Neural Radiance Fields for View Synthesis. In Andrea Vedaldi, Horst  
338 Bischof, Thomas Brox, and Jan-Michael Frahm, editors, *Computer Vision – ECCV 2020*, Lecture Notes in  
339 Computer Science, pages 405–421, Cham, 2020. Springer International Publishing.
- 340 [16] Matthias Mueller, Neil Smith, and Bernard Ghanem. A benchmark and simulator for uav tracking. In  
341 *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14,*  
342 *2016, Proceedings, Part I 14*, pages 445–461. Springer, 2016.
- 343 [17] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics prim-  
344 itives with a multiresolution hash encoding. *ACM Transactions on Graphics*, 41(4):1–15, July 2022.  
345 titleTranslation:.

- 346 [18] Stefano Pellegrini, Andreas Ess, Konrad Schindler, and Luc Van Gool. You'll never walk alone: Modeling  
347 social behavior for multi-target tracking. In *2009 IEEE 12th international conference on computer vision*,  
348 pages 261–268. IEEE, 2009.
- 349 [19] Amir M Rahimi, Raphael Ruschel, and BS Manjunath. Uav sensor fusion with latent-dynamic conditional  
350 random fields in coronal plane estimation. In *Proceedings of the IEEE Conference on Computer Vision  
351 and Pattern Recognition*, pages 4527–4534, 2016.
- 352 [20] Konstantinos Rematas, Andrew Liu, Pratul P Srinivasan, Jonathan T Barron, Andrea Tagliasacchi, Thomas  
353 Funkhouser, and Vittorio Ferrari. Urban radiance fields. In *Proceedings of the IEEE/CVF Conference on  
354 Computer Vision and Pattern Recognition*, pages 12932–12942, 2022.
- 355 [21] Barbara Roessle, Jonathan T Barron, Ben Mildenhall, Pratul P Srinivasan, and Matthias Nießner. Dense  
356 depth priors for neural radiance fields from sparse input views. In *Proceedings of the IEEE/CVF Conference  
357 on Computer Vision and Pattern Recognition*, pages 12892–12901, 2022.
- 358 [22] Johannes L. Schönberger and Jan-Michael Frahm. Structure-from-Motion Revisited. In *2016 IEEE  
359 Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4104–4113, June 2016. ISSN:  
360 1063-6919.
- 361 [23] Shital Shah, Debadeepta Dey, Chris Lovett, and Ashish Kapoor. Airsim: High-fidelity visual and physical  
362 simulation for autonomous vehicles. In *Field and Service Robotics: Results of the 11th International  
363 Conference*, pages 621–635. Springer, 2018.
- 364 [24] Kirill Sibiriakov. Artstation page <https://www.artstation.com/vegaart>, 2022.
- 365 [25] Hao Su, Charles R Qi, Yangyan Li, and Leonidas J Guibas. Render for cnn: Viewpoint estimation in images  
366 using cnns trained with rendered 3d model views. In *Proceedings of the IEEE international conference on  
367 computer vision*, pages 2686–2694, 2015.
- 368 [26] Matthew Tancik, Vincent Casser, Xinchun Yan, Sabeek Pradhan, Ben Mildenhall, Pratul P Srinivasan,  
369 Jonathan T Barron, and Henrik Kretzschmar. Block-nerf: Scalable large scene neural view synthesis. In  
370 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8248–8258,  
371 2022.
- 372 [27] Matthew Tancik, Pratul P. Srinivasan, Ben Mildenhall, Sara Fridovich-Keil, Nithin Raghavan, Utkarsh  
373 Singhal, Ravi Ramamoorthi, Jonathan T. Barron, and Ren Ng. Fourier Features Let Networks Learn High  
374 Frequency Functions in Low Dimensional Domains, June 2020. arXiv:2006.10739 [cs].
- 375 [28] Shubham Tulsiani and Jitendra Malik. Viewpoints and keypoints. In *Proceedings of the IEEE Conference  
376 on Computer Vision and Pattern Recognition*, pages 1510–1519, 2015.
