

# Neural-PBIR Reconstruction of Shape, Material, and Illumination

Cheng Sun<sup>1,2\*</sup>    Guangyan Cai<sup>1,3\*</sup>    Zhengqin Li<sup>1</sup>    Kai Yan<sup>3</sup>    Cheng Zhang<sup>1</sup>  
 Carl Marshall<sup>1</sup>    Jia-Bin Huang<sup>1,4</sup>    Shuang Zhao<sup>3</sup>    Zhao Dong<sup>1</sup>

<sup>1</sup>Meta RLR    <sup>2</sup>National Tsing Hua University    <sup>3</sup>University of California, Irvine    <sup>4</sup>University of Maryland, College Park

\*The authors contribute equally to this paper

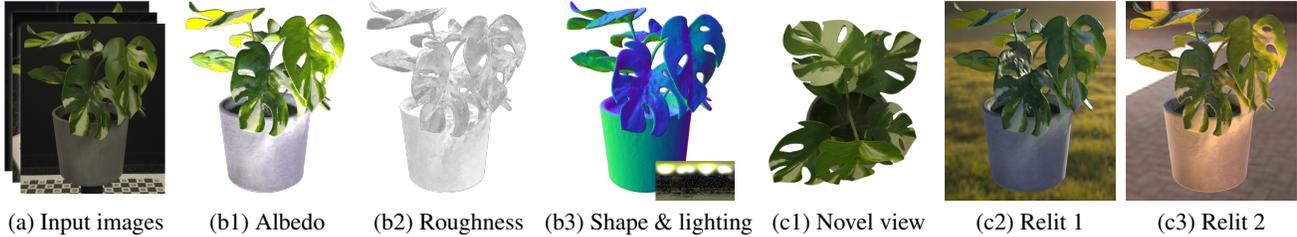


Figure 1: **Neural-PBIR** recovers high-fidelity material (b1,2), shape and lighting (b3), enabling realistic re-rendering (c1-3).

## Abstract

*Reconstructing the shape and spatially varying surface appearances of a physical-world object as well as its surrounding illumination based on 2D images (e.g., photographs) of the object has been a long-standing problem in computer vision and graphics. In this paper, we introduce an accurate and highly efficient object reconstruction pipeline combining neural based object reconstruction and physics-based inverse rendering (PBIR). Our pipeline firstly leverages a neural SDF based shape reconstruction to produce high-quality but potentially imperfect object shape. Then, we introduce a neural material and lighting distillation stage to achieve high-quality predictions for material and illumination. In the last stage, initialized by the neural predictions, we perform PBIR to refine the initial results and obtain the final high-quality reconstruction of object shape, material, and illumination. Experimental results demonstrate our pipeline significantly outperforms existing methods quality-wise and performance-wise. Code: <https://neural-pbir.github.io/>*

## 1. Introduction

Reconstructing geometry, material reflectance, and lighting from images, also known as inverse rendering, is a long-standing challenge in computer vision and graphics. Conventionally, the acquisition of the three intrinsic components has been mainly studied independently. For instance, multiview-stereo (MVS) [29, 30, 10] and time-of-flight [48] methods only focus on recovering object geometry, usually

based on diffuse reflectance assumption. Classical material acquisition methods typically assume known or simple geometries (e.g., a planar surface) with highly controlled illuminations [25, 38, 47], usually created with a light stage or gantry. This significantly limits their practicality when such capturing conditions are unavailable.

Recently, the advent of novel techniques enables us to jointly reconstruct shape, material, and lighting from 2D images of an object. At a high level, these techniques can be classified into two categories. Neural reconstruction methods encode the appearance of objects into a multi-layer perceptron (MLP) and optimize the network by minimizing the rendering errors from different views through differentiable volume ray tracing. NeRF [21] reconstructs a density field-based radiance field that allows high-quality view synthesis but not relighting. A series of methods [34, 39, 27] compute the density field from the signed distance function to achieve high-quality geometry reconstruction. Recent works [44, 45, 4, 5, 23, 12, 46] seek to fully decompose shapes, materials, and lighting from input images. However, due to the high computational cost of volume ray tracing and neural rendering, those methods take hours [23, 12] or a day [4, 46] to run and usually cannot model more complex indirect illumination [44, 45, 4, 5, 23, 12], causing shadows and color bleeding to be baked into the material reflectance. Several new methods [31, 22, 8] significantly reduce the computational cost of radiance field reconstruction using hybrid volume representations and efficient MLPs. We are among the first that adopt these novel techniques for efficient joint recovery of geometry, material, and lighting.

On the contrary, *physics-based inverse rendering* (PBIR)

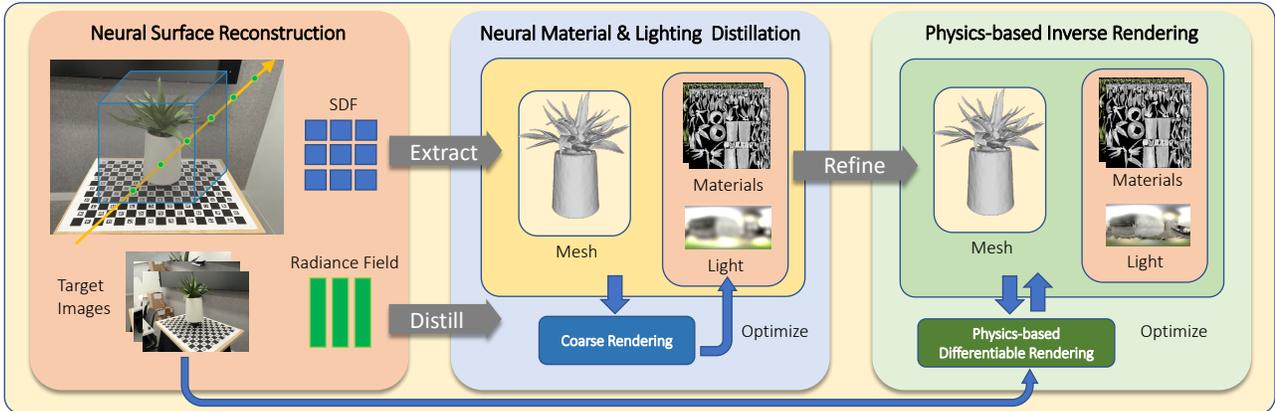


Figure 2: Our pipeline for joint shape, material, and lighting estimation.

[7, 19, 26, 1] optimizes shape, material, and lighting by computing unbiased gradients of image appearance with respect to scene parameters. Leveraging physics-based differentiable renderers [42, 13, 18], state-of-the-art PBIR pipelines can efficiently handle complex light transport effects such as soft shadow and interreflection. Such complex light transport effects cannot be easily handled through volume-based neural rendering. On the other hand, since PBIR methods rely on gradient-based optimization to refine intrinsic components, they can be prone to local minima and overfitting. Therefore, they may require a good initialization to achieve optimal reconstruction quality.

In this paper, we present a highly efficient and accurate inverse rendering pipeline with the advantages of both neural reconstruction and PBIR methods. Our pipeline attempt to estimate geometry, spatially varying material reflectance, and an HDR environment map from multiple images of an object captured under static but arbitrary lighting. As shown in Fig. 2, our pipeline consists of three stages. In the first stage, we propose a hybrid neural volume-based method for fast neural SDF and radiance field reconstruction, which achieves state-of-the-art geometry accuracy and efficiency. In the next stage, based on the reconstructed geometry and radiance field, we design an efficient optimization method to distill materials and lighting by fitting the surface light field. Our method relies on a radiance field to handle visibility and indirect illumination but avoids expensive volume ray tracing to significant computational cost compared to some recent works. Finally, we use an advanced PBIR framework [42, 13] to jointly refine the geometry, materials, and lighting. Note that our PBIR framework models complex light transport effects such as visibility, occlusion, soft shadows and interreflection in a physically correct and unbiased way while still being much faster than recent inverse rendering methods.

Concretely, our contributions include the following:

- A hybrid volume representation for fast and accurate geometry reconstruction.

- A efficient optimization scheme to distill high-quality initial material and lighting estimation from the reconstructed geometry and radiance field.
- An advanced PBIR framework that jointly optimizes materials, lighting and geometry with visibility and inter-reflection handled in a physically unbiased way.
- A end-to-end pipeline that achieves state-of-the-art geometry, material and lighting estimation that enables realistic view synthesis and relighting.

## 2. Related Works

**Volumetric surface reconstruction.** Recent progress in volumetric-based surface reconstruction of a static scene shows high-quality and robust results. NeuS [34] and VolSDF [39] replace NeRF [21]’s density prediction with signed distance values and proposes an unbiased and occlusion-aware volume rendering, achieving promising results. Subsequent works improve quality by introducing regularization losses to encourage surface smoothness [43], multiview consistency [9], and manhattan alignment [11]. These methods use a large MLP as their volumetric representation, which is however slow (many hours) to optimize per scene. Recent advances [31, 40, 22, 8] show great optimization acceleration without loss of quality by using an explicit grid. Unfortunately, using explicit volume to model SDF leads to bumpy and noisy surfaces [41]. Supervisions from SfM sparse point cloud [43, 9] or monocular depth/normal [41] may mitigate the difficulty but it’s out of our scope. We propose simple and effective regularization losses for explicit SDF grid optimization, achieving fast and high-quality surfaces without using external priors.

**Material and lighting estimation.** Several neural reconstruction methods [44, 45, 4, 5, 23, 12, 46] adopt the same setting as ours to simultaneously reconstruct geometry, material, and lighting from multiple images. An earlier work [45] directly optimizes a low-resolution environment map and materials from fixed geometry, leading to noisy light-

ing reconstruction caused by the highly ill-posed nature of this problem. More recent methods model lighting with a mixture of spherical Gaussians [44, 4, 46] or pre-filtered approximation [5, 23], and constrain material reflection with a spherical Gaussian lobe [44] or a low dimensional latent code [4, 5, 46] to improve the reconstruction quality. Most of them only consider direct illumination without occlusion [44, 4, 23] or use a fixed shadow map [45, 5] due to the high computational cost of MLP-based neural rendering. Nvdiffrmc [12] models shadows using a differentiable ray tracer with a denoiser that significantly reduces the number of samples but still cannot model interreflection. MII [46] trains another network to predict visibility and interreflection from the radiance field, which takes several hours. On the contrary, our hybrid, highly optimized neural SDF framework enables us to efficiently extract geometry, materials, lighting, visibility, and interreflection in less than 10 minutes. Our advanced PBIR framework for the first time allows holistic optimization of lighting, material, and geometry with direct and indirect illumination modeled in a physically unbiased way.

**Physics-based inverse rendering.** Different from neural-based methods, several recent works [19, 12, 7] utilize classic image formulation models developed in the graphics community and try to inverse this process using gradient-based optimizations, with the gradients computed using physics-based differentiable renderers. These approaches can handle complex light transport effects during optimization but are prone to local minima and overfitting. In our work, instead of optimizing from scratch, we perform PBIR as a refinement stage for better robustness.

### 3. Our Method

Provided multi-view images of an opaque object under fixed unknown illumination (with known camera parameters), our technique reconstructs the shape and reflectance of the object as well as the illumination condition.

As illustrated in Fig. 2, our pipeline is comprised of three main stages. The first stage (Sec. 3.1) is a *fast and precise surface reconstruction step* that brings direct SDF grid optimization into NeuS [34]. Associated with this surface is an overfitted radiance field that does not fully model the surface reflectance of the object. Our second stage (Sec. 3.2) is an *efficient neural distillation method* that converts the radiance fields to physics-based reflectance [6] and illumination models. Lastly, our third stage (Sec. 3.3) utilizes *physics-based inverse rendering* (PBIR) to further refine the object geometry and reflectance reconstructed by the first two stages. This stage leverages physics-based differentiable rendering that captures global illumination (GI) effects such as soft shadows and interreflection.

### 3.1. Neural Surface Reconstruction

Taking as input multiple images of an opaque object under fixed illumination (and with known camera parameters), the first stage of our method produces a detailed reconstruction of the object’s surface. To this end, we optimize the object surface and an outgoing radiance field that best describes the input images. In this stage, for efficiency and robustness, we express the object surface as the zero-level set  $\{\mathbf{x} \in \mathbb{R}^3 \mid S(\mathbf{x}) = 0\}$  of a signed distance field (SDF)  $S(\mathbf{x})$ , and the radiance field  $L_o(\mathbf{x}, -\mathbf{v})$  (for any position  $\mathbf{x} \in \mathbb{R}^3$  and viewing direction  $\mathbf{v} \in \mathbb{S}^2$ ) as a non-physics-based general function approximator.

