

# NeRO: Neural Geometry and BRDF Reconstruction of Reflective Objects from Multiview Images

YUAN LIU, The University of Hong Kong, China

PENG WANG, The University of Hong Kong, China

CHENG LIN, Tencent Games, China

XIAOXIAO LONG, The University of Hong Kong, China

JIEPENG WANG, The University of Hong Kong, China

LINGJIE LIU, Max Planck Institute for Informatics, Germany and University of Pennsylvania, USA

TAKU KOMURA, The University of Hong Kong, China

WENPING WANG, Texas A&M University, U.S.A



Fig. 1. **NeRO**. We present NeRO for reconstructing the geometry and the BRDF of reflective objects with strong reflective appearances. NeRO only requires multiview input images of the reflective object under an unknown illumination condition. The output of NeRO is a triangular mesh with material parameters, which can easily be used in rendering software for relighting and other applications.

Authors' addresses: Yuan Liu, The University of Hong Kong, Hong Kong, China; Peng Wang, The University of Hong Kong, Hong Kong, China; Cheng Lin, Tencent Games, Shen Zhen, China; Xiaoxiao Long, The University of Hong Kong, Hong Kong, China; Jiepeng Wang, The University of Hong Kong, Hong Kong, China; Lingjie Liu, Max Planck Institute for Informatics, Germany, University of Pennsylvania, USA; Taku

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed

Komura, The University of Hong Kong, Hong Kong, China; Wenping Wang, Texas A&M University, Texas, U.S.A.

for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM

We present a neural rendering-based method called NeRO for reconstructing the geometry and the BRDF of reflective objects from multiview images captured in an unknown environment. Multiview reconstruction of reflective objects is extremely challenging because specular reflections are view-dependent and thus violate the multiview consistency, which is the cornerstone for most multiview reconstruction methods. Recent neural rendering techniques can model the interaction between environment lights and the object surfaces to fit the view-dependent reflections, thus making it possible to reconstruct reflective objects from multiview images. However, accurately modeling environment lights in the neural rendering is intractable, especially when the geometry is unknown. Most existing neural rendering methods, which can model environment lights, only consider direct lights and rely on object masks to reconstruct objects with weak specular reflections. Therefore, these methods fail to reconstruct reflective objects, especially when the object mask is not available and the object is illuminated by indirect lights. We propose a two-step approach to tackle this problem. First, by applying the split-sum approximation and the integrated directional encoding to approximate the shading effects of both direct and indirect lights, we are able to accurately reconstruct the geometry of reflective objects without any object masks. Then, with the object geometry fixed, we use more accurate sampling to recover the environment lights and the BRDF of the object. Extensive experiments demonstrate that our method is capable of accurately reconstructing the geometry and the BRDF of reflective objects from only posed RGB images without knowing the environment lights and the object masks. Codes and datasets are available at <https://github.com/liuyuan-pal/NeRO>.

CCS Concepts: • **Computing methodologies** → **Mesh geometry models**.

Additional Key Words and Phrases: neural representation, neural rendering, multiview reconstruction

#### ACM Reference Format:

Yuan Liu, Peng Wang, Cheng Lin, Xiaoxiao Long, Jiepeng Wang, Lingjie Liu, Taku Komura, and Wenping Wang. 2023. NeRO: Neural Geometry and BRDF Reconstruction of Reflective Objects from Multiview Images. *ACM Trans. Graph.* 1, 1 (May 2023), 26 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

## 1 INTRODUCTION

Multiview 3D reconstruction, a fundamental task in computer graphics and vision [Hartley and Zisserman 2003], has witnessed tremendous progress in recent years [Oechsle et al. 2021; Schönberger et al. 2016; Wang et al. 2021a,b; Yao et al. 2018; Yariv et al. 2021, 2020]. Despite the compelling results achieved, the reconstruction of reflective objects, which are frequently seen in the real-world environment, remains a challenging and outstanding problem. Reflective objects usually have glossy surfaces on which some or all of the lights that strike the object are reflected. The reflection leads to inconsistent colors when observing the objects from different views. However, most multiview reconstruction methods rely heavily on view consistency for stereo matching. This constitutes a significant barrier to the reconstruction quality of existing techniques. Fig. 2 (b) shows the reconstructions of widely-used COLMAP [Schönberger et al. 2016] on reflective objects.

As an emerging trend for multi-view reconstruction, modeling surfaces based on neural rendering exhibits a powerful ability for tackling complex objects [Oechsle et al. 2021; Wang et al. 2021b; Yariv et al. 2021, 2020]. In these so-called neural reconstruction methods, the underlying surface geometry is represented as an implicit function, e.g., a signed distance function (SDF) encoded by a multi-layer perception (MLP). To reconstruct the geometry, these methods optimize the neural implicit function by modeling the view-dependent colors and minimizing the difference between the rendered and the input images. However, neural reconstruction methods still struggle to reconstruct reflective objects. Examples are provided in Fig. 2 (c). The reason is that the color function used in these methods only correlates the color with the view direction and surface geometry, rather than explicitly considering the underlying shading mechanism for reflections. Consequently, fitting the specular color variations in different view directions on the surface leads to erroneous geometry, even with higher frequency in positional encoding, or deeper and wider MLP networks.

To address the challenging surface reflections, we propose to explicitly incorporate the formulation of the rendering equation [Kajiya 1986] into the neural reconstruction framework. The rendering equation enables us to consider the interaction between the surface Bidirectional Reflectance Distribution Function (BRDF) [Nicodemus 1965] and the environment lights. Since the appearances of reflective objects are strongly affected by the environment lights, the view-dependent specular reflection can be well-explained by the rendering equation. With the explicit rendering function, the representation ability of the existing neural reconstruction framework is substantially enhanced to capture the high-frequency specular color variations, which significantly benefits the geometry reconstruction of reflective objects.

Explicitly incorporating the rendering equation in a neural reconstruction framework is not trivial. Accurately evaluating the rendering equation on a surface point requires computing the integral of environment lights, which is intractable with unknown surface locations and unknown environment lights. In order to tractably evaluate the rendering equation, existing material estimation methods [Boss et al. 2021a,b; Hasselgren et al. 2022; Munkberg et al. 2022; Verbin et al. 2022; Zhang et al. 2021a,b, 2022b] strongly rely on object masks to obtain a correct surface reconstruction and are mainly designed for material estimation of objects without strong specular reflections, which perform much worse on reflective objects as shown in Fig. 2 (d,e). Moreover, most of these methods further simplify the rendering process to only consider the lights from distant regions (*direct lights*) [Boss et al. 2021a,b; Munkberg et al. 2022; Verbin et al. 2022; Zhang et al. 2021a], which thus struggle to reconstruct surfaces illuminated by reflected lights from the object itself or nearby regions (*indirect lights*). Although there are methods [Hasselgren et al. 2022; Zhang et al. 2021b, 2022b] considering indirect lights in the rendering, they either require a reconstructed radiance field with known geometry [Zhang et al. 2021b, 2022b] or only use very few ray samples to compute the lights [Hasselgren et al. 2022], which results in unstable convergence on reflective objects or additional dependence on object masks. Thus, considering both direct and indirect lights to correctly reconstruct the unknown surfaces of reflective objects is still challenging.

must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

© 2023 Association for Computing Machinery.  
0730-0301/2023/5-ART \$15.00  
<https://doi.org/10.1145/nnnnnnn.nnnnnnn>



Fig. 2. We apply different multiview reconstruction methods to reconstructing reflective objects. (a) The input images. Reconstruction results of (b) COLMAP [Schönberger et al. 2016], (c) NeuS [Wang et al. 2021b], (d) Ref-NeRF [Verbin et al. 2022], (e) NDRMC [Hasselgren et al. 2022] and (f) our method. In images of column (a), we use red bounding boxes to indicate regions illuminated by indirect lights. \*NDRMC uses ground-truth object masks for training while all the other methods are trained without object masks.

By incorporating the rendering equation in a neural reconstruction framework, we propose a method called NeRO for reconstructing both the geometry and the BRDF of reflective objects from only posed RGB images. The key component of NeRO is a novel light representation. In this light representation, we use two individual MLPs to encode the radiance of direct and indirect lights respectively and compute an occlusion probability to determine whether direct or indirect lights should be used in the rendering. Such a light representation efficiently accommodates both direct lights and indirect lights for accurate surface reconstruction of reflective objects. Based on the proposed light representation, NeRO adopts a two-stage strategy for a tractable evaluation of the rendering equation in the neural reconstruction. The first stage of NeRO employs a split-sum approximation and the integrated directional encoding [Verbin et al. 2022] to evaluate the rendering equation, which produces accurate geometry reconstruction with compromised environment lights and surface BRDF estimation. Then, with the reconstructed geometry fixed, the second stage of NeRO improves the estimated BRDF by more accurately evaluating the rendering equation with Monte Carlo sampling. With the light representation and the two-stage design, the proposed method essentially extends the representational power of neural rendering methods on reflective objects, making it achieve the full potential of learning geometric surfaces.

To evaluate the performance of NeRO, we introduce a synthetic dataset and a real dataset, both of which contain reflective objects

illuminated by complex environment lights. On both datasets, NeRO successfully reconstructs both the geometry and the surface BRDF of reflective objects, on which the baseline MVS methods and neural reconstruction methods fail. The output of our method is a triangular mesh with the estimated BRDF parameters, which can be easily used in downstream applications such as relighting.

## 2 RELATED WORKS

### 2.1 Multiview 3D reconstruction

Multiview 3D reconstruction or Multiview Stereo (MVS) has been studied for decades [Campbell et al. 2008; Furukawa and Ponce 2009; Strecha et al. 2006]. Traditional multiview reconstruction methods mainly rely on the multiview consistency of 3D points to build correspondences and estimate the depth values on different views [Barron and Poole 2016; Bleyer et al. 2011; Campbell et al. 2008; Furukawa and Ponce 2009; Gallup et al. 2007; Hosni et al. 2012; Richardt et al. 2010; Schönberger et al. 2016; Strecha et al. 2006]. With the advances of deep learning techniques, many recent works [Cheng et al. 2020; Wang et al. 2021a; Yan et al. 2020; Yang et al. 2020; Yao et al. 2018] try to introduce neural networks to estimate correspondences for the MVS task, which demonstrates impressive reconstruction quality on widely-used benchmarks [Geiger et al. 2013; Jensen et al. 2014; Scharstein and Szeliski 2002]. In this paper, we aim to reconstruct reflective objects with strong specular reflections. The strong specular reflections violate the multiview

consistency so these correspondence-based methods do not perform well on reflective objects.

**Neural surface reconstruction.** Neural rendering and neural representations [Mildenhall et al. 2020; Park et al. 2019; Sitzmann et al. 2019; Tewari et al. 2020, 2022] have attracted much attention due to their strong representation ability and impressive improvements on the novel-view-synthesis task. DVR [Niemeyer et al. 2020] first introduces the neural rendering and neural surface representation in the multiview reconstruction. IDR [Yariv et al. 2020] improves the reconstruction quality with the differentiable sphere tracing and the Eikonal regularization [Gropp et al. 2020]. UNISURF [Oechsle et al. 2021], VolSDF [Yariv et al. 2021] and NeuS [Wang et al. 2021b] introduce the differentiable volume rendering in the multiview surface reconstruction with improved robustness and quality. Subsequent works improve the volume-rendering-based multiview reconstruction framework in various aspects, such as introducing Manhattan or normal priors [Guo et al. 2022b; Wang et al. 2022c], utilizing symmetry [Insafutdinov et al. 2022; Zhang et al. 2021c], extracting image features [Darmon et al. 2022; Long et al. 2022], improving fidelity [Fu et al. 2022; Wang et al. 2022b] and efficiency [Li et al. 2022; Sun et al. 2022; Wang et al. 2022a; Wu et al. 2022; Zhao et al. 2022a]. Similar to these works, we also follow the volume rendering framework for surface reconstruction but we focus on reconstructing reflective objects with strong specular reflections, an outstanding problem that has not been explored by existing neural reconstruction methods.

## 2.2 Reflective object reconstruction

Only a few works try to reconstruct reflective objects in the multiview stereo setting by using additional object masks [Godard et al. 2015] or removing the reflections [Wu et al. 2018]. Other than the uncontrolled multiview reconstruction, some works [Han et al. 2016; Roth and Black 2006] resort to constrained settings with known specular flows [Roth and Black 2006] or known environments [Han et al. 2016] for the reconstruction of ideal mirror-like objects. Some other works utilize additional ray information by encoding rays [Tin et al. 2016] or utilizing polarization images [Dave et al. 2022; Kadambi et al. 2015; Rahmann and Canterakis 2001] to reconstruct the objects with specular reflections. [Whelan et al. 2018] reconstructs mirror planes in a scene by utilizing reflected images of the scanner. These methods are limited to a relatively strict setting with specially-designed capturing devices. In contrast, we aim to directly reconstruct the reflective objects from posed multiview images, which can be easily captured using a cellphone camera.

Some image-based rendering methods [Rodriguez et al. 2020; Sinha et al. 2012] are specially designed for glossy or reflective objects for the NVS task. NeRFren [Guo et al. 2022a] reconstructs a neural density field of a scene with the existence of mirror-like planes. Neural Point Catacaustics [Kopanas et al. 2022] applies a warp field to improve the rendering quality on reflective objects. Ref-NeRF [Verbin et al. 2022] proposes integrated direction encoding (IDE) to improve the NVS quality on reflective materials. Our method incorporates the IDE in reconstructing reflective objects with a neural SDF for the surface reconstruction. A concurrent work ORCA [Tiwary et al. 2022] extends to reconstruct the radiance field

of a scene from the reflections on a glossy object, which also reconstructs the object in the pipeline. Since the target of ORCA is mainly to reconstruct the radiance field of the scene, it relies on object masks for the reconstruction of the reflective objects. In comparison, our method does not require object masks and our main target is to reconstruct the geometry and BRDF of the object.

## 2.3 BRDF estimation

Estimating the surface BRDF from images is mainly based on the inverse rendering techniques [Barron and Malik 2014; Nimier-David et al. 2019]. Some methods [Gao et al. 2019; Guo et al. 2020; Li et al. 2020, 2018; Wimbauer et al. 2022; Ye et al. 2022] rely on an object prior or a scene prior to directly estimating BRDF and lighting. Differentiable renderers [Chen et al. 2019, 2021; Kato et al. 2018; Liu et al. 2019; Nimier-David et al. 2019] allow direct optimization of the BRDF from image losses. To enable more accurate BRDF estimation, most methods [Bi et al. 2020, [n.d.]; Cheng et al. 2021; Kuang et al. 2022; Li and Li 2022a,b; Nam et al. 2018; Schmitt et al. 2020; Yang et al. 2022a,b; Zhang et al. 2022a] require multiple images of the object to be illuminated by different collocated flashlights. In this paper, we estimate the BRDF in a static scene with moving cameras, which is also the setting adopted by [Boss et al. 2021a, 2022, 2021b; Deng et al. 2022; Hasselgren et al. 2022; Munkberg et al. 2022; Zhang et al. 2021a,b, 2022b]. Among these works, PhySG [Zhang et al. 2021a], NeRD [Boss et al. 2021a], Neural-PIL [Boss et al. 2021b] and NDR [Munkberg et al. 2022] consider the interaction between the direct environment lights and the surfaces for the BRDF estimation. Subsequent works MII [Zhang et al. 2022b], NDRMC [Hasselgren et al. 2022], DIP [Deng et al. 2022] and NeLF [Yao et al. 2022] add indirect lights, which improves the quality of estimated BRDF. These methods mainly aim to reconstruct the BRDF of common objects without too many specular reflections, which produces low-quality BRDF on reflective objects. Some other methods [Chen and Liu 2022; Duchêne et al. 2015; Gao et al. 2020; Liu et al. 2021; Lyu et al. 2022; Nestmeyer et al. 2020; Philip et al. 2019, 2021; Rudnev et al. 2022; Shih et al. 2013; Yu et al. 2020; Yu and Smith 2019; Zhao et al. 2022b; Zheng et al. 2021] are mainly targeted to the relighting task but not designed for reconstructing the surface geometry or the BRDF. NeLF [Yao et al. 2022] is the most similar work to the Stage II of our method, both of which fix the geometry to optimize BRDF with MC sampling. However, NeLF does not use importance sampling on the specular lobe and simply predicts lights from a position and a direction without considering the occlusions. In comparison, our method explicitly distinguishes direct and indirect lights and uses importance sampling on diffuse and specular lobes for a better BRDF estimation on reflective objects.

# 3 METHOD

## 3.1 Overview

Given a set of RGB images with known camera poses as input, our target is to reconstruct the surface and BRDF of the reflective object in the images. Note that our method does not require knowing the object masks or environment lights. The pipeline of NeRO consists of two stages. In Stage I (Sec. 3.3), we reconstruct the geometry of the reflective object by optimizing a neural SDF with the volume

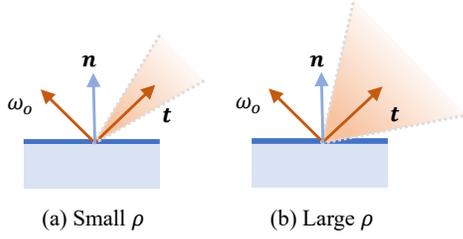


Fig. 3. **Specular lobe** in Eq. 7 is determined by the roughness  $\rho$  and the reflective direction  $\mathbf{t}$ . (a) A smooth surface with small  $\rho$  has a smaller specular lobe while (b) a rougher surface with large  $\rho$  has a larger specular lobe.

rendering, in which approximate direct and indirect lights are estimated to model the view-dependent specular colors. In Stage II (Sec. 3.4), we fix the object geometry and finetune the direct and indirect lights to compute an accurate BRDF of the reflective object. In the following, we begin with a brief review of NeuS [Wang et al. 2021b] and a micro-facet BRDF model [Cook and Torrance 1982; Torrance and Sparrow 1967].

### 3.2 Preliminaries

**NeuS [Wang et al. 2021b] for surface reconstruction.** We follow NeuS to represent the object surface by an SDF encoded by an MLP network  $g_{\text{sdf}}(\mathbf{x})$ . The surface is the zero-level set with  $\{\mathbf{x} \in \mathbb{R}^3 | g_{\text{sdf}}(\mathbf{x}) = 0\}$ . Then, volume rendering [Mildenhall et al. 2020] is applied to render images from the neural SDF. Given a camera ray  $\mathbf{o} + t\mathbf{v}$  emitting from the camera center  $\mathbf{o}$  to the space along the direction  $\mathbf{v}$ , we sample  $n$  points on the ray  $\{\mathbf{p}_j = \mathbf{o} + t_j\mathbf{v} | t_j > 0, t_{j-1} < t_j\}$ . Then, the rendered color for this camera ray is computed by

$$\hat{\mathbf{c}} = \sum_n w_j \mathbf{c}_j, \quad (1)$$

where  $w_j$  is the weight for the  $j$ -th point, which is derived from the SDF value via the opaque density proposed in [Wang et al. 2021b].  $\mathbf{c}_j$  is the color for this point, which is decoded from an MLP network by  $\mathbf{c}_j = g_{\text{color}}(\mathbf{p}_j, \mathbf{v})$  in NeuS. Then, by minimizing the difference between the rendered color  $\hat{\mathbf{c}}$  and the input ground-truth color  $\mathbf{c}$ , the parameters of two MLP networks  $g_{\text{sdf}}$  and  $g_{\text{color}}$  are learned. The reconstructed surface is extracted from the zero-level set of  $g_{\text{sdf}}$ . To enable the color function to correctly represent the specular colors on the reflective surfaces, NeRO replaces the color function of NeuS with the shading function using a Micro-facet BRDF.

**Micro-facet BRDF [Cook and Torrance 1982].** On the point  $\mathbf{p}_j$ , we compute its colors  $\mathbf{c}_j$  by the rendering equation (the subscript  $j$  is omitted in the following discussion for simplicity)

$$\mathbf{c}(\omega_o) = \int_{\Omega} L(\omega_i) f(\omega_i, \omega_o) (\omega_i \cdot \mathbf{n}) d\omega_i, \quad (2)$$

where  $\omega_o = -\mathbf{v}$  is the outgoing view direction,  $\mathbf{c}(\omega_o)$  is the output color  $\mathbf{c}_j$  for this point  $\mathbf{p}_j$  in the view direction  $\omega_o$ ,  $\mathbf{n}$  is the surface normal,  $\omega_i$  is the input light direction on the upper half sphere  $\Omega$ ,  $f(\omega_i, \omega_o) \in [0, 1]^3$  is the BRDF function,  $L(\omega_i) \in [0, +\infty)^3$  is the radiance of incoming lights. In NeRO, the normal  $\mathbf{n}$  is computed from the gradient of the SDF. The BRDF function consists of a diffuse

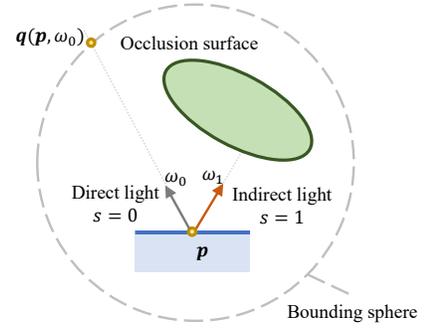


Fig. 4. **Direct and indirect lights for a point  $\mathbf{p}$ .** The direct light in the direction  $\omega_0$  is not occluded while the indirect light in the direction  $\omega_1$  is occluded by surfaces inside the unit sphere.  $s$  is the occlusion probability.  $\mathbf{q}(\mathbf{p}, \omega)$  is the intersection point on the bounding sphere of the ray emitting from  $\mathbf{p}$  along the direction  $\omega$ .

and a specular part

$$f(\omega_i, \omega_o) = \underbrace{(1-m)\frac{\mathbf{a}}{\pi}}_{\text{diffuse}} + \underbrace{\frac{DFG}{4(\omega_i \cdot \mathbf{n})(\omega_o \cdot \mathbf{n})}}_{\text{specular}}, \quad (3)$$

where  $m \in [0, 1]$  is the metalness of the point,  $1-m$  is the weight for the diffuse part,  $\mathbf{a} \in [0, 1]^3$  is the albedo color of the point,  $D$  is the normal distribution function,  $F$  is the Fresnel term and  $G$  is the geometry term.  $D$ ,  $F$  and  $G$  are all determined by the metalness  $m$ , the roughness  $\rho \in [0, 1]$  and the albedo  $\mathbf{a}$ . The expressions of  $D$ ,  $F$ , and  $G$  can be found in Sec. A.1 of the Appendix. In summary, the BRDF of the point is specified by the metalness  $m$ , the roughness  $\rho$ , and the albedo  $\mathbf{a}$ , all of which are predicted by a material MLP  $g_{\text{material}}$  in NeRO, i.e.,  $[m, \rho, \mathbf{a}] = g_{\text{material}}(\mathbf{p})$ .

Combining Eq. 2 and Eq. 3, we have

$$\mathbf{c}(\omega_o) = \mathbf{c}_{\text{diffuse}} + \mathbf{c}_{\text{specular}}, \quad (4)$$

$$\mathbf{c}_{\text{diffuse}} = \int_{\Omega} (1-m)\frac{\mathbf{a}}{\pi} L(\omega_i) (\omega_i \cdot \mathbf{n}) d\omega_i, \quad (5)$$

$$\mathbf{c}_{\text{specular}} = \int_{\Omega} \frac{DFG}{4(\omega_i \cdot \mathbf{n})(\omega_o \cdot \mathbf{n})} L(\omega_i) (\omega_i \cdot \mathbf{n}) d\omega_i. \quad (6)$$

As explained before, accurately evaluating the integrals of Eq. 5 and Eq. 6 for every sample point in the volume rendering is intractable. Therefore, we propose a two-step framework to approximately compute these two integrals. In the first stage, our priority is to faithfully reconstruct the geometric surface.

### 3.3 Stage I: Geometry reconstruction

In order to reconstruct surfaces of reflective objects, we adopt the same neural SDF representation and the volume rendering algorithm (Eq. 1) as NeuS [Wang et al. 2021b] but with a different color function. In NeRO, we predict a metalness  $m$ , a roughness  $\rho$  and an albedo  $\mathbf{a}$  to compute a color  $\mathbf{c}_j$  (i.e.,  $\mathbf{c}(\omega_o)$ ) using the micro-facet BRDF (Eq. 4-6). To make the computation tractable in the volume rendering of NeuS, we adopt the split-sum approximation [Karis and Games 2013], which separates the integral of the product of lights and BRDFs into two individual integrals.

**Split-sum approximation.** For the specular color in Eq. 6, we follow [Karis and Games 2013] to approximate it with

$$\mathbf{c}_{\text{specular}} \approx \underbrace{\int_{\Omega} L(\omega_i) D(\rho, \mathbf{t}) d\omega_i}_{L_{\text{specular}}} \cdot \underbrace{\int_{\Omega} \frac{DFG}{4(\omega_o \cdot \mathbf{n})} d\omega_i}_{M_{\text{specular}}}, \quad (7)$$

where  $L_{\text{specular}}$  is the integral of lights on the normal distribution function  $D(\rho, \mathbf{t}) \in [0, 1]$  (also called the *specular lobe*) as shown in Fig. 3,  $\mathbf{t}$  is the reflective direction,  $M_{\text{specular}}$  denotes the integral of BRDF. Note that a rougher surface has a larger specular lobe while a smoother surface has a smaller lobe. The integral of BRDF can be directly computed by  $M_{\text{specular}} = ((1 - m) * 0.04 + m * \mathbf{a}) * F_1 + F_2$ , where  $F_1$  and  $F_2$  are two pre-computed scalars depending on the roughness  $\rho$ , the view direction  $\omega_o$  and the normal  $\mathbf{n}$  as introduced in more details in Sec. A.1 in Appendix. The diffuse color in Eq. 5 can be written as

$$\mathbf{c}_{\text{diffuse}} = \mathbf{a}(1 - m) \underbrace{\int_{\Omega} L(\omega_i) \frac{\omega_i \cdot \mathbf{n}}{\pi} d\omega_i}_{L_{\text{diffuse}}}, \quad (8)$$

where we use  $L_{\text{diffuse}}$  to denote the diffuse light integral.

Split-sum has already been proven to be a good approximation and is widely used in real-time rendering. With predicted material parameters albedo  $\mathbf{a}$ , roughness  $\rho$  and metalness  $m$  from the material MLP, the only two unknowns in Eq. 8 and Eq. 7 are light integrals  $L_{\text{diffuse}}$  and  $L_{\text{specular}}$ . However, to compute the light integrals, we do not prefilter environment lights as previous methods [Boss et al. 2021b; Munkberg et al. 2022] but use the integrated directional encoding [Verbin et al. 2022]. In the following, we first introduce our light representation  $L(\omega_i)$ .

**Light representation.** In NeRO, we define a *bounding sphere* around the object to build the neural SDF. Since we only aim to reconstruct the surfaces inside the bounding sphere, we call all lights coming from the outside of the bounding sphere as *direct lights* while we name lights reflected by the surfaces inside the bounding sphere as *indirect lights*, which is shown in Fig. 4. Then, we represent the light  $L(\omega_i)$  by

$$L(\omega_i) = [1 - s(\omega_i)]g_{\text{direct}}(SH(\omega_i)) + s(\omega_i)g_{\text{indirect}}(SH(\omega_i), \mathbf{p}), \quad (9)$$

where  $g_{\text{direct}}$  and  $g_{\text{indirect}}$  are two MLPs for the direct lights and the indirect lights respectively,  $s(\omega_i) \in [0, 1]$  is the occlusion probability that the ray from the point  $\mathbf{p}$  to the direction  $\omega_i$  is occluded by the surfaces inside the bounding sphere,  $SH$  is the directional encoding using spherical harmonics (SH) as basis functions. Note that  $s(\omega_i) = g_{\text{occ}}(SH(\omega_i), \mathbf{p})$  is also predicted by an MLP  $g_{\text{occ}}$ .

**Motivation of the light representation design.** The direct light  $g_{\text{direct}}(\omega_i)$  only depends on the direction  $\omega_i$  so that all points are illuminated by the same direct environment light. This provides a strong global prior to explaining the view-dependent colors of the reflective object. The indirect light  $g_{\text{indirect}}(\omega_i, \mathbf{p})$  additionally takes the point position  $\mathbf{p}$  as input because indirect lights vary in the space. Additional discussion about the light representation is provided in Sec. A.2 of the Appendix.

**Light integral approximation.** We use the integrated directional encoding to approximate the light integrals  $L_{\text{specular}}$  and

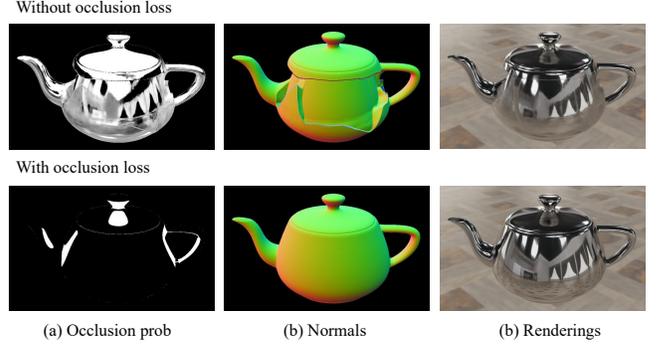


Fig. 5. **Effects of occlusion loss.** (Top) Without  $\ell_{\text{occ}}$ , the occlusion probability predicted by  $g_{\text{occ}}$  will be completely inconsistent with the reconstructed geometry and causes incorrect reconstruction. (Bottom) With  $\ell_{\text{occ}}$ , the predicted occlusion probability is accurate and the reconstruction is correct.

$L_{\text{diffuse}}$ . The specular light integral is

$$\begin{aligned} L_{\text{specular}} &\approx [1 - s(\mathbf{t})] \int_{\Omega} g_{\text{direct}}(SH(\omega_i)) D(\rho, \mathbf{t}) d\omega_i + \\ &\quad s(\mathbf{t}) \int_{\Omega} g_{\text{indirect}}(SH(\omega_i), \mathbf{p}) D(\rho, \mathbf{t}) d\omega_i \\ &\approx [1 - s(\mathbf{t})] g_{\text{direct}} \left( \int_{\Omega} SH(\omega_i) D(\rho, \mathbf{t}) d\omega_i \right) + \\ &\quad s(\mathbf{t}) g_{\text{indirect}} \left( \int_{\Omega} SH(\omega_i) D(\rho, \mathbf{t}) d\omega_i, \mathbf{p} \right). \end{aligned} \quad (10)$$

In the first approximation, we use the occlusion probability  $s(\mathbf{t})$  of the reflective direction  $\mathbf{t}$  to replace occlusion probabilities  $s(\omega_i)$  of different rays. In the second approximation, we exchange the order of the MLP and the integral. Discussion about the rationale of two approximations is provided in Sec. A.3. With Eq. 10, we only need to evaluate the MLP networks  $g_{\text{direct}}$  and  $g_{\text{indirect}}$  once on the integrated directional encoding  $\int_{\Omega} SH(\omega_i) D(\rho, \mathbf{t}) d\omega_i$ . By choosing the normal distribution function  $D$  to be a von Mises–Fisher (vMF) distribution (Gaussian distribution on a sphere), Ref-NeRF [Verbin et al. 2022] has shown that  $\int_{\Omega} SH(\omega_i) D(\rho, \mathbf{t}) d\omega_i$  has an approximated closed-form solution. In this case, we use this closed-form solution here to approximate the integral of lights.

Similarly, for  $L_{\text{diffuse}}$ , the cosine lobe  $\frac{\omega_i \cdot \mathbf{n}}{\pi}$  is also a probability distribution, which can be approximated by

$$\frac{\omega_i \cdot \mathbf{n}}{\pi} \approx D(1.0, \mathbf{n}). \quad (11)$$

Thus, we can also compute the diffuse light integral as Eq. 10. With the integrals of lights, we are able to compute the diffuse colors Eq. 5 and specular colors Eq. 6 and composite them as the color for each sample point. Note that both the split-sum approximation and the light integral approximation are only used in Stage I to enable tractable computation and will be replaced by more accurate Monte Carlo sampling in Stage II.

**Occlusion loss.** In the light representation, we use an occlusion probability  $s$  predicted by an MLP  $g_{\text{occ}}$  to determine whether direct lights or indirect lights will be used in rendering. However, as shown in Fig. 5, if we put no constraint on the occlusion probability  $s$  and

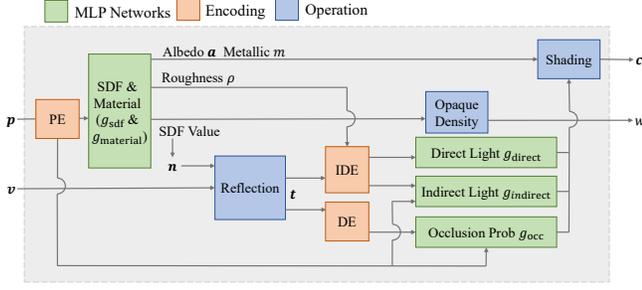


Fig. 6. **Architecture of networks in Stage I.** “PE” is positional encoding [Mildenhall et al. 2020] while “IDE” and “DE” are integrated direction encoding [Verbin et al. 2022] and vanilla directional encoding, respectively.

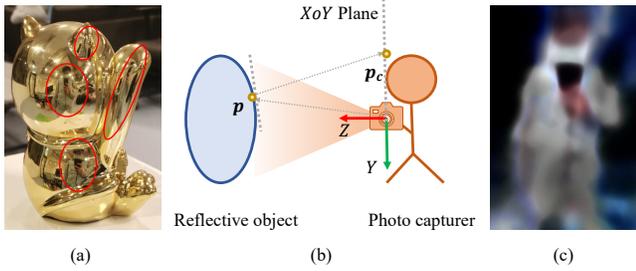


Fig. 7. (a) A reflective object shows the reflections of the photo capturer. (b) We build a 2D NeRF on the  $XoY$  plane to model the lights occluded by the photo capturer. (c) The estimated 2D radiance field of the capturer.

just let the network  $g_{occ}$  learn  $s$  from rendering loss, the predicted occlusion probability will be completely inconsistent with the reconstructed geometry and cause unstable convergence. Thus, we use the neural SDF to constrain the predicted occlusion probability. Given the ray emitting from a sample point  $\mathbf{p}$  to its reflective direction  $\mathbf{t}$ , we compute its occlusion probability  $s_{march}$  by ray-marching in the neural SDF  $g_{SDF}$  and enforce the consistency between the computed probability  $s_{march}$  and the predicted probability  $s$  with

$$\ell_{occ} = \|s_{march} - s\|_1, \quad (12)$$

where  $\ell_{occ}$  is the loss for this occlusion probability regularization.

**Training Losses.** Based on the volume rendering Eq. 1, we compute the color for the camera ray and compute the Charbonier loss [Barron et al. 2022; Charbonnier et al. 1994] between the rendered color and the input ground-truth color as the rendering loss  $\ell_{render}$ . Meanwhile, we observe that the first few training steps for the SDF are not stable, which either extremely enlarges the surface or squashes the surface too small. A stabilization regularization loss  $\ell_{stable}$  is applied for the first 1k steps. (This is discussed in detail in Sec. A.4.) In summary, the final loss is

$$\ell = \ell_{render} + \lambda_{eikonal} \ell_{eikonal} + \lambda_{occ} \ell_{occ} + \mathbb{1}(\text{step} < 1000) \ell_{stable}, \quad (13)$$

where we also adopt the Eikonal loss [Gropp et al. 2020] to regularize the norms of SDF gradients to be 1,  $\mathbb{1}$  is the indicator function,  $\lambda_{eikonal}$  and  $\lambda_{occ}$  are two predefined scalars. The overview of the network architecture is shown in Fig. 6.

**Reflection of the capturer.** We assume a static illumination environment in our model. However, in reality, there is always a person holding a camera to capture images around a reflective object. The moving person will be visible in the reflection of the object, thus violating the assumption of static illumination, as shown in the red circle of Fig. 7 (a). Since the photo capturer is relatively static to the camera, we build a 2D NeRF on the  $XoY$  plane in the camera coordinate system as shown in Fig. 7 (b). In the computation of the direct light  $g_{direct}$  in Eq. 9, we additionally check whether the ray hits the  $XoY$  plane of the camera coordinate system or not. If a hit point  $\mathbf{p}_c$  exists, we will use an MLP  $g_{camera}$  to compute an alpha value  $\alpha_{camera}$  and a color  $\mathbf{c}_{camera}$  by

$$[\alpha_{camera}, \mathbf{c}_{camera}] = g_{camera}(\mathbf{p}_c), \quad (14)$$

where  $\alpha_{camera}$  indicates whether the ray is occluded by the capturer or not while the  $\mathbf{c}_{camera}$  represents the color of the capturer on this point. Then, the direct light is  $(1 - \alpha_{camera})g_{direct}(\omega_i) + \alpha_{camera}\mathbf{c}_{camera}$ .

### 3.4 Stage II: BRDF estimation

So far, after Stage I, we have faithfully reconstructed the geometry of the reflective object but obtained only a rough BRDF estimation, which needs to be further refined. In Stage II, we aim to accurately evaluate the rendering equation so as to precisely estimate the surface BRDF i.e. metalness  $m$ , albedo  $\mathbf{a}$ , and  $\rho$ . With the fixed geometry from Stage I, we only need to evaluate the rendering equation on surface points. Thus, it is now feasible to apply Monte Carlo sampling to compute the diffuse colors in Eq. 5 and the specular colors in Eq. 6. In the MC sampling, we conduct importance sampling on both the diffuse lobe and the specular lobe as follows.

**Importance sampling.** In Monte Carlo sampling, the diffuse color  $\mathbf{c}_{diffuse}$  is computed by sampling  $N_d$  rays with a cosine-weighted hemisphere probability

$$\mathbf{c}_{diffuse} = \frac{1}{N_d} \sum_i^{N_d} (1 - m) \mathbf{a} L(\omega_i), \quad (15)$$

where  $i$  is  $i$ -th sample ray and  $\omega_i$  is the direction of this sample ray. For the specular color  $\mathbf{c}_{specular}$ , we apply the GGX distribution as normal distribution  $D$ . Then, the specular color  $\mathbf{c}_{specular}$  is computed by sampling  $N_s$  rays with DDX distribution [Cook and Torrance 1982]

$$\mathbf{c}_{specular} = \frac{1}{N_s} \sum_i^{N_s} \frac{FG(\omega_o \cdot \mathbf{h})}{(\mathbf{n} \cdot \mathbf{h})(\mathbf{n} \cdot \omega_o)} L(\omega_i), \quad (16)$$

where  $\mathbf{h}$  is the half-way vector between  $\omega_i$  and  $\omega_o$ . To evaluate Eq. 15 and Eq. 16, we still use the same material MLP  $[m, \rho, \mathbf{a}] = g_{material}$  as Stage I to compute the metalness  $m$ , roughness  $\rho$  and albedo  $\mathbf{a}$ . The light representation  $L(\omega_i)$  in Stage II is also Eq. 9 with the reflected lights of the capturer in Eq. 14, which is the same as Stage I. Since the geometry is fixed, we directly compute the occlusion probability  $s$  by tracing rays in the given geometry instead of predicting it from the MLP  $g_{occ}$ . Meanwhile, for real data, we add the intersection point on the bounding sphere  $\mathbf{q}_{p,\omega}$  of the ray emitting from  $\mathbf{p}$  along the direction  $\omega$ , as shown in Fig. 4, as an additional input to the direct light MLP  $g_{direct}$ . More details are discussed in Sec. A.2

**Regularization terms.** We follow previous works [Hasselgren et al. 2022; Munkberg et al. 2022] to impose two regularization terms on the final loss. The first is a smoothness regularization  $\ell_{\text{smooth}}$

$$\ell_{\text{smooth}} = \|g_{\text{material}}(\mathbf{p}) - g_{\text{material}}(\mathbf{p} + \epsilon)\|_2, \quad (17)$$

where  $\epsilon = 5e - 3$ . The  $\ell_{\text{smooth}}$  makes the predicted materials (roughness, metallic, and albedo) more smooth in the space. Furthermore, a light regularization  $\ell_{\text{light}}$  is imposed to make the diffuse lights  $L_{\text{diffuse}} = \frac{1}{N_d} \sum_i L(\omega_i)$  to be neutral white lighting. We compute the mean of RGB values of the diffuse light and minimize the difference between the RGB values and their mean,

$$\ell_{\text{light}} = \sum_c ([L_{\text{diffuse}}]_c - \frac{1}{3} \sum_c [L_{\text{diffuse}}]_c)^2, \quad (18)$$

where  $c$  is the index of RGB channel of  $L_{\text{diffuse}}$ . Along with the rendering loss, the final loss for Stage II is

$$\ell = \ell_{\text{render}} + \lambda_{\text{smooth}} \ell_{\text{smooth}} + \lambda_{\text{light}} \ell_{\text{light}}, \quad (19)$$

where  $\lambda_{\text{smooth}}$  and  $\lambda_{\text{light}}$  are two predefined scalars.

## 4 EXPERIMENTS

### 4.1 Experiment protocol

**Datasets.** To evaluate the performance of NeRO, we propose a synthetic dataset called the Glossy-Blender dataset and a real dataset called the Glossy-Real dataset. The Glossy-Blender dataset consists of 8 objects with low roughness and strong reflective appearances. The copyright information of these objects are included in Sec. A.12. For each object, we uniformly rendered 128 images of resolution  $800 \times 800$  around the object with the Cycles renderer in Blender. In rendering, we randomly select indoor HDR images from PolyHeaven [Poly Heaven 2022] as sources of environment lights. Among 8 objects, Bell, Cat, Teapot, Potion, and TBell have large smooth surfaces with strong specular effects while Angel, Luyu, and Horse contain more complex geometry with small reflective fragments, as shown in Fig. 8. The Glossy-Real dataset contains 5 objects (Coral, Maneki, Bear, Bunny and Vase), as shown in Fig. 9 (a). We capture about 100-130 images around each object and use COLMAP [Schönberger and Frahm 2016] to track the camera poses for all images. To enable robust camera tracking for COLMAP, we place the object on a calibration board with strong textures. All images are captured by a cellphone camera with a resolution of  $1024 \times 768$ . To get the ground-truth surfaces, we use the structure light-based RGBD sensor EinScan Pro 2X to scan these objects. Since the scanner is unable to reconstruct reflective objects, we manually paint non-reflective substances on these objects before scanning, as shown in Fig. 9 (b). Sec. A.6 includes the detailed dataset statistics.

**Geometry evaluation.** For the evaluation of geometry, we adopt the Chamfer distance (CD) between the reconstructed surfaces and the ground-truth surfaces as the metric. Note that we only evaluate the reconstructed surfaces that are visible in images. To achieve this, we first select 16 input cameras using the farthest point sampling on the camera locations. Then, we project the reconstructed surface or the ground-truth surface onto each sampled camera to get a depth map. Finally, we fuse all the depth maps into a complete point cloud. The Chamfer distance is computed between the point cloud of the



Fig. 8. **Objects from the Glossy-Blender dataset.** “TBell” means “Table Bell” and “Luyu” is a Chinese tea master worshiped as the Sage of Tea.

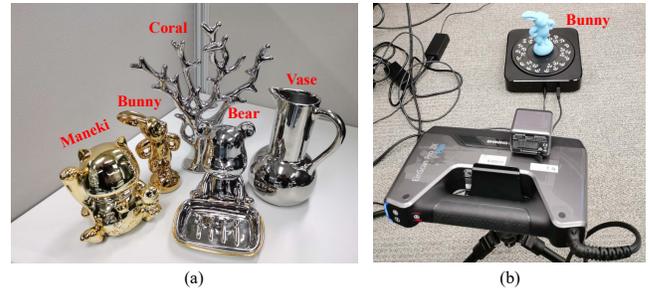


Fig. 9. **The Glossy-Real dataset.** (a) Objects and their names. (b) We paint non-reflective substances on objects and scan objects with an EinScan Pro 2X scanner.

reconstructed surface and the one of the ground-truth surface. On the real data, NeRO reconstructs all surfaces inside the bounding sphere rather than only reconstructing the surfaces of the reflective object, as shown in Fig. 10 (b). From the reconstructed mesh, we manually crop the mesh of the object for the evaluation, as shown in Fig. 10 (c).

**BRDF evaluation.** Directly evaluating the quality of BRDF is difficult because different baseline methods adopt different BRDF models. Instead, we evaluate the quality of relighted images to reveal the quality of the estimated BRDF. For each object in the Glossy-Blender dataset, we use three new HDR images as the environment lights to illuminate the object. For each environment light, we render 16 evenly-distributed relighted images around the object. Finally, we compute PSNR, SSIM [Wang et al. 2004], and LPIPS [Zhang et al. 2018] between the relighted images of different methods and the images rendered by Blender. There is an indeterminable scale factor between the albedo and the lights. The indeterminable light-albedo scales may be different for different baselines. For a fair comparison, we normalize the relighted images to match the average colors of ground-truth images before computing the metrics. On all the real data, we do not adopt such normalization and provide a visual comparison.

**Baselines for geometry estimation.** For geometry reconstruction, we compare our method with COLMAP [Schönberger et al.

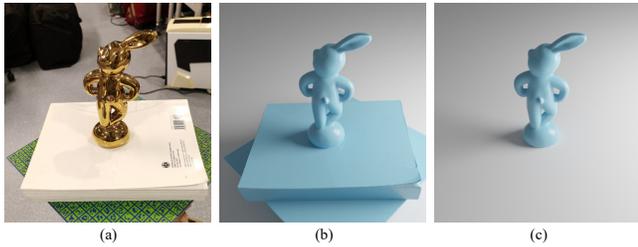


Fig. 10. (a) An input image. (b) The completed reconstruction of our method. (c) The cropped mesh for the evaluation.

2016], Ref-NeRF [Verbin et al. 2022], NeuS [Wang et al. 2021b], NDR [Munkberg et al. 2022], and NDRMC [Hasselgren et al. 2022]. COLMAP [Schönberger et al. 2016] adopts the traditional MVS algorithm patch match stereo [Bleyer et al. 2011] for the reconstruction. Since COLMAP fails to reconstruct complete meshes on reflective objects, we only provide a qualitative comparison with COLMAP. Ref-NeRF represents the surface with a density field and uses the reflective direction and IDE in its color function. For Ref-NeRF, we use grid search to find the best density threshold in  $[0, 15]$  to extract a level set as the output surface. NDR adopts the differentiable marching tetrahedral (DMTet) [Shen et al. 2021] as the surface representation and uses prefiltered split-sum to approximate direct illuminations. NDRMC replaces the prefiltered split-sum of NDR with a differentiable MC sampling. Note that both NDR and NDRMC use ground-truth object masks for training. NeuS uses volume rendering on SDF to extract the surfaces. For all baseline methods, we adopt their official implementations.

**Baselines for BRDF estimation.** We choose baseline methods NDR [Munkberg et al. 2022], NDRMC [Hasselgren et al. 2022], MII [Zhang et al. 2022b] and NeILF [Yao et al. 2022] for comparison. NDR only uses split-sum to consider direct lights while NDRMC, NeILF, and MII all consider indirect lights in the BRDF estimation. NDRMC uses the differentiable MC sampling with a denoiser. NeILF also adopts the MC sampling but with an MLP-based light model. MII uses the Spherical Gaussian to represent direct or indirect lights and the BRDF. In order to reconstruct reasonable geometry on reflective objects, MII, NDR and NDRMC all rely on the object masks in training. NeILF does not involve surface reconstruction in its pipeline and needs a reconstructed mesh as input. Thus, we use the mesh reconstructed by our method as the input to NeILF. MII can be regarded as a combined version of NeRFactor [Zhang et al. 2021b] and PhySG [Zhang et al. 2021a]. NDR is similar to Neural-PIL [Boss et al. 2021b], both of which are based on the prefiltered split-sum. In this case, we do not exhaustively include a comparison with NeRFactor, PhySG, and Neural-PIL.

**Experimental details.** We train NeRO on each object for 300k steps for the surface reconstruction and 100k steps for the BRDF estimation. On each training step, 512 camera rays are sampled for both stages. We adopt the Adam optimizer [Kingma and Ba 2014] ( $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$ ) with a warm-up learning rate schedule, in which the learning rate first increases from  $1e-5$  to  $5e-4$  in 5k (or 1k) steps for the surface reconstruction (or the BRDF estimation) and then decreases to  $1e-5$  in the subsequent steps. To compute

	NDR*	NDRMC*	NeuS	Ref-NeRF	Ours
Bell	0.0122	<u>0.0045</u>	0.0146	0.0137	<b>0.0032</b>
Cat	0.0344	<u>0.1299</u>	0.0278	<u>0.0201</u>	<b>0.0044</b>
Teapot	0.0530	<u>0.0052</u>	0.0546	0.0143	<b>0.0037</b>
Potion	0.0554	0.0415	0.0393	<u>0.0131</u>	<b>0.0053</b>
TBell	0.0821	<u>0.0046</u>	0.0348	0.0216	<b>0.0035</b>
Angel	0.0056	<b>0.0034</b>	<u>0.0035</u>	0.0291	<b>0.0034</b>
Horse	0.0077	<u>0.0052</u>	0.0053	0.0071	<b>0.0049</b>
Luyu	0.0131	0.0082	<u>0.0066</u>	0.0141	<b>0.0054</b>
Avg.	0.0329	0.0253	<u>0.0233</u>	0.0241	<b>0.0042</b>

Table 1. **Reconstruction quality in Chamfer Distance (CD $\downarrow$ ) on the Glossy-Synthetic dataset.** We compare our method with NeuS [Wang et al. 2021b], Ref-NeRF [Verbin et al. 2022], NDR [Munkberg et al. 2022] and NDRMC [Hasselgren et al. 2022]. \*NDRMC and NDR use ground-truth object masks while the other methods do not use object masks. **Bold** means the best performance and underline means the second best performance.

the MC sampling in Stage II,  $N_d = 512$  ray samples and  $N_s = 256$  ray samples are used on the diffuse lobe and the specular lobe respectively. The whole training process takes about 20 hours for the surface reconstruction and 5 hours for the BRDF estimation on a 2080Ti GPU. Note that the most time-consuming part in training is that we sample a massive number of points in the volume rendering. Recent voxel-based representations could largely reduce sample points and thus could speed up the NeRO training, which we leave for future works. On the output mesh, we assign BRDF parameters to vertices and interpolate the parameters on each face. More details about the architecture can be found in Sec. A.5.

## 4.2 Results on Glossy-Blender

**Geometry.** The quantitative results in CD are reported in Table 1. Fig. 2 and Fig. 11 show the visualization of reconstructed surfaces. We interpret the results in the following aspects:

- (1) On “Luyu”, “Horse” and “Angel” which have complex geometry but only small reflective fragments, all methods achieve good reconstruction. The reason is that the surface reconstruction is mainly dominated by the silhouettes of these objects and the appearances of small reflective fragments are relatively less informative. However, on the other objects “Bell”, “Cat”, “Teapot”, “Potion”, and “TBell” which contain a large reflective surface, baseline methods suffer from view-dependent reflections and struggle to reconstruct correct surfaces. In comparison, our method can accurately reconstruct the surfaces.
- (2) As a traditional MVS method, COLMAP fails to reconstruct the reflective objects due to the strong reflections, which either cannot reconstruct any points due to the lack of multi-view consistency or only reconstructs some erroneous points inside the object surfaces.
- (3) Ref-NeRF explicitly considers the direct environment lights by encoding colors in reflective directions. However, a density

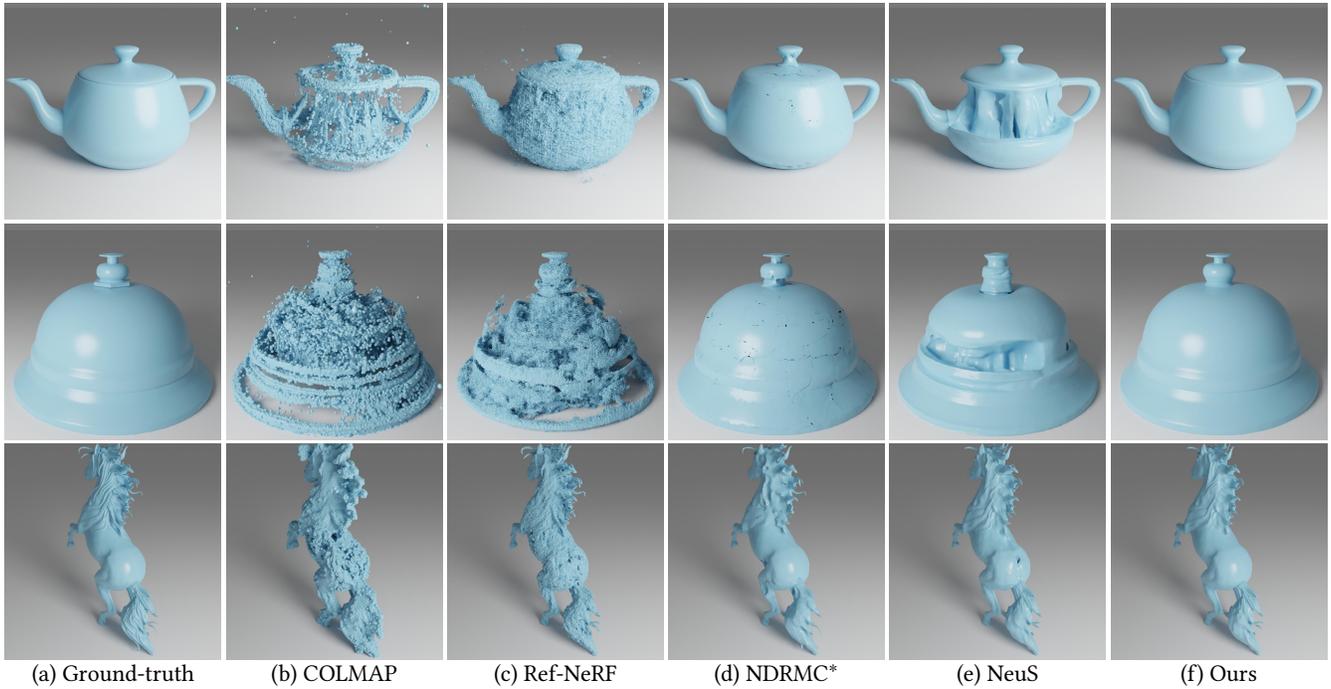


Fig. 11. **Ground-truth and reconstructed surfaces of the Glossy-Blender dataset.** We compare our results with COLMAP [Schönberger et al. 2016], Ref-NeRF [Verbin et al. 2022], NDRMC [Hasselgren et al. 2022], and NeuS [Wang et al. 2021b]. \*NDRMC is trained with ground-truth object masks while the other methods do not use object masks.



Fig. 12. **Relighting objects in the Glossy-Blender dataset.** We compare our method with NDR [Munkberg et al. 2022], NDRMC [Hasselgren et al. 2022], MII [Zhang et al. 2022b] and NeLF [Yao et al. 2022]. All the relighted images are normalized to match the average colors of the ground-truth images.

	NDR	NDRMC	MII	NeLF	Ours
Bell	25.05	24.30	21.58	<u>25.40</u>	<b>31.05</b>
Cat	<u>24.65</u>	23.88	23.46	23.04	<b>30.03</b>
Teapot	19.19	22.45	21.76	<u>24.49</u>	<b>30.95</b>
Potion	21.95	22.07	<u>25.62</u>	24.86	<b>31.17</b>
TBell	16.10	22.60	19.27	<u>22.97</u>	<b>27.48</b>
Angel	22.90	22.89	22.47	<u>24.56</u>	<b>25.49</b>
Horse	25.56	<u>26.42</u>	20.46	25.97	<b>27.41</b>
Luyu	23.72	23.60	23.42	<u>24.62</u>	<b>26.61</b>
Avg.	22.39	23.53	22.28	<u>24.49</u>	<b>28.77</b>

Table 2. **Relighting quality in PSNR $\uparrow$  on the Glossy-Blender dataset.** We compare our method with NDR [Munkberg et al. 2022], NDRMC [Hasselgren et al. 2022], MII [Zhang et al. 2022b] and NeLF [Yao et al. 2022]. **Bold** means the best performance while underline means the second best.

field is not a good representation for the surface reconstruction, which not only results in noisy surfaces but also causes errors in estimating surface normals when the reflections are strong, as shown by the “TBell” Row 2 in Fig. 11. Meanwhile, for regions illuminated by indirect lights, Ref-NeRF also reconstructs incorrect surfaces, as shown by “Bell” Row 1 and “Cat” Row 2 in Fig. 2.

- (4) NDRMC produces good general shape reconstruction but the reconstructed surfaces are not smooth and contain holes or cracks. The CD of NDRMC in Table 1 is large because these small cracks or holes will produce points that are far away from the ground-truth surfaces. Moreover, it is still intractable to apply MC-based shading in volume rendering with multiple sample points on a ray. Thus, similar to IDR [Yariv et al. 2020], NDRMC computes shading on one point per ray and strongly relies on object masks for shape reconstruction.
- (5) NeuS incorrectly distorts the large reflective surfaces to fit the reflection colors, leading to erroneous reconstructions on the objects with large reflective surfaces.
- (6) Our method can correctly reconstruct all these reflective objects by considering both direct and indirect lights to explain the reflective appearance. Especially on “Teapot”, “TBell”, “Cat” and “Potion” with large reflective surfaces, our method accurately reconstructs surfaces no matter whether the regions are illuminated by direct or indirect lights.

**BRDF.** For the evaluation of BRDF estimation, we report the quality of relighted images of different methods. The quantitative results in PSNR are reported in Table 2 and some qualitative results are shown in Fig. 12. The results of SSIM and LPIPS are included in Sec. A.8. NDR [Munkberg et al. 2022] mainly suffers from incorrectly reconstructed surfaces, which prevents it from accurately estimating BRDF. NDRMC tends to produce rough materials with blurred reflections for these glossy objects. MII represents both lights and BRDF with Spherical Gaussian, which saves computations but is unable to represent high-frequency environment lights and smooth reflective BRDFs. NeLF also produces rough materials because it only uses a uniform sampling on the upper half sphere and has difficulty in converging to a smooth surface. In comparison

	NDR*	NDRMC*	NeuS	Ref-NeRF	Ours
Bunny	0.0047	0.0042	<u>0.0022</u>	0.0054	<b>0.0012</b>
Coral	0.0025	0.0022	<u>0.0016</u>	0.0052	<b>0.0014</b>
Maneki	0.0148	0.0117	0.0091	<u>0.0084</u>	<b>0.0024</b>
Bear	0.0104	0.0118	<u>0.0074</u>	0.0109	<b>0.0033</b>
Vase	0.0201	<u>0.0058</u>	0.0101	0.0091	<b>0.0011</b>
Avg.	0.0105	0.0071	<u>0.0061</u>	0.0078	<b>0.0019</b>

Table 3. **Reconstruction quality in Chamfer Distance (CD $\downarrow$ ) on the Glossy-Real dataset.** We compare our method with NeuS [Wang et al. 2021b], Ref-NeRF [Verbin et al. 2022], NDR [Munkberg et al. 2022] and NDRMC [Hasselgren et al. 2022]. \*NDR and NDRMC use ground-truth object masks while the other methods do not use object masks. **Bold** means the best performance and underline means the second best performance.

with NDRMC using only 16 ray samples, we directly use MLP to decode indirect lights and thus are able to use 512 ray samples for the cosine lobe and 256 ray samples for the specular lobe, which makes our method accurately capture the materials. In comparison with NeLF which only uses a Fibonacci ray sampling on the upper half-sphere, we adopt the importance sampling on both the cosine diffuse lobe and the specular lobe. This is essential for recovering the smooth material because the specular lobe is small but vital for the reflective appearance and Fibonacci ray sampling often fails to sample rays in the small specular lobe. Note that our method does not simply predict smooth materials but also correctly recovers the rough materials on the “Potion” row 1 of Fig. 12.

### 4.3 Results on Glossy-Real

**Geometry.** The quantitative results in CD on the Glossy-Real dataset are reported in Table 3 and qualitative results are shown in Fig. 13. In general, the results on the Glossy-Real dataset are similar to the results on the Glossy-Blender dataset, where our method can accurately reconstruct all reflective surfaces while each baseline method fails to reconstruct several objects. On all the objects of the Glossy-Real dataset, surfaces reconstructed by Ref-NeRF are very noisy. On “Vase”, “Bear” and “Maneki”, NeuS incorrectly distorts the large reflective surfaces and on “Coral” and “Bunny”, the reconstruction of NeuS is generally correct but still contains non-smooth artifacts on the regions highlighted by the red bounding boxes in Fig. 13. NDRMC can reconstruct the general shape with the help of object masks but also produces holes on the surfaces of “Vase”, “Bear” and “Maneki”. In comparison, our method can correctly reconstruct all objects. The main artifact of our method is that some details on the “Maneki” are missing.

**BRDF.** We provide qualitative relighting results in Fig. 14 to compare our method with MII [Zhang et al. 2022b], NeLF [Yao et al. 2022], NDR [Munkberg et al. 2022] and NDRMC [Hasselgren et al. 2022]. Similar to the relighting results on the synthetic dataset, our method produces more realistic relighting results than the two baselines. NDRMC tends to produce rough materials while MII estimates smooth but plastic-like materials for these smooth metallic objects. Our method accurately estimates the smooth metallic materials for all reflective surfaces and the rough materials for the back of “Bear”.

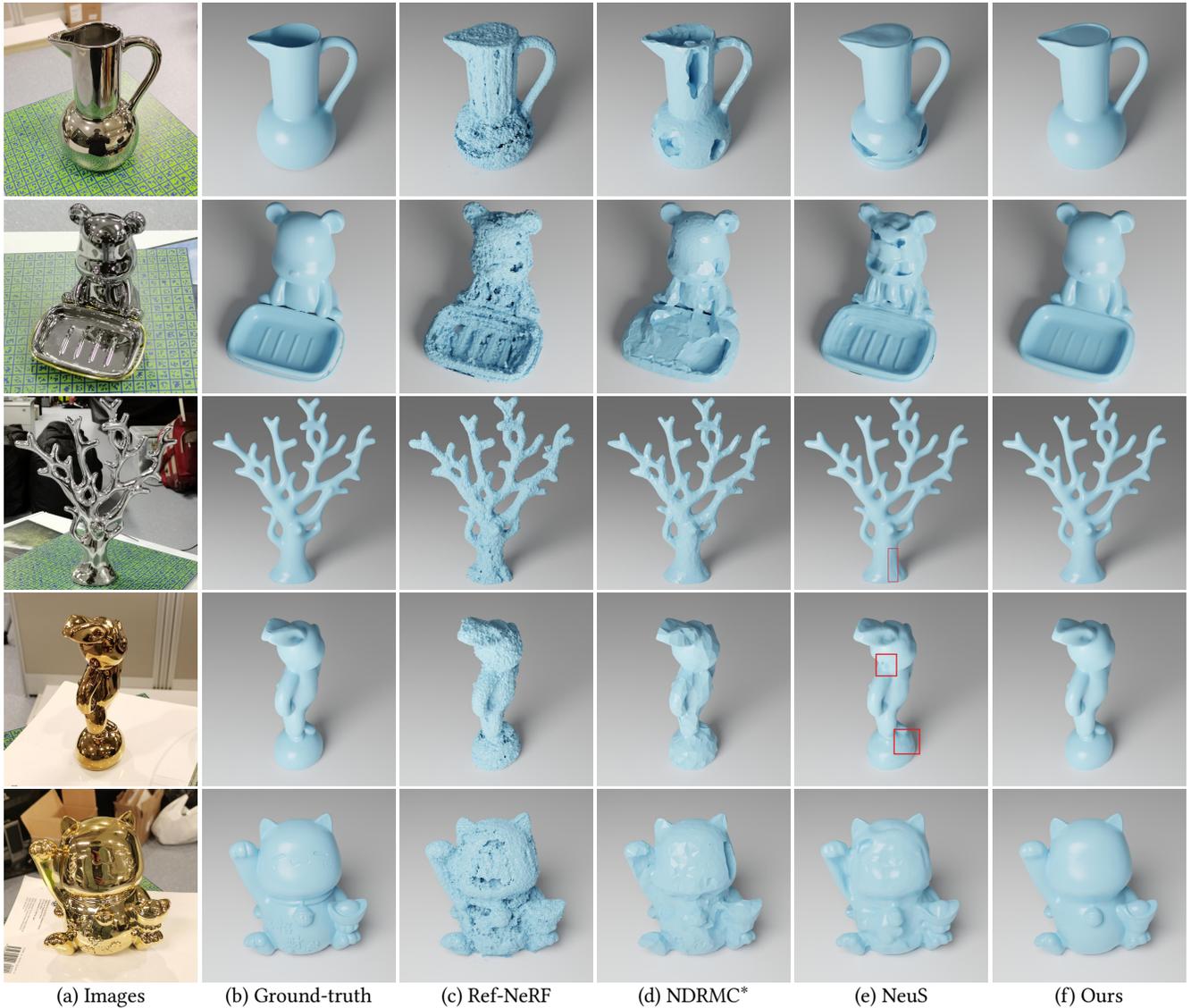


Fig. 13. **Images, ground-truth and reconstructed surfaces of the Glossy-Real dataset.** We compare our results with Ref-NeRF [Verbin et al. 2022], NDRMC [Hasselgren et al. 2022], and NeuS [Wang et al. 2021b]. \*NDRMC is trained with ground-truth object masks while the other methods do not use object masks.

#### 4.4 Ablation studies

In this section, we conduct ablation studies on geometry reconstruction and BRDF estimation.

**Ablation studies on the geometry reconstruction.** The key component of our method for surface reconstruction is to consider the environmental lights for the reconstruction. Thus, we provide ablation studies on different environment light representations as shown in Table 4 and Fig. 15.

ID	Description	Angel	Bell	Cat	Teapot	Avg.
0	NeuS	0.0035	0.0146	0.0278	0.0546	0.0251
1	Only direct lights	0.0033	0.0051	0.0217	0.0076	0.0094
2	Only indirect lights	0.0037	0.0043	0.0210	0.0064	0.0089
3	No occlusion loss $\ell_{occ}$	0.0033	0.0036	0.0212	0.0287	0.0142
4	No Eikonal loss $\ell_{eikonal}$	0.0027	0.0030	0.0212	0.0184	0.0113
5	Full model	0.0034	0.0032	0.0044	0.0037	0.0037

Table 4. **Ablation studies on the geometry reconstruction** using objects from the Glossy-Synthetic dataset. Chamfer Distance (CD $\downarrow$ ) is reported.



Fig. 14. **Relighting objects from the Glossy-Real dataset.** We provide a visual comparison with NDR [Munkberg et al. 2022], NDRMC [Hasselgren et al. 2022], MII [Zhang et al. 2022b], and NeLF [Yao et al. 2022]. We provide the input image with the nearest viewpoint and the relighting HDR map as a reference.

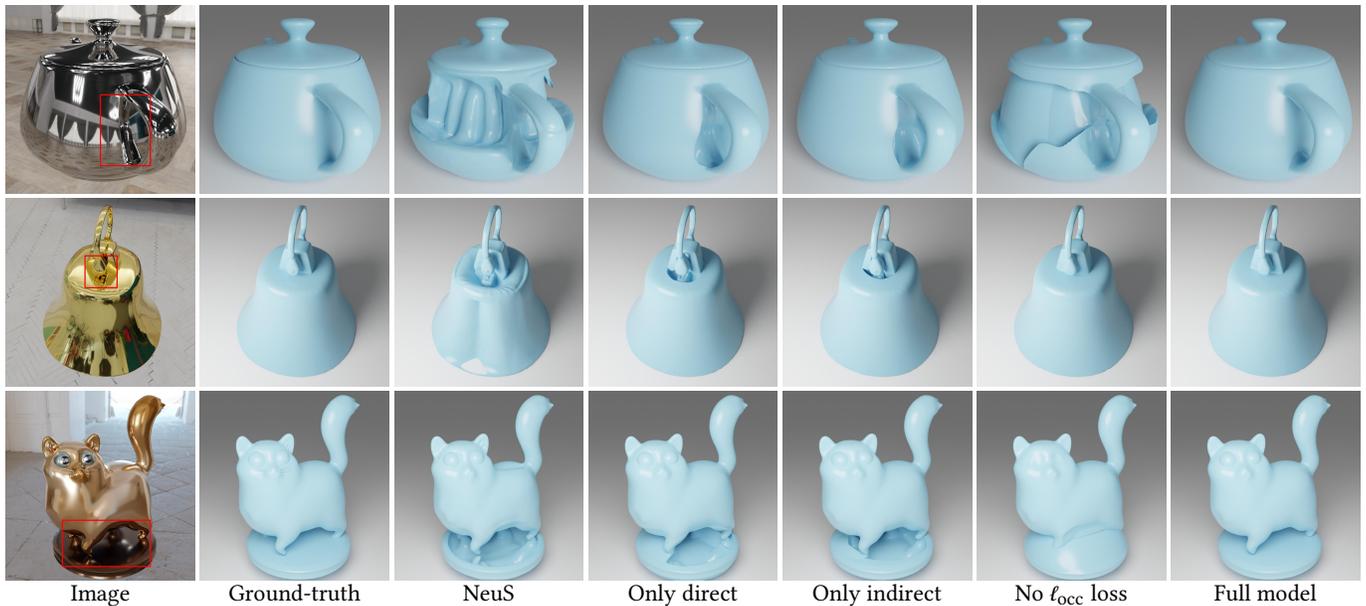


Fig. 15. **Ablation studies on the surface reconstruction.** We use red bounding boxes to highlight regions illuminated by indirect lights in the images. Results of NeuS [Wang et al. 2021b] are provided for comparison. “Only direct” means the model using only the direct lights (Model 1 in Table 4). “Only indirect” means the model using only indirect lights (Model 2 in Table 4). “No  $\ell_{occ}$  loss” means the model uses both direct and indirect lights but not use  $\ell_{occ}$  on the occlusion probability (Model 3 in Table 4). “Full model” contains all components (Model 4 in Table 4).

Fig. 16. Ablation studies on the Eikonal loss  $\ell_{\text{eikonal}}$ .

Index	Description	Potion	Bell	Cat	Teapot	Avg.
0	Stage I	22.44	28.83	21.54	28.53	25.34
1	No specular lobe sampling	25.94	22.84	24.69	24.25	24.43
2	Only direct	29.74	30.03	29.21	29.84	29.71
3	Only Indirect	29.95	30.94	28.46	30.78	30.03
4	Without $\ell_{\text{light}}$	30.70	31.07	29.76	30.80	30.58
5	Without $\ell_{\text{smooth}}$	30.90	30.45	29.41	31.88	30.66
6	Full model	31.17	31.05	30.03	30.95	30.80

Table 5. Ablation studies on the BRDF estimation using objects from the Glossy-Synthetic dataset. PSNR $\uparrow$  on of the relighting images is reported.

- (1) In Model 1, we simply only consider the direct lights with  $L(\omega_i) = g_{\text{direct}}(\omega_i)$  in shading. In comparison with the original color function of NeuS, Model 1 is already able to reconstruct the surfaces illuminated by direct lights and avoids the undesirable surface distortion of NeuS. However, Model 1 is unable to correctly reconstruct the regions illuminated by indirect lights as shown in red boxes of Fig. 15.
- (2) Model 2 is designed only to consider the indirect lights with  $L(\omega_i) = g_{\text{indirect}}(\mathbf{p}, \omega_i)$ , where  $g_{\text{indirect}}$  is supposed to predict location-dependent lights and is not limited to the direct environment lights. However, as shown in Fig. 15, Model 2 still struggles to reconstruct the surfaces illuminated by the indirect lights. The reason is that when the surfaces are mostly illuminated by the direct environment lights, Model 2 tends to neglect the input position  $\mathbf{p}$  and only predicts the position-independent lights on all surface points.
- (3) Model 3 applies the combination of direct and indirect lights as Eq. 9 but does not use the occlusion loss in Eq. 12 to constrain the predicted occlusion probability. In this case, Model 3 learns the occlusion probability from only the rendering loss. As shown in Fig. 15, the convergence of Model 3 is not stable due to the lack of clear supervision on the occlusion probability, which succeeds in reconstructing all surfaces of “Bell” but fails on “Teapot” and “Cat”.
- (4) Model 4 does not use the Eikonal loss [Gropp et al. 2020]. Without Eikonal loss, the reconstructed surfaces are slightly better on “Angel” and “Bell”. However, on “Cat” and “Teapot”, Model 4 tends to reconstruct double-layer surfaces on regions illuminated by indirect lights as shown in Fig. 16.
- (5) Our full model successfully reconstructs all surfaces of the objects. On the object “Angel”, all models reconstruct its surface successfully as discussed in Sec. 4.2.

**Ablation studies on the BRDF estimation.** We conduct ablation studies on the importance sampling strategy and the proposed

light representation. The relighting results are shown in Table 5 and Fig. 17. First, the BRDF estimated in stage I of our method usually has very small roughness because our method in stage I tends to use a smooth material to improve the ability to fit specular colors. Second, about the importance sampling, if we do not apply importance sampling on the specular lobe, the estimated surface roughness will be large because it is unable to capture the high-frequency specular colors, which is similar to the results of NeILF. Third, we compare different light representations with the same setting as the ablation studies on geometry reconstruction. Only using direct lights prevents the model from estimating a correct BRDF on surfaces illuminated by indirect lights. Only using indirect lights results in over-smooth surfaces because, on rough materials showing low-frequency color changes, the model with only indirect lights will predict a smooth material with low-frequency lights rather than predicting a rough material. We further conduct ablation studies on the neutral light loss  $\ell_{\text{light}}$  and the smoothness loss  $\ell_{\text{smooth}}$  [Hasselgren et al. 2022; Munkberg et al. 2022]. Neutral light loss makes the estimated albedo more accurate and reasonable while smoothness loss avoids noisy material prediction, both of which improve the relighting quality as shown in Models 4, 5, and 6 in Table 5.

**Visualization of image decomposition.** To show the quality of the estimated BRDF and lights, we provide two qualitative examples of image decomposition by our method in Fig. 18. Though the appearances of reflective objects contain strong specular effects, our method successfully separates the high-frequency lights from the low-frequency albedo. The predicted metalness and roughness are also very reasonable. Our method accurately predicts a rougher material for the base of the table bell and a smooth material for the table bell lid in Fig. 18 (Row 1). Meanwhile, our method also distinguishes the metallic vase from the non-metallic calibration board in Fig. 18 (Row 2).

#### 4.5 Limitations

**Geometry.** Though we successfully reconstruct the shape of reflective objects, our method still fails to capture some subtle details, as shown in Fig. 19. The main reason is that the rendering function strongly relies on the surface normals estimated by the neural SDF but a neural SDF tends to produce smooth surface normals. Thus, it is hard for the neural SDF to produce abrupt normal changes to reconstruct subtle details like the cloth textures of “Angel”, the beards of “Cat”, and the textures of “Maneki”.

**BRDF.** In the experiments, we observe that our BRDF estimation mainly suffers from incorrect geometry, especially on “Angel” as shown in Fig. 20. Since the appearance of reflective objects strongly relies on the surface normals to compute the reflective directions, the incorrectness of surface normals will make our method struggle to fit correct colors, which leads to inaccurate BRDF estimation. Meanwhile, the BRDF in NeRO does not support advanced reflections such as anisotropic reflections.

**Pose estimation.** Another limitation is that our method relies on accurate input camera poses and estimating camera poses on reflective objects usually requires stable textures like calibration boards for image matching. Without calibration boards, we may



Fig. 17. **Ablation studies on the BRDF estimation.** “Stage 1” means the estimated BRDF in Stage I. “No sp. lobe” means not using the importance sampling on the specular lobe. “Only direct” means only using direct lights in the light representation while “Only indirect” means only using indirect lights. Note that all relighted images are normalized to match the average colors of the ground-truth images.



Fig. 18. **Decomposition of the input image.** After the estimation of surface BRDF, our method is able to automatically decompose the input image into albedo, lights, metalness, and roughness.

recover poses from other co-visible non-reflective objects or with the help of devices like IMU.

## 5 CONCLUSION

We have presented NeRO, a neural reconstruction method for accurately reconstructing the geometry and the BRDF of reflective objects, without knowing the environment light conditions and the object masks. The key idea of NeRO is to explicitly incorporate the rendering equation in a neural reconstruction framework. NeRO achieves this challenging goal by proposing a novel light representation and adopting a two-stage approach. In the first stage, by applying tractable approximations, we model both the direct and indirect lights with shading MLPs and learn the surface geometry faithfully. In the second stage, we fix the geometry and reconstruct

a more accurate surface BRDF as well as the environment light by Monte Carlo sampling. Experiments have demonstrated that NeRO achieves better surface reconstruction quality and BRDF estimation of reflective objects compared to the state-of-the-art.

**Acknowledgement.** Lingjie Liu has been supported by the ERC Consolidator Grant 4DReply (77078). We sincerely thank the Tam Wing Fan Innovation Wing of HKU for providing the EinScan scanner to make the Glossy-Real dataset.

## REFERENCES

- Matan Atzmon and Yaron Lipman. 2020. SAL: Sign agnostic learning of shapes from raw data. In *CVPR*.
- Jonathan T Barron and Jitendra Malik. 2014. Shape, illumination, and reflectance from shading. *TPAMI* 37, 8 (2014), 1670–1687.
- Jonathan T Barron, Ben Mildenhall, Dor Verbin, Pratul P Srinivasan, and Peter Hedman. 2022. Mip-NeRF 360: Unbounded anti-aliased neural radiance fields. In *CVPR*.

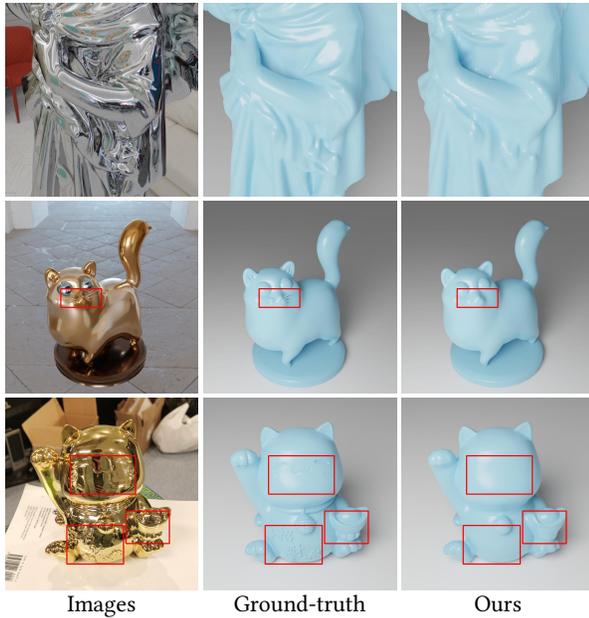


Fig. 19. **Failure cases for the geometry reconstruction.** Though our method is able to correctly reconstruct the overall shapes of glossy objects, some subtle details are missing or incorrect because the neural SDF field has difficulty in producing abrupt normal changes on these subtle regions.



Fig. 20. **Failure cases for the BRDF estimation.** Reflective appearances are very sensitive to the reflective direction. With inaccurate surface geometry, our method fails to find a correct reflective direction and thus predicts inaccurate BRDFs.

Jonathan T Barron and Ben Poole. 2016. The fast bilateral solver. In *ECCV*.  
 Sai Bi, Zexiang Xu, Pratul Srinivasan, Ben Mildenhall, Kalyan Sunkavalli, Miloš Hašan, Yannick Hold-Geoffroy, David Kriegman, and Ravi Ramamoorthi. 2020. Neural reflectance fields for appearance acquisition. *arXiv preprint arXiv:2008.03824* (2020).  
 Sai Bi, Zexiang Xu, Kalyan Sunkavalli, Miloš Hašan, Yannick Hold-Geoffroy, David Kriegman, and Ravi Ramamoorthi. [n.d.]. Deep reflectance volumes: Relightable reconstructions from multi-view photometric images. In *ECCV*. Springer.  
 Michael Bleyer, Christoph Rhemann, and Carsten Rother. 2011. Patchmatch stereo-stereo matching with slanted support windows. In *BMVC*.  
 Mark Boss, Raphael Braun, Varun Jampani, Jonathan T Barron, Ce Liu, and Hendrik Lensch. 2021a. NerD: Neural reflectance decomposition from image collections. In *CVPR*.  
 Mark Boss, Andreas Engelhardt, Abhishek Kar, Yuanzhen Li, Deqing Sun, Jonathan T Barron, Hendrik Lensch, and Varun Jampani. 2022. SAMURAI: Shape And Material from Unconstrained Real-world Arbitrary Image collections. In *NeurIPS*.  
 Mark Boss, Varun Jampani, Raphael Braun, Ce Liu, Jonathan Barron, and Hendrik Lensch. 2021b. Neural-PIL: Neural pre-integrated lighting for reflectance decomposition. In *NeurIPS*.  
 Neill DF Campbell, George Vogiatzis, Carlos Hernández, and Roberto Cipolla. 2008. Using multiple hypotheses to improve depth-maps for multi-view stereo. In *ECCV*.  
 Pierre Charbonnier, Laure Blanc-Feraud, Gilles Aubert, and Michel Barlaud. 1994. Two deterministic half-quadratic regularization algorithms for computed imaging. In

*ICIP*.  
 Wenzheng Chen, Huan Ling, Jun Gao, Edward Smith, Jaakko Lehtinen, Alec Jacobson, and Sanja Fidler. 2019. Learning to predict 3d objects with an interpolation-based differentiable renderer. *NeurIPS* 32 (2019).  
 Wenzheng Chen, Joey Litalien, Jun Gao, Zian Wang, Clement Fuji Tsang, Sameh Khamis, Or Litany, and Sanja Fidler. 2021. DIB-R++: Learning to predict lighting and material with a hybrid differentiable renderer. *NeurIPS* (2021).  
 Zhaoxi Chen and Ziwei Liu. 2022. Relighting4d: Neural relightable human from videos. In *ECCV*.  
 Shuo Cheng, Zexiang Xu, Shilin Zhu, Zhuwen Li, Li Erran Li, Ravi Ramamoorthi, and Hao Su. 2020. Deep stereo using adaptive thin volume representation with uncertainty awareness. In *CVPR*.  
 Ziang Cheng, Hongdong Li, Yuta Asano, Yinqiang Zheng, and Imari Sato. 2021. Multi-view 3D Reconstruction of a Texture-less Smooth Surface of Unknown Generic Reflectance. In *CVPR*.  
 Robert L Cook and Kenneth E. Torrance. 1982. A reflectance model for computer graphics. *ACM Transactions on Graphics (TOG)* 1, 1 (1982), 7–24.  
 François Darmon, Bénédicte Bascle, Jean-Clément Devaux, Pascal Monasse, and Mathieu Aubry. 2022. Improving neural implicit surfaces geometry with patch warping. In *CVPR*.  
 Akshat Dave, Yongyi Zhao, and Ashok Veeraraghavan. 2022. PANDORA: Polarization-Aided Neural Decomposition Of Radiance. In *ECCV*.  
 Youming Deng, Xueting Li, Sifei Liu, and Ming-Hsuan Yang. 2022. DIP: Differentiable Interreflection-aware Physics-based Inverse Rendering. *arXiv preprint arXiv:2212.04705* (2022).  
 Sylvain Duchêne, Clement Riant, Gaurav Chaurasia, Jorge Lopez-Moreno, Pierre-Yves Laffont, Stefan Popov, Adrien Bousseau, and George Drettakis. 2015. Multi-view intrinsic images of outdoors scenes with an application to relighting. *ACM Transactions on Graphics (TOG)* (2015).  
 Qiancheng Fu, Qingshan Xu, Yew-Soon Ong, and Wenbing Tao. 2022. Geo-Neus: Geometry-Consistent Neural Implicit Surfaces Learning for Multi-view Reconstruction. In *NeurIPS*.  
 Yasutaka Furukawa and Jean Ponce. 2009. Accurate, dense, and robust multiview stereopsis. *TPAMI* 32, 8 (2009), 1362–1376.  
 David Gallup, Jan-Michael Frahm, Philippos Mordohai, Qingxiang Yang, and Marc Pollefeys. 2007. Real-time plane-sweeping stereo with multiple sweeping directions. In *CVPR*.  
 Duan Gao, Guojun Chen, Yue Dong, Pieter Peers, Kun Xu, and Xin Tong. 2020. Deferred neural lighting: free-viewpoint relighting from unstructured photographs. *ACM Transactions on Graphics (TOG)* 39, 6 (2020), 1–15.  
 Duan Gao, Xiao Li, Yue Dong, Pieter Peers, Kun Xu, and Xin Tong. 2019. Deep inverse rendering for high-resolution SVBRDF estimation from an arbitrary number of images. *ToG* 38, 4 (2019), 134–1.  
 Andreas Geiger, Philip Lenz, Christoph Stiller, and Raquel Urtasun. 2013. Vision meets robotics: The KITTI dataset. *The International Journal of Robotics Research* 32, 11 (2013), 1231–1237.  
 Clement Godard, Peter Hedman, Wenbin Li, and Gabriel J Brostow. 2015. Multi-view reconstruction of highly specular surfaces in uncontrolled environments. In *3DV*.  
 Amos Gropp, Lior Yariv, Niv Haim, Matan Atzmon, and Yaroslav Lