

LightOctree: Lightweight 3D Spatially-Coherent Indoor Lighting Estimation

Supplementary Material

001 1. Details in Rendering Layer

002 In this section, we delve into the rendering layer employed
003 within our framework. Leveraging a multi-level octree for
004 sparse illumination representation, we have devised a fast
005 differentiable cone tracing method to render 2D panoramic
006 environment maps from lighting octree.

007 In classical volume rendering, we calculate the light ra-
008 diation $L(s)$ from a viewpoint \mathbf{x} in direction ω using the
009 rendering equation (Equation 1). The radiance function is
010 defined as the integral $L(s)$ of transmittance $T(t)$, density
011 $\sigma(t)$, and radiance intensity $L_e(t)$ over the path length s .

$$012 \quad L(s) = \int_{t=0}^s T(t)\sigma(t)L_e(t)dt \quad (1)$$

013 To solve this equation, we discretize it with a sampling dis-
014 tance δ_k between consecutive sampling points:

$$015 \quad L(s) = \sum_{n=1}^N C_n T_n (1 - e^{-\sigma_n \delta_n}), \quad T_n = e^{-\sum_{k=1}^{n-1} \sigma_k \delta_k} \quad (2)$$

016 We optimize the rendering process with cone tracing to
017 reduce samples and enhance sampling efficiency, inspired by
018 prior works [1, 3]. This approach also effectively integrates
019 multi-scale outputs from the lighting estimation network,
020 adapting to the characteristics of octree-based networks.

021 Given the angle θ of the cone, the sampling position s_n
022 and sampling step distance δ_n of each sampling point are
023 computed based on the distance traveled and the angle of
024 the previous point (Equation 3). This approach allows for
025 more frequent sampling of nearby light sources, enhancing
026 sensitivity to critical lighting information.

$$027 \quad s_n = \sum_{k=0}^n \delta_k, \quad \delta_n = \begin{cases} c_0, & n = 0 \\ c \cdot s_{n-1} \tan \theta, & n \geq 1 \end{cases} \quad (3)$$

028 Then, by calculating the cone radius at the sampling point
029 location, we sample octree nodes at depth d_n :

$$030 \quad d_n = \lceil \log_2 \frac{\delta_n}{l_0} \rceil \quad (4)$$

031 where l_0 illustrate the minimum side length of leaf node
032 in the octree. Based on these processes, the new rendering
033 equation can be formulated as equation 5.

$$034 \quad L(s) = \sum_{n=1}^N w_n C_n^{d_n} T_n (1 - e^{-\sigma_n^{d_n} \delta_n}) \quad (5)$$

035 where $C_n^{d_n}$ and $\sigma_n^{d_n}$ represents the radiance and density of
036 sampling point n at a depth of d_n . And $w_n = \delta_n/s_n$ is the
037 weight for each sampling point.

038 To ensure the differentiability of the above process, we
039 derive the derivatives of the color and opacity with respect
040 to the forward rendering process, as shown in equation 6,
041 where $E_n = 1 - e^{-\sigma_n^{d_n} \delta_n}$ to simplify the formula.

$$\frac{\partial L(\mathbf{x}, \omega)}{\partial \sigma_t^{d_t}} = w_t C_t^{d_t} T_{t+1} \delta_t - \left[L(\mathbf{x}, \omega) - \sum_{n=1}^t w_n C_n^{d_n} T_n E_n \right] \delta_t \quad (6) \quad 042$$

$$\frac{\partial L(\mathbf{x}, \omega)}{\partial C_t^{d_t}} = \frac{\partial \left(\sum_{n=1}^N C_n^{d_n} T_n (1 - e^{-\sigma_n^{d_n} \delta_n}) \right)}{\partial C_t^{d_t}} = T_t E_t \quad (7) \quad 044$$

045 The above equations describe how to calculate derivatives
046 during rendering. The derivatives of color can be computed
047 during the sampling process, while the derivatives of density
048 require the use of the final rendering results. Building upon
049 the derived formulas, we have implemented a differentiable
050 renderer using Taichi[2], enabling rapid forward rendering
051 and gradient calculations.

052 2. Details in Object Insertion

053 Given that common object representations are typically
054 meshes, while point clouds are prevalent in AR applica-
055 tions, and our lighting representation employs a voxel oc-
056 tree structure, we've devised a fusion rendering method that
057 seamlessly combines point clouds, meshes, and octrees.

