

I²-SDF: Intrinsic Indoor Scene Reconstruction and Editing via Raytracing in Neural SDFs

Jingsen Zhu¹ Yuchi Huo^{2,1} Qi Ye^{3,6} Fujun Luan⁴ Jifan Li¹ Dianbing Xi¹
 Lisha Wang¹ Rui Tang⁵ Wei Hua² Hujun Bao¹ Rui Wang¹

¹State Key Lab of CAD&CG, Zhejiang University ²Zhejiang Lab ³Zhejiang University

⁴Adobe Research ⁵KooLab, Manycore ⁶Key Lab of CS&AUS of Zhejiang Province

<https://jingsenzhu.github.io/i2-sdf>

Abstract

In this work, we present I²-SDF, a new method for intrinsic indoor scene reconstruction and editing using differentiable Monte Carlo raytracing on neural signed distance fields (SDFs). Our holistic neural SDF-based framework jointly recovers the underlying shapes, incident radiance and materials from multi-view images. We introduce a novel bubble loss for fine-grained small objects and error-guided adaptive sampling scheme to largely improve the reconstruction quality on large-scale indoor scenes. Further, we propose to decompose the neural radiance field into spatially-varying material of the scene as a neural field through surface-based, differentiable Monte Carlo raytracing and emitter semantic segmentations, which enables physically based and photorealistic scene relighting and editing applications. Through a number of qualitative and quantitative experiments, we demonstrate the superior quality of our method on indoor scene reconstruction, novel view synthesis, and scene editing compared to state-of-the-art baselines.

1. Introduction

Reconstructing 3D scenes from multi-view images is a fundamental task in computer graphics and vision. Neural Radiance Field (NeRF) [18] and its follow-up research leverage multi-layer perceptions (MLPs) as implicit functions, taking as input the positional and directional coordinates, to approximate the underlying geometry and appearance of a 3D scene. Such methods have shown compelling and high-fidelity results in novel view synthesis. However, we argue that novel view synthesis itself is insufficient for scene editing applications such as inserting virtual objects,



Figure 1. I²-SDF. Left: State-of-the-art neural implicit surface representation method [43] fails in reconstructing small objects inside an indoor scene (e.g. lamps and chandeliers), which is resolved by our bubbling method. Middle and Right: Our intrinsic decomposition and raytracing method enable photo-realistic scene editing and relighting applications.

relighting and editing surface materials with global illumination.

On the other hand, inverse rendering or *intrinsic decomposition*, which reconstructs and decomposes the scene into shape, shading and surface reflectance from single or multiple images, enables photorealistic scene editing possibilities. It is a long-term challenge especially for large-scale indoor scenes because they typically exhibit complex geometry and spatially-varying global illumination appearance. As intrinsic decomposition is an extremely ill-posed task, a physically based shading model will crucially affect the decomposition quality. Existing neural rendering methods [2, 20, 44, 47] rely on simple rendering algorithms (such as pre-filtered shading) for the decomposition and use a global lighting representation (e.g., spherical Gaussians). Although these methods have demonstrated the effective-

ness on object-level inverse rendering, they are inapplicable to complex indoor scenes. Moreover, indoor scene images are usually captured from the inside out and most lighting information has already presented inside the room. As a result, the reconstructed radiance field already provides sufficient lighting information without the need of active, external capture lighting setup.

To tackle the above challenges, we propose I^2 -SDF, a new method to decompose a 3D scene into its underlying shape, material, and incident radiance components using implicit neural representations. We design a robust two-stage training scheme that first reconstructs a neural SDF with radiance field, and then conducts raytracing in the SDF to decompose the radiance field into material and emission fields. As complex indoor scenes typically contain many fine-grained, thin or small structures with high-frequency details that are difficult for an implicit SDF function to fit, we propose a novel bubble loss and an error-guided adaptive sampling scheme that greatly improve the reconstruction quality on small objects in the scene. As a result, our approach achieves higher reconstruction quality in both geometry and novel view synthesis, outperforming previous state-of-the-art neural rendering methods in complex indoor scenes. Further, we present an efficient intrinsic decomposition method that decomposes the radiance field into spatially-varying material and emission fields using surface-based, differentiable Monte Carlo raytracing, enabling various scene editing applications.

In summary, our contributions include:

- We introduce I^2 -SDF¹, a holistic neural SDF-based framework for complex indoor scenes that jointly recovers the underlying shape, radiance, and material fields from multi-view images.
- We propose a novel bubble loss and error-guided adaptive sampling strategy to effectively reconstruct fine-grained small objects inside the scene.
- We are the first that introduce Monte Carlo raytracing technique in scene-level neural SDF to enable photo-realistic indoor scene relighting and editing.
- We provide a high-quality synthetic indoor scene multi-view dataset, with ground truth camera pose and geometry annotations.

2. Related Work

Neural implicit scene representations (or *neural fields*) have recently received extensive attention from the research community for representing 3D geometry and radiance information. Neural radiance field (NeRF) [18] uses a single MLP to encode a scene as a continuous volumetric field of

RGB radiance and density, giving promising results in novel view synthesis. Follow-up works accelerate reconstruction speed using voxels [27, 31], hashgrids [19] or deep image features [4, 42]. Neural fields can also be applied to represent 3D geometric functions [23, 24, 36, 41]. Despite their success in reconstructing small-scaled and textured objects, they have difficulties in handling shape-radiance ambiguity on texture-less surfaces. In this paper, we adopt one of the state-of-the-art implicit SDF methods, volSDF [41], as our neural implicit representation backbone and overcome the difficulties in indoor scene reconstruction task.

Neural 3D reconstruction for indoor scenes. Traditional multi-view stereo methods [28, 29] can produce plausible geometry of textured surfaces, but struggle with texture-less regions such as white walls commonly seen in indoor scenes. Recently, learning-based MVS methods have been widely studied, which can be divided into two categories: depth-based methods and TSDF (truncated signed distance function) based methods. NeuralRecon [32] proposes a coarse-to-fine framework to regress input images to TSDF incrementally. NerfingMVS [39] leverages depth priors to guide the point sampling in NeRF to reduce shape-radiance ambiguity. NeRFusion [46] combines the advantages of NeRF and TSDF-based fusion techniques to achieve reconstruction and rendering. Neural implicit SDF methods [36, 41], which succeed in object-level 3D reconstruction, have also been applied to indoor scene reconstruction. To tackle with texture-less regions, additional priors are exploited to guide the network optimization, including semantic priors [8], normal priors [35, 43] and depth priors [43]. The usage of priors results in improved reconstruction quality.

Inverse rendering (also known as *intrinsic decomposition*) is a long-term challenging and ill-posed problem in computer graphics and vision, which attempts to reconstruct and factorize the scene with geometry, material and lighting from single or multiple images. Monocular methods [6, 9, 15–17, 38, 48] rely on strong priors from large-scale datasets. Recent methods exploit on physically-based techniques such as differentiable rendering [15] or raytracing [48] to achieve high-fidelity predictions, but cannot recover full 3D reconstructions that can be viewed from arbitrary viewpoints. Multi-view methods holistically recover factorized 3D models for relighting and novel view synthesis from additional observations instead of strong priors. Neural implicit representations have been widely researched to estimate BRDF and lighting from image collections. NeRV [30] model light transport to support lighting effects such as shadows with high computational costs. Recent methods jointly estimate 3D geometry, BRDF and lighting from images Illumination is represented as spheri-

¹ I^2 meaning “Intrinsics and Indoor”

cal Gaussians (NeRD [2], PhysSG [44]), low-resolution environment maps (NeRFactor [47]), split-sum lighting model (Neural-PIL [3], NVDIFFREC [20]), or Monte-Carlo estimator (NVDIFFREC-MC [11]). However, these recent methods mainly focus on single object reconstruction and do not handle spatially-varying lighting conditions, which is not applicable to indoor scenes with complex geometry and lighting variations. In contrast, our method handles spatially-varying lighting and achieves indoor relighting with high fidelity.

3. Overview

Our goal is to jointly decompose the underlying shape, incident radiance and material of the indoor scene according to multi-view input images. We use implicit representations [18, 41] to model the scene geometry, radiance and material, parameterizing each factor as a single MLP. Fig. 2 shows an overview of our pipeline, which consists of the neural SDF field (F_d), the neural radiance field (F_c), the neural material fields (F_a and F_ρ) and emission field (F_e with $\mathbf{L}[\cdot]$), and finally the Monte Carlo rendering layer which uses the decomposed factors to re-render the scene image.

To avoid the training ambiguities, we adopt a two-stage training scheme: We first train the geometry (F_d), radiance (F_c) and emitter semantic (F_e) fields, and then train the material (F_a , F_ρ) and emission ($\mathbf{L}[\cdot]$) fields. During the optimization of material and emission fields, F_d , F_c , F_e are fixed and detached from the gradient descents.

In the following sections, we first review the concepts of neural SDF field and volume rendering in Sec. 3.1. Next, we will respectively introduce the design details of decomposed components in Sec. 4, Monte Carlo rendering technique with raytracing in Sec. 5 and the training strategy in Sec. 6.

