
UE4-NeRF: Neural Radiant Field for Real-Time Rendering of Large-Scale Scene

Anonymous Author(s)

Abstract

1 Neural Radiance Field (NeRF) is an implicit 3D reconstruction method that has
2 shown immense potential and has gained significant attention for its ability to
3 reconstruct 3D scenes solely from a set of photographs. However, its real-time
4 rendering capability, especially for interactive real-time rendering of large-scale
5 scenes, has significant limitations. To address this challenge, we propose a novel
6 neural rendering system called UE4-NeRF that is designed for real-time rendering
7 of large-scale scenes. Our proposed approach partitions large scenes into sub-
8 NeRFs, and uses polygonal meshes to represent them. In order to represent the
9 partitioned independent scene, we initialize polygonal meshes by constructing
10 multiple regular octahedra within the scene and the vertices of the polygonal faces
11 are continuously optimized during the training process. Drawing inspiration from
12 the Level of Detail (LOD) techniques, we train meshes with varying levels of detail
13 for different observation levels. Our approach combines with the rasterization
14 pipeline in Unreal Engine 4 (UE4), achieving real-time rendering of large-scale
15 scenes at 4K resolution with a frame rate of up to 43 FPS. Our experimental results
16 demonstrate that our method attains rendering quality on par with state-of-the-art
17 approaches, while additionally offering the advantage of real-time performance.
18 Furthermore, rendering within UE4 facilitates scene editing in subsequent stages.
19 Project page: <https://jamchaos.github.io/UE4-NeRF/>.

20

1 Introduction

21 Real-time interactive rendering capability for 3D large-scale scenes is a crucial tool in creating digital
22 worlds for applications such as virtual reality (VR), computer games, and the Metaverse. Neural
23 Radiance Fields (NeRF) [23] is a pioneering method for 3D reconstruction and synthesis of novel
24 views, which is different from traditional 3D reconstruction techniques that represent scenes using
25 explicit expressions such as point clouds, grids, and voxels [5, 4, 14, 1]. NeRF samples each ray and
26 retrieves the 3D location of each sampling point and the 2D viewing direction of the ray. Subsequently,
27 these 5D vector values are fed into a neural network to determine the color and volume density of
28 each sampling point. NeRF constructs a field parameterized by an MLP neural network [16, 34] to
29 continuously optimize parameters and reconstruct the scene. NeRF has demonstrated remarkable
30 performance for 3D reconstruction when provided with a set of posed camera images [29, 2, 23], but
31 early NeRF works [2, 23] had only focused on small-scale scenes and object-centric reconstruction
32 due to the model computational complexity. Although Block-NeRF [37] and BungeeNeRF [43]
33 are capable of reconstructing large-scale environments, they are still unable to achieve real-time
34 rendering. Traditional NeRF obtains the color of each pixel through volume rendering, but this
35 rendering process is far too slow for interactive visualization and compromises high fidelity. Recently,
36 Mobile-NeRF [6] combined with classic polygon rasterization pipeline has shown to achieve real-time
37 rendering on various devices. However, training just a small-scale scene incurs significant overheads,
38 such as training time and GPU resources. It requires 8 v100 GPUs to train the model successfully
39 over several days. Additionally, due to opacity binarization, Mobile-NeRF cannot handle outdoor
40 scenes with a variety of materials, including transparent or semi-transparent objects such as water,

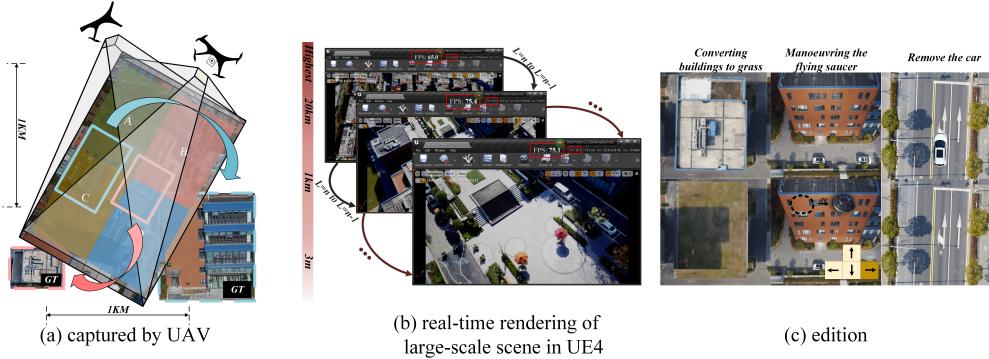


Figure 1: (a) Large-scale scenes characterized by complex textures and significant height differences are captured with monocular or multi-camera drone footage. We divide scene into multiple individual sub-NeRFs.(b) Real-time interactive rendering in UE4 with distinct levels of detail, along with frame rate disparities among them. (c) Defferent editing capabilities on the scenes.

41 glass and plastic shed. Performing subsequent editing on the scenes generated by NeRF is another
 42 meaningful but challenging task, and previous works [13, 44, 19] required training complex neural
 43 networks to achieve this capability. These methods are time-consuming and are limited in their
 44 editing capabilities.

45 This paper introduces UE4-NeRF to address the aforementioned issues regarding large-scale scene
 46 rendering and editing. UE4-NeRF combines NeRF with the Unreal Engine 4 (UE4) [17, 30, 17]
 47 to achieve real-time interactive rendering of large-scale scenes with the scene editing ability. As
 48 shown in Figure 1(a) , to achieve this, large-scale scene such as industrial parks with large height
 49 differences, roads and farmlands are partitioned into separately trained blocks. The size of individual
 50 scenes captured by drones is in the range of several square kilometers. The multi-resolution hash
 51 encoding proposed in Instant-NGP [25] is employed for training acceleration. In order to efficiently
 52 utilize the parallel computing capabilities of GPUs, the computationally intensive portions are
 53 implemented using CUDA. Therefore, the training speed of UE4-NeRF is $1000 \times$ faster than Mobile-
 54 NeRF. For each block, we represent the scene using textured polygonal meshes obtained from
 55 regular tetrahedrons. Similar to Mobile-NeRF, the texture atlas stores feature vectors and opacity
 56 information. We utilize an MLP as the High-Level Shading Language(HLSL) fragment shader in
 57 UE4, which takes in the feature vectors and outputs colors. Furthermore, to ensure real-time rendering
 58 performance for different observation heights viewpoints, we adopted a method akin to the Level of
 59 Detail (LOD) [8, 21] approach commonly used in computer graphics, where the complexity of 3D
 60 model representation is quantified in terms of metrics such as geometric detail and texture resolution.
 61 Figure 1(b) demonstrates that UE4-NeRF is capable of achieving real-time rendering at different
 62 relative observation heights by training different levels of polygon meshes, where coarser meshes are
 63 used for relatively higher observation heights. One of the other advantages of combining NeRF with
 64 UE4 is the powerful development capabilities of the rendering engine, allowing further processing of
 65 the final scene rendered in UE4. Users can freely combine different sub-NeRFs and add objects to
 66 the scene as needed. Moreover, object manipulation, as illustrated in Figure 1(c) , is also supported.