058 2.1. Scene Representation

059 **Background Representation** Differential rendering en-
060 compasses the rendering of virtual objects and the shading
061 computation, which requires multiple environment lighting
062 queries. Each query involves casting a cone from the surface
063 position \mathbf{x} into the environment to find radiance incident
064 from the ray direction (solid angle direction w_i). We employ
065 a point cloud to represent the original scene, including depth
066 and color information from the viewpoint to the surfaces.
067 This point cloud is per-pixel, allowing us to skip ray colli-
068 sions from the viewpoint to the scene surfaces and directly
069 use the projected depth from the point cloud as the position.
070 We utilize this information to determine the spatial relation-
071 ship between the scene and virtual objects. We perform cone
072 tracing to illuminate the rendered virtual objects.

073 **Virtual Object Representation** To efficiently perform ray
074 intersections and integrate it into the Taichi framework, we

075 normalizes the bounding boxes of virtual objects to the range
076 of 0 to 1. For each triangular face, a multi-level regular grid
077 space division is performed in this three-dimensional space,
078 which is managed using the SNode structure.

079 2.2. Differential Rendering

080 Due to the fact that virtual objects in the scene typically oc-
081 cupy only a small portion, yet they may consist of thousands
082 or even tens of thousands of triangles, achieving real-time
083 rendering with three different data structures in AR applica-
084 tions requires a new approach to index lookup. During the
085 ray traversal through the bounding box, we can calculate the
086 index of the first intersected object grid based on the ray's
087 origin R_o and direction R_d , as shown below:

$$H_{min} = \max \left(\min \left(\frac{Box_{min} - R_o}{R_d}, \frac{Box_{max} - R_o}{R_d} \right) \right)$$

$$I_0 = floor \left(\frac{R_o + R_d H_{min} - Box_{min}}{Box_{size}} \right) \quad (8)$$

088 where Box_{min} and Box_{max} represent the coordinates of
089 the bounding box's vertices, and Box_{size} represents the size
090 of the bounding box.
091

092 In subsequent tracing steps, the index of the next grid
093 can be obtained based on the ray's current position R_t and
094 normalized forward direction R'_d as follows:

$$R'_t = (1 - R_t) * (R'_d \geq 0) + R_t * (R'_d < 0)$$

$$I_{t+1} = I_t + G \left(\frac{R'_t}{R'_d} \right) * sign(R'_d) \quad (9)$$

096 where $G(\cdot)$ calculates the axis with the smallest advancement
097 in the axis directions. For example, if $G((0.5, 0.2, 0.7)) =$
098 $(0, 1, 0)$. During rendering, we emit a ray from the camera
099 viewpoint for each pixel on the screen and employ ray tracing
100 to trace it until it intersects with an object. Upon detecting an
101 intersection, we emit a cone of rays from the object's surface
102 in various directions based on the surface normal. This cone
103 of rays is then traced and rendered using ray marching, with
104 the results obtained from volume rendering serving as the
105 illumination for surface shading. Same to the rendering layer,
106 we begin with a minimum step size δ_{min} and later perform
107 long-distance mip-map sampling based on the cone's radius.

108 As the scene is managed using an octree, obtaining the
109 voxel index directly from the sampling point's position is
110 not feasible. Therefore, we trace the voxel recursively based
111 on the current sampling point's location within the volume.
112 To improve performance, we utilize a shared recursive stack
113 for all cone rays, modifying the sampling order based on
114 the ray's direction. Instead of clearing the stack after each
115 sampling, we continue the recursion in the next sampling, en-
116 hancing efficiency. Once we determine the sampling point's
117 position and its associated node, we perform trilinear inter-

polation using the values within the node block, accumulate
the results along the ray and apply for shading.