3.1. Implicit Neural Surface Representation and Volume Rendering

We represent the scene geometry as an implicit signed distance function (SDF). A signed distance function is a continuous function d that maps a 3D point \mathbf{x} to the closest distance of \mathbf{x} to surface:

$$d : \mathbb{R}^3 \rightarrow \mathbb{R} \quad d(\mathbf{x}) = (-1)^{1_{\Omega}(\mathbf{x})} \min_{\mathbf{y} \in \mathcal{M}} \|\mathbf{x} - \mathbf{y}\|, \quad (1)$$

where Ω is the scene space and $\mathcal{M} = \partial\Omega$ is the scene surface. The sign of $d(\mathbf{x})$ indicates whether the point \mathbf{x} is inside or outside the scene. In this work, we parameterize the SDF function as a single MLP F_d . Inspired by NeRF [18], we also parameterize the scene appearance as a view-dependent radiance field F_c :

$$(d(\mathbf{x}), \mathbf{z}(\mathbf{x})) = F_d(\mathbf{x}), \quad (2)$$

$$\mathbf{c}(\mathbf{x}) = F_c(\mathbf{z}(\mathbf{x}), \mathbf{v}), \quad (3)$$

where $\mathbf{z}(\mathbf{x})$ is a neural feature output by F_d providing deep geometric cues to the radiance field F_c , and \mathbf{v} is the view direction to model view-dependent visual effects such as specular reflections.

Following NeRF [18], we adopt differentiable volume rendering to learn scene implicit representation network from images. Specifically, to render a pixel, we cast a ray \mathbf{r} from camera position \mathbf{o} through the pixel along view direction \mathbf{v} . M points $\mathbf{x}_r^i = \mathbf{o} + t_r^i \mathbf{v}$ are sampled along the ray and fed into F_d and F_c to predict their SDF value and radiance. In order to apply color accumulation, we transform SDF value $d_r^i = d(\mathbf{x}_r^i)$ to volume density $\sigma_r^i = \sigma(\mathbf{x}_r^i)$ [41]:

$$\sigma(\mathbf{x}) = \begin{cases} \frac{1}{\beta} \left(1 - \frac{1}{2} \exp\left(-\frac{d(\mathbf{x})}{\beta}\right) \right) & \text{if } d(\mathbf{x}) < 0, \\ \frac{1}{2\beta} \exp\left(-\frac{d(\mathbf{x})}{\beta}\right) & \text{if } d(\mathbf{x}) \geq 0, \end{cases} \quad (4)$$

where β is a learnable parameter to control the sparsity near the surface. The pixel color $\hat{\mathbf{C}}(\mathbf{r})$ for ray \mathbf{r} is rendered via numerical integration [18]:

$$\hat{\mathbf{C}}(\mathbf{r}) = \sum_{i=1}^M T_r^i \alpha_r^i \hat{\mathbf{c}}_r^i, \quad (5)$$

where δ_r^i is the distance between adjacent sampled points \mathbf{x}_r^i and \mathbf{x}_r^{i+1} , $\alpha_r^i = 1 - \exp(-\sigma_r^i \delta_r^i)$ is the alpha value of each point, and $T_r^i = \prod_{j=1}^{i-1} (1 - \alpha_r^j)$ is the accumulated transmittance. Similarly, depth $\hat{D}(\mathbf{r})$, normal $\hat{\mathbf{N}}(\mathbf{r})$ of the surface can be accumulated as:

$$\hat{D}(\mathbf{r}) = \sum_{i=1}^M T_r^i \alpha_r^i t_r^i \cos \theta_v, \quad \hat{\mathbf{N}}(\mathbf{r}) = \sum_{i=1}^M T_r^i \alpha_r^i \hat{\mathbf{n}}_r^i, \quad (6)$$

where θ_v is the angle between ray direction \mathbf{v} and the principle view direction of the camera, and normal values $\mathbf{n}(\mathbf{x})$ can be estimated by computing the gradient of SDF function at point \mathbf{x} .

4. Intrinsics Decomposition

4.1. Geometry Field

4.1.1 Bubbling for Small Objects

Reconstructing small objects inside an indoor scene remains a challenging problem. Unlike single-object scenarios, indoor scenes contain objects of different scales and different visibility levels, some of which appear rarely in the input views or restrict in view poses due to their location (e.g. in the corner). We observe that existing indoor reconstruction methods [35, 43] frequently fail in recognizing and reconstructing thin (e.g. chair legs) or suspended (e.g. chandeliers) objects in the room, *even with dense geometry priors*. Fig. 1 shows a reconstruction error on the chandelier.

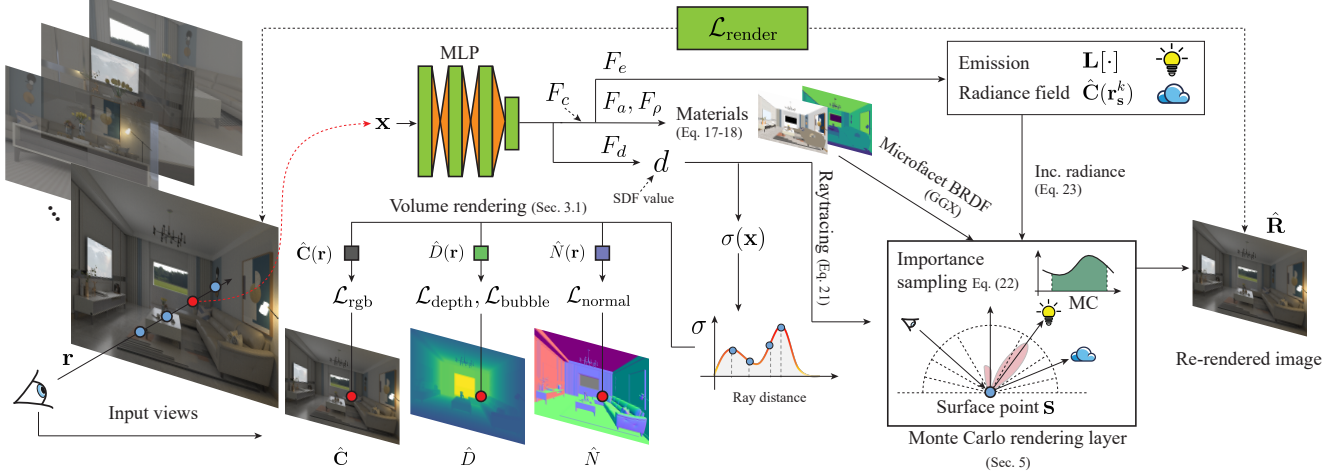


Figure 2. **An overview of our pipeline.** Multi-view images are used to learn the underlying neural SDF field (F_d), radiance field (F_c), material fields (F_a and F_ρ), and emission field (F_e with $\mathbf{L}[\cdot]$), producing an intrinsic neural scene re-renderable for various applications.

The reason for the failure of recovering small objects in the scenes can be attributed to the inherent nature of the neural network as elaborated in [33]: the low-frequency information in a neural network tends to converge faster than the higher-frequency information. In an indoor scene, the SDF for large objects like walls and tables (corresponding to low frequency information) converges in earlier iterations. In some tasks the high-frequency details can be recovered in later iterations during training, while methods based on SDF as geometry representations can hardly recover the fine details due to vanishing gradients for small objects.

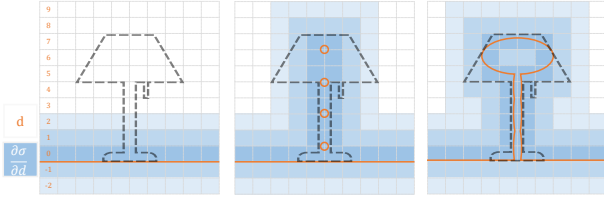


Figure 3. Concept image of bubble. Left: $\frac{\partial \sigma}{\partial d}$ rapidly vanishes as d increases, and thus the SDF cannot learn the thin object. Middle: inserting bubbles (creating zero-value surfaces) recovers the gradient flow around the missing object. Right: the bubbles grow with the introduced gradients to recover the thin object.

An example of the vanishing gradients is illustrated in Fig. 3. Suppose the loss for the radiance \mathcal{L} is a function $c(\sigma(d(\mathbf{x}; \theta)))$. The derivative of \mathcal{L} w.r.t the network parameters θ for SDF can be written as $\frac{\partial \mathcal{L}}{\partial \theta} = \frac{\partial \mathcal{L}}{\partial c} \frac{\partial c}{\partial \sigma} \frac{\partial \sigma}{\partial d} \frac{\partial d}{\partial \theta}$. For $d(\mathbf{x}) \geq 0$, $\frac{\partial \sigma}{\partial d} = \frac{1}{2\beta^3} \exp\left(-\frac{d(\mathbf{x})}{\beta}\right)$ (from Eq. (4)). Usually the learned β is large to make σ fall off the target surface rapidly, indicating $\frac{\partial \sigma}{\partial d}$ vanishes fast with $d(\mathbf{x})$ increases. When the SDF for a surface with a thin object converges

to the status as shown in the left of Fig. 3, the gradients $\frac{\partial \sigma}{\partial d}$ for points near the object off the surface are almost zero and therefore $\frac{\partial \mathcal{L}}{\partial \theta}$ remains near zero (see the gradient fields $\frac{\partial \sigma}{\partial d}$ Fig. 3); in other words, no gradients can be acquired to recover the thin object. To address the problem, we propose to insert “bubbles” for the missing surface points to create gradients for SDF near small or thin objects. The middle figure of Fig. 3 exemplifies inserted bubbles. The inserted bubble creates many local surface islands and gradient fields near them, which then allows the growing up of new objects that are previously ignored.