**Unbiased volume rendering with SDF.** To compute the color of a pixel, we first sample  $N$  query points  $\{\mathbf{x}_i = \mathbf{o} + t_i \mathbf{v}\}_{i=1}^N$  (with  $0 < t_1 < t_2 < \dots < t_N$ ) along the corresponding camera ray originated at the camera’s location  $\mathbf{o}$  with viewing direction  $\mathbf{v}$ . We then query the scene representation for points sign distance  $S(\mathbf{x}_i)$  and radiance  $L_o(\mathbf{x}_i, \mathbf{v})$ . Following NeuS [34]’s unbiased rendering, we activate the queried signed distance into alpha for all  $1 \leq i \leq N$  via

$$\alpha_i = \max\left(0, \frac{\sigma(S(\mathbf{x}_i)) - \sigma(S(\mathbf{x}_{i+1}))}{\sigma(S(\mathbf{x}_i))}\right), \quad (1)$$

where  $\sigma$  denotes the Sigmoid function. Then, the pixel color  $C$  is computed by the alpha blending of the queried alphas and radiance values

$$C = \sum_{i=1}^N T_i \alpha_i L_o(\mathbf{x}_i, -\mathbf{v}) \quad \text{where } T_i = \prod_{j=1}^{i-1} (1 - \alpha_j). \quad (2)$$

**Voxel-based scene representation.** NeuS uses a large MLP to model the SDF  $S$  and radiance field  $L_o$ . Unfortunately, since this MLP is expensive to query, NeuS’ optimization processes can be very time-consuming.

For better performance, we adapt the dense-grid-based scene representation from DVGO [31] by setting

$$\begin{aligned} S(\mathbf{x}) &= \text{interp}(\mathbf{x}, V^{(\text{sdf})}), \\ L_o(\mathbf{x}, -\mathbf{v}) &= \text{MLP}\left(\text{interp}(\mathbf{x}, V^{(\text{feat})}), \mathbf{v}\right), \end{aligned} \quad (3)$$

where  $\text{interp}()$  indicates trilinear interpolation,  $V^{(\text{sdf})}$  is a dense SDF grid,  $V^{(\text{feat})}$  is a dense feature grid, and MLP is a single-hidden-layer MLP with ReLU activation.

In practice, we use two sets of  $V^{(\text{sdf})}$  and  $V^{(\text{feat})}$  grids to model the foreground (i.e., the object of interest) and the background with varying resolutions. Please refer to the supplement for more details.

**Adaptive Huber loss.** We incorporate a Huber loss to reduce the impact of specular highlights which can cause bumpy artifacts in reconstructed surfaces:

$$\mathcal{L}_{\text{photo}} = \sum_r \begin{cases} (\mathcal{I}[r] - C[r])^2, & (|\mathcal{I}[r] - C[r]| \leq t) \\ 2t|\mathcal{I}[r] - C[r]| - t^2, & (\text{otherwise}) \end{cases} \quad (4)$$

where:  $\mathcal{I}[r]$  and  $C[r]$  denote, respectively, the colors of the  $r$ -th pixels of the target and rendered images;  $t \in \mathbb{R}_{>0}$  is a hyper-parameter.

Intuitively, the infinity norm of the gradient of  $\mathcal{L}_{\text{photo}}$  is clamped by  $2t$  to prevent bright pixels (e.g., those exhibiting specular highlights) dominating the optimization. In practice, instead of tuning a constant  $t$ , we adaptively update  $t$  by the running mean of the median of  $|\mathcal{I}[r] - C[r]|$  in each iteration to clamp gradient for about half of the pixels.

**Laplacian regularization.** To further improve the robustness of our approach and reduce potential artifacts, we regularize our SDF grid  $V^{(\text{sdf})}$  using a Laplacian loss:

$$\mathcal{L}_{\text{lap}} = \sum_u \left( \sum_{u' \in N(u)} \left( V^{(\text{sdf})}[u'] - V^{(\text{sdf})}[u] \right) \right)^2, \quad (5)$$

where  $V^{(\text{sdf})}[u]$  indicates the SDF value stored at the grid point  $u$ , and  $N(u)$  denotes the six direct neighbors of  $u$ .

**Training.** Following DVGO [31], we use progressive grid scaling for efficiency and more coherent results. Our experiments also indicate that using a per-point RGB loss  $\mathcal{L}_{\text{pp-rgb}}$  improves convergence speed and quality:

$$\mathcal{L}_{\text{pp-rgb}} = \sum_{r,i} T_i^{(r)} \alpha_i^{(r)} \left| L_o(\mathbf{x}_i^{(r)}, -\mathbf{v}^{(r)}) - \mathcal{I}[r] \right|, \quad (6)$$

where  $\alpha_i^{(r)}$  and  $T_i^{(r)}$  are computed using Eqs. (1) and (2) for each pixel  $r$ . In summary, to train our SDF and feature grids  $V^{(\text{sdf})}$  and  $V^{(\text{feat})}$  in Eq. (3), we minimize the objective:

$$\mathcal{L}_{\text{surf}} = \mathcal{L}_{\text{photo}} + w_{\text{lap}} \mathcal{L}_{\text{lap}} + w_{\text{pp-rgb}} \mathcal{L}_{\text{pp-rgb}}. \quad (7)$$

**Postprocessing.** Optionally, our surface reconstruction stage generates for each input image an anti-aliased mask that separates the object from the background using mesh rasterisation and alpha matting [28]. The generated masks can be used by our physics-based inverse rendering stage (Sec. 3.3) to facilitate the refinement of object shapes.

### 3.2. Neural material and lighting distillation

Provided the optimized object surface and outgoing radiance field from the surface reconstruction stage (Sec. 3.1), the goal of the second stage of our pipeline is to obtain an initial estimation of surface reflectance and illumination condition. To this end, we leverage the radiance field  $L_o$  as the teacher model to distill the learned surface color into physics-based material and illumination models (that can be rendered to reproduce  $L_o$ ) as follows.

As preprocessing for this stage, we extract a triangle mesh  $M^{(0)}$  from the optimized SDF  $S(\mathbf{x})$  given by Eq. (3) using (non-differentiable) marching cube. Additionally, for each mesh vertex  $v$ , we compute its normal  $M_n[v] \in \mathbb{S}^2$  based on the gradient  $\nabla S$  of the SDF.

**Material and illumination models.** To model spatially varying surface reflectance, we use the widely-adopted Disney microfacet BRDF [15] parameterized by surface albedo and roughness. In practice, instead of using a large MLP to model the spatially varying BRDF parameters [44, 45, 23, 46], we store the albedo  $M_a[v] \in [0, 1]$  and roughness  $M_r[v] \in \mathbb{R}_{>0}$  for each vertex  $v$  of the extracted mesh  $M^{(0)}$  (and interpolate them in the interior of triangle faces) for better performance. To model environmental illumination, we use mixtures of spherical Gaussians [33, 20, 44, 46] with 256 lobes.

**Coarse rendering.** For fast training, we opt to use numerical integration on stratified pre-sampled light directions  $\Omega$  for efficiency (we use 256 samples in practice). Specifically, for each vertex  $v$  of the extracted mesh  $M^{(0)}$  and light direction  $\omega_i \in \Omega$ , we precompute the visibility  $\text{Vis}[v, \omega_i]$  and indirect illumination  $L_i^{(\text{ind})}[v, \omega_i]$  by tracing a ray from the vertex in the direction  $\omega_i$ . If this ray intersects the object surface at some  $\mathbf{x}$ , we set  $\text{Vis}[v, \omega_i] = 0$  and  $L_i^{(\text{ind})}[v, \omega_i] = L_o(\mathbf{x}, -\omega_i)$  [46]. Otherwise, we set  $\text{Vis}[v, \omega_i] = 1$ . Then, the incident radiance at each vertex can be expressed as

$$L_i(\omega_i) = L_{\text{env}}^{\text{SG}}(\omega_i) \text{Vis}[\omega_i] + L_i^{(\text{ind})}[\omega_i] (1 - \text{Vis}[\omega_i]), \quad (8)$$

where  $L_{\text{env}}^{\text{SG}}$  is the SG-based environmental illumination. Lastly, the outgoing radiance is computed by

$$\hat{L}_o(\omega_o) = \frac{1}{Z} \sum_{\omega_i \in \Omega} L_i(\omega_i) f(\omega_i, \omega_o, M_n) (M_n \cdot \omega_i), \quad (9)$$

where  $Z$  is a normalization factor,  $f$  denotes the BRDF function parameterized by  $M_a$  and  $M_r$ . In Eqs. (8) and (9), we omit dependencies on vertex  $v$  for brevity.

**Training.** In each iteration, we randomly sample a outgoing direction  $\omega_o$  (with  $M_n[v] \cdot \omega_o[v] > 0$ ) for each vertex  $v$  and minimize the absolute difference with the teacher model:

$$\mathcal{L}_{\text{distill}} = \sum_v \left| \hat{L}_o(v, \omega_o[v]) - L_o(v, \omega_o[v]) \right|. \quad (10)$$

To encourage smoothness, we apply total variation loss on the per-vertex albedo  $M_a$  and roughness  $M_r$ :

$$\mathcal{L}_{\text{v-reg}} = \sum_{(v_1, v_2) \in \text{edge}} \sum_{* \in \{a, r\}} |M_*[v_1] - M_*[v_2]|. \quad (11)$$

Also, we regularize the SG-based illumination  $L_{\text{env}}^{\text{SG}}$  using

$$\mathcal{L}_{\text{bg}} = \sum_{\omega_i \in \Omega} \left| L_{\text{env}}^{\text{SG}}(\omega_i) - (L_{\text{env}}^{\text{SG}})'(\omega_i) \right|, \quad (12)$$

where  $(L_{\text{env}}^{\text{SG}})'$  indicates the averaged background observation (see the supplement for more details).

In summary, we optimize per-vertex albedo  $M_a$ , roughness  $M_r$ , and the environmental illumination  $L_{\text{env}}^{\text{SG}}$  by minimizing

$$\mathcal{L}_{\text{distill.total}} = \mathcal{L}_{\text{distill}} + w_{\text{v-reg}} \mathcal{L}_{\text{v-reg}} + w_{\text{bg}} \mathcal{L}_{\text{bg}}. \quad (13)$$

**Postprocessing.** With the per-vertex attributes  $M_a$  and  $M_r$  obtained, we UV-parameterize the extracted mesh  $M^{(0)}$  and pixelize  $M_a$  and  $M_r$ , respectively, into an albedo map  $T_a^{(0)}$  and a roughness map  $T_r^{(0)}$ . These maps will be refined by our following inverse-rendering stage (Sec. 3.3).

### 3.3. Physics-Based Inverse Rendering

Reconstructions produced by our shape reconstruction and neural material distillation stages (Secs. 3.1 and 3.2) comprise three components: (i) environmental illumination  $L_{\text{env}}^{(0)} := L_{\text{env}}^{\text{SG}}$  expressed as spherical Gaussians; (ii) object reflectance parameterized with an albedo map  $T_a^{(0)}$  and a roughness map  $T_r^{(0)}$ ; and (iii) object surface mesh  $M^{(0)}$ .

Despite being high-fidelity, these reconstructions can still lack sharp details and are not immune to artifacts (see Fig. 7). To further improve reconstruction quality, we utilize a physics-based inverse rendering stage. Specifically, initialized our neural reconstructions (i.e.,  $L_{\text{env}}^{(0)}$ ,  $T_a^{(0)}$ ,  $T_r^{(0)}$ , and  $M^{(0)}$ ), we apply gradient-based optimization to minimize the *inverse-rendering loss*:

$$\mathcal{L}_{\text{IR}} = \mathcal{L}_{\text{img}} + w_{\text{mask}} \mathcal{L}_{\text{mask}} + w_{\text{reg}} \mathcal{L}_{\text{reg}}, \quad (14)$$

where

$$\mathcal{L}_{\text{img}}(L_{\text{env}}, T_a, T_r, M) = \sum_j \|\mathcal{I}_j - \mathcal{R}_j(L_{\text{env}}, T_a, T_r, M)\|_1, \quad (15)$$

is the *image loss* given by the sum of L1 losses between the  $j$ -th target image  $\mathcal{I}_j$  and the corresponding rendered image  $\mathcal{R}_j(L_{\text{env}}, T_a, T_r, M)$ . Additionally,  $\mathcal{L}_{\text{reg}}$  and  $\mathcal{L}_{\text{mask}}$  in Eq. (14) denote, respectively, the *regularization* and the *mask losses*—which we will discuss in the following.

**Environment map optimization.** We recall that the initial illumination  $L_{\text{env}}^{(0)} = L_{\text{env}}^{\text{SG}}$  is a coarse reconstruction expressed as a set of spherical Gaussians (SGs). To further refine it, we employ a two-step process as follows. In the first step, we directly optimize the SG parameters (i.e., per-lobe means and variances). In the second step, we pixelize the optimized SG representation into an environment map (using the latitude-longitude parameterization) and perform per-pixel optimization.

**SVBRDF optimization.** Since our initializations  $T_a^{(0)}$  and  $T_r^{(0)}$  are already high-quality, we perform per-vertex optimization for the albedo and roughness maps  $T_a$  and  $T_r$ . Further, to make our optimization less prone to Monte Carlo noises produced by our physics-based renderer (discussed later), we regularize  $T_r$  using a total variation loss:

$$\mathcal{L}_{\text{reg}}(T_r) = \sum_{(x,y)} \sum_{(x',y')} |T_r[x',y'] - T_r[x,y]|, \quad (16)$$

where  $T_r[x,y]$  denotes the  $(x,y)$ -th texel of the roughness map  $T_r$ , and  $(x',y') \in \{(x+1,y), (x,y+1)\}$  are two direct neighbors of  $(x,y)$ .