67 In summary, our paper presents the following contributions:(1) We propose a novel UE4-NeRF
 68 technique which has 1000 times faster training times than MobileNeRF, while maintaining higher
 69 quality and smaller storage overhead;(2) The demonstration of the capability of UE4-NeRF to
 70 reconstruct large-scale scenes in real-time and showcase the practical viability of NeRF;(3) The
 71 integration of NeRF and UE4 facilitates editing and manipulation of 3D scenes.

72 2 Related work

73 **Large Scale 3D Reconstruction.** The original NeRF algorithm was limited to object-centric
 74 scenes or scenes with boundaries. NeRF++ [46] expanded upon NeRF by introducing an inverted
 75 sphere parameterization to handle unbounded scenarios. Mip-NeRF 360 [3] addressed the challenges
 76 associated with unbounded scenes through the use of non-linear scene parameterization, online

77 distillation and a novel distortion-based regularizer. However, it was observed that simply expanding
78 the scene without adjusting the model capacity, i.e., using a fixed-size MLP, can lead to unstable
79 training and loss of high-fidelity. Meanwhile, naively increasing the model capacity can significantly
80 increase the training cost. BungeeNeRF [43] proposed a progressive approach to build and train
81 large-scale scene models, which promoted layer specialization and unleashed the power of positional
82 encoding over the full frequency range. Block-NeRF [37] and Mega-NeRF [40] spatially partitioned
83 a large scene into several individually trained sub-NeRFs. They also incorporated appearance
84 embeddings like NeRF-W [22] to handle appearance variations in the scene. Similar to Block-NeRF,
85 our approach partitions the large-scale scene into sub-NeRFs, and trains a coarse scene model to
86 assist in accurately delineating each sub-NeRF. *However, our approach is not only applicable to*
87 *the grid structure like Block-NeRF, but is also generally suitable for various types of scenes, e.g.*
88 *mountainous terrain. Additionally, in contrast to the aforementioned NeRFs designed for large-scale*
89 *scenes, our method accomplishes real-time rendering.*

90 **Rendering of Neural Fields.** NeRF has tremendous potential for rendering and has gained attention
91 for 3D reconstruction from a set of posed camera images [29, 2, 23]. However, the computational
92 overhead of NeRF limits its practical applications. Some prior works [25, 11] have addressed training
93 costs, while others have focused on rendering [9, 36, 27, 31]. SNeRG [20] and PlenOctrees [45]
94 improved performance by baking the components of NeRF. However, these methods do not take
95 sufficient advantage of the texture-based rasterization pipeline and require significant GPU memory
96 to store volumetric textures. MobileNeRF [6] is a variation of NeRF that can run in real-time on a
97 variety of common mobile devices. To render an image, Mobile-NeRF utilized the classic polygon
98 rasterization pipeline with Z-buffering to produce a feature vector for each pixel and then passed it to
99 a lightweight MLP running in a GLSL fragment shader to produce the output color. *In this work, we*
100 *accelerate the rendering of NeRF by converting the learned knowledge into a mesh representation for*
101 *efficient mesh rasterization pipelines. Moreover, we achieve a significant milestone by successfully*
102 *implementing NeRF for real-time interactive rendering of large-scale scenes, which was previously*
103 *unexplored.*

104 **Manipulation and Composition of NeRF.** Practical applications of 3D representations require
105 effective manipulation and composition capabilities. While, explicit 3D representations, such as
106 meshes and voxels, inherently support editing and composition, it becomes challenging for neural
107 network-based implicit representations such as vanilla NeRF [23]. NSVF [20] uses a sparse voxel
108 octree for modeling, employing voxel embeddings for each scene while sharing the same MLP for
109 predicting density and color. However, the flexibility of this approach is somewhat constrained.
110 GIRAFFE [28], combined with GAN [49, 18], enabled the addition of multiple objects in a scene,
111 expanding from single-object to multi-object generation, even in the absence of such objects in the
112 training data. Plenoxels [11] utilizes explicit sparse voxel representation for direct object composition
113 of different objects, but suffers from large storage requirements of the dense index matrix. None of
114 these methods support object manipulation in a 3D scene. *In contrast, our approach leverages UE4*
115 *to achieve large-scale scene rendering, enabling users to freely add and manipulate objects within*
116 *the UE4-NeRF environment. For instance, users can use a broomstick to fly or drive a car within the*
117 *rendered environment powered by UE4-NeRF.*

118 3 Proposed Method

119 Our novel method involves representing large scenes as polygon meshes, where the positions of mesh
120 vertices are iteratively optimized during the training phase. The encoding network generates texture
121 maps that store opacity and feature vectors, essential for rendering. To accomplish real-time rendering
122 of large scenes, we combine the pre-rendered polygon meshes with the rasterization pipeline available
123 in Unreal Engine 4 (UE4). In the subsequent subsections, we present a detailed and comprehensive
124 description of entire procedure.

125 3.1 Training

126 3.1.1 Block Partition

127 As depicted in Figure 2, a large scene is partitioned into multiple blocks. To ensure smooth transitions
128 at the boundaries of the model, each training area, indicated by dashed lines, is slightly larger than