Specifically, we obtain the bubbles from depth images. Given a depth image $D(u, v)$ with corresponding camera pose $[\mathbf{R}|\mathbf{t}]$ and intrinsics \mathbf{K} , the 3D point $\mathbf{x}(\mathbf{p})$ associated with a pixel $\mathbf{p} = (u, v)$ is

$$\mathbf{x}(\mathbf{p}, D) = \mathbf{t} + D(u, v) \left(\mathbf{R} \mathbf{K}^{-1} \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} \right). \quad (7)$$

The 3D points from depth images, *i.e.* the bubbles, can be regarded as an approximation of surface points during ray casting and its SDF value should be zero. In order to enforce precise surface reconstructions, we define bubble loss $\mathcal{L}_{\text{bubble}}$ to minimize the absolute SDF value of these surface points:

$$\mathcal{L}_{\text{bubble}} = \sum_{\mathbf{p} \in \mathcal{P}} |d(\mathbf{x}(\mathbf{p}, D))|, \quad (8)$$

where $d(\mathbf{x})$ is the predicted SDF value of a point \mathbf{x} (Eq. (2)), D is the depth image and \mathcal{P} denotes the minibatch of sampled pixels in each iteration.

4.1.2 Error-Guided Adaptive Sampling Strategy

Applying the bubble loss on the implicit SDF field can improve the 3D reconstruction quality. However, the geometry of large planar areas can already be reconstructed very well with the image-space depth loss. Applying $\mathcal{L}_{\text{bubble}}$ on 3D points within these areas is a waste of computation. Furthermore, as analyzed in Sec. 1, small objects (which is our target) make up only a small percentage of all pixels, so it will be favorable if the bubble loss is applied more frequently in these areas.

Determining the areas of “small objects” is a perceptual task, which cannot be easily accomplished without semantic segmentation or manual marking. However, due to the feature of neural networks that low-frequency signals are always easier to fit than high-frequency signal [33], we can naturally filter out the uninterested large planar areas (low frequency) and preserve small-object areas (high frequency) according to the reconstruction error of the neural network. Specifically, given an error metric E (depth loss is adopted in our case), we leverage *importance sampling* algorithm according to a probability distribution determined by E . The PDF (Probability Density Function) of a point $\mathbf{x}(\mathbf{p}; D)$ (Eq. (7)) is proportional to the reconstruction error $E(\mathbf{p})$ at the pixel \mathbf{p} . We also prune the pixels with errors lower than a particular threshold P_{\min} , indicating that those pixels are not of our interest. In this way, the network can pay more attention to the erroneous areas and converge faster. Please refer to our supplementary material for visualizations of PDF map.

The PDF values are updated dynamically: We maintain a PDF map for each training image. After each training iteration, the PDF values for pixels in the current batch are updated with their error metrics E .

Geometry loss. Together with the bubble loss, our geometry loss is as follows to approximate the geometry field:

$$\mathcal{L}_{\text{geo}} = \lambda_1 \mathcal{L}_{\text{eikonal}} + \lambda_2 \mathcal{L}_{\text{depth}} + \lambda_3 \mathcal{L}_{\text{normal}} + \lambda_4 \mathcal{L}_{\text{smooth}} + \lambda_5 \mathcal{L}_{\text{bubble}}. \quad (9)$$

As suggested by previous work [41], we apply Eikonal term [7] to regularize SDF values

$$\mathcal{L}_{\text{eikonal}} = \sum_{\mathbf{x} \in \mathcal{X}} (\|\nabla d(\mathbf{x})\|_2 - 1)^2, \quad (10)$$

where \mathcal{X} is the minibatch of 3D points uniformly sampled in 3D space and nearby surface.

We use depth and normal priors to supervise our network to handle shape-radiance ambiguity. We compute the surface depth and normal by Eq. (6) and use L_1 loss for depth

and angular L_1 loss for normal:

$$\mathcal{L}_{\text{depth}} = \sum_{\mathbf{r} \in \mathcal{R}} \|\hat{D}(\mathbf{r}) - D(\mathbf{r})\|_1, \quad (11)$$

$$\mathcal{L}_{\text{normal}} = \sum_{\mathbf{r} \in \mathcal{R}} \|1 - \hat{N}(\mathbf{r}) \cdot N(\mathbf{r})\|_1. \quad (12)$$

We also use smoothness loss on the gradient of the SDF field suggested by UNISURF [23] to encourage smooth surface reconstruction:

$$\mathcal{L}_{\text{smooth}} = \sum_{\mathbf{x} \in \mathcal{S}} \|\nabla d(\mathbf{x}) - \nabla d(\mathbf{x} + \epsilon)\|_2, \quad (13)$$

where \mathcal{S} is the minibatch of points sampled near the surface.

4.2. Emitter Semantic Field

Considering that our radiance field F_c is trained from LDR images, the light intensity from emitters (e.g. lamps/windows) is usually under-estimated. This will cause an over-dark estimation in the re-render stage (Sec. 5). We resolve this by introducing semantic labels of emitters to optimize the radiance value emitted from light sources. We add a neural emitter semantic field F_e into our module, which determines whether the input 3D point \mathbf{x} is on an emitter or not. The estimated emitter mask $\hat{M}_e(\mathbf{r})$ for a ray can also be evaluated by volume accumulation

$$\hat{m}(\mathbf{x}) = F_e(\mathbf{z}(\mathbf{x})), \quad (14)$$

$$\hat{M}_e(\mathbf{r}) = \sum_{i=1}^M T_{\mathbf{r}}^i \alpha_{\mathbf{r}}^i \hat{m}_{\mathbf{r}}^i, \quad (15)$$

where $\mathbf{z}(\mathbf{x})$ is the latent code output by SDF field F_d (Eq. (2)). We choose $\mathbf{z}(\mathbf{x})$ as the MLP input because it can provide latent information about the scene.

Emitter segmentation loss. Given ground truth emitter masks M_e , we optimize F_e by a binary cross-entropy loss:

$$\mathcal{L}_{\text{emi}} = \sum_{\mathbf{r} \in \mathcal{R}} M_e(\mathbf{r}) \log \hat{M}_e(\mathbf{r}) + (1 - M_e(\mathbf{r})) \log (1 - \hat{M}_e(\mathbf{r})). \quad (16)$$

After F_e is trained, it can be used by our raytracing stage to indicate if a ray hits an emitter. We use K-Means algorithm [10] to cluster emitter points as K emitters. To model HDR emissions, we define an array $\mathbf{L}[\cdot]$ with size K as a learnable parameter that corresponds to the emission values of each emitter. $\mathbf{L}[\cdot]$ will be queried in the raytracing and re-rendering stage, which will be described in detail in Sec. 5.

4.3. Material Field

Similar to the geometry and radiance field, we parameterize the spatially-varying material of the scene as a neural field. We use physically-based GGX microfacet BRDF

model [34] to present scene material and introduce two MLPs to model the albedo and roughness of the scene, respectively:

$$K_d(\mathbf{x}), K_s(\mathbf{x}) = F_a(\mathbf{z}(\mathbf{x})), \quad (17)$$

$$\rho(\mathbf{x}) = F_\rho(\mathbf{z}(\mathbf{x})), \quad (18)$$

where $K_d(\mathbf{x})$ and $K_s(\mathbf{x})$ are the diffuse and specular albedo at 3D point \mathbf{x} , $\rho(\mathbf{x})$ is the surface roughness at \mathbf{x} . Note that the SDF field F_d and radiance field F_c have been pretrained and fixed. The estimated material parameter associated with a ray \mathbf{r} can be calculated by volumetric accumulation similar to Eq. (5) and Eq. (6).

Material regularizations. To enforce physical correctness for predicted material parameters, we define regularizations as

$$\mathcal{L}_{\text{mreg}} = \sum_{\mathbf{x} \in \mathcal{S}} \|\hat{M}(\mathbf{x}) - \hat{M}(\mathbf{x} + \epsilon)\|_2 \quad (19)$$

$$+ \sum_{\mathbf{x} \in \mathcal{S}} \left(\hat{K}_d(\mathbf{x}) + \hat{K}_s(\mathbf{x}) - 1 \right)_+, \quad (20)$$

where $\hat{M} \in \{\hat{K}_d, \hat{K}_s, \hat{\rho}\}$ denotes 3 material parameters. Similar to Eq. (13), we encourage smooth estimation of materials in the first loss term. According to energy conservation law, the sum of diffuse and specular albedo should not exceed 1 (the second loss term), where $(\cdot)_+$ is the ReLU function [21].

5. Differentiable Monte Carlo Raytracing

Given material, geometry, and lighting components, the scene appearance can be re-rendered using surface rendering algorithms. Inspired by [48], we use differentiable Monte Carlo rendering algorithm with raytracing to recover scene appearance from shape, material and illumination. The difference between our work and [48] is that their raytracing is performed on screen space while ours is in a 3D volumetric space. For a ray $\mathbf{r} : \mathbf{x} = \mathbf{o} + t\mathbf{v}_s$, we can trace and intersect it with the neural SDF field. The intersection point \mathbf{s} can be estimated by

$$\mathbf{s} = \text{trace}(\mathbf{r}) = \mathbf{o} + \left(\sum_{i=1}^M T_{\mathbf{r}}^i \alpha_{\mathbf{r}}^i t_{\mathbf{r}}^i \right) \mathbf{v}_s. \quad (21)$$

We leverage Monte Carlo rendering technique to perform the scene re-rendering. We first cast the rays from camera view to obtain the surface points associated with each pixel by Eq. (21), as well as their corresponding surface normal by Eq. (6). Then, given a sample rate N , we use GGX importance sampling to generate N outgoing rays $\{\mathbf{r}_s^k : \mathbf{x} = \mathbf{s} + t\mathbf{d}_s^k\}_{k=1}^N$ starting from a surface point \mathbf{s} according to the surface normal and material parameters

$\hat{N}(\mathbf{s}), \hat{K}_d(\mathbf{s}), \hat{K}_s(\mathbf{s}), \hat{\rho}(\mathbf{s})$. The surface color can be rendered by Monte Carlo integration:

$$\hat{\mathbf{R}}(\mathbf{s}) = \frac{1}{N} \sum_{k=1}^N \frac{f_r(\mathbf{v}_s, \mathbf{d}_s^k; \hat{N}, \hat{K}_d, \hat{K}_s, \hat{\rho}) L_s^k \cos \theta_k}{p(\mathbf{v}_s, \mathbf{d}_s^k)}, \quad (22)$$

where f_r and p is the evaluation and PDF value of GGX microfacet BRDF model determined by the material parameters. L_s^k is the predicted radiance of ray \mathbf{r}_s^k :

$$L_s^k = \begin{cases} \hat{\mathbf{C}}(\mathbf{r}_s^k) & \text{if not } \hat{M}_e(\mathbf{r}_s^k), \\ \mathbf{L}[\text{index}(\text{trace}(\mathbf{r}_s^k))] & \text{if } \hat{M}_e(\mathbf{r}_s^k), \end{cases} \quad (23)$$

which can be divided into two cases: we use F_e to determine if \mathbf{r}_s^k hits an emitter. If so, we obtain the emitter index by K-means and retrieve its emission from emission field $\mathbf{L}[\cdot]$ (defined in Sec. 4.2). Otherwise, we use F_c and volume rendering (Eq. (5)) to predict the radiance of the ray.