3. Details in Dataset Construction

120 We train our model using photorealistic renderings of in-
121 door scenes from the FutureHouse synthetic dataset[4]. This
122 dataset consists of artist-designed indoor panoramas with
123 high-quality geometry and HDR environment maps. We can
124 extract photographs from this dataset to obtain input/output
125 pairs for training. It is worth mentioning that we chose to
126 use panorama datasets instead of InteriorNet dataset used
127 by [5] and [12], or the OpenRoom dataset used by [9] and
128 [10]. This is because the images in InteriorNet have low
129 dynamic range, and the resolution of the panoramic data in
130 OpenRoom is too low. These shortcomings make these two
131 datasets less ideal for lighting estimation tasks. Although
132 InteriorVerse proposed by [11] is a better choice, at the time
133 of writing this paper, the lighting data of this dataset are
134 still unavailable. Fortunately, with the HDR panoramas and
135 related geometry information provided by FutureHouse, we
136 can construct the training data that meets our needs.
137



Figure 1. Schematic visualization of the constructed dataset.

Each training sample $\{I, D, O_s, \{I_{nv}, P_{nv}\}_N\}$ contains
one LDR perspective image I as inputs and paired data as
ground truth to supervise depth and lighting voxel octree.
The construction process of a single sample is as follows:
For a perfect 360-degree camera, we can treat its imaging
plane as a sphere, and any point in world coordinates can be
projected onto this sphere as a spherical projection. There-
fore, for each RGB panoramic image and its corresponding
depth panoramic image in the dataset, we can reproject all
points onto the world coordinate system using the inverse
transformation of spherical projection, thus obtaining the 3D
point cloud of the scene \mathcal{P} . Based on the 3D point cloud
 \mathcal{P} , we can construct all the required ground truth. For the
input image I and its corresponding depth map D , we use
a perspective camera model with a field of view (FOV) of
90 degrees and an aspect ratio of 1 (similar to what [6] did)
to reproject points on the sphere onto a cube face, obtaining
four perspective images that cover a 360-degree horizontal
view. For the ground truth of the scene's voxel octree, we
can directly construct it from the point cloud \mathcal{P} using the
method introduced by [7, 8]. For the new viewpoint images

159 $\{I_{nv}, P_{nv}\}_N$, we construct a uniform voxel grid with a res-
 160 olution of 256^3 from the point cloud \mathcal{P} , and then render N
 161 HDR panoramic images at locations that are visible in the
 162 perspective input image I using a rendering method similar
 163 to NeRFs. The corresponding camera positions P_{nv} are also
 164 recorded during rendering. Based on the 28,579 panoramic
 165 views from 1,752 house-scale scenes provided by Future-
 166 House, we construct and select 113,232 pairs of data for
 167 model training. We use 90% (1,570) of the scenes to train
 168 our model and reserve 10% (180) for evaluation.

169 4. Additional Details

170 **Bad cases.** Depth estimation inaccuracies may lead to un-
 171 foreseen outcomes in environments with complex lighting or
 172 geometry, as depicted in Figure 2. Artifacts observed in both
 173 Wang et al. and ours are presumably due to erroneous initial
 174 depth compromising alpha prediction accuracy. However,
 175 the distant environmental illumination assumption may also
 176 fail in these scenarios with intense spatially-varying lighting.

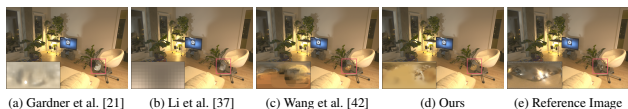


Figure 2. Cases of failure in certain special circumstances.

177 **Data used in user study.** Images are showcased in Figure
 178 3, where video frames in virtual object insertion results were
 179 extracted and amalgamated into a single composite image.

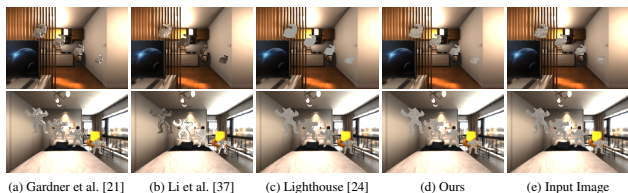


Figure 3. Virtual object insertion examples used in user study.