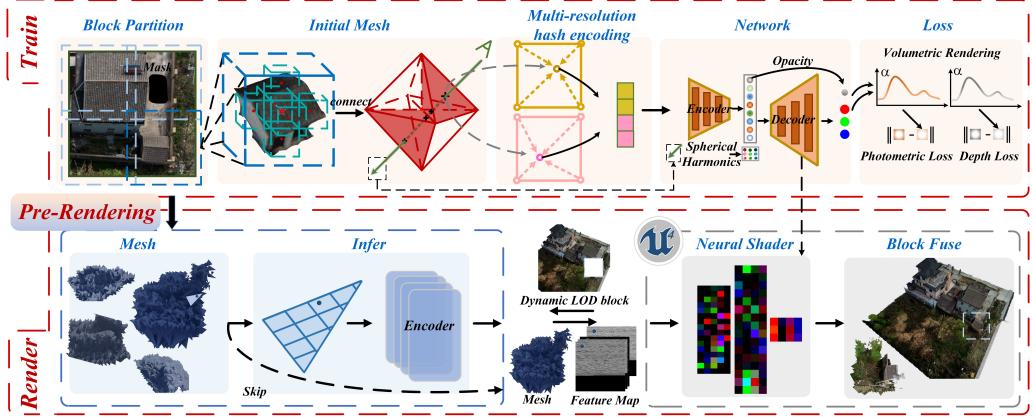


Figure 2: **Overview of UE4-NeRF.** UE4-NeRF consists of three modules: (1) In the training module, we partition the sub-NeRFs and initialize the grids in each partitioned scene. During the training process, we continuously optimize the parameters of the Encoder-Decoding network and the vertex positions of the meshes; (2) Through the pre-rendering module, we extract the polygonal meshes at different levels of detail that will be used for the final rendering; (3) The rendering module consists of an inference submodule and an UE4 submodule, which communicate with each other to achieve the final real-time rendering.

129 its corresponding scene area, with overlapping borders with neighboring blocks. We train a coarse
 130 model for sub-region segmentation, focusing on capturing the overall terrain and building outlines.
 131 We then filter the images associated with the sub-regions and use a mask to exclude pixels outside the
 132 scene area.

133 3.1.2 Optimization

134 We begin by initializing a $128 \times 128 \times 128$ grid, and selecting the center point of each grid along with
 135 its six neighboring grids (front, back, left, right, top, and bottom) to create polygonal meshes. **To**
 136 **address the slow and unstable convergence issue encountered by Mobile-NeRF when dealing with**
 137 **tilted surfaces, we utilize a regular octahedron with 20 faces.** This is illustrated in the "Initial Mesh"
 138 section of Figure 2. The ray emitted from the camera origin to the pixel is defined as $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$
 139 and $\mathbf{d} = (\theta, \phi)$ is the viewing direction of the ray. For each ray emitted from the camera origin to the
 140 pixel, we calculate its intersection with the polygonal meshes. The point of intersection serves as the
 141 sampling point for that ray. For each block, we define an Encoder-Decoding network:

$$\mathcal{E}(\text{hash}(\mathbf{p}_i); \theta_{\mathcal{E}}) \rightarrow \mathbf{M}_i, \alpha_i \quad (1)$$

$$\mathcal{D}(\mathbf{M}_i, \mathcal{SH}(\mathbf{d}_i); \theta_{\mathcal{D}}) \rightarrow \mathbf{c}_i \quad (2)$$

142 The encoder network, denoted as MLP \mathcal{E} , takes positional information encoded by multi-resolution
 143 hash functions as input. It generates an 8D feature vector \mathbf{M}_i , which incorporates opacity information.
 144 On the other hand, the decoder network \mathcal{D} utilizes \mathbf{M}_i and the ray direction as inputs and predicts the
 145 color of the sampling point. The direction of the ray is encoded using Spherical Harmonics.

146 We use the traditional NeRF volume rendering method (Eq. 3) to obtain the predicted color $\hat{\mathbf{C}}(\mathbf{r})$.
 147 Due to the complexity of calculating the distance between sampling points, we simplify the process
 148 by fixing the distance value. The opacity α is then directly predicted instead of the volume density.

$$\hat{\mathbf{C}}(\mathbf{r}) = \sum_{i=1}^N T_i \alpha_i \mathbf{c}_i, \quad T_i = \prod_{j=0}^{i-1} (1 - \alpha_j) \quad (3)$$

149 Unlike the traditional photometric loss, our proposed photometric loss consists of two components.
 150 The first component referred to as *level-1* photometric loss is calculated using Eq. 4. It quantifies the
 151 mean squared error between the predicted colors and the ground truth colors of the corresponding
 152 pixel.

$$\mathcal{L}_{rgb}^{part1}(\theta, V_p) = \sum_{r \in \mathcal{R}} \left\| \hat{\mathbf{C}}(\mathbf{r}) - \mathbf{C}(\mathbf{r}) \right\|_2^2 \quad (4)$$

153 where V_p is the position of the vertices of the polygonal mesh under the corresponding level
 154 and \mathcal{R} is the set of rays in each batch. $\mathbf{C}(\mathbf{r})$
 155 represents the ground truth RGB colors for ray
 156 r . In addition, to improve the proximity of the
 157 triangle meshes to the object surface, we design
 158 a second component of the photometric loss, de-
 159 noted as $\mathcal{L}_{rgb}^{part2}$. As depicted in Figure 3, in
 160 contrast to the previous sampling rule, we only
 161 select sampling points with opacity greater than
 162 a threshold f . Furthermore, the volumetric ren-
 163 dering process is halted when the accumulated
 164 opacity exceeds 0.8, indicating that the remain-
 165 ing opacity T' is less than 0.2. The threshold
 166 value f is computed using the formula provided
 167 in Eq. 5. Initially, during the first 10,000 epochs,
 168 f is fixed at 0.3. Subsequently, f increases grad-
 169 ually with the progression of epochs.

$$f = \min\{0.3, \left\lfloor \frac{\text{epoch} - 10000}{10000} \right\rfloor^2 \times 0.012\} \quad (5)$$

170 The photometric loss for the second part of the *level-1* can be obtained using the following formula:

$$\mathcal{L}_{rgb}^{part2}(\theta, V_p) = \sum_{r \in \mathcal{R}} \left\| \hat{\mathbf{C}}'(\mathbf{r}) - \mathbf{C}(\mathbf{r}) \right\|_2^2, \quad \hat{\mathbf{C}}'(\mathbf{r}) = \frac{\sum_{\phi=1}^N T_\phi \alpha_\phi \mathbf{c}_\phi}{1 - T'_\phi} \quad (6)$$

171 Where $\hat{\mathbf{C}}'(\mathbf{r})$ is obtained through volume rendering under the aforementioned sampling scheme. To
 172 cover the entire color range, we divide the expression by $1 - T'_\phi$ as shown in Eq. 6, since our volume
 173 rendering stops when the remaining opacity is less than 0.2. We obtain the final photometric loss of
 174 the *level-1* as follows:

$$\beta \mathcal{L}_{rgb}^{part1}(\theta, V_p) + (1 - \beta) \mathcal{L}_{rgb}^{part2}(\theta, V_p), \quad \beta = 0.8^{\lfloor \frac{\text{epoch}}{10000} \rfloor} \quad (7)$$