- 377 [29] Haithem Turki, Deva Ramanan, and Mahadev Satyanarayanan. Mega-NeRF: Scalable Construction of  
378 Large-Scale NeRFs for Virtual Fly- Throughs. In *2022 IEEE/CVF Conference on Computer Vision and  
379 Pattern Recognition (CVPR)*, pages 12912–12921, June 2022. ISSN: 2575-7075.
- 380 [30] Zirui Wang, Shangzhe Wu, Weidi Xie, Min Chen, and Victor Adrian Prisacariu. Nerf-: Neural radiance  
381 fields without known camera parameters. *arXiv preprint arXiv:2102.07064*, 2021.
- 382 [31] Longyin Wen, Dawei Du, Pengfei Zhu, Qinghua Hu, Qilong Wang, Liefeng Bo, and Siwei Lyu. Detection,  
383 tracking, and counting meets drones in crowds: A benchmark. In *Proceedings of the IEEE/CVF Conference  
384 on Computer Vision and Pattern Recognition*, pages 7812–7821, 2021.
- 385 [32] Yuanbo Xiangli, Linning Xu, Xingang Pan, Nanxuan Zhao, Anyi Rao, Christian Theobalt, Bo Dai, and  
386 Dahua Lin. Bungeenerf: Progressive neural radiance field for extreme multi-scale scene rendering. In *The  
387 European Conference on Computer Vision (ECCV)*, 2022.
- 388 [33] Dan Xie, Sinisa Todorovic, and Song-Chun Zhu. Inferring "dark matter" and "dark energy" from videos.  
389 In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2224–2231, 2013.
- 390 [34] Ziyang Xie, Junge Zhang, Wenye Li, Feihu Zhang, and Li Zhang. S-nerf: Neural radiance fields for street  
391 views. In *The Eleventh International Conference on Learning Representations*, 2023.
- 392 [35] Linning Xu, Yuanbo Xiangli, Sida Peng, Xingang Pan, Nanxuan Zhao, Christian Theobalt, Bo Dai, and  
393 Dahua Lin. Grid-guided neural radiance fields for large urban scenes, 2023.
- 394 [36] Anqi Joyce Yang, Can Cui, Ioan Andrei Bârsan, Raquel Urtasun, and Shenlong Wang. Asynchronous  
395 Multi-View SLAM. In *2021 IEEE International Conference on Robotics and Automation (ICRA)*, pages  
396 5669–5676, May 2021. ISSN: 2577-087X.

- 397 [37] Lin Yen-Chen, Pete Florence, Jonathan T. Barron, Alberto Rodriguez, Phillip Isola, and Tsung-Yi Lin.  
398 iNeRF: Inverting Neural Radiance Fields for Pose Estimation. In *2021 IEEE/RSJ International Conference*  
399 *on Intelligent Robots and Systems (IROS)*, pages 1323–1330, September 2021. ISSN: 2153-0866.
- 400 [38] Zehao Yu, Songyou Peng, Michael Niemeyer, Torsten Sattler, and Andreas Geiger. Monosdf: Exploring  
401 monocular geometric cues for neural implicit surface reconstruction. *Advances in Neural Information*  
402 *Processing Systems (NeurIPS)*, 2022.
- 403 [39] Kuo-Hao Zeng, Roozbeh Mottaghi, Luca Weihs, and Ali Farhadi. Visual reaction: Learning to play catch  
404 with your drone. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*,  
405 pages 11573–11582, 2020.
- 406 [40] Richard Zhang, Phillip Isola, Alexei A. Efros, Eli Shechtman, and Oliver Wang. The Unreasonable  
407 Effectiveness of Deep Features as a Perceptual Metric. In *2018 IEEE/CVF Conference on Computer Vision*  
408 *and Pattern Recognition*, pages 586–595, June 2018. ISSN: 2575-7075.
- 409 [41] Runze Zhang, Siyu Zhu, Tian Fang, and Long Quan. Distributed very large scale bundle adjustment by  
410 global camera consensus. In *Proceedings of the IEEE International Conference on Computer Vision*, pages  
411 29–38, 2017.
- 412 [42] Yi Zhou, Connelly Barnes, Jingwan Lu, Jimei Yang, and Hao Li. On the continuity of rotation representa-  
413 tions in neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*  
414 *Recognition*, pages 5745–5753, 2019.