180 References

181 [1] Cyril Crassin, Fabrice Neyret, Miguel Sainz, Simon Green,
 182 Elmar Eisemann, and MLeal Llaguno. Interactive indirect
 183 illumination using voxel cone tracing. *Computer Graphics*
 184 *Forum*, page 1921–1930, 2011. 1

185 [2] Yuanming Hu, Tzu-Mao Li, Luke Anderson, Jonathan Ragan-
 186 Kelley, and Frédo Durand. Taichi: a language for high-
 187 performance computation on spatially sparse data structures.
 188 *ACM Transactions on Graphics (TOG)*, 38(6):201, 2019. 1

189 [3] Samuli Laine and Tero Karras. Efficient sparse voxel octrees.
 190 In *Proceedings of the 2010 ACM SIGGRAPH symposium on*
 191 *Interactive 3D Graphics and Games*, pages 55–63, 2010. 1

192 [4] Zhen Li, Lingli Wang, Xiang Huang, Cihui Pan, and Jiaqi
 193 Yang. Phyr: Physics-based inverse rendering for panoramic
 194 indoor images. In *IEEE/CVF Conference on Computer Vision*
 195 *and Pattern Recognition, CVPR 2022, New Orleans, LA, USA,*
 196 *June 18-24, 2022*, pages 12703–12713. IEEE, 2022. 2

[5] Pratul P. Srinivasan, Ben Mildenhall, Matthew Tancik,
 Jonathan T. Barron, Richard Tucker, and Noah Snavely. Light-
 house: Predicting lighting volumes for spatially-coherent illu-
 mination. In *2020 IEEE/CVF Conference on Computer Vision*
and Pattern Recognition, CVPR 2020, Seattle, WA, USA, June
13-19, 2020, pages 8077–8086. Computer Vision Foundation
 / IEEE, 2020. 2

[6] Fu-En Wang, Hou-Ning Hu, Hsien-Tzu Cheng, Juan-Ting
 Lin, Shang-Ta Yang, Meng-Li Shih, Hung-Kuo Chu, and Min
 Sun. Self-supervised learning of depth and camera motion
 from 360 videos. In *Asian Conference on Computer Vision*,
 pages 53–68. Springer, 2018. 2

[7] Peng-Shuai Wang, Chun-Yu Sun, Yang Liu, and Xin Tong.
 Adaptive o-cnn: A patch-based deep representation of 3d
 shapes. *ACM Transactions on Graphics*, page 1–11, 2018. 2

[8] Peng-Shuai Wang, Yang Liu, and Xin Tong. Dual octree
 graph networks for learning adaptive volumetric shape repre-
 sentations. *ACM Transactions on Graphics*, page 1–15, 2022.
 2

[9] Zhengqin Li, Mohammad Shafiei, Ravi Ramamoorthi, Kalyan
 Sunkavalli, and Manmohan Chandraker. Inverse rendering
 for complex indoor scenes: Shape, spatially-varying lighting
 and svbrdf from a single image. In *2020 IEEE/CVF Con-*
ference on Computer Vision and Pattern Recognition, CVPR
2020, Seattle, WA, USA, June 13-19, 2020, pages 2472–2481.
 Computer Vision Foundation / IEEE, 2020. 2

[10] Zhengqin Li, Ting-Wei Yu, Shen Sang, Sarah Wang, Meng
 Song, Yuhan Liu, Yu-Ying Yeh, Rui Zhu, Nitesh B. Gun-
 davarapu, Jia Shi, Sai Bi, Hong-Xing Yu, Zexiang Xu, Kalyan
 Sunkavalli, Milos Hasan, Ravi Ramamoorthi, and Manmohan
 Chandraker. Openrooms: An open framework for photoreal-
 istic indoor scene datasets. In *IEEE Conference on Computer*
Vision and Pattern Recognition, CVPR 2021, virtual, June
19-25, 2021, pages 7190–7199. Computer Vision Foundation
 / IEEE, 2021. 2

[11] Jingsen Zhu, Fujun Luan, Yuchi Huo, Zihao Lin, Zhihua
 Zhong, Dianbing Xi, Jiayang Zheng, Rui Tang, Hujun Bao,
 and Rui Wang. Learning-based inverse rendering of complex
 indoor scenes with differentiable monte carlo raytracing. In
SIGGRAPH Asia, 2022. 2

[12] Zian Wang, Jonah Philion, Sanja Fidler, and Jan Kautz. Learn-
 ing indoor inverse rendering with 3d spatially-varying light-
 ing. In *2021 IEEE/CVF International Conference on Com-*
puter Vision, ICCV 2021, Montreal, QC, Canada, October
10-17, 2021, pages 12518–12527. IEEE, 2021. 2