175 During the initial 10,000 epochs, the second component of the photometric loss does not have
 176 any impact. This is done to align with the traditional volume rendering approach and adhere to
 177 the modeling principles of NeRF. **However, as the training progresses, the weight for the second
 178 component gradually increases. This enables the preservation of opacity only on specific meshes
 179 and ensures consistency with the later stages, where only meshes with opacity greater than 0.3 are
 180 exported for rasterization.**

181 After 10000 epochs, a process is initiated where we randomly select sampling points within each
 182 grid, calculate the opacity of each grid, and select the vertex of the grid with the highest opacity
 183 among the eight grids as the vertex of the merged grid. This allows us to generate grids of $64 \times 64 \times 64$,
 184 $32 \times 32 \times 32$, $16 \times 16 \times 16$, and $8 \times 8 \times 8$ sequentially. Each grid level is tailored for specific observation
 185 heights, enabling us to minimize the rendering cost of the scene. It is important to note that for the
 186 same block, different levels always use the same network. The final loss function of the photometric
 187 part is calculated as follows:

$$\mathcal{L}_{rgb} = \sum_{l=1}^5 \beta \mathcal{L}_{rgb}^{part1}(\theta, V_p) + (1 - \beta) \mathcal{L}_{rgb}^{part2}(\theta, V_p) \quad (8)$$

188 **Pseudo-depth.** Utilizing depth information [10, 33] as a supervisory signal has demonstrated
 189 the ability to enhance the convergence rate and improve the quality of reconstructions. However,
 190 when operating under the influence of ambient light in outdoor environments, RGB-D cameras are
 191 susceptible to produce inaccurate depth measurements and have limited measurement ranges. In
 192 our approach, we propose to use the pseudo-depth as the supervision information for NeRF. In
 193 conventional NeRF, structure-from-motion (SFM) [35] solvers such as COLMAP [15, 39, 12] are
 194 typically utilized to estimate the camera poses. As part of this process, a sparse point cloud can be
 195 obtained. By aligning the camera and point cloud in the same coordinate system, we calculate the
 196 pseudo-depth by measuring the distance between the sparse point cloud and the camera position. For

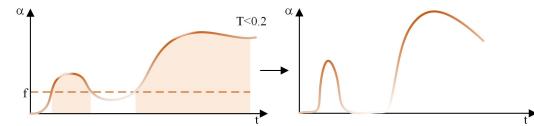


Figure 3: **Illustration of sampling details.** Our approach involves selectively sampling intersection points that are located above a threshold value f and terminating the volume rendering process when the remaining opacity drops below 0.2. This strategy effectively concentrates opacity on a limited number of specific meshes and compresses it towards the object, leading to improved rendering results.

197 a detailed explanation of the pseudo-depth loss (\mathcal{L}_D) calculation, please refer to the supplementary
 198 material. The final loss \mathcal{L} of the model is calculated as follows:

$$\mathcal{L} = \mathcal{L}_{rgb} + \mathcal{L}_D \quad (9)$$

199 **Transient Objects.** Moving/dynamic objects, such as pedestrians and cars, have the potential to
 200 cause inconsistency in the observed scene across different viewpoints and can also create gaps or
 201 "holes" in the depth map. To address this challenge, we integrate a semantic segmentation model
 202 [7, 48] into our training process. By leveraging this model, we generate masks that identify the
 203 regions occupied by the moving objects, effectively excluding them from the training of the scene.
 204 This strategy mitigates the impact of moving objects on the scene's depth map, and thus improves the
 205 overall fidelity and generalization performance of the model.

206 3.2 Pre-rendering

207 To optimize computational efficiency, we leverage the predicted values from the acceleration
 208 grid to identify and remove meshes that do not possess significant geometric surfaces. Once
 209 the training phase is complete, we obtain the network weights as well as the meshes with
 210 optimized vertex positions. Subsequently, we perform the pre-rendering process. The overall
 211 workflow is illustrated in Figure 4. **First**, for each block, in addition to the camera view used
 212 for training, we incorporate parallel lights from various angles are added above. We compute
 213 the intersection of each ray with the polygon mesh, and the traversal of light is halted when
 214 the cumulative opacity surpasses 0.8. We filter out meshes that have an intersection opacity less
 215 0.3, resulting in the retention of only those meshes that meet this criterion. From this, we obtain the
 216 final polygon mesh representations of the scene. Based on our statistical analysis, the final extracted
 217 meshes account for only 5% of the original mesh count. This significant reduction effectively im-
 218 proves mesh utilization, inference speed during rendering, and reduces storage costs. **Second**, we
 219 perform coordinate transformation to ensure the obtained mesh aligns with real-world scale. **Third**,
 220 we map the vertex coordinates of the triangles to their corresponding positions in UV coordinates.
 221

222 3.3 Rendering

223 After training, we obtain the weights of the Encoder-Decoder network, and refined polygonal meshes
 224 at different levels through vertex optimization for each individual block. Subsequently, we perform
 225 pre-rendering and extract the final meshes. For each polygonal mesh, we perform rasterization using a
 226 texture map of size 32×32 . However, storing texture maps for numerous polygonal meshes in a large-
 227 scale scene can lead to significant storage overhead. To address this, we utilize dynamic inference
 228 instead of directly storing all texture maps, as done in Mobile-NeRF [6]. Real-time rendering is
 229 achieved through continuous and efficient communication between the inference submodule and the
 230 UE4 submodule. During the block composition process, since the large-scale scene is accurately
 231 divided into multiple blocks and scene fusion is also considered during training, transitions between
 232 different blocks in a large-scale scene are smooth and interpolation is not required at the junction.
 233

234 **Texture map.** Prior to rasterizing the faces of the polygonal mesh, we perform sampling on each
 235 mesh and use the encoder network to obtain 8 channels feature vectors. These 8 channels along
 236 with the opacity α obtained are processed to generate nine single-channel texture maps using BC4
 237 compression [26]. Compared to the original texture map, this reduces 2/3 of the GPU memory
 238 consumption during rendering without significantly compromising the image quality. *Please note*
 239 *that texture maps do not store colors, but rather spatial features, which also contain transparency*
 240 *information.* Taking into account the principle of locality, we employ a strategy where we collectively
 241 infer and store information from nearby regions into texture maps. This enables us to directly access
 242 the texture maps without the need for re-inference when traversing these regions again.
 243

Table 1: **Details of the real-world scene** Table 2: **Quantitative Comparisons with Existing captured by UAV.** The memory requirement **Methods.** Our method achieves remarkable results in increases with the size of the scene, while terms of both rendering quality and speed. Among the the variation in FPS is primarily influenced compared methods, UE4-NeRF requires the least GPU by the materials present in the scene itself. memory. The training costs are also provided.