**Shape refinement.** In all our experiments, object geometries predicted by our neural stages accurately recover object topology. Thus, although it is possible to directly optimize SDFs in inverse rendering [2, 32], we opt to use explicit mesh-based representations for the object surface  $M$  in this stage for better efficiency. To make our per-vertex optimization more robust to Monte Carlo noises, we utilize Nicolet et al.’s AdamUniform optimizer [24].

To further improve efficiency, we leverage object masks produced either as input or by our surface reconstruction (Sec. 3.1) and introduce a mask loss:

$$\mathcal{L}_{\text{mask}}(M) = \sum_j \|S_j - \mathcal{R}_j^{\text{mask}}(M)\|_1, \quad (17)$$

where  $S_j$  is our predicted mask for the  $j$ -th target image and  $\mathcal{R}_j^{\text{mask}}$  the corresponding rendered mask.

**Differentiable rendering.** To differentiate the image and the mask losses defined in Eqs. (15) and (17), we develop a physics-based Monte Carlo differentiable renderer that implements path-space differentiable rendering [42] and utilizes state-of-the-art numerical backend [13].

Most, if not all, previous techniques including MII [46] and nvdiffrmc [12] rely on highly simplified differentiable rendering processes that typically neglect global-illumination (GI) effects and produce biased gradients with respect to geometry. On the contrary, our differentiable renderer offers unbiased gradients and the generality of differentiating GI (and anti-aliased masks) with respect to surface geometries. Consequently, our pipeline is capable of generating higher-quality reconstructions than state-of-the-art methods, which we will demonstrate in Sec. 4.

## 4. Results

To demonstrate the effectiveness of our method, we show our reconstructions on synthetic input images and captured photographs in Secs. 4.1 and 4.2, respectively. We compare reconstructions obtained with our technique with two state-of-the-art baselines: MII [46] and nvdiffrmc [12]. Additionally, we conduct ablation studies to evaluate several components of our pipeline in Sec. 4.3. Please refer to the supplement for more results.

### 4.1. Synthetic Data

We assess the effectiveness of our proposed method using the synthetic dataset made available by MII [46]. This dataset comprises four virtual scenes—each of which includes multi-view renderings (with object masks) of an object under some natural environmental illumination (with the ground truth environment map provided). For each scene, the dataset provides 100 target images accompanied by camera poses and object masks for training. Additionally, the testing set consists of renderings of ground truth

Method	Speed	Relighting			Aligned albedo			Albedo			Rough.
	Time↓	PSNR↑	SSIM↑	LPIPS↓	PSNR↑	SSIM↑	LPIPS↓	PSNR↑	SSIM↑	LPIPS↓	MSE↓
nvdifrec-mc [12]	~2 h	23.93	0.946	0.074	<b>29.72</b>	<b>0.959</b>	<b>0.057</b>	18.25	0.899	0.103	<u>0.009</u>
MII [46]	~10 h	27.53	0.947	0.087	25.77	0.935	<u>0.066</u>	24.62	0.931	<b>0.064</b>	<b>0.008</b>
Ours - Distilled only	<15 m	30.26	0.961	0.059	27.67	0.933	0.079	26.20	0.931	0.093	<u>0.009</u>
Ours - Const. init.	~37 m	30.25	<b>0.970</b>	0.050	28.55	0.940	0.070	25.83	0.940	0.080	<u>0.010</u>
Ours - w/o GI	~45 m	30.57	0.960	0.050	27.71	0.940	0.070	26.38	0.940	0.082	<u>0.009</u>
Ours - w/o shape ref.	~45 m	<u>30.61</u>	0.965	<u>0.049</u>	28.74	0.944	0.067	<u>26.84</u>	<u>0.941</u>	0.082	<b>0.008</b>
Ours - Full	~1 h	<b>30.73</b>	<u>0.966</u>	<b>0.047</b>	<u>29.06</u>	<u>0.946</u>	0.067	<b>26.85</b>	<b>0.944</b>	<u>0.080</u>	<b>0.008</b>

Table 1: **Relighting, material reconstruction, and view-interpolation quality on MII dataset [46].** We compare our method with MII and Nvdifrec-mc. The highest performing number is presented in bold, while the second best is under-scored. For “Ours - Distilled only”, we did not run the PBIR stage. For “Ours - Const. init.” and “Ours - w/o GI”, we did not run the shape refinement. Please refer to Secs. 4.1 and 4.3 for more details.

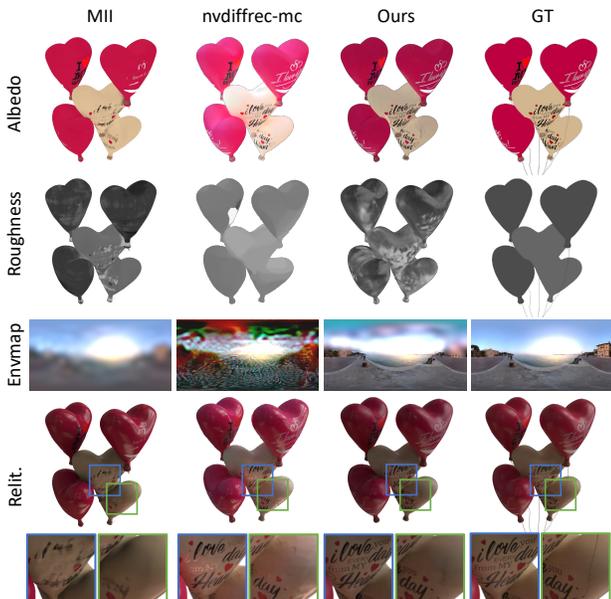


Figure 3: **Qualitative comparisons on the MII data.**

albedo, roughness, and the object under two novel lighting conditions in 200 poses. We apply our technique and the baselines to the posed training target images and masks (but not the GT environment maps).

Tab. 1 presents the quantitative comparisons where we compare our method’s reconstructions with the baselines by rendering them in the testing poses and comparing them with the ground truth images. In line with MII, we report PSNR, SSIM, and LPIPS for the relit and albedo images, and the MSE for the roughness images averaged across all scenes. MII also reports error metrics on aligned albedo to reduce the impact of albedo-light ambiguity [46]. Specifically, we compute an RGB scale that minimizes the difference between the reconstructed albedo images and the ground truth albedo images, and reconstructed albedo images are scaled before evaluation. Our method outperforms the baselines and takes less time.

Methods	Novel-view		Captured light re-rendering			
			raw		aligned	
	PSNR↑	SSIM↑	PSNR↑	SSIM↑	PSNR↑	SSIM↑
Nvdifrecmc	30.3	0.94	21.8	0.92	28.0	0.94
MII	28.9	0.94	27.5	0.94	28.6	0.94
Ours	<b>31.6</b>	<b>0.96</b>	<b>28.8</b>	<b>0.95</b>	<b>30.7</b>	<b>0.95</b>

Table 2: **Quantitative comparison on Our Real Dataset.** To inspect material quality for relighting, we capture 360 images for our dataset and evaluate the rendering results under the captured lighting.

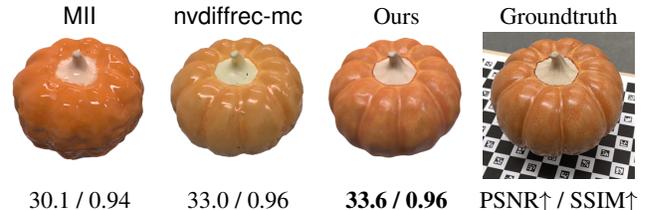


Figure 5: **Rerendering of reconstruction results under captured (GT) illumination.** We rescale all renderings to match the overall brightness of the GT image.

We show qualitative results in Fig. 3. Since MII imposes a sparsity constraint to reconstruct a sparse set of materials, the results tend to be over-blurred. Despite being able to produce sharp texture maps, nvdifrec-mc suffers from inaccurate illumination reconstructions, causing their reconstructed albedo maps to be color-shifted. Compared with these two baselines, our method produces significantly more accurate material and illumination reconstructions.

## 4.2. Real Data

To further demonstrate the robustness of our method, we captured five real scenes under indoor lighting and again compared the reconstruction quality with MII and nvdifrec-mc. For each scene, we use around 200 measured images for training and 10–20 for testing. We place a checkerboard beneath the object for the geometric calibration of camera

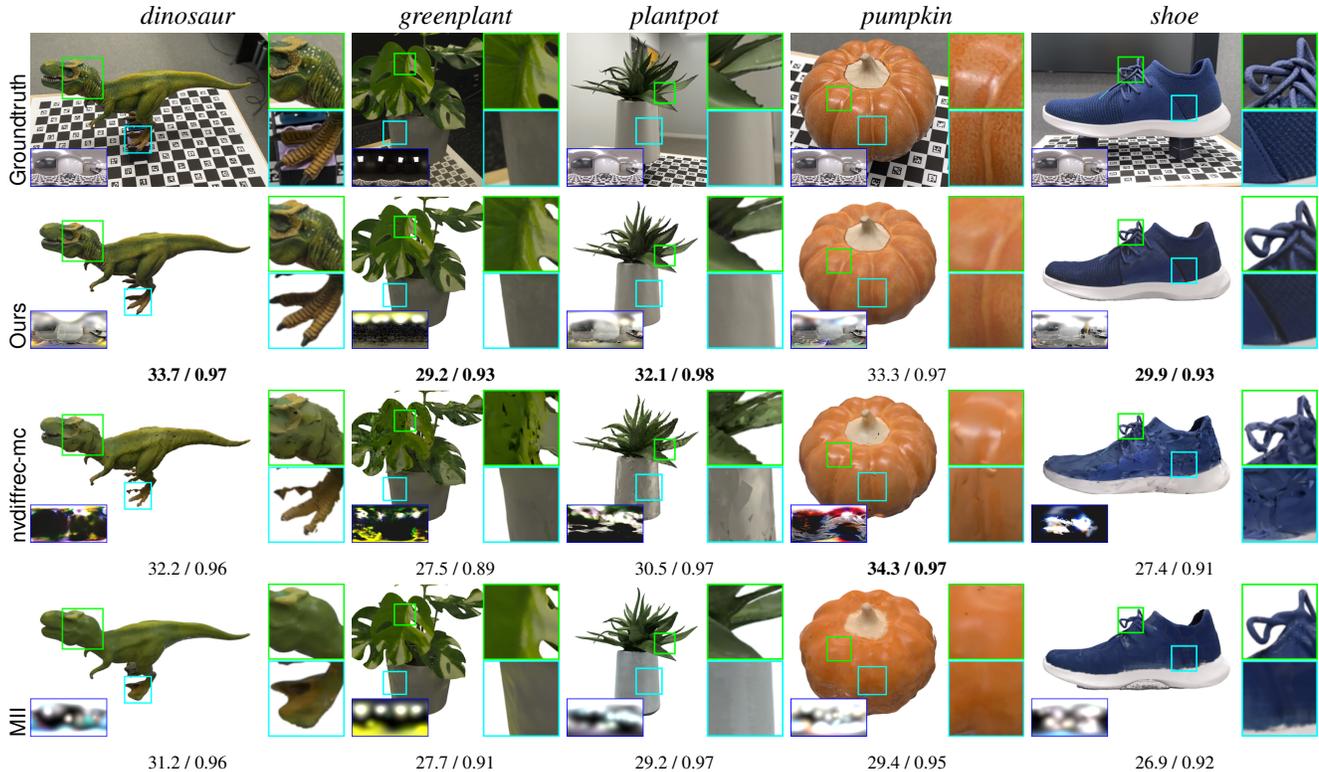


Figure 4: **Novel-view interpolation on our real dataset.** Our technique produces high-fidelity reconstructions with minimal artifacts. We report the average PSNR $\uparrow$  and SSIM $\uparrow$  below each image.

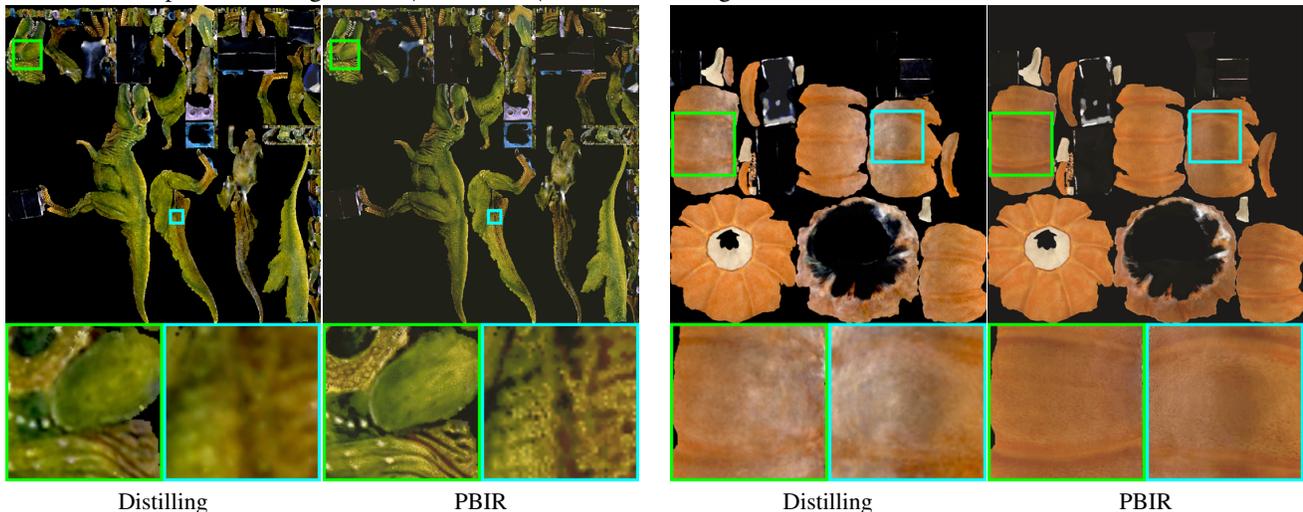


Figure 7: **Reconstructed albedo before and after PBIR.** Despite our neural distilling outperform previous arts by a large margin, we still observe blurriness and light baking in the reconstructed materials. Our PBIR stage can provide sharper details and remove light baking efficiently.

parameters. For each object, we manually annotate an object 3D bounding box to crop the reconstructed meshes to focus the comparisons on the object of interest. In addition, we measure ground truth environment map (which we use for evaluation only) for each scene using an HDR 360 camera. As there are no ground truth foreground masks while

our baselines require them, we use masks generated in our surface reconstruction stage (§3.1).

We compare the qualities of reconstructions produced by our method and the baselines using novel-view renderings (since we do not have the GTs). As shown in Fig. 4, our reconstructions produce detailed renderings that closely re-

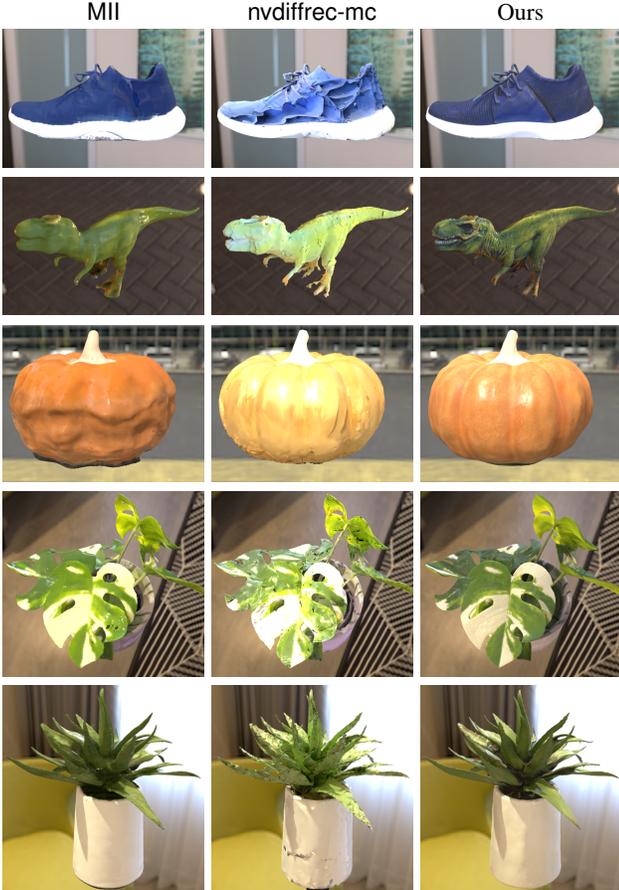


Figure 6: **Novel-illumination renderings of the reconstructed models.** Our results offer the highest overall quality with minimal artifacts.

semble the GTs, while MII and nvdiffrrec-mc’s results suffer from various artifacts (*e.g.* lacking details and bumpy surfaces). In addition, to evaluate the quality of reconstructed object shape and reflectance, we compare renderings under captured GT illumination in Fig. 5. The quantitative comparison is summarized in Tab. 2 where we achieve the best novel-view synthesis results under both the reconstructed and the captured environment maps. We show qualitative comparison under novel-illumination in Fig. 6.

### 4.3. Evaluations and Ablations

**Surface reconstruction.** To evaluate the quality of reconstructed geometries produced by our surface reconstruction stage (Sec. 3.1), we use the subsampled 15 scenes of DTU dataset [14]. Each scene has 49 or 64 images with camera parameters and masks provided. The quantitative results are concluded in Tab. 3. Our direct SDF grid optimization uses significantly less time than NeuS [34] while still outperforming NeuS surface accuracy measured in chamfer distance (CD). Our simple surface reconstruction with dense grid also achieves similar performance comparing to

Methods	COLMAP	NeuS	Voxurf	NeuS2	Our
Runtime↓	1 hrs	5.5 hrs	16 mins	<b>5 mins</b>	<b>5 mins</b>
CD (mm)↓	1.36	0.77	0.72	0.70	<b>0.66</b>

Table 3: **Surface reconstruction quality on DTU dataset [14].** The results are averaged across the 15 scenes. We include two recent works, Voxurf [37] and NeuS2 [35], for reference. See supp. for results breakdown.

	COLMAP	NeuS	Voxurf	NeuS2	Our
$\mathcal{L}_{\text{lap}}$			✓	✓	✓
$\mathcal{L}_{\text{pp,rgb}}$		✓		✓	✓
ada. huber	✓		✓	✓	✓
CD (mm)↓	1.50	1.37	1.24	1.11	1.03
					1.00
					0.89
					<b>0.68</b>

Table 4: **Ablation study of SDF grid regularizations.** The results are averaged over the 15 objects on DTU dataset.

the most recent Voxurf [37] and NeuS2 [35].

We conduct a comprehensive ablation study for the three losses that help us achieve direct SDF grid optimization. As shown in Tab. 4, naively adapting DVGO without regularization leads poor quality. Among the regularization,  $\mathcal{L}_{\text{lap}}$  is the most significant, which enforces the SDF grid to evolve smoothly during optimization. Adding  $\mathcal{L}_{\text{pp,rgb}}$  and using adaptive Huber loss also shows good improvement. Combining all three losses offers the best results.

We also show the chamfer distance (CD) in the first 15k iterations **with** and **without**  $\mathcal{L}_{\text{pp,rgb}}$  here. The solid lines are CD, while the dashed lines are the difference to the final CD. The results are averaged over the 15 DTU scenes. Using  $\mathcal{L}_{\text{pp,rgb}}$  improves final quality and speeds up convergence to the final CD in fewer iterations.

Please note that our results on DTU dataset are directly from the shape reconstruction stage. We skip evaluating the shape refinement of our physics-based inverse rendering as DTU exhibit vary light occlusion from robot arms. Please see the supplement for more results.

**Neural material distilling.** We now demonstrate the usefulness of our neural material distilling stage (Sec. 3.2)—which provides high-quality material predictions used to initialize our PBIR stage (Sec. 3.3). Specifically, when using constant initializations for surface reflectance (with  $T_a = 0.5$ ,  $T_r = 0.5$ ) and illumination (with a gray environment map), it takes our PBIR pipeline 37 minutes (see the row labeled as “Ours - const. init.” in Tab. 1) to generate reconstructions with a similar level of accuracy as the neural distilling stage does in less than 15 minutes (see the row labeled as “Ours - distilled only”).

**Importance of PBIR.** The distilling method is efficient by sidestepping ray tracing and approximating global illu-



Figure 8: **Reconstructed albedo with/without GI.** We show the impact of modeling global illumination (GI) on material reconstruction. The above images are the reconstructed albedo maps with GI on/off, with their errors listed below (PSNR $\uparrow$  / SSIM $\uparrow$ ). Without GI, the optimization tends to “bake” the indirect lighting into the albedo map.

mination but it sacrifices accuracy and tends to “bake” inter-reflection into the materials. We use our physically-based differentiable renderer to further refine the material. In 7, we show the albedo maps before and after PBIR on our real dataset. We achieve much higher fidelity with less baking after PBIR. The result is consistent with the quantitative results on synthetic dataset (refer to “Ours - Distilled only” versus “Ours - Full”) in 1.

**Importance of GI.** Our differentiable renderer used in the PBIR stage (Sec. 3.3) is capable of handling global-illumination (GI) effects such as interreflection. To demonstrate the importance of GI, we run PBIR optimizations with and without GI and compare the material reconstruction qualities. As shown in rows labeled as “Ours - w/o GI” and “Ours - w/o shape ref.” in Tab. 1, enabling GI improves the accuracy of material and lighting reconstructions.

A main reason for the usefulness of GI is that, without GI, inverse-rendering optimizations tend to “bake” effects like interreflections into reflectance (*e.g.* albedo) maps, limiting their overall accuracy. We demonstrate this by comparing albedo maps optimized with and without GI in Fig. 8 (using the balloons data from the MII dataset).

**Shape refinement.** Lastly, we demonstrate the usefulness of shape refinement (by optimizing the vertex positions of the extracted mesh) using rows labeled as “Ours - w/o shape ref.” and “Ours - Full” in Tab. 1 and Fig. 9. Please refer to the supplement for more examples.

## 5. Discussion and Conclusion

**Limitations.** Due to the fundamentally under-constrained nature of inverse rendering, our technique is not immune to “baking” artifacts—especially when predictions provided by our neural stages are far from the groundtruth. Also, although our technique handles global-illumination effects including interreflection, it assumes for opaque materials and does not currently support the reconstruction of transparent or translucent objects.

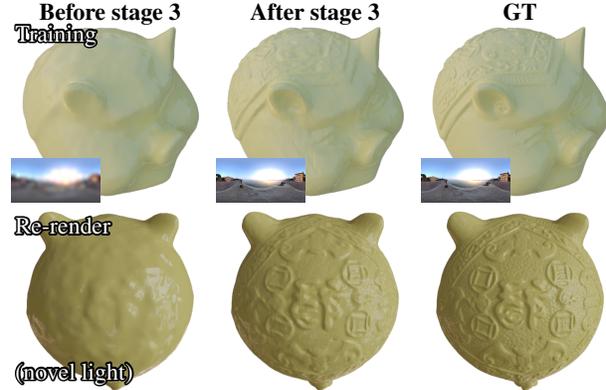


Figure 9: **An example showcasing the usefulness of our PBIR shape refinement.** Our stage 1 can sometimes lack geometry details while our PBIR in stage 3 can recover these details.

**Conclusion.** We introduced a new method for reconstructing an object’s shape and reflectance under unknown environmental illumination. Our technique is comprised of three stages: (i) a surface reconstruction stage fitting an implicit surface and a radiance field from input images; (ii) a neural distilling stage decomposing material and lighting from the radiance field; and (iii) a PBIR stage jointly refining the geometry, material, and lighting. Broad experiments show our effectiveness.

**Acknowledgement.** This project has been partially supported by NSF grant 1900927.

## A. Results on Stanford ORB dataset

We evaluate Neural-PBIR on the recently released Stanford-ORB dataset [17]<sup>1</sup>.

**Stanford-ORB.** The dataset scans 14 real-world objects each under 3 different lighting condition, resulting in 42 capturing sequence. The ground-truth mesh and material are also scanned from studio. The relative camera poses of an object captured from different lighting are also provided for evaluating the equality of re-lighting quality.

**Evaluation protocol.** We follow official guideline to train our method on each of the 42 scenes separately under benchmark resolution of  $512 \times 512$ . The same set hyperparameters as the MII dataset are applied for all the Stanford-ORB scenes. We follow official train-test split and use official script (<https://github.com/StanfordORB/Stanford-ORB>) for authentic evaluation scores.

**Results.** The quantitative comparison with previous arts is provided in Tab. 5, where our Neural-PBIR outperforms previous methods on most metrics. We show some qualitative results in Fig. 10.

<sup>1</sup>We add Neural-PBIR’s result in Dec 2023.

Method	Geometry			Novel Scene Relighting				Novel View Synthesis			
	Depth↓	Normal↓	Shape↓	PSNR-H↑	PSNR-L↑	SSIM↑	LPIPS↓	PSNR-H↑	PSNR-L↑	SSIM↑	LPIPS↓
PhySG [44]	1.90	0.17	9.28	21.81	28.11	0.960	0.055	24.24	32.15	0.974	0.047
NVDiffRec [23]	<u>0.31</u>	0.06	0.62	22.91	29.72	0.963	0.039	21.94	28.44	0.969	0.030
NeRD [4]	1.39	0.28	13.7	23.29	29.65	0.957	0.059	25.83	32.61	0.963	0.054
NeRFactor [45]	0.87	0.29	9.53	23.54	30.38	0.969	0.048	26.06	33.47	0.973	0.046
InvRender [36]	0.59	<u>0.06</u>	<u>0.44</u>	23.76	30.83	0.970	0.046	25.91	34.01	0.977	0.042
NVDiffRecMC [12]	0.32	<b>0.04</b>	0.51	<u>24.43</u>	<u>31.60</u>	<u>0.972</u>	<u>0.036</u>	28.03	36.40	<u>0.982</u>	<u>0.028</u>
Neural-PBIR	<b>0.30</b>	<u>0.06</u>	<b>0.43</b>	<b>26.01</b>	<b>33.26</b>	<b>0.979</b>	<b>0.023</b>	<b>28.83</b>	<b>36.80</b>	<b>0.986</b>	<b>0.019</b>

Table 5: Geometry, relighting, and view-interpolation quality on Stanford-ORB dataset [17].