6. Training

The training of intrinsic decomposition and the reconstruction of the indoor scenes are conducted in two stages: 1) the training of geometry and radiance fields, 2) the training of the material and emission fields.

Training of geometry and radiance fields. The training scheme of the SDF network F_d and radiance network F_c is end-to-end. We firstly use geometric initialization [1] to initialize F_d , and then optimize the networks with following loss:

$$\mathcal{L}_1 = \mathcal{L}_{\text{rgb}} + \lambda_{\text{geo}} \mathcal{L}_{\text{geo}} + \lambda_{\text{emi}} \mathcal{L}_{\text{emi}}, \quad (24)$$

where $\mathcal{L}_{\text{rgb}} = \sum_{\mathbf{r} \in \mathcal{R}} \|\hat{\mathbf{C}}(\mathbf{r}) - \mathbf{C}(\mathbf{r})\|_1$. Eq. (5) provides volume rendering results from 3D neural representation to 2D images. \mathcal{R} denotes the set of pixels/rays sampled in the minibatch and $\mathbf{C}(\mathbf{r})$ is the ground truth pixel color.

The weight hyperparameters of some losses vary during our 3-step training:

1. (Warm-up step) In the early stage of training, bubble loss is not applied (*i.e.* $\lambda_5 = 0$) until the network can reconstruct a coarse scene geometry. At the end of this step, the error (PDF) maps are initialized per image.
2. (Bubble step) The adaptive sampling and bubble loss are enabled in this step to reconstruct missing small structures. Note that the bubble loss breaks the stable status of the converged SDF field so far, which makes the Eikonal and smoothness regularization increase. Therefore, we disable them (*i.e.* $\lambda_1, \lambda_4 = 0$) to prevent potential contradictions.

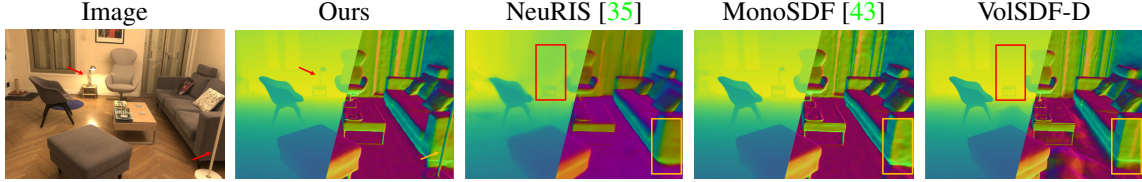


Figure 4. Qualitative comparisons of reconstructed depth map and normal map.

Table 1. Quantitative comparisons of novel view synthesis results on synthetic data and real data. Data in brackets denote metrics in training views.

Method	Synthetic Data			Real Data		
	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
NeuRIS [35]	25.02 (26.01)	0.77 (0.77)	0.46 (0.42)	22.33 (24.45)	0.68 (0.71)	0.53 (0.50)
MonoSDF [43]	25.74 (26.73)	0.79 (0.78)	0.38 (0.43)	22.38 (24.71)	0.69 (0.73)	0.49 (0.45)
NeRF [18]	26.81 (27.74)	0.85 (0.85)	0.18 (0.22)	24.06 (25.58)	0.77 (0.78)	0.27 (0.27)
Instant-NGP [19]	23.89 (31.59)	0.78 (0.90)	0.25 (0.12)	22.95 (28.51)	0.76 (0.88)	0.20 (0.09)
Ours	28.42 (29.70)	0.87 (0.87)	0.15 (0.19)	24.66 (26.33)	0.78 (0.79)	0.25 (0.26)

3. (Smooth step) Since the bubble loss will affect the smoothness and stability of the SDF field, which will have negative effects such as imprecision in normal estimation. Therefore, the bubble loss is disabled (*i.e.* $\lambda_5 = 0$) in this step. To restore the SDF field smoothness, the training continues with the Eikonal and smoothness loss enabled again.

Training of material and emission fields. We jointly train the material network F_a , F_ρ and the emission array $\mathbf{L}[\cdot]$ after the geometry and radiance networks (F_d , F_c , F_e) have been pretrained. During the optimization of intrinsic decomposition, the parameters of the geometry and radiance networks are fixed. Since the ground truths of material are impossible to capture from images, we weakly supervise the network by re-render results instead of direct supervision from strong material priors. The overall loss for this stage is

$$\mathcal{L}_2 = \mathcal{L}_{\text{render}} + \lambda_{\text{mat}} \mathcal{L}_{\text{mat}}, \quad (25)$$

where $\mathcal{L}_{\text{render}} = \sum_{\mathbf{r} \in \mathcal{R}} \|\hat{\mathbf{R}}(\mathbf{s}(\mathbf{r})) - \mathbf{C}(\mathbf{r})\|_1$, minimizing the L_1 error between the re-rendered result (Eq. (22)) and input image.

7. Experiments

We first analyze and compare our method with the state-of-the-art methods in terms of geometry reconstruction and novel view synthesis. We then demonstrate qualitative scene editing and relighting results based on our intrinsic decomposition method. Finally, we also perform ablation studies to prove the effectiveness of our design.

Datasets. We propose a new synthetic multi-view indoor scene dataset, which includes well-designed scenes

by artists and provides high quality rendered images with ground truth camera poses and geometry annotations (depth and normal maps). Existing datasets (such as ScanNet [5]) suffers from inaccurate camera calibration, erroneous depth capture and low image quality (such as motion blur), which will crucially affect the reconstruction quality. Our dataset provides well-designed indoor scenes with ground truth camera poses, normal and depth maps, with superior image quality to existing datasets. We also test our method on some real-world scenes, specifically, a living room scene from [26] and 3 scenes from Scalable-NISR [40]. All the real data contains calibrated camera poses and depth maps. To train the emitter field, we annotate emitter masks for the real data.

Metrics. For the 3D geometry reconstruction, following previous works [43], we compare on mesh-based metrics including accuracy, precision, recall and F-score. We also present image-space geometry errors including depth error and normal angular error in the supplementary. For the novel view synthesis, we use widely-used image metrics including PSNR, SSIM [37] and LPIPS [45].

Baselines. We compare against state-of-the-art neural reconstruction, multi-view stereo and novel view synthesis methods. In particular, for the geometry reconstruction, VolSDF [41], NeuRIS [35] and MonoSDF [43] are compared. Since the original VolSDF method suffers from shape-radiance ambiguity and usually fails in reconstructing plausible scene structure, we add an additional depth loss when optimizing it. We mark it as “VolSDF-D” in the following figures and tables. For novel view synthesis, we also compare with NeRF [18] and Instant-NGP [19].

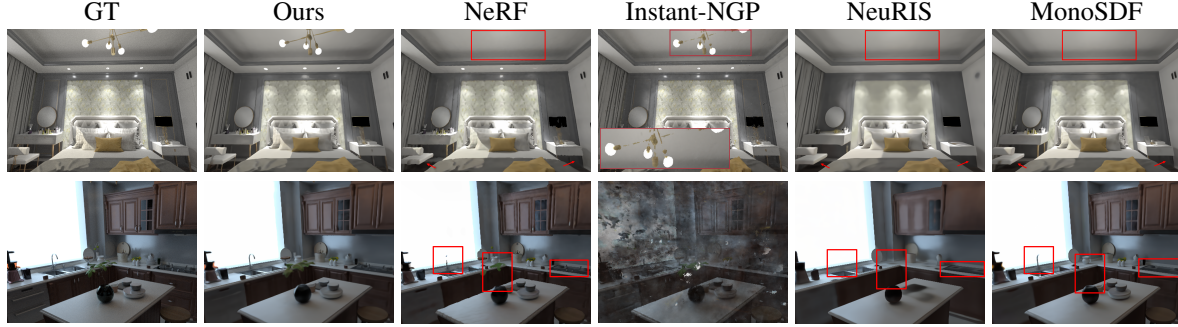


Figure 5. **Qualitative comparisons of novel view synthesis.** More results are presented in supplementary material.

Table 2. **Quantitative comparisons of geometric reconstruction results** on synthetic data.

Method	Ours	VolSDF-D	NeuRIS	MonoSDF
Acc.↓	0.035	0.041	0.036	0.052
Prec.↑	0.87	0.76	0.74	0.82
Recall↑	0.79	0.64	0.65	0.72
F-Score↑	0.83	0.68	0.66	0.77

7.1. Comparisons with state-of-the-art methods

Geometry reconstruction. We evaluate 3D geometry metrics on our synthetic dataset with ground truth meshes. Tab. 2 shows quantitative results on mesh evaluation. Our method outperforms all baselines due to the precise reconstructions on small objects. Fig. 4 shows qualitative results on the reconstructed depth and normal maps. We observe significant reconstruction losses for small objects in the baselines. In contrast, our method can faithfully reconstruct these small objects on account of the usage of bubble loss and adaptive sampling. Please refer to our supplementary material for per-scene detailed results and quantitative results in depth and normal errors.

Novel view synthesis. We evaluate the quality of novel view synthesis on both synthetic and real data. Tab. 1 and Fig. 5 give quantitative results and qualitative results. NeRF and SDF-baselines struggle in recognizing small objects, leading to poor results. Instant-NGP performs best in training views. However, this is achieved by color over-fitting instead of accurate geometry understanding, leading to poor quality in test views. Benefiting from our high-quality geometry, our method provides good results in novel views. We also provide view interpolation results, which will be displayed in the supplementary video.

7.2. Scene Editing

With the decomposition results of shape, material and lighting, we can enable photo-realistic scene editing tasks such as material editing and relighting, as shown in Figs. 1 and 6. In Fig. 6, we change the hue of the light in the room

(1st column), increase the emission intensity of the lamp (2nd column), and change the material of the closet door into a mirror (3rd column), respectively. Note that the reflections in the specular mirror is consistent to the surroundings. With our physically-based rendering algorithm, our method can produce photo-realistic lighting effects such as specular reflections. More results and videos are presented in the supplementary.



Figure 6. **Qualitative results in material editing and relighting.**

7.3. Ablation Studies

Robustness on inaccurate depth information. In real-world scenarios, the captured depth always contains noises and errors. To simulate depth noise, all rendered depth images are added with noise scaling approximately quadratically with depth z . As suggested by [14, 22], the noise model is $\epsilon = N(\mu(z), \sigma(z))$, where $\mu(z) = 0.0001125z^2 + 0.0048875$, $\sigma(z) = 0.002925z^2 + 0.003325$. We use the noisy depth to supervise our bubble and depth loss. Fig. 7 shows that noisy depths produce negligible impacts on the prediction, demonstrating the robustness of our method.

Effectiveness of adaptive sampling strategy. We ablate between the sampling strategies on bubble points sampling.

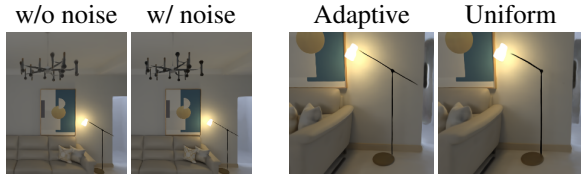


Figure 7. Noisy depth.

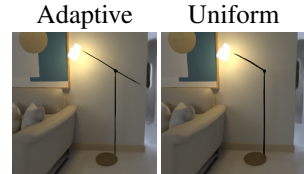


Figure 8. Sampling strategy.

We compare our error-guided adaptive sampling with uniform sampling. In Fig. 8, uniform sampling is unable to reconstruct the complete lamp pole. This is because insufficient bubble points are sampled from the missing pole.

8. Conclusion and Limitations

This work proposes I²-SDF that reconstructs an intrinsic neural scene from multi-view images, enabling physically-realistic novel view synthesis of editable indoor scenes. With the novel bubbling strategy, we are able to recover the small objects in large-scale scenes and obtain SOTA geometry and novel view synthesis results. This work’s limitations are: Firstly, the MLP-based network backbone is not powerful enough to capture high-frequency textures. Secondly, the time-consuming MC raytracing increases the total reconstruction time. We consider them as problems to be solved in the future.

References

- [1] Matan Atzmon and Yaron Lipman. Sal: Sign agnostic learning of shapes from raw data. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020. 6
- [2] Mark Boss, Raphael Braun, Varun Jampani, Jonathan T. Barron, Ce Liu, and Hendrik P.A. Lensch. Nerd: Neural reflectance decomposition from image collections. In *IEEE International Conference on Computer Vision (ICCV)*, 2021. 1, 3
- [3] Mark Boss, Varun Jampani, Raphael Braun, Ce Liu, Jonathan T. Barron, and Hendrik P.A. Lensch. Neural-pil: Neural pre-integrated lighting for reflectance decomposition. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2021. 3
- [4] Anpei Chen, Zexiang Xu, Fuqiang Zhao, Xiaoshuai Zhang, Fanbo Xiang, Jingyi Yu, and Hao Su. Mvnerf: Fast generalizable radiance field reconstruction from multi-view stereo. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14124–14133, 2021. 2
- [5] Angela Dai, Angel X. Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *Proc. Computer Vision and Pattern Recognition (CVPR)*, IEEE, 2017. 7
- [6] Mathieu Garon, Kalyan Sunkavalli, Sunil Hadap, Nathan Carr, and Jean-François Lalonde. Fast spatially-varying indoor lighting estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6908–6917, 2019. 2
- [7] Amos Gropp, Lior Yariv, Niv Haim, Matan Atzmon, and Yaron Lipman. Implicit geometric regularization for learning shapes. In *International Conference on Machine Learning*, pages 3789–3799. PMLR, 2020. 5
- [8] Haoyu Guo, Sida Peng, Haotong Lin, Qianqian Wang, Guofeng Zhang, Hujun Bao, and Xiaowei Zhou. Neural 3d scene reconstruction with the manhattan-world assumption. In *CVPR*, 2022. 2
- [9] Jie Guo, Shuichang Lai, Chengzhi Tao, Yuelong Cai, Lei Wang, Yanwen Guo, and Ling-Qi Yan. Highlight-aware two-stream network for single-image svbrdf acquisition. *ACM Transactions on Graphics (TOG)*, 40(4):1–14, 2021. 2
- [10] John A Hartigan and Manchek A Wong. Algorithm as 136: A k-means clustering algorithm. *Journal of the royal statistical society. series c (applied statistics)*, 28(1):100–108, 1979. 5
- [11] Jon Hasselgren, Nikolai Hofmann, and Jacob Munkberg. Shape, Light, and Material Decomposition from Images using Monte Carlo Rendering and Denoising. *arXiv:2206.03380*, 2022. 3
- [12] James T Kajiya. The rendering equation. In *Proceedings of the 13th annual conference on Computer graphics and interactive techniques*, pages 143–150, 1986. 11
- [13] Brian Karis and Epic Games. Real shading in unreal engine 4. *Proc. Physically Based Shading Theory Practice*, 4(3):1, 2013. 11
- [14] Leonid Keselman, John Iselin Woodfill, Anders Grunnet-Jepsen, and Achintya Bhowmik. Intel realsense stereoscopic depth cameras. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, pages 1–10, 2017. 8
- [15] 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 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2475–2484, 2020. 2
- [16] Zhengqin Li, Jia Shi, Sai Bi, Rui Zhu, Kalyan Sunkavalli, Miloš Hašan, Zexiang Xu, Ravi Ramamoorthi, and Manmohan Chandraker. Physically-based editing of indoor scene lighting from a single image. *arXiv preprint arXiv:2205.09343*, 2022. 2
- [17] Zhengqin Li, Zexiang Xu, Ravi Ramamoorthi, Kalyan Sunkavalli, and Manmohan Chandraker. Learning to reconstruct shape and spatially-varying reflectance from a single image. In *SIGGRAPH Asia 2018 Technical Papers*, page 269. ACM, 2018. 2
- [18] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *ECCV*, 2020. 1, 2, 3, 7, 11
- [19] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM Trans. Graph.*, 41(4):102:1–102:15, July 2022. 2, 7
- [20] Jacob Munkberg, Jon Hasselgren, Tianchang Shen, Jun Gao, Wenzheng Chen, Alex Evans, Thomas Müller, and Sanja Fidler. Extracting triangular 3d models, materials, and lighting from images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8280–8290, 2022. 1, 3
- [21] Vinod Nair and Geoffrey E Hinton. Rectified linear units improve restricted boltzmann machines. In *Icml*, 2010. 6

- [22] Chuong V Nguyen, Shahram Izadi, and David Lovell. Modeling kinect sensor noise for improved 3d reconstruction and tracking. In *2012 second international conference on 3D imaging, modeling, processing, visualization & transmission*, pages 524–530. IEEE, 2012. 8
- [23] Michael Oechsle, Songyou Peng, and Andreas Geiger. Unisurf: Unifying neural implicit surfaces and radiance fields for multi-view reconstruction. In *International Conference on Computer Vision (ICCV)*, 2021. 2, 5
- [24] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. Deepsdf: Learning continuous signed distance functions for shape representation. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019. 2
- [25] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32, 2019. 11
- [26] Julien Philip, Sébastien Morgenthaler, Michaël Gharbi, and George Drettakis. Free-viewpoint indoor neural relighting from multi-view stereo. *ACM Transactions on Graphics*, 2021. 7
- [27] Sara Fridovich-Keil and Alex Yu, Matthew Tancik, Qinhong Chen, Benjamin Recht, and Angjoo Kanazawa. Plenoxels: Radiance fields without neural networks. In *CVPR*, 2022. 2
- [28] Johannes Lutz Schönberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 2
- [29] Johannes Lutz Schönberger, Enliang Zheng, Marc Pollefeys, and Jan-Michael Frahm. Pixelwise view selection for unstructured multi-view stereo. In *European Conference on Computer Vision (ECCV)*, 2016. 2
- [30] Pratul P. Srinivasan, Boyang Deng, Xiuming Zhang, Matthew Tancik, Ben Mildenhall, and Jonathan T. Barron. Nerv: Neural reflectance and visibility fields for relighting and view synthesis. In *CVPR*, 2021. 2
- [31] Cheng Sun, Min Sun, and Hwann-Tzong Chen. Direct voxel grid optimization: Super-fast convergence for radiance fields reconstruction. In *CVPR*, 2022. 2
- [32] Jiaming Sun, Yiming Xie, Linghao Chen, Xiaowei Zhou, and Hujun Bao. NeuralRecon: Real-time coherent 3D reconstruction from monocular video. In *CVPR*, 2021. 2
- [33] Matthew Tancik, Pratul P. Srinivasan, Ben Mildenhall, Sara Fridovich-Keil, Nithin Raghavan, Utkarsh Singhal, Ravi Ramamoorthi, Jonathan T. Barron, and Ren Ng. Fourier features let networks learn high frequency functions in low dimensional domains. *NeurIPS*, 2020. 4, 5
- [34] Bruce Walter, Stephen R Marschner, Hongsong Li, and Kenneth E Torrance. Microfacet models for refraction through rough surfaces. *Rendering techniques*, 2007:18th, 2007. 6, 11
- [35] Jiepeng Wang, Peng Wang, Xiaoxiao Long, Christian Theobalt, Taku Komura, Lingjie Liu, and Wenping Wang. Neuris: Neural reconstruction of indoor scenes using normal priors. *arXiv preprint arXiv:2211.03017*, 2022. 2, 3, 7
- [36] Peng Wang, Lingjie Liu, Yuan Liu, Christian Theobalt, Taku Komura, and Wenping Wang. Neus: Learning neural implicit surfaces by volume rendering for multi-view reconstruction. *NeurIPS*, 2021. 2
- [37] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612, 2004. 7
- [38] Zian Wang, Jonah Philion, Sanja Fidler, and Jan Kautz. Learning indoor inverse rendering with 3d spatially-varying lighting. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 12538–12547, 2021. 2
- [39] Yi Wei, Shaohui Liu, Yongming Rao, Wang Zhao, Jiwen Lu, and Jie Zhou. Nerfingmvs: Guided optimization of neural radiance fields for indoor multi-view stereo. In *ICCV*, 2021. 2
- [40] Xiuchao Wu, Jiamin Xu, Zihan Zhu, Hujun Bao, Qixing Huang, James Tompkin, and Weiwei Xu. Scalable neural indoor scene rendering. *ACM Transactions on Graphics (TOG)*, 2022. 7
- [41] Lior Yariv, Jiatao Gu, Yoni Kasten, and Yaron Lipman. Volume rendering of neural implicit surfaces. *Advances in Neural Information Processing Systems*, 34:4805–4815, 2021. 2, 3, 5, 7, 11
- [42] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelNeRF: Neural radiance fields from one or few images. In *CVPR*, 2021. 2
- [43] Zehao Yu, Songyou Peng, Michael Niemeyer, Torsten Sattler, and Andreas Geiger. Monosdf: Exploring monocular geometric cues for neural implicit surface reconstruction. *arXiv:2022.00665*, 2022. 1, 2, 3, 7
- [44] Kai Zhang, Fujun Luan, Qianqian Wang, Kavita Bala, and Noah Snavely. PhysSG: Inverse rendering with spherical gaussians for physics-based material editing and relighting. In *The IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 1, 3
- [45] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*, 2018. 7
- [46] Xiaoshuai Zhang, Sai Bi, Kalyan Sunkavalli, Hao Su, and Zexiang Xu. Nerfusion: Fusing radiance fields for large-scale scene reconstruction. In *CVPR*, 2022. 2
- [47] Xiuming Zhang, Pratul P Srinivasan, Boyang Deng, Paul Debevec, William T Freeman, and Jonathan T Barron. Nerfactor: Neural factorization of shape and reflectance under an unknown illumination. *ACM Transactions on Graphics (TOG)*, 40(6):1–18, 2021. 1, 3
- [48] Jingsen Zhu, Fujun Luan, Yuchi Huo, Zihao Lin, Zhihua Zhong, Dianbing Xi, Jiaxiang Zheng, Rui Tang, Hujun Bao, and Rui Wang. Learning-based inverse rendering of complex indoor scenes with differentiable monte carlo raytracing. *arXiv preprint arXiv:2211.03017*, 2022. 2, 6

A. Network Details

A.1. Architecture

All neural fields in our network are implemented by multi-layer perceptron (MLP). For neural SDF field F_d and radiance field F_c , we follow VolSDF [41]’s default setting, where F_d is a 8-layer MLP with hidden dimension 256 and F_c is a 4-layer MLP with hidden dimension 256. The dimension of the latent code output by F_d (*i.e.* $\mathbf{z}(\mathbf{x})$) is 256, and a skip connection is used in the 4th layer in F_d . The input position \mathbf{x} and view direction \mathbf{v} are encoded by positional encoding, the same as in NeRF [18]. The emitter semantic field F_e is a 2-layer MLP with hidden dimension 128, while material fields F_ρ, F_a are 3-layer MLPs with hidden dimension 64.

A.2. Implementation Details

The network model, as well as the training and evaluation scripts, are implemented with Pytorch [25]. The network is trained per-scene on a single NVIDIA Tesla V100 GPU. We adopt two stage training scheme, the training details of the 2 stages are as follows:

Training of geometry and radiance fields. We jointly optimize SDF network F_d , radiance network F_c and emitter semantic network F_e in this stage. We optimize our model for 200k iterations in this stage, which takes about 15 hours for a scene. The training loss

$$\mathcal{L}_{\text{geo}} = \lambda_1 \mathcal{L}_{\text{eikonal}} + \lambda_2 \mathcal{L}_{\text{depth}} + \lambda_3 \mathcal{L}_{\text{normal}} + \lambda_4 \mathcal{L}_{\text{smooth}} + \lambda_5 \mathcal{L}_{\text{bubble}} \quad (26)$$

$$\mathcal{L}_1 = \mathcal{L}_{\text{rgb}} + \lambda_{\text{geo}} \mathcal{L}_{\text{geo}} + \lambda_{\text{emi}} \mathcal{L}_{\text{emi}} \quad (27)$$

where the weight hyperparameters are $\lambda_{\text{geo}} = 1$, $\lambda_{\text{emi}} = 0.5$, and $\lambda_2 = 0.1$, $\lambda_3 = 0.05$, respectively. For $\lambda_1, \lambda_4, \lambda_5$, since the training process is further divided into 3 steps (*i.e.* warm-up, bubble and smooth), their values are adjusted during the training accordingly:

1. In warm-up step, $\lambda_5 = 0$, $\lambda_1 = 0.1$, $\lambda_4 = 0$.
2. In bubble step, $\lambda_1 = \lambda_4 = 0$, $\lambda_5 = 0.5$.
3. In smooth step, $\lambda_5 = 0$, $\lambda_1 = 0.1$, $\lambda_4 = 0.01$.

The number of iterations assigned to the three steps are 50k, 100k and 50k in sequence.

We use error-guided adaptive sampling in bubble step, where the pruning threshold $P_{\min} = 0.05$.

Training of material and emission fields. In this stage, we use importance sampling and Monte Carlo estimation to compute the rendering result. We generate N outgoing rays to perform Monte Carlo integration. In practice, the sample

rate N is set to 16, which is a trade-off between quality and performance.

We jointly optimize $F_a, F_\rho, \mathbf{L}[\cdot]$ for 100k iterations. The bottleneck of computational cost lies in the prediction of incident radiance L_s^k , which grows proportional to the sample rate N . With $N = 16$, the training lasts for about 2-3 days.

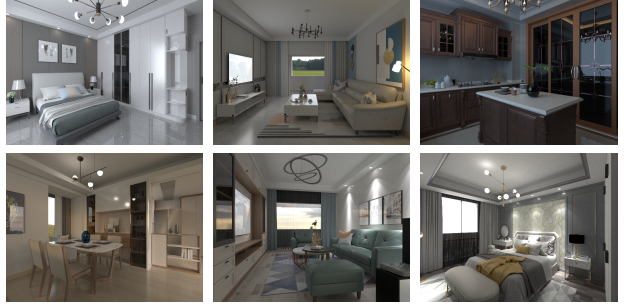


Figure 9. Display of indoor scenes in our dataset.

B. Dataset Details

Our synthetic dataset contains 12 scenes in total: 6 bedrooms, 2 living rooms, 2 dining rooms and 2 kitchens. All the scenes are well-designed by artists with detailed geometry and fine textures. We use GPU-accelerated path tracing algorithm [12] to render the images, which can create photo-realistic rendering results with global illumination. All data are rendered on a NVIDIA RTX 3090 GPU, with 4096 samples per pixel (spp). The rendering time is roughly 15 seconds per image. Fig. 9 displays some of the indoor scenes in our dataset.

All images in our dataset are annotated by ground truth camera intrinsics and poses, normal maps, depth maps and emitter semantic masks.

C. BRDF Model

We use GGX microfacet BRDF model [34] to approximate the surface reflection properties by a set of material parameters, including diffuse albedo K_d , specular albedo K_s and roughness ρ . In our implementation, we refer to Unreal Engine [13]’s implementation of microfacet BRDF model. The BRDF (bidirectional reflectance distribution function) $f_r(\mathbf{v}, \mathbf{d}; N, K_d, K_s, \rho)$ (where \mathbf{v} and \mathbf{d} are view and lighting directions, and N is the surface normal) can be decomposed into diffuse and specular components and computed by

$$f_r(\mathbf{v}, \mathbf{d}) = f_d(K_d) + f_s(\mathbf{v}, \mathbf{d}; N, K_s, \rho) \quad (28)$$

$$f_d(K_d) = \frac{K_d}{\pi} \quad \alpha = \rho^2 \quad (29)$$

$$f_s(\mathbf{v}, \mathbf{d}; N, K_s, \rho) = D(\alpha, N, h) G_2(\alpha, N, \mathbf{v}, \mathbf{d}) F(K_s, \mathbf{d}, h) \quad (30)$$

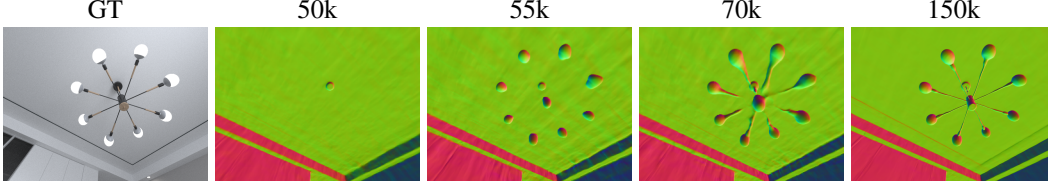


Figure 10. **Intermediate results in bubble step.**

where f_d and f_s are the diffuse and specular components respectively, while D, G_2, F are the distribution, Fresnel and geometric terms, defined as

$$D(\alpha, \mathbf{N}, \mathbf{h}) = \frac{\alpha^2}{\pi((\alpha^2 - 1)(\mathbf{N} \cdot \mathbf{h})^2 + 1)^2} \quad (31)$$

$$S(\alpha, \mathbf{N}, \mathbf{v}, \mathbf{d}) = (\mathbf{N} \cdot \mathbf{d}) \sqrt{\alpha^2 + (\mathbf{N} \cdot \mathbf{v})^2(1 - \alpha^2)} \quad (32)$$

$$G_2(\alpha, \mathbf{N}, \mathbf{v}, \mathbf{d}) = \frac{1}{2(S(\alpha, \mathbf{N}, \mathbf{v}, \mathbf{d}) + S(\alpha, \mathbf{N}, \mathbf{d}, \mathbf{v}))} \quad (33)$$

$$\text{lum}(C) = 0.213C.r + 0.715C.g + 0.072C.b \quad (34)$$

$$F_{90}(K_s) = \min\left(\frac{\text{lum}(K_s)}{0.04}, 1\right) \quad (35)$$

$$F(K_s, \mathbf{N}, \mathbf{d}) = K_s + (F_{90}(K_s) - K_s)(1 - (\mathbf{N} \cdot \mathbf{d}))^5 \quad (36)$$

In importance sampling and Monte Carlo integration, we also need to calculate the PDF value $p(\mathbf{v}, \mathbf{d})$ corresponding to the view and lighting direction:

$$w_d = \frac{\text{lum}(K_d)}{\text{lum}(K_d) + \text{lum}(K_s)} \quad (37)$$

$$p(\mathbf{v}, \mathbf{d}) = w_d p_d(\mathbf{v}, \mathbf{d}) + (1 - w_d) p_s(\mathbf{v}, \mathbf{d}) \quad (38)$$

$$p_d(\mathbf{v}, \mathbf{d}) = \frac{\mathbf{N} \cdot \mathbf{d}}{\pi} \quad (39)$$

$$G_1(\alpha, N, \mathbf{v}) = \frac{2}{\sqrt{1 + \frac{\alpha^2(1 - (\mathbf{N} \cdot \mathbf{v})^2)}{(\mathbf{N} \cdot \mathbf{v})^2}} + 1} \quad (40)$$

$$p_s(\mathbf{v}, \mathbf{d}) = \frac{D(\alpha, N, h) G_1(\alpha, N, \mathbf{v})}{4(\mathbf{N} \cdot \mathbf{v})} \quad (41)$$

where p_d and p_s are the diffuse and specular components, which are mixed according to the luminance of K_d and K_s .

D. Details of Bubbling and Adaptive Sampling

Intermediate results in bubble step. Fig. 10 presents the process of how the missing objects are reconstructed by our bubbling method. The chandelier is missing initially at 50k iterations. In the early stage of bubble step, the light balls are reconstructed rapidly (55k), since they are relatively large inside the chandelier. On the other hand, thin components (e.g. poles) grow slowly (70k). Eventually (150k) the entire chandelier is successfully reconstructed.

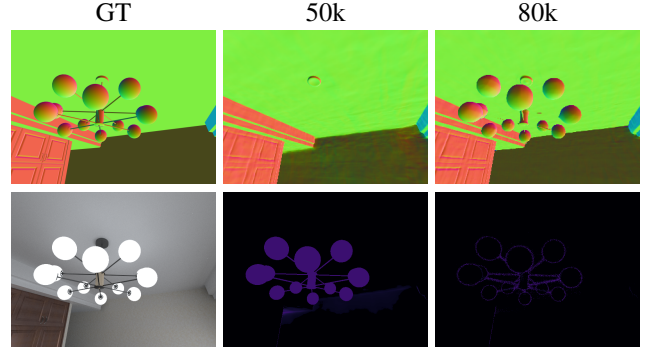


Figure 11. **Visualization of error-guided sampling map during training.** The first row displays intermediate normal reconstruction results, while the second row displays the corresponding error PDF map.

Visualization of error-guided sampling map during training. As shown in Fig. 11, at the training iteration of 50k (i.e. the start of bubble step), the chandelier is completely ignored and thus the corresponding pixels in the PDF map have high values. As the training proceeds to 80k iterations, the light balls have already been well-reconstructed, with chandelier poles still missing. Therefore, the value of pixels corresponding to light balls are reduced to 0, whereas chandelier pole pixels still need to be further sampled.

E. Additional Experimental Results

Novel view synthesis. Tab. 3 displays quantitative results of per-scene novel view PSNR. It turns out that our method outperforms all of the baselines, benefiting from our precise reconstruction of small objects and proper handling of shape-radiance ambiguity. Qualitative results are presented in Fig. 12. NeRF, NeuRIS and MonoSDF fails to reconstruct small objects such as chandeliers and lamp poles, while NeRF and Instant-NGP also suffers from fractured reconstruction results. While Instant-NGP usually capture most high-frequency details, it likely overfits to single-view radiance and fails to ensure multi-view geometry consistency in indoor scenes, leading to poor novel view synthesis results with floating artifacts.

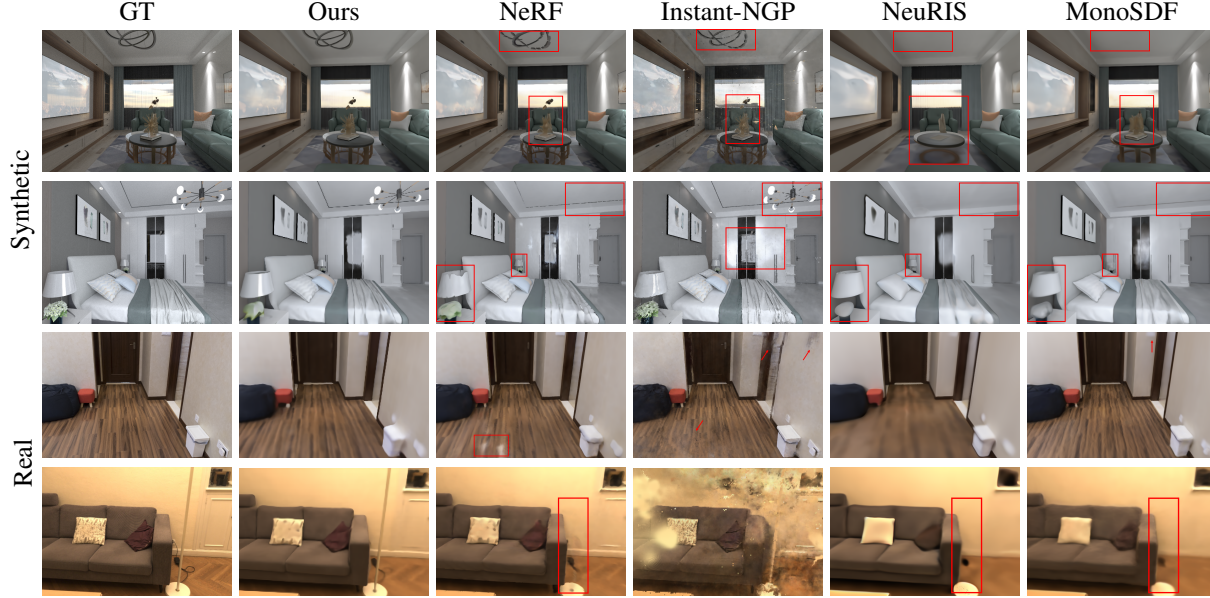


Figure 12. **Qualitative comparisons of novel view synthesis** on synthetic data and real data. Zoom in for details.

Table 3. **Comparisons of per-scene novel view PSNR.**

	Ours	NeRF	Instant-NGP	NeuRIS	MonoSDF
syn_1	28.03	26.24	25.8	25.53	27.37
syn_2	30.09	29.27	27.34	24.30	24.83
syn_3	27.46	25.73	26.58	24.30	24.48
syn_4	29.64	27.99	25.95	26.93	26.67
syn_5	27.71	26.93	16.82	24.38	26.03
syn_6	28.55	27.27	16.27	25.21	26.49
syn_7	28.04	27.65	24.67	25.18	24.67
syn_8	27.83	25.31	27.72	24.36	25.39
mean	29.70	27.09	23.89	25.02	25.74
real_1	26.63	26.48	19.10	25.91	26.21
real_2	28.01	27.21	26.32	23.87	24.38
real_3	24.58	24.11	23.29	23.31	22.72
real_4	21.39	20.86	19.23	19.82	20.05
mean	25.15	24.66	21.99	23.22	23.34

Geometry Reconstruction. Tab. 4 displays quantitative comparisons of per-scene normal angular L_1 error and depth L_1 error between our method and baselines. The definition of normal angular L_1 error is

$$\mathcal{L}_{\text{normal}} = \|1 - \hat{N} \cdot N\|_1 \quad (42)$$

Our method also outperforms NeuRIS and MonoSDF, indicating superior 3D reconstruction quality. Fig. 13 presents qualitative comparisons of the reconstructed normal maps. Our method can even recover high-frequency details on geometry, such as the spikes on the ball.

Material Decomposition. Fig. 14 presents qualitative results of the decomposed diffuse albedo K_d , specular albedo K_s and roughness ρ .

Scene editing. With the intrinsic decomposition results, we can enable photo-realistic scene editing tasks such as material editing and relighting. Fig. 15 shows qualitative results of scene editing results in both real and synthetic data. We explore mirror insertion (top-left and bottom-right), texture editing (top-right and mid-left), object insertion (mid-right) and relighting (bottom-left). Note that the edited specular reflections (on mirrors and inserted metal ball) are consistent with the surroundings. On account of

Table 4. Comparisons of per-scene normal angular error and depth L_1 loss.

	Normal-Angular- $L_1 \downarrow$			Depth- $L_1 \downarrow$		
	Ours	NeuRIS	MonoSDF	Ours	NeuRIS	MonoSDF
syn_1	0.040	0.051	0.036	0.014	0.240	0.010
syn_2	0.030	0.041	0.035	0.019	0.331	0.048
syn_3	0.054	0.080	0.073	0.021	0.319	0.061
syn_4	0.053	0.071	0.065	0.068	0.312	0.103
syn_5	0.064	0.096	0.058	0.025	0.227	0.034
syn_6	0.057	0.082	0.054	0.051	0.355	0.043
syn_7	0.057	0.071	0.076	0.065	0.335	0.099
syn_8	0.072	0.071	0.064	0.016	0.274	0.033
mean	0.053	0.070	0.058	0.035	0.299	0.054

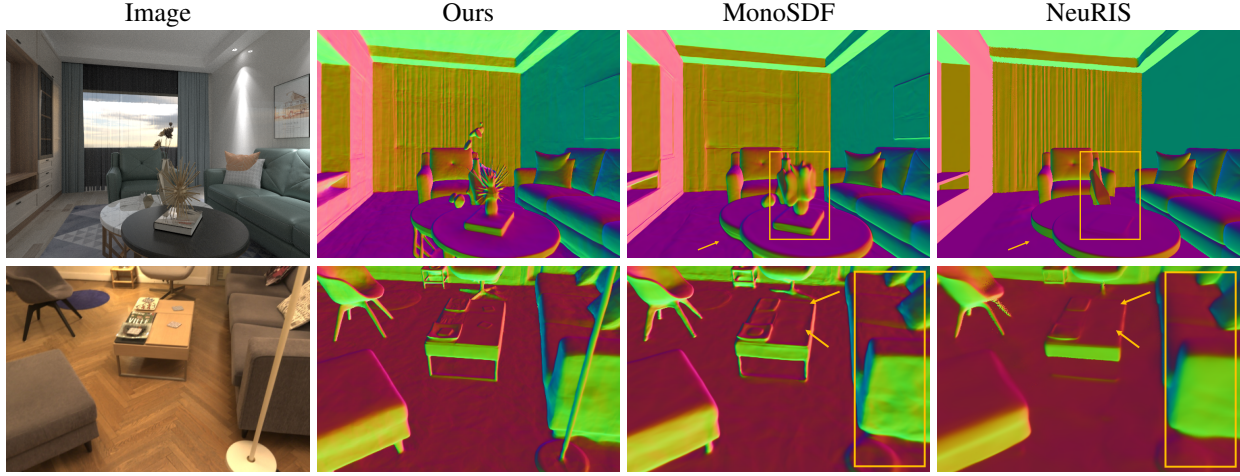


Figure 13. Qualitative comparisons of normal estimation on synthetic data and real data.

our raytracing algorithm, our method is capable of casting shadows of the inserted object (see the shadows of the inserted ball on the sofa).

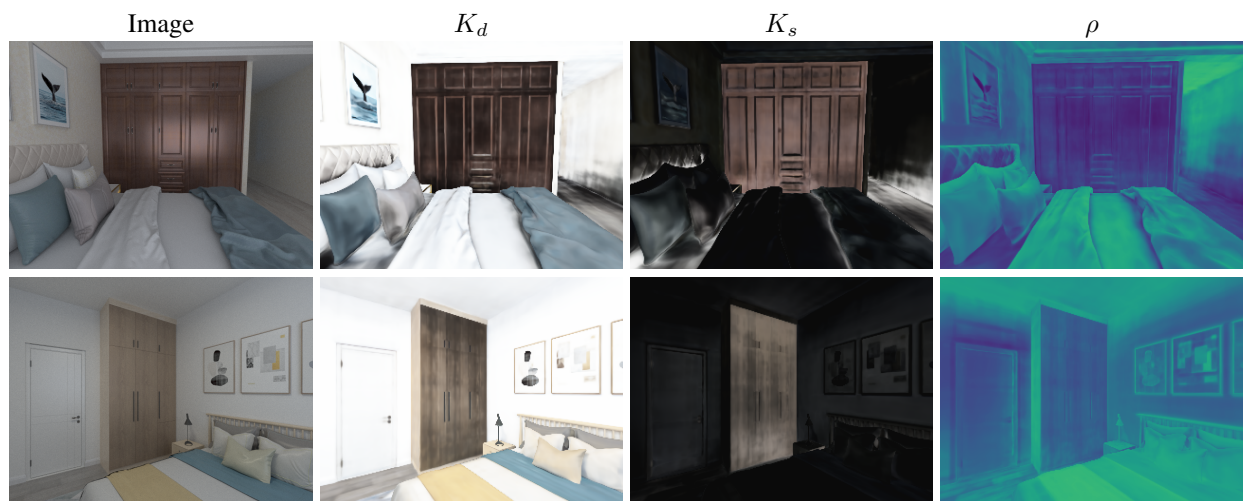


Figure 14. Qualitative results of decomposed materials.

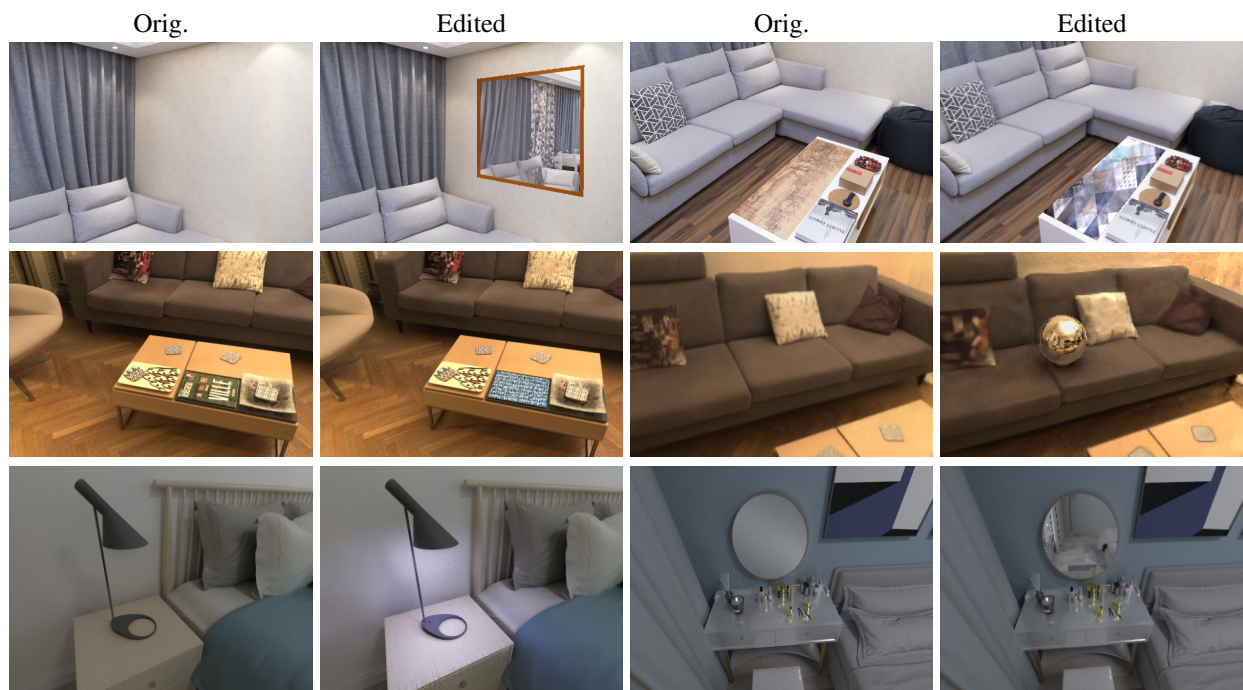


Figure 15. Qualitative results of scene editing and relighting.