	Memory(4k)		Area(m^2)	FPS		CA	Method	PSNR↑	SSIM↑	LPIPS↓	FPS(2K)	Memory(GB) Host GPU	time(mins)	
	Host	GPU		2k	4k									
FL	~ 12GB	~ 5GB	180×240	55	36	120m	NeRFacto [38]	20.99	0.663	0.389	0.5	~ 14	~ 3	~ 45
VL	~ 16GB	~ 7GB	180×240	46	30	70m	Mip-NeRF [2]	16.99	0.516	0.458	X	~ 19	~ 6	~ 120
CS	~ 18GB	~ 8GB	240×240	66	43	70m	Urban-NeRF [32]	20.42	0.563	0.409	X	-	-	-
IP	~ 40GB	~ 14GB	600×600	58	35	120m	NeRF-w [22]	17.16	0.455	0.620	X	-	-	-
							Instant-NGP [25]	22.00	0.631	0.426	11	~ 10	~ 32	~ 15
							UE4-NeRF	25.03	0.704	0.287	55	~ 19	~ 3	~ 40

248 **Rasterization.** We utilize the traditional rasterization pipeline in UE4 to rasterize each face of the
 249 polygonal mesh. Subsequently, the decoder network is applied to convert the 17D features of each
 250 pixel into RGB colors (see UE4 submodule in Figure 2). The 17D features comprise 8 channels
 251 of learned features and a 9D view direction, which is encoded using spherical harmonic functions
 252 [11, 24]. The decoder network is implemented as a HLSL fragment shader of the UE4, which is
 253 called Neural Shader. To achieve translucent effects without the need for volume rendering, we
 254 employ alpha-dithered [42], which offers a relatively low-overhead approach.

255 **Dynamic inference.** The interaction between the UE4 submodule and the inference submodule is
 256 depicted in Figure 2. When a block requiring rendering is encountered, the UE4 submodule supplies
 257 the block and level index information, while the inference submodule provides the corresponding
 258 regional meshes and feature maps. In cases where feature maps are not readily for a specific region,
 259 the inference operation is carried out to generate them. This communication mechanism ensures
 260 efficient and prompt processing between these two submodules.

261 4 Experiments

262 4.1 Implementation Details

263 Similar to the Instant-NGP approach, we employ multi-resolution hash encoding and construct an
 264 Encoder network comprising a 4-layer MLP ($32 \times 64, 64 \times 64, 64 \times 64, 64 \times 8$) to estimate the opacity
 265 and an 8-dimensional feature vector. Additionally, we develop a Decoder network with 3-layer MLP
 266 ($17 \times 16, 16 \times 16, 16 \times 3$) to predict the final color of the sampled points. For each block, we train
 267 the model using one Nvidia RTX 3090 GPU for 80,000 epochs to achieve convergence, with an
 268 approximate duration of 40 minutes.

269 **Datasets and Metric.** Our objective is to enable real-time rendering of large-scale scenes with
 270 complex textures and materials. Table 1 presents four distinct scenarios, namely Farmland(FL),
 271 Village(VL), Construction site(CS) and Industrial parks(IP) which were captured using drones to
 272 address the challenges that UE4-NeRF was designed to tackle. Each captured image has a resolution
 273 of 5472×3648 pixels. At *level-1*, UE4-NeRF can achieve an maximum frame rate of around 43 FPS
 274 at 4K resolution. Due to the limited battery life of the drone, the captured areas were constrained in
 275 size. Nevertheless, subsequent experiments have revealed that limitation does not represent the upper
 276 limit of UE4-NeRF. In fact, UE4-NeRF can render even larger-scale scenes in real-time. To evaluate
 277 the reconstruction quality and fidelity of our model, we use evaluation metrics such as PSNR/SSIM
 278 [41] and LPIPS [47]. We also compare the real-time rendering speed of UE4-NeRF with other models
 279 when rendering vast scenes.

280 4.2 Comparisons with Existing Methods

281 We performed qualitative comparisons between UE4-NeRF and other NeRF implementations on
 282 four scenes captured by UAV, and the results are depicted in Figure 5. Existing tools like NeRFacto
 283 [38] and Instant-NGP [25] aim to achieve real-time rendering, but tend to procude subpar results
 284 when applied to large-scale scenes, as evident in Figure 5. Mip-NeRF, despite increasing the number
 285 of training iterations, still exhibits noticeable blurriness when rendering natural scenes. In contrast,
 286 UE4-NeRF showcases remarkable accuracy in rendering small objects, such as power lines, in the

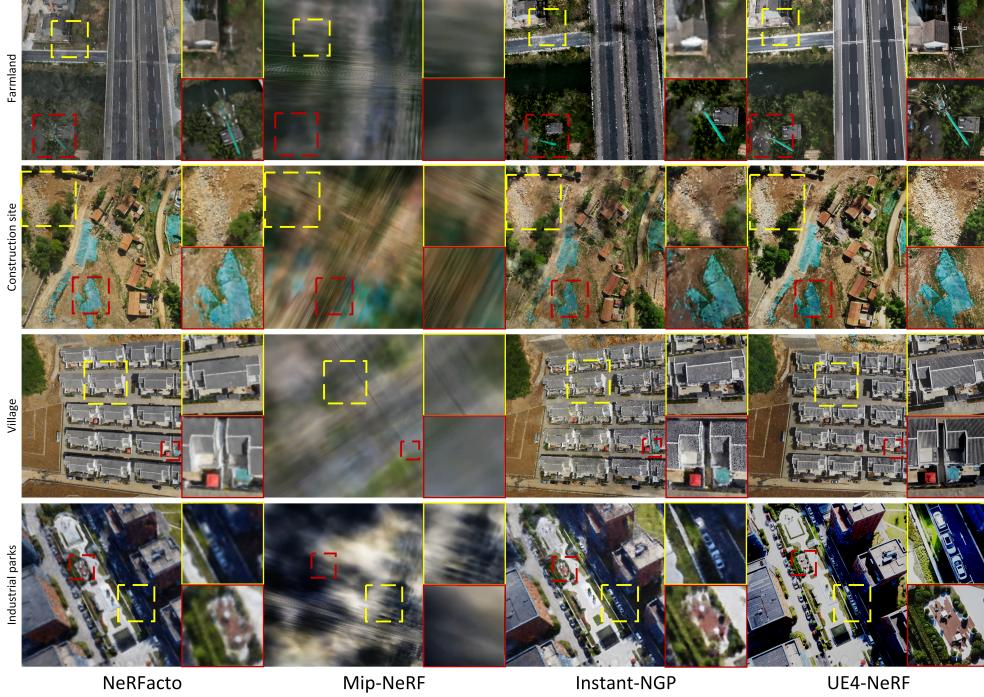


Figure 5: **Qualitative Comparisons with Existing Methods.** We visually compare our proposed technique with recently implemented real-time rendering NeRF models [25, 38] and Mip-NeRF [3]. All these methods are implemented within nerfstudio [38]. It is evident that UE4-NeRF yields the best reconstruction for all the four scenes.



Figure 6: **Rendering result at different LOD.** The red boxes indicate the presentation of the same content at different levels. In addition to achieving high FPS rendering at different LOD in 2K resolution, UE4-NeRF also demonstrates a remarkable photo-realistic.

287 "farmland" scene. It also excels in handling dense textures, as demonstrated in the "village" scene.
 288 The "construction site" scene showcases the rendering effects of translucent objects. In Table 2,
 289 we provide a quantitative comparison of our proposed technique with other existing methods using
 290 different metrics. It is important to note that the actual rendering quality of UE4-NeRF is expected
 291 to surpass the metric's performance, as achieving consistent exposure matching with the original
 292 image is challenging within the Unreal Engine environment. Despite this limitation, UE4-NeRF still
 293 outperforms other methods by a significant margin and its real-time rendering capability surpasses
 294 that of the current state-of-the-art real-time rendering models.

295 4.3 LOD render results

296 Figure 6 showcases the qualitative rendering results at different levels of detail, along with the
 297 corresponding real-time frame rates, at a resolution of 2K. As the level of hierarchy decreases, the

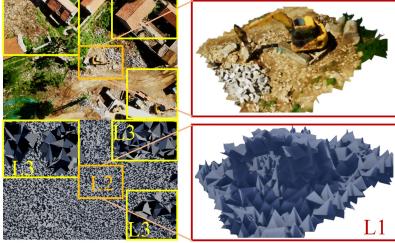


Figure 7: **Shading mesh**– not textured. Polygonal meshes at different levels are displayed, revealing the varying complexity of the scene.

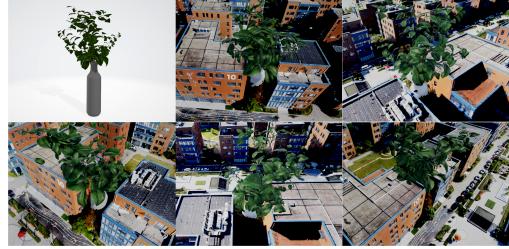


Figure 8: **Scene edition**. The results of multi-view scene editing are presented. One can carefully observe the occlusion relationship between objects and the scene.

298 rendering quality and clarity progressively improve. UE4-NeRF demonstrates a high level of fidelity
 299 when compared to the ground truth. In addition, it can be observed that the LOD approach effectively
 300 addresses the challenge of real-time rendering at higher levels although it may sacrifice some rendering
 301 quality at higher levels for the sake of real-time performance. When comparing materials and frame
 302 rates across different scenes, it is worth noting that the presence of semi-transparent objects, such as
 303 vegetation(as seen in the last scene of Figure 6), can potentially lead to a decrease in frame rate due
 304 to increased computational overhead.

305 **Shading mesh.** Figure 7 showcases the shading meshes extracted without textures. Despite the
 306 coarse appearance of the meshes and the limited level of detail even in the closest *level-1*, such
 307 as excavators, the final rendering quality remains impressive. This can be attributed to the vertex
 308 optimization carried out during the training process, the texture maps encoded in the network outputs,
 309 and the deployment of the final neural shader during rendering.

310 4.4 Scene edition

311 Editing scenes in UE4-NeRF, especially the composition of different scenes, is a seamless process.
 312 UE4 supports multiple 3D model file formats, allowing us to import and edit the rendered scenes. The
 313 Unreal Engine automatically handles occlusion between the original scene and added objects, making
 314 the editing process effortless. We have complete freedom to manipulate the newly added objects
 315 according to our needs. In Figure 8, when an object obstructs the NeRF-rendered scene, the occluded
 316 parts do not require shading by the neural shader. Interestingly, we observed that after adding and
 317 manipulating objects, the frame rate of the scene actually slightly improves. This is attributed to the
 318 reduction in NeRF computation, resulting in improved frame rates.

319 5 Conclusion and Limitations

320 In this research paper, we introduce UE4-NeRF, a real-time interactive rendering system designed
 321 for large-scale scenes. UE4-NeRF partitions the scene into blocks and trains a separate NeRF
 322 model for each block. Inspired by LOD techniques, different levels of mesh precision are used for
 323 different layers, and the NeRF model is integrated with the rasterization pipeline in UE4. Currently,
 324 UE4-NeRF stands as the only NeRF model capable of achieving real-time interactive rendering of
 325 large-scale realistic scenes while maintaining high fidelity. Moreover, its integration with UE4 opens
 326 up possibilities of scene editing and manipulation, allowing for the flexible addition and manipulation
 327 of traditional 3D models such as *.obj* and *.fbx* within the scene. Overall, UE4-NeRF presents a
 328 comprehensive neural rendering system that combines large-scale, real-time interactive rendering
 329 with photo-realism and editability. **However**, there are still some issues that we aim to address in
 330 the future: 1) UE4-NeRF requires the direct use of CUDA and therefore relies on NVIDIA GPUs.
 331 Rendering large scenes incurs significant memory overhead, and increasing rendering accuracy further
 332 amplifies memory consumption. 2) During the pre-rendering process, it is difficult to capture rays
 333 from any viewing angles and extract meshes for each viewpoint, resulting in gaps or holes in the final
 334 rendered scene. The strategy of pre-rendering can be optimized.

335 **References**

- 336 [1] John Ashburner and Karl J Friston. Voxel-based morphometry—the methods. *Neuroimage*, 11(6):805–821,
337 2000.
- 338 [2] Jonathan T Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and Pratul P
339 Srinivasan. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields. In *Proceedings
340 of the IEEE/CVF International Conference on Computer Vision*, pages 5855–5864, 2021.
- 341 [3] Jonathan T Barron, Ben Mildenhall, Dor Verbin, Pratul P Srinivasan, and Peter Hedman. Mip-nerf 360:
342 Unbounded anti-aliased neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer
343 Vision and Pattern Recognition*, pages 5470–5479, 2022.
- 344 [4] Matthew Berger, Andrea Tagliasacchi, Lee M Seversky, Pierre Alliez, Gael Guennebaud, Joshua A Levine,
345 Andrei Sharf, and Claudio T Silva. A survey of surface reconstruction from point clouds. In *Computer
346 graphics forum*, volume 36, pages 301–329. Wiley Online Library, 2017.
- 347 [5] Mario Botsch, Leif Kobbelt, Mark Pauly, Pierre Alliez, and Bruno Lévy. *Polygon mesh processing*. CRC
348 press, 2010.
- 349 [6] Zhiqin Chen, Thomas Funkhouser, Peter Hedman, and Andrea Tagliasacchi. Mobilenerf: Exploiting the
350 polygon rasterization pipeline for efficient neural field rendering on mobile architectures. *arXiv preprint
351 arXiv:2208.00277*, 2022.
- 352 [7] Bowen Cheng, Maxwell D Collins, Yukun Zhu, Ting Liu, Thomas S Huang, Hartwig Adam, and Liang-
353 Chieh Chen. Panoptic-deeplab: A simple, strong, and fast baseline for bottom-up panoptic segmentation. In
354 *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 12475–12485,
355 2020.
- 356 [8] James H Clark. Hierarchical geometric models for visible surface algorithms. *Communications of the
357 ACM*, 19(10):547–554, 1976.
- 358 [9] Alvaro Collet, Ming Chuang, Pat Sweeney, Don Gillett, Dennis Evseev, David Calabrese, Hugues Hoppe,
359 Adam Kirk, and Steve Sullivan. High-quality streamable free-viewpoint video. *ACM Transactions on
360 Graphics (ToG)*, 34(4):1–13, 2015.
- 361 [10] Kangle Deng, Andrew Liu, Jun-Yan Zhu, and Deva Ramanan. Depth-supervised nerf: Fewer views and
362 faster training for free. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern
363 Recognition*, pages 12882–12891, 2022.
- 364 [11] Sara Fridovich-Keil, Alex Yu, Matthew Tancik, Qinhong Chen, Benjamin Recht, and Angjoo Kanazawa.
365 Plenoxels: Radiance fields without neural networks. In *Proceedings of the IEEE/CVF Conference on
366 Computer Vision and Pattern Recognition*, pages 5501–5510, 2022.
- 367 [12] Yasutaka Furukawa and Jean Ponce. Accurate, dense, and robust multi-view stereopsis (pmvs). In *IEEE
368 Computer Society Conference on Computer Vision and Pattern Recognition*, volume 2, 2007.
- 369 [13] Michelle Guo, Alireza Fathi, Jiajun Wu, and Thomas Funkhouser. Object-centric neural scene rendering.
370 *arXiv preprint arXiv:2012.08503*, 2020.
- 371 [14] Yulan Guo, Hanyun Wang, Qingyong Hu, Hao Liu, Li Liu, and Mohammed Bennamoun. Deep learning
372 for 3d point clouds: A survey. *IEEE transactions on pattern analysis and machine intelligence*, 43(12):
373 4338–4364, 2020.
- 374 [15] Richard Hartley and Andrew Zisserman. *Multiple view geometry in computer vision*. Cambridge university
375 press, 2003.
- 376 [16] Kurt Hornik, Maxwell Stinchcombe, and Halbert White. Multilayer feedforward networks are universal
377 approximators. *Neural networks*, 2(5):359–366, 1989.
- 378 [17] Brian Karis and Epic Games. Real shading in unreal engine 4. *Proc. Physically Based Shading Theory
379 Practice*, 4(3):1, 2013.
- 380 [18] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial
381 networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages
382 4401–4410, 2019.
- 383 [19] Verica Lazova, Vladimir Guzov, Kyle Olszewski, Sergey Tulyakov, and Gerard Pons-Moll. Control-nerf:
384 Editable feature volumes for scene rendering and manipulation. In *Proceedings of the IEEE/CVF Winter
385 Conference on Applications of Computer Vision*, pages 4340–4350, 2023.

- 386 [20] Lingjie Liu, Jiatao Gu, Kyaw Zaw Lin, Tat-Seng Chua, and Christian Theobalt. Neural sparse voxel fields.
 387 *Advances in Neural Information Processing Systems*, 33:15651–15663, 2020.
- 388 [21] David Luebke, Martin Reddy, Jonathan D Cohen, Amitabh Varshney, Benjamin Watson, and Robert
 389 Huebner. *Level of detail for 3D graphics*. Morgan Kaufmann, 2003.
- 390 [22] Ricardo Martin-Brualla, Noha Radwan, Mehdi SM Sajjadi, Jonathan T Barron, Alexey Dosovitskiy, and
 391 Daniel Duckworth. Nerf in the wild: Neural radiance fields for unconstrained photo collections. In
 392 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7210–7219,
 393 2021.
- 394 [23] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng.
 395 Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM*, 65
 396 (1):99–106, 2021.
- 397 [24] Claus Müller. *Spherical harmonics*, volume 17. Springer, 2006.
- 398 [25] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives
 399 with a multiresolution hash encoding. *ACM Transactions on Graphics (ToG)*, 41(4):1–15, 2022.
- 400 [26] Jae-Ho Nah. Quicketc2: Fast etc2 texture compression using luma differences. *ACM Transactions on
 401 Graphics (TOG)*, 39(6):1–10, 2020.
- 402 [27] Thomas Neff, Pascal Stadlbauer, Mathias Parger, Andreas Kurz, Joerg H Mueller, Chakravarty R Alla
 403 Chaitanya, Anton Kaplanyan, and Markus Steinberger. Donerf: Towards real-time rendering of compact
 404 neural radiance fields using depth oracle networks. In *Computer Graphics Forum*, volume 40, pages 45–59.
 405 Wiley Online Library, 2021.
- 406 [28] Michael Niemeyer and Andreas Geiger. Giraffe: Representing scenes as compositional generative neural
 407 feature fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*,
 408 pages 11453–11464, 2021.
- 409 [29] Keunhong Park, Utkarsh Sinha, Jonathan T Barron, Sofien Bouaziz, Dan B Goldman, Steven M Seitz, and
 410 Ricardo Martin-Brualla. Nerfies: Deformable neural radiance fields. In *Proceedings of the IEEE/CVF
 411 International Conference on Computer Vision*, pages 5865–5874, 2021.
- 412 [30] Weichao Qiu and Alan Yuille. Unrealcv: Connecting computer vision to unreal engine. In *Computer
 413 Vision–ECCV 2016 Workshops: Amsterdam, The Netherlands, October 8–10 and 15–16, 2016, Proceedings,
 414 Part III 14*, pages 909–916. Springer, 2016.
- 415 [31] Christian Reiser, Songyou Peng, Yiyi Liao, and Andreas Geiger. Kilonerf: Speeding up neural radiance
 416 fields with thousands of tiny mlps. In *Proceedings of the IEEE/CVF International Conference on Computer
 417 Vision*, pages 14335–14345, 2021.
- 418 [32] Konstantinos Rematas, Andrew Liu, Pratul P Srinivasan, Jonathan T Barron, Andrea Tagliasacchi, Thomas
 419 Funkhouser, and Vittorio Ferrari. Urban radiance fields. In *Proceedings of the IEEE/CVF Conference on
 420 Computer Vision and Pattern Recognition*, pages 12932–12942, 2022.
- 421 [33] Antoni Rosinol, John J Leonard, and Luca Carlone. Nerf-slam: Real-time dense monocular slam with
 422 neural radiance fields. *arXiv preprint arXiv:2210.13641*, 2022.
- 423 [34] Franco Scarselli and Ah Chung Tsoi. Universal approximation using feedforward neural networks: A
 424 survey of some existing methods, and some new results. *Neural networks*, 11(1):15–37, 1998.
- 425 [35] Johannes L Schonberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Proceedings of the
 426 IEEE conference on computer vision and pattern recognition*, pages 4104–4113, 2016.
- 427 [36] Vincent Sitzmann, Semon Reznichenko, Bill Freeman, Josh Tenenbaum, and Fredo Durand. Light field
 428 networks: Neural scene representations with single-evaluation rendering. *Advances in Neural Information
 429 Processing Systems*, 34:19313–19325, 2021.
- 430 [37] Matthew Tancik, Vincent Casser, Xinchen Yan, Sabeek Pradhan, Ben Mildenhall, Pratul P Srinivasan,
 431 Jonathan T Barron, and Henrik Kretzschmar. Block-nerf: Scalable large scene neural view synthesis. In
 432 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8248–8258,
 433 2022.
- 434 [38] Matthew Tancik, Ethan Weber, Evonne Ng, Ruilong Li, Brent Yi, Justin Kerr, Terrance Wang, Alexander
 435 Kristoffersen, Jake Austin, Kamyar Salehi, et al. Nerfstudio: A modular framework for neural radiance
 436 field development. *arXiv preprint arXiv:2302.04264*, 2023.

- 437 [39] Bill Triggs, Philip F McLauchlan, Richard I Hartley, and Andrew W Fitzgibbon. Bundle adjustment—a
 438 modern synthesis. In *Vision Algorithms: Theory and Practice: International Workshop on Vision Algo-*
 439 *rithms Corfu, Greece, September 21–22, 1999 Proceedings*, pages 298–372. Springer, 2000.
- 440 [40] Haithem Turki, Deva Ramanan, and Mahadev Satyanarayanan. Mega-nerf: Scalable construction of
 441 large-scale nerfs for virtual fly-throughs. In *Proceedings of the IEEE/CVF Conference on Computer Vision*
 442 *and Pattern Recognition*, pages 12922–12931, 2022.
- 443 [41] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error
 444 visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612, 2004.
- 445 [42] Peter L Williams, Randall J Frank, and Eric C LaMar. Alpha dithering to correct low-opacity 8 bit
 446 compositing errors. *Lawrence Livermore National Laboratory Technical Report UCRL-ID-153185*, 2003.
- 447 [43] Yuanbo Xiangli, Lining Xu, Xingang Pan, Nanxuan Zhao, Anyi Rao, Christian Theobalt, Bo Dai,
 448 and Dahua Lin. Bungeenerf: Progressive neural radiance field for extreme multi-scale scene rendering.
 449 In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022,*
 450 *Proceedings, Part XXXII*, pages 106–122. Springer, 2022.
- 451 [44] Bangbang Yang, Yinda Zhang, Yinghao Xu, Yijin Li, Han Zhou, Hujun Bao, Guofeng Zhang, and Zhaopeng
 452 Cui. Learning object-compositional neural radiance field for editable scene rendering. In *Proceedings of*
 453 *the IEEE/CVF International Conference on Computer Vision*, pages 13779–13788, 2021.
- 454 [45] Alex Yu, Ruilong Li, Matthew Tancik, Hao Li, Ren Ng, and Angjoo Kanazawa. Plenoctrees for real-time
 455 rendering of neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer*
 456 *Vision*, pages 5752–5761, 2021.
- 457 [46] Kai Zhang, Gernot Riegler, Noah Snavely, and Vladlen Koltun. Nerf++: Analyzing and improving neural
 458 radiance fields. *arXiv preprint arXiv:2010.07492*, 2020.
- 459 [47] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable
 460 effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer*
 461 *vision and pattern recognition*, pages 586–595, 2018.
- 462 [48] Hengshuang Zhao, Jianping Shi, Xiaojuan Qi, Xiaogang Wang, and Jiaya Jia. Pyramid scene parsing
 463 network. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages
 464 2881–2890, 2017.
- 465 [49] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using
 466 cycle-consistent adversarial networks. In *Proceedings of the IEEE international conference on computer*
 467 *vision*, pages 2223–2232, 2017.

468 **Checklist**

- 469 1. For all authors...
- 470 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
 471 contributions and scope? **[Yes]**
- 472 (b) Did you describe the limitations of your work? **[Yes]** See Section 5 Limitation.
- 473 (c) Did you discuss any potential negative societal impacts of your work? **[Yes]** We discuss
 474 it in supplementary materials.
- 475 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
 476 them? **[Yes]**
- 477 2. If you are including theoretical results...
- 478 (a) Did you state the full set of assumptions of all theoretical results? **[N/A]**
- 479 (b) Did you include complete proofs of all theoretical results? **[N/A]**
- 480 3. If you ran experiments...
- 481 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
 482 mental results (either in the supplemental material or as a URL)? **[Yes]**
- 483 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
 484 were chosen)? **[Yes]** See Section 4.1.

- (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [N/A]

(d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes] See Section 4.1.

4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...

 - (a) If your work uses existing assets, did you cite the creators? [Yes]
 - (b) Did you mention the license of the assets? [Yes] See supplementary materials.
 - (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
 - (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [N/A]
 - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [N/A]

5. If you used crowdsourcing or conducted research with human subjects...

 - (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
 - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
 - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]