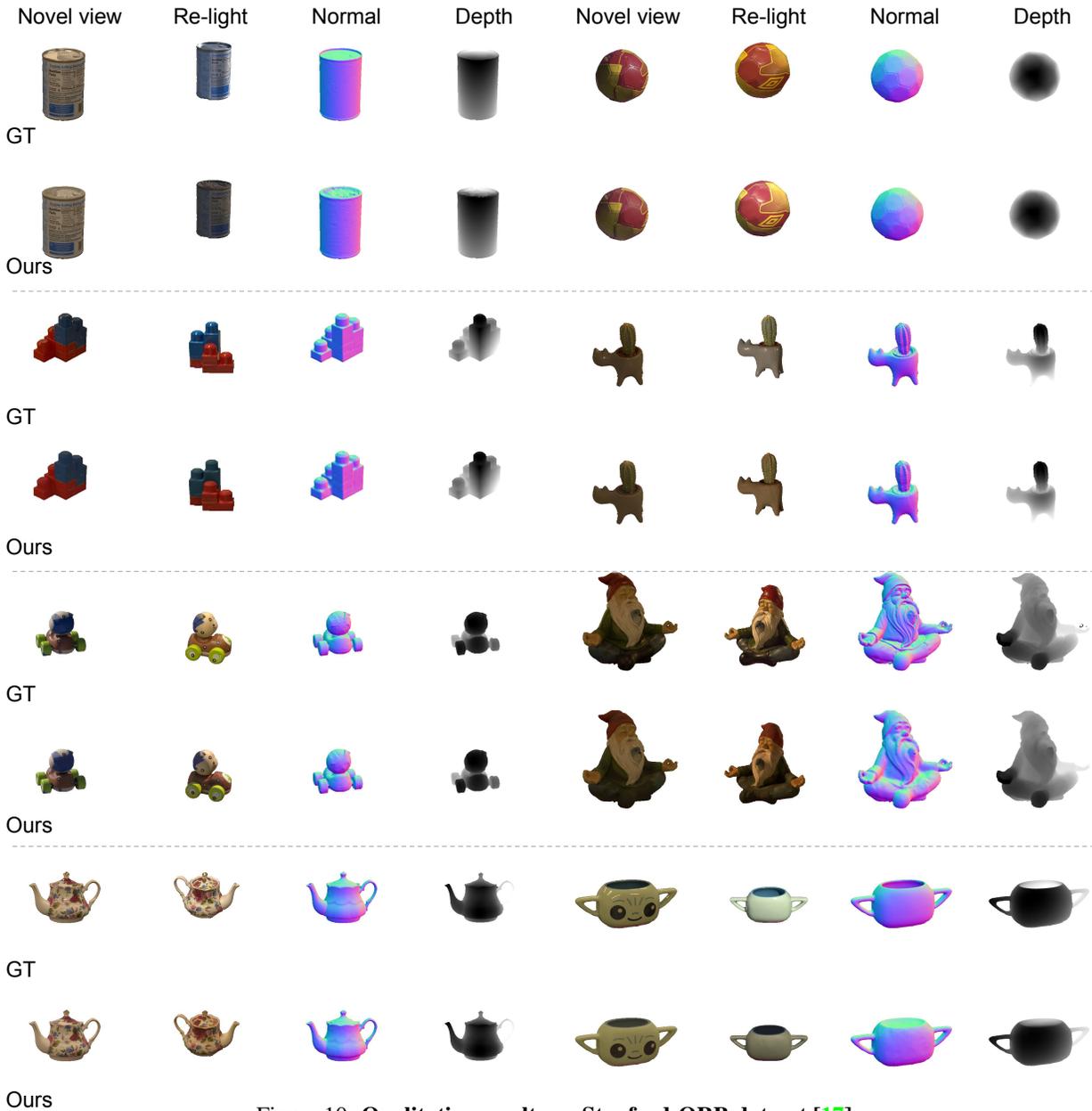


Figure 10: Qualitative results on Stanford-ORB dataset [17].

## B. Technical Details

In what follows, we elaborate technical details of our Neural-PBIR pipeline’s three main stages. Code will be released upon internal approval for future extension and reproduction.

### B.1. Neural Surface Reconstruction

**Sharpness term in unbiased volume rendering.** Following NeuS [34], we use a scaled sigmoid  $\sigma_s$  function in the SDF for alpha activation:

$$\alpha_i = \max \left( 0, \frac{\sigma_s(S(\mathbf{x}_i)) - \sigma_s(S(\mathbf{x}_{i+1}))}{\sigma_s(S(\mathbf{x}_i))} \right), \quad (18)$$

where:

- $S$  is the signed-distance function;
- $\sigma_s(y) = (1 + \exp(-sy))^{-1}$  with  $s > 0$  being the *sharpness term*;
- $\{\mathbf{x}_i = \mathbf{o} + t_i \mathbf{v}\}_{i=1}^N$  (with  $0 < t_1 < t_2 < \dots < t_N$ ) are the  $N$  sampled points on the camera ray originated at the camera’s location  $\mathbf{o}$  with viewing direction  $\mathbf{v}$ .

Specifically, we start with  $s = 30$  if foreground masks are provided (e.g., for the synthetic and DTU datasets) and  $s = 5$  otherwise (e.g. for our measured real-world dataset). In practice, we use a *scheduled* sharpness  $s$  (instead of updating  $s$  with gradient descent) as we find it more stable. Then, we update the sharpness  $s$  by setting  $s \leftarrow \min(s + 0.02, 300)$  after each iteration.

**Background modeling.** As stated in Sec. 3.1 of the main paper, we use two sets of  $V^{(\text{sdf})}$  and  $V^{(\text{feat})}$  grids to model the foreground and the background (via Eq. (3) in the main paper), respectively. Specifically, the foreground region is defined as the volume inside a (predetermined) small bounding box. The background, on the other hand, is the volume inside a much larger bounding box.<sup>2</sup> Given a camera ray, we categorize the sample points  $\{\mathbf{x}_i = \mathbf{o} + t_i \mathbf{v}\}_{i=1}^N$  along the ray as foreground or background and then evaluate signed distance  $S(\mathbf{x}_i)$  and radiance  $L_o(\mathbf{x}_i, -\mathbf{v})$  for each  $\mathbf{x}_i$  using the corresponding grids.

Thanks to the background scene volume, our method can work without external mask supervision (e.g., our own real-world dataset).

**Points sampling on rays.** When sampling 3D points  $\{\mathbf{x}_i = \mathbf{o} + t_i \mathbf{v}\}_{i=1}^N$  along a camera ray, we use  $t_i = i\Delta t$  with  $\Delta t$  being half the size of a grid voxel for all  $i = 1, 2, \dots, N$ .

<sup>2</sup>In practice, we use background bounding boxes that are  $16\times$  as large as the foreground ones.



Figure 11: Examples of the averaged background observation. The black region indicates missing observation.

**Coarse-to-fine optimization.** For better efficiency and more coherent results, when optimizing the  $V^{(\text{sdf})}$  and  $V^{(\text{feat})}$  grids, we leverage a coarse-to-fine scheme by doubling the number of voxels every 1k iterations for the first 10k iterations. The final voxel resolutions are  $300^3$  for the foreground grids (which contain the object of interest) and  $160^3$  for the background ones.

**Optimization details.** We optimize the  $V^{(\text{sdf})}$  and  $V^{(\text{feat})}$  grids for the foreground and the background jointly using the Adam [16] method with  $\beta = (0.9, 0.99)$  and  $\epsilon = 10^{-12}$  in 20k iterations. When computing the loss, we use weights  $w_{\text{lap}} = 10^{-8}$  and  $w_{\text{pp\_rgb}} = 0.01$ . Also, when using the running means to update the threshold  $t$  in the adaptive Huber loss, we set the momentum to 0.99 and clamp  $t$  to a minimum of 0.01.

When training the SDF grids  $V^{(\text{sdf})}$ , we use an initial learning rate of 0.01 that then decays to 0.001 at 10k iterations. When training the outgoing radiance field  $L_o$ , we use a learning rate of 0.001 for the MLPs and 0.1 for the feature grids  $V^{(\text{feat})}$ .

### B.2. Neural Distillation of Material and Lighting

**Initialization.** We initialize the roughnesses to  $M_r[v] = 0.25$  for each vertex  $v$ . For per-vertex albedo, we initialize  $M_a[v]$  to the median of the outgoing radiance from the teacher model  $L_o$ :

$$M_a[v] = \text{Median} \left\{ L_o(\mathbf{x}[v], \boldsymbol{\omega}_o) \mid \boldsymbol{\omega}_o \in \Omega, (\boldsymbol{\omega}_o \cdot M_n[v]) > 0 \right\}, \quad (19)$$

where  $\mathbf{x}[v]$  and  $M_n[v]$  indicate, respectively, the position and the normal of vertex  $v$ , and  $\Omega$  is the predetermined set of outgoing directions.

**Fresnel term.** In addition to albedo and roughness, we also need the Fresnel term  $F_0$  [15] to model specular reflection. Following MII, we assume the object to be reconstructed is dielectric and make  $F_0$  constant. We set  $F_0 = 0.02$  for all synthetic data since it is used by MII’s open-source implementation, and  $F_0 = 0.04$  for real-world data since it is the industrial standard.

**Averaged background constraint.** Recap that we regularize our SG-based illumination  $L_{\text{env}}^{\text{SG}}$  to be similar to the averaged background observation  $(L_{\text{env}}^{\text{SG}})'$ . We now detail how the latter is obtained.

First, we gather all “background” training pixels (toward which the camera rays miss our reconstructed mesh). Then, we compute  $(L_{\text{env}}^{\text{SG}})'$  as an environment map under the latitude-and-longitude representation as follows. For each background pixel with intensity  $I$  and viewing direction  $\mathbf{v}$ , we set the value of the corresponding pixel  $j$  in the environment map  $(L_{\text{env}}^{\text{SG}})'$ —based on latitude and longitude coordinates of  $\mathbf{v}$ —as  $(L_{\text{env}}^{\text{SG}})'[j] = I$ . When multiple pixels (from different camera locations) contribute to one pixel  $j$  of  $(L_{\text{env}}^{\text{SG}})'$ , we set  $(L_{\text{env}}^{\text{SG}})'[j]$  using the average intensity of all such pixels.

We show some examples of the averaged background observations  $(L_{\text{env}}^{\text{SG}})'$  in Fig. 11. We only compute the regularization loss for the observed viewing directions.

**Optimization details.** To optimize per-vertex appearance parameters, we use the Adam method with  $\beta = (0.9, 0.999)$  and  $\epsilon = 10^{-8}$  in 2k iterations. When computing losses, we use the weights  $w_{\text{v.reg}} = 0.1$  and  $w_{\text{bg}} = 10$ . We use a learning rate 0.01 for per-vertex attributes and 0.001 for the spherical Gaussian (SG) parameters (representing the illumination  $L_{\text{env}}^{\text{SG}}$ ).

### B.3. Physics-Based Inverse Rendering

**Optimization details.** Initialized using the mesh  $M_0$  predicted by the surface reconstruction stage as well as albedo/roughness maps  $T_{\text{a}}^{(0)}, T_{\text{r}}^{(0)}$  (for surface reflectance) and SG-based illumination  $L_{\text{env}}^{\text{SG}}$  produced by the neural distillation stage, our physics-based inverse rendering (PBIR) stage involves the following three steps:

1. We jointly optimize (using 1k iterations) the albedo/roughness maps  $T_{\text{a}}, T_{\text{r}}$  and the SG parameters  $L_{\text{env}}^{\text{SG}}$  while keeping the mesh geometry fixed.
2. We first pixelize the SG-based  $L_{\text{env}}^{\text{SG}}$  into an environment map  $L_{\text{env}}$  and then perform joint per-pixel optimizations (using 1k iterations) for the albedo, roughness, and environment maps  $T_{\text{a}}, T_{\text{r}}$ , and  $L_{\text{env}}$ .
3. We jointly optimize (using 500 iterations) all maps and the mesh geometry (per-vertex).

In practice, when optimizing albedo and roughness maps  $T_{\text{a}}$  and  $T_{\text{r}}$  in all three steps, we use the Adam optimizer with  $\beta = (0.9, 0.999)$ ,  $\epsilon = 10^{-8}$ , and the learning rates  $10^{-2}$  for  $T_{\text{a}}$  and  $5 \times 10^{-3}$  for  $T_{\text{r}}$ . When computing losses, we use  $w_{\text{mask}} \approx 10$  and  $w_{\text{reg}} \approx 0.1$  (which we slightly adjust for each example).

Additionally, in the first step, we use the Adam optimizer [16] for the SG parameters with  $\beta = (0.9, 0.999)$ ,  $\epsilon = 10^{-8}$ , and learning rates around 0.001 (which we slightly adjust per example). In the second step, to suppress the impact of Monte Carlo noises during environment map optimization, we utilize the AdamUniform optimizer [24] with  $\lambda = 1$  and a learning rate of 0.01. In the last step, when

optimizing the mesh geometry, we again use the AdamUniform optimizer with  $\lambda = 100$ .

## C. Additional Results and Evaluations

### C.1. Additional Results

**Video for view synthesis and relighting.** Since results of novel-view synthesis and relighting are best viewed animated, we encourage readers to see our supplementary video ([video.mp4](#)) for a more convincing comparison on our five real-world objects.

Similar to the results shown in Fig. 6 of the main paper, our method significantly outperforms nvdiffrmc [12] and MII [46]. nvdiffrmc’s geometry and material reconstructions contain heavy artifacts. Despite nvdiffrmc showing better novel-view results than MII (main paper’s Fig. 4), the artifacts become visually prominent under novel illuminations as can be seen in the video and main paper’s Fig. 6. MII offers better overall albedo than nvdiffrmc but suffer from over-blurring in both geometry and material reconstructions. Overall, our results show significant better quality in both geometry and material.

**Outdoor illumination.** The five real-world objects presented the main paper are captured under indoor lighting. In Fig. 12, we showcase the results of two of these objects re-captured under outdoor illumination. Same as the results under indoor lighting, our reconstructions are more detailed, allowing their rerenderings (under novel views) to achieve better PSNR and SSIM.

**Synthetic MII dataset.** The authors of MII have kindly shared their rendered results for us to compare. As their evaluation scripts are unavailable, we use our own implementation for all the quantitative results. Due to the different implementation of the evaluation metrics, MII’s quantitative results presented in our main paper differ slightly from those reported in their paper.

In Figs. 13 to 16, we show more qualitative results on the synthetic MII dataset. Overall, our method offers more detailed albedo reconstructions than the baseline methods. On the other hand, none of the methods performs well on roughness estimation—likely due to the lack of robust priors. The qualitative results are consistent with the quantitative comparison in Tab. 1 of the main paper.

**Our synthetic dataset.** Since the MII dataset does not contain groundtruth meshes, it is difficult to evaluate the accuracy of reconstructed shapes. To address this, we create two extra synthetic scenes—*buddha* and *lion*—with groundtruth meshes for evaluation. For each scene, the training set includes 190 posed images with masks. The testing set consists of visualizations of groundtruth albedo, roughness, and renderings of the object under seven novel lighting conditions in 10 poses.

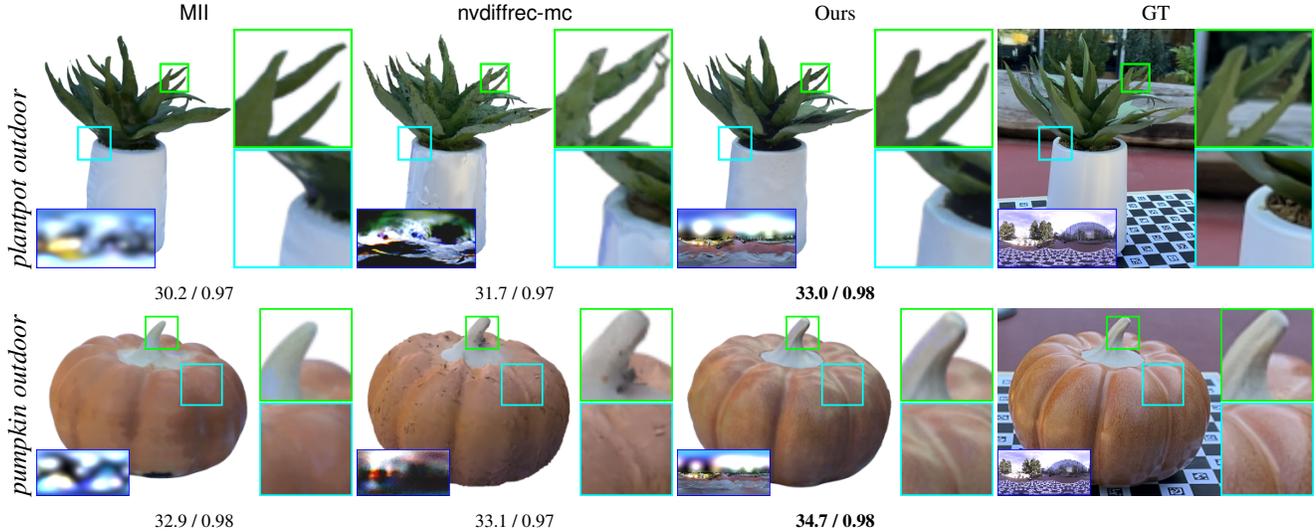


Figure 12: **Novel-view interpolation on our additional two real outdoor data.** We report the average PSNR $\uparrow$  and SSIM $\uparrow$  below each image. The results show that our method achieves good quality and outperforms previous arts under outdoor lighting as well.

Table 6 shows quantitative comparisons between our method and the baselines. In addition to metrics used in the MII dataset, we also measure Chamfer distances [3] between optimized and groundtruth shapes (normalized so that the groundtruth has unit bounding boxes). Our method again outperforms the baselines.

As shown in Figs. 17 and 18, since the background is fully visible (i.e., each pixel of  $L_{env}$  is visible as the background of at least one input image), our method is capable of reconstructing the environment map almost perfectly. Because of this, our albedo reconstructions are not hindered by the albedo-light ambiguity—as demonstrated in Tab. 6 where the error metrics barely change with or without albedo alignment. We note that this might not apply to all scenarios, for instance, the background might not be fully visible, as shown in the MII dataset. Reconstructing indoor lighting perfectly is also challenging even if the background is completely visible, because it breaks the assumption of environmental (i.e., distant) lighting.

## C.2. Additional Evaluations

**Surface quality on the DTU dataset.** We show quantitative results breakdown for the 15 scenes from DTU dataset [14] in Tab. 7. We use the official evaluation script to measure Chamfer distances. Please note that our results evaluated here are directly from the shape reconstruction stage. We skip evaluating the shape refinement of our physics-based inverse rendering on DTU dataset as DTU exhibit vary light occlusion from robot arms.

**Usefulness of shape refinement.** Lastly, we demonstrate the usefulness of our shape refinement (as the last step of the

physics-based inverse rendering stage) via an ablation. As shown in Fig. 19 and Tab. 6, our shape refinement improves the accuracy of reconstructed object geometries.



Figure 13: Qualitative comparisons of *air\_balloons* from the MII dataset.



Figure 14: Qualitative comparisons of *chair* from the MII dataset.

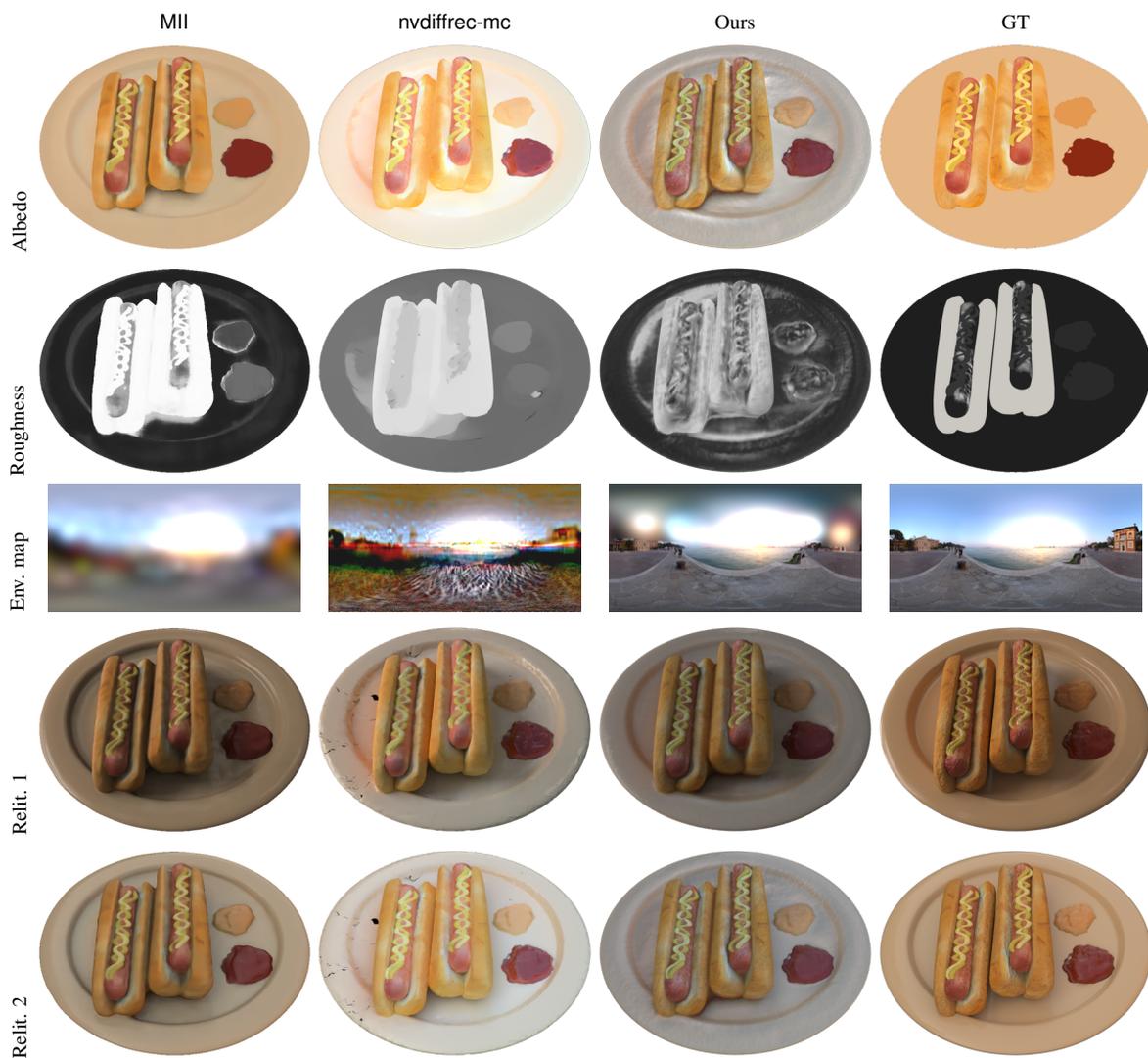


Figure 15: **Qualitative comparisons of hotdog from the MII dataset.**

Method	Speed	Relighting			Aligned albedo			Albedo			Rough.	Shape
	Time↓	PSNR↑	SSIM↑	LPIPS↓	PSNR↑	SSIM↑	LPIPS↓	PSNR↑	SSIM↑	LPIPS↓	MSE↓	CD↓
nvdiffrmc [12]	~ 2 h	21.77	0.936	0.071	33.29	0.964	0.037	16.14	0.910	0.068	0.013	9.42e-5
MI [46]	~ 10 h	24.94	0.952	0.051	30.92	0.962	0.044	19.80	0.923	0.065	<u>0.003</u>	5.92e-5
Ours - Distilled only	< 15 m	33.90	0.976	0.034	34.09	0.971	0.034	34.09	0.972	0.034	0.005	<u>2.61e-5</u>
Ours - w/o shape ref.	~ 45 m	<u>34.18</u>	<u>0.980</u>	<u>0.028</u>	<u>35.57</u>	<u>0.983</u>	<u>0.026</u>	<u>35.57</u>	<u>0.983</u>	<u>0.026</u>	<u>0.003</u>	<u>2.61e-5</u>
Ours - Full	~ 1 h	<b>35.30</b>	<b>0.982</b>	<b>0.026</b>	<b>37.69</b>	<b>0.985</b>	<b>0.023</b>	<b>37.68</b>	<b>0.985</b>	<b>0.023</b>	<b>0.002</b>	<b>2.56e-5</b>

Table 6: **Relighting, material reconstruction, and mesh quality on our synthetic dataset.** We compare our method with MII and nvdiffrmc. The highest performing number is presented in bold, while the second best is underscored. We measure the shape quality using Chamfer distances (CD).

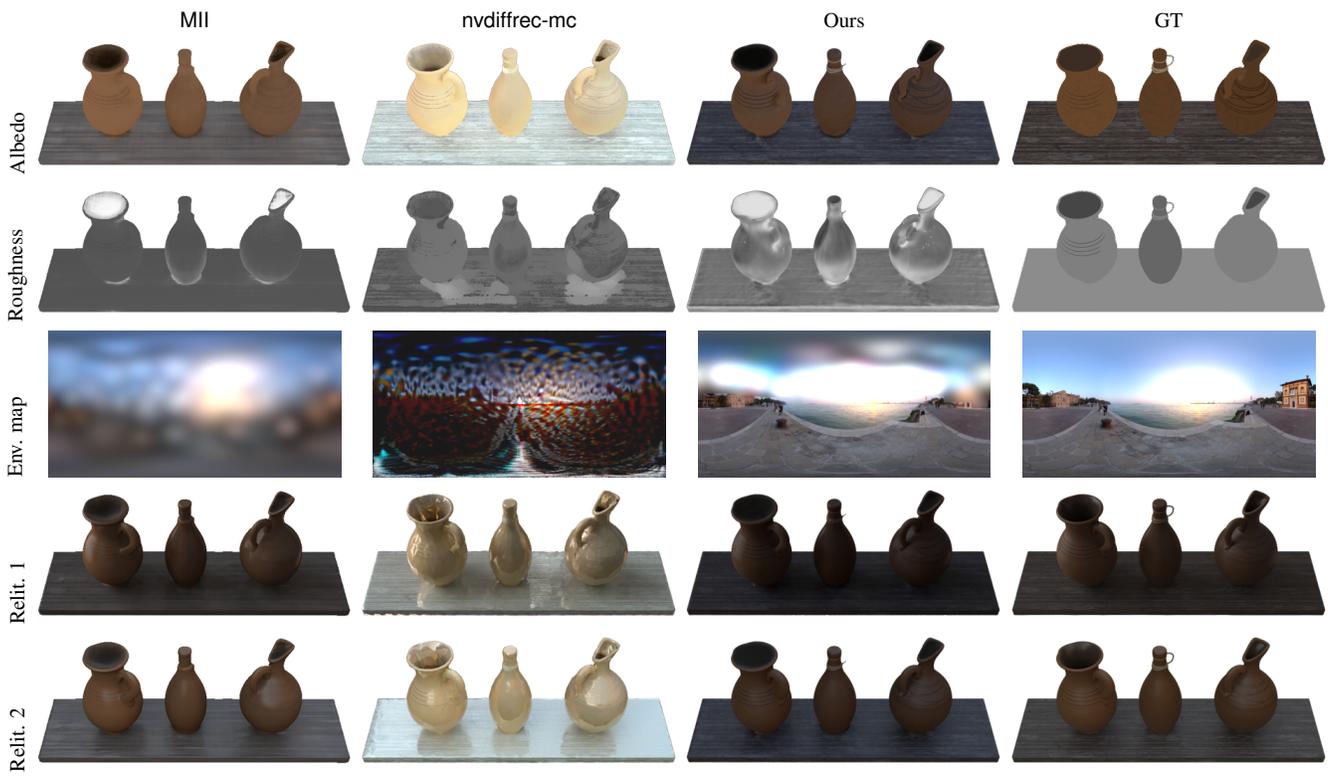


Figure 16: Qualitative comparisons of *jugs* from the MII dataset.

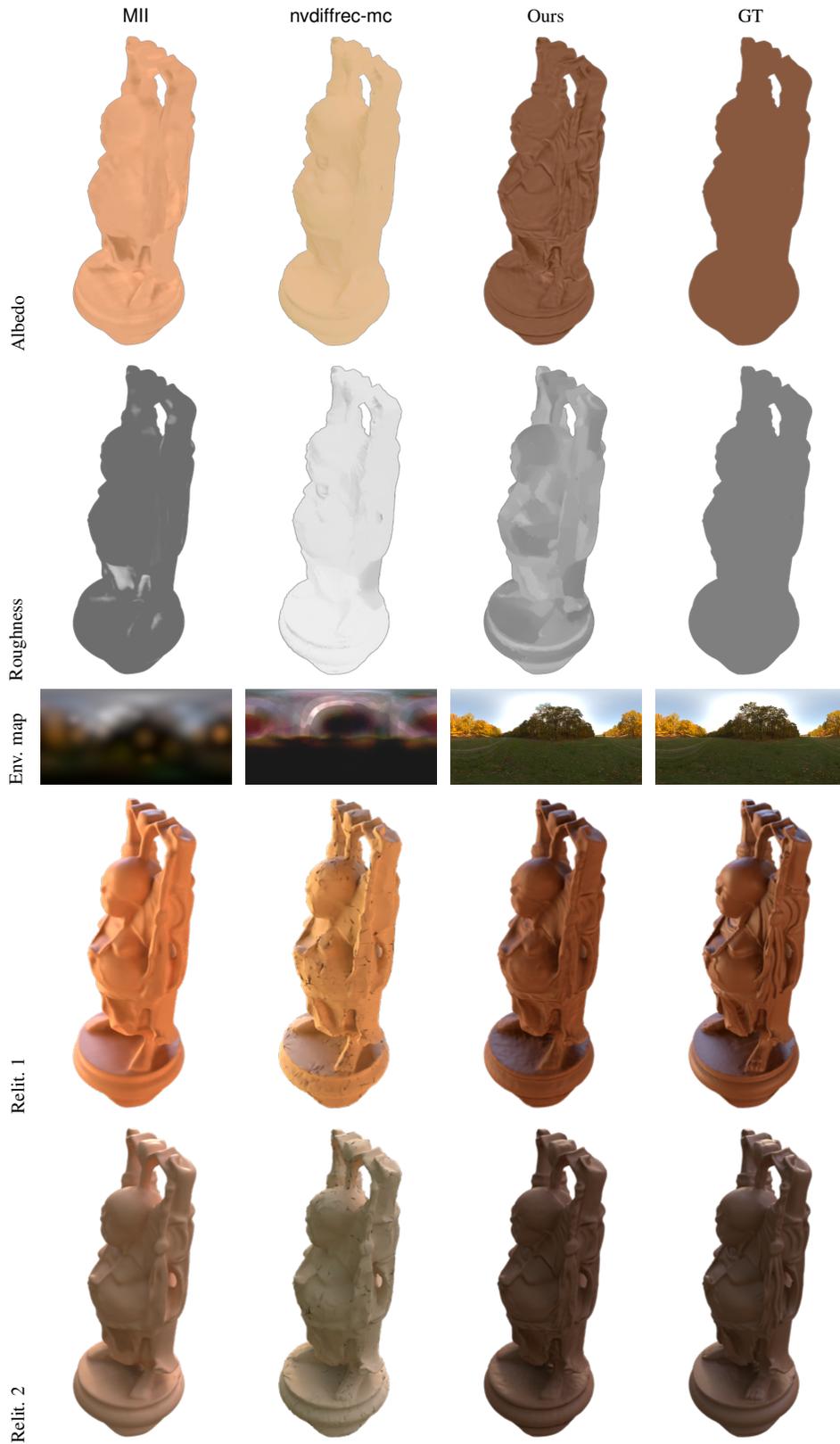


Figure 17: Qualitative comparisons of *buddha* from our dataset.



Figure 18: Qualitative comparisons of *lion* from our dataset.

Method	Time	avg.	24	37	40	55	63	65	69	83	97	105	106	110	114	118	122
COLMAP	1 h	1.36	0.81	2.05	0.73	1.22	1.79	1.58	1.02	3.05	1.40	2.05	1.00	1.32	0.49	0.78	1.17
NeuS	5.5 h	0.77	0.83	0.98	0.56	0.37	1.13	<b>0.59</b>	<b>0.60</b>	1.45	<b>0.95</b>	0.78	<b>0.52</b>	1.43	<b>0.36</b>	0.45	<b>0.45</b>
Ours	<b>5 m</b>	<b>0.66</b>	<b>0.52</b>	<b>0.72</b>	<b>0.36</b>	<b>0.35</b>	<b>0.97</b>	0.68	0.61	<b>1.27</b>	1.06	<b>0.71</b>	<b>0.52</b>	<b>0.78</b>	<b>0.36</b>	<b>0.43</b>	0.56

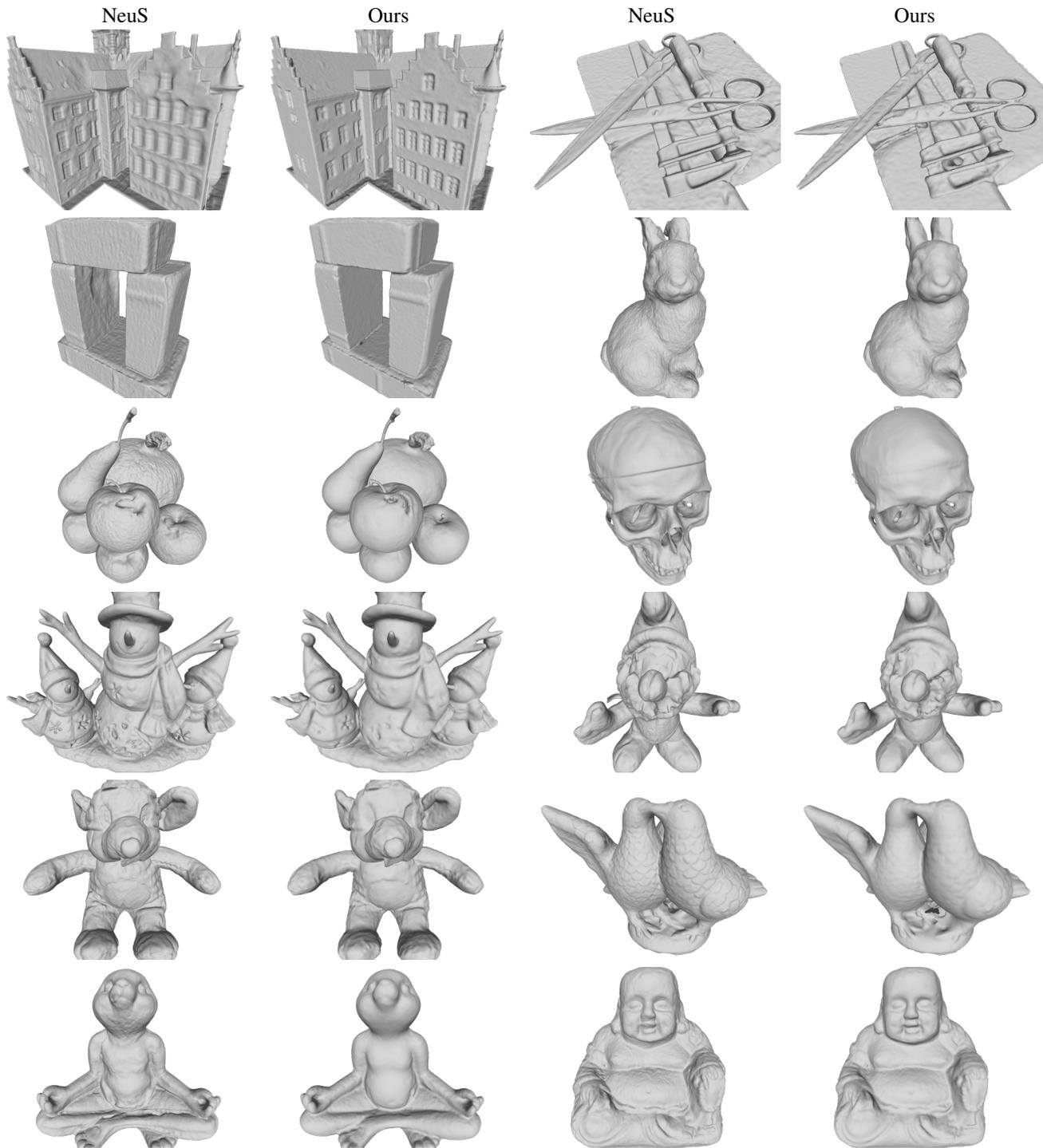


Table 7: **Quantitative results breakdown and visualization on the DTU MVS dataset [14].** We use official evaluation script to measure Chamfer distances (in mm). Our results are typically smoother with some details missing. We do not apply PBIR shape refinement as DTU exhibits significant lighting variation. See Fig. 19 and the main paper for the experiments about shape refinement.

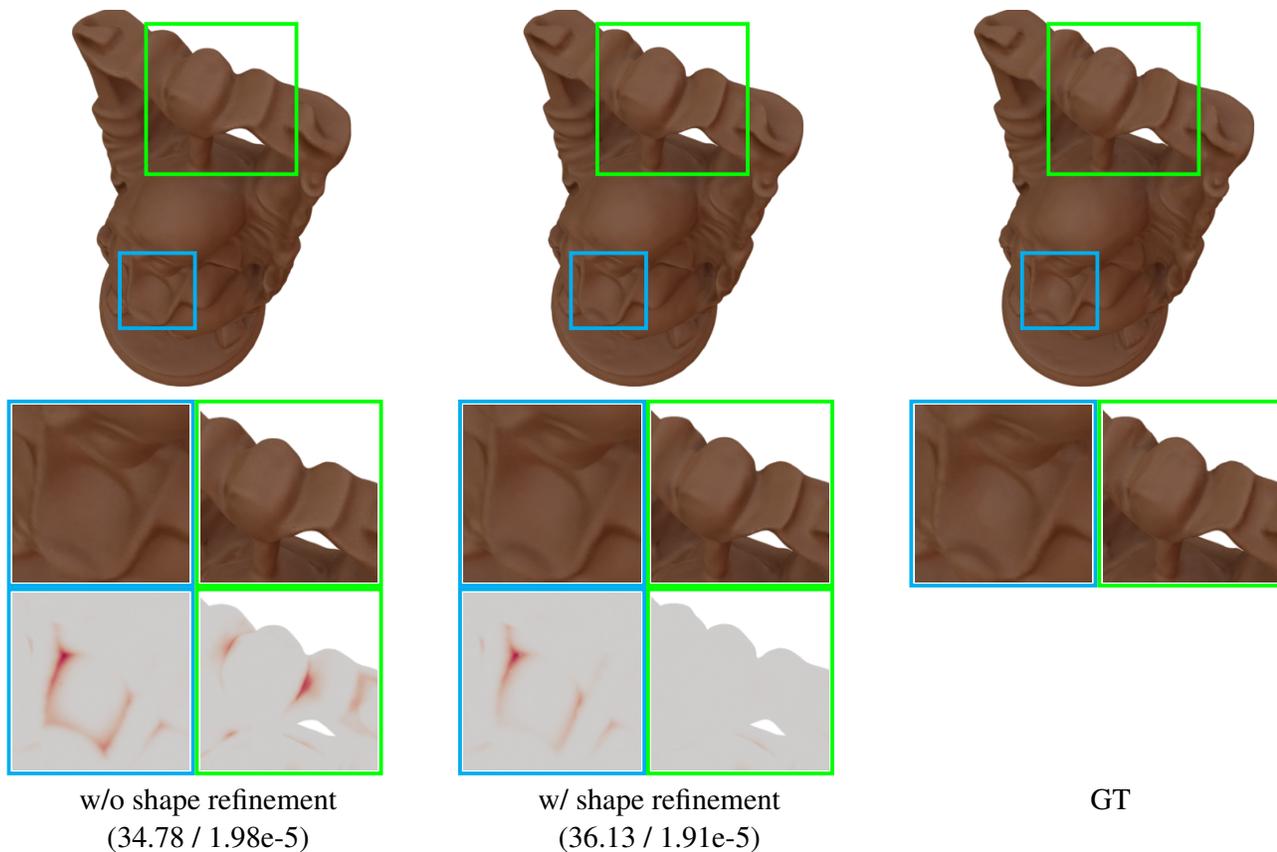


Figure 19: **Usefulness of our shape refinement in the physics-based inverse rendering (PBIR) stage.** To showcase the effectiveness of our shape refinement, we employ the *buddha* scene and present zoom-in renderings along with Chamfer distance visualizations where darker colors indicate higher errors. Additionally, we report the PSNR for relighting and the Chamfer distance, presented at the bottom of our results.

## References

- [1] Dejan Azinovic, Tzu-Mao Li, Anton Kaplanyan, and Matthias Nießner. Inverse path tracing for joint material and lighting estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2447–2456, 2019. [2](#)
- [2] Sai Bangaru, Michael Gharbi, Tzu-Mao Li, Fujun Luan, Kalyan Sunkavalli, Milos Hasan, Sai Bi, Zexiang Xu, Gilbert Bernstein, and Fredo Durand. Differentiable rendering of neural sdfs through reparameterization. In *ACM SIGGRAPH Asia 2022 Conference Proceedings*, 2022. [5](#)
- [3] Harry G Barrow, Jay M Tenenbaum, Robert C Bolles, and Helen Cf Wolf. Parametric correspondence and chamfer matching: Two new techniques for image matching. In *Proceedings: Image Understanding Workshop*, pages 21–27. Science Applications, Inc, 1977. [13](#)
- [4] Mark Boss, Raphael Braun, Varun Jampani, Jonathan T. Barron, Ce Liu, and Hendrik P. A. Lensch. Nerf: Neural reflectance decomposition from image collections. In *ICCV*, 2021. [1](#), [2](#), [3](#), [10](#)
- [5] 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 *NeurIPS*, 2021. [1](#), [2](#), [3](#)
- [6] Brent Burley and Walt Disney Animation Studios. Physically-based shading at Disney. In *Acm Siggraph*, volume 2012, pages 1–7. vol. 2012, 2012. [3](#)
- [7] Guangyan Cai, Kai Yan, Zhao Dong, Ioannis Gkioulekas, and Shuang Zhao. Physics-based inverse rendering using combined implicit and explicit geometries. In *Computer Graphics Forum*, volume 41, pages 129–138. Wiley Online Library, 2022. [2](#), [3](#)
- [8] Anpei Chen, Zexiang Xu, Andreas Geiger, Jingyi Yu, and Hao Su. Tensorf: Tensorial radiance fields. In *ECCV*, 2022. [1](#), [2](#)
- [9] Qiancheng Fu, Qingshan Xu, Yew-Soon Ong, and Wenbing Tao. Geo-neus: Geometry-consistent neural implicit surfaces learning for multi-view reconstruction. In *NeurIPS*, 2022. [2](#)
- [10] Yasutaka Furukawa and Jean Ponce. Accurate, dense, and robust multiview stereopsis. *IEEE transactions on pattern analysis and machine intelligence*, 32(8):1362–1376, 2009. [1](#)
- [11] 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](#)
- [12] Jon Hasselgren, Nikolai Hofmann, and Jacob Munkberg. Shape, light & material decomposition from images using monte carlo rendering and denoising. In *NeurIPS*, 2022. [1](#), [2](#), [3](#), [5](#), [6](#), [10](#), [12](#), [16](#)
- [13] Wenzel Jakob, Sébastien Speierer, Nicolas Roussel, and Delio Vicini. Dr.jit: A just-in-time compiler for differentiable rendering. *Transactions on Graphics (Proceedings of SIGGRAPH)*, 41(4), July 2022. [2](#), [5](#)
- [14] Rasmus Ramsbøl Jensen, Anders Lindbjerg Dahl, George Vogiatis, Engin Tola, and Henrik Aanæs. Large scale multi-view stereopsis evaluation. In *CVPR*, 2014. [8](#), [13](#), [20](#)
- [15] Brian Karis and Epic Games. Real shading in unreal engine 4. *Proc. Physically Based Shading Theory Practice*, 4(3):1, 2013. [4](#), [11](#)
- [16] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015. [11](#), [12](#)
- [17] Zhengfei Kuang, Yunzhi Zhang, Hong-Xing Yu, Samir Agarwala, Shangzhe Wu, and Jiajun Wu. Stanford-orb: A real-world 3d object inverse rendering benchmark, 2023. [9](#), [10](#)
- [18] Tzu-Mao Li, Miika Aittala, Frédo Durand, and Jaakko Lehtinen. Differentiable monte carlo ray tracing through edge sampling. *ACM Transactions on Graphics (TOG)*, 37(6):1–11, 2018. [2](#)
- [19] Fujun Luan, Shuang Zhao, Kavita Bala, and Zhao Dong. Unified shape and svbrdf recovery using differentiable monte carlo rendering. In *Computer Graphics Forum*, volume 40, pages 101–113. Wiley Online Library, 2021. [2](#), [3](#)
- [20] Julian Meder and Beat D. Brüderlin. Hemispherical gaussians for accurate light integration. In *ICCVG*, 2018. [4](#)
- [21] 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](#)
- [22] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM Trans. Graph.*, 2022. [1](#), [2](#)
- [23] Jacob Munkberg, Wenzheng Chen, Jon Hasselgren, Alex Evans, Tianchang Shen, Thomas Müller, Jun Gao, and Sanja Fidler. Extracting triangular 3d models, materials, and lighting from images. In *CVPR*, 2022. [1](#), [2](#), [3](#), [4](#), [10](#)
- [24] Baptiste Nicolet, Alec Jacobson, and Wenzel Jakob. Large steps in inverse rendering of geometry. *ACM Transactions on Graphics (Proceedings of SIGGRAPH Asia)*, 40(6), Dec. 2021. [5](#), [12](#)
- [25] Jannik Boll Nielsen, Henrik Wann Jensen, and Ravi Ramamoorthi. On optimal, minimal brdf sampling for reflectance acquisition. *ACM Transactions on Graphics (TOG)*, 34(6):1–11, 2015. [1](#)
- [26] Merlin Nimier-David, Zhao Dong, Wenzel Jakob, and Anton Kaplanyan. Material and lighting reconstruction for complex indoor scenes with texture-space differentiable rendering. 2021. [2](#)
- [27] Michael Oechsle, Songyou Peng, and Andreas Geiger. UNISURF: unifying neural implicit surfaces and radiance fields for multi-view reconstruction. In *ICCV*, 2021. [1](#)
- [28] Gyutae Park, Sungjoon Son, Jaeyoung Yoo, Seho Kim, and Nojun Kwak. Matteformer: Transformer-based image matting via prior-tokens. In *CVPR*, 2022. [4](#)
- [29] Johannes Lutz Schönberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. [1](#)
- [30] 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. [1](#)

- [31] Cheng Sun, Min Sun, and Hwann-Tzong Chen. Direct voxel grid optimization: Super-fast convergence for radiance fields reconstruction. In *CVPR*, 2022. 1, 2, 3, 4
- [32] Delio Vicini, Sébastien Speierer, and Wenzel Jakob. Differentiable signed distance function rendering. *Transactions on Graphics (Proceedings of SIGGRAPH)*, 41(4):125:1–125:18, July 2022. 5
- [33] Jiaping Wang, Peiran Ren, Minmin Gong, John M. Snyder, and Baining Guo. All-frequency rendering of dynamic, spatially-varying reflectance. *ACM Trans. Graph.*, 2009. 4
- [34] 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. In *NeurIPS*, 2021. 1, 2, 3, 8, 11
- [35] Yiming Wang, Qin Han, Marc Habermann, Kostas Daniilidis, Christian Theobalt, and Lingjie Liu. Neus2: Fast learning of neural implicit surfaces for multi-view reconstruction. *CoRR*, abs/2212.05231, 2022. 8
- [36] Haoqian Wu, Zhipeng Hu, Lincheng Li, Yongqiang Zhang, Changjie Fan, and Xin Yu. Nefii: Inverse rendering for reflectance decomposition with near-field indirect illumination. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4295–4304, June 2023. 10
- [37] Tong Wu, Jiaqi Wang, Xingang Pan, Xudong Xu, Christian Theobalt, Ziwei Liu, and Dahua Lin. Voxurf: Voxel-based efficient and accurate neural surface reconstruction. *CoRR*, abs/2208.12697, 2022. 8
- [38] Zexiang Xu, Jannik Boll Nielsen, Jiyang Yu, Henrik Wann Jensen, and Ravi Ramamoorthi. Minimal brdf sampling for two-shot near-field reflectance acquisition. *ACM Transactions on Graphics (TOG)*, 35(6):1–12, 2016. 1
- [39] Lior Yariv, Jiatao Gu, Yoni Kasten, and Yaron Lipman. Volume rendering of neural implicit surfaces. In *NeurIPS*, 2021. 1, 2
- [40] Alex Yu, Sara Fridovich-Keil, Matthew Tancik, Qinhong Chen, Benjamin Recht, and Angjoo Kanazawa. Plenoxels: Radiance fields without neural networks. In *CVPR*, 2022. 2
- [41] Zehao Yu, Songyou Peng, Michael Niemeyer, Torsten Sattler, and Andreas Geiger. Monosdf: Exploring monocular geometric cues for neural implicit surface reconstruction. In *NeurIPS*, 2022. 2
- [42] Cheng Zhang, Bailey Miller, Kai Yan, Ioannis Gkioulekas, and Shuang Zhao. Path-space differentiable rendering. *ACM Trans. Graph.*, 39(4):143:1–143:19, 2020. 2, 5
- [43] Jingyang Zhang, Yao Yao, Shiwei Li, Tian Fang, David McKeinnon, Yanghai Tsin, and Long Quan. Critical regularizations for neural surface reconstruction in the wild. In *CVPR*, 2022. 2
- [44] Kai Zhang, Fujun Luan, Qianqian Wang, Kavita Bala, and Noah Snavely. Physg: Inverse rendering with spherical gaussians for physics-based material editing and relighting. In *CVPR*, 2021. 1, 2, 3, 4, 10
- [45] Xiuming Zhang, Pratul P. Srinivasan, Boyang Deng, Paul E. Debevec, William T. Freeman, and Jonathan T. Barron. Nerfactor: neural factorization of shape and reflectance under an unknown illumination. *ACM Trans. Graph.*, 2021. 1, 2, 3, 4, 10
- [46] Yuanqing Zhang, Jiaming Sun, Xingyi He, Huan Fu, Rongfei Jia, and Xiaowei Zhou. Modeling indirect illumination for inverse rendering. In *CVPR*, 2022. 1, 2, 3, 4, 5, 6, 12, 16
- [47] Zhiming Zhou, Guojun Chen, Yue Dong, David Wipf, Yong Yu, John Snyder, and Xin Tong. Sparse-as-possible svbrdf acquisition. *ACM Transactions on Graphics (TOG)*, 35(6):1–12, 2016. 1
- [48] Jiejie Zhu, Liang Wang, Ruigang Yang, James E Davis, et al. Reliability fusion of time-of-flight depth and stereo geometry for high quality depth maps. *IEEE transactions on pattern analysis and machine intelligence*, 33(7):1400–1414, 2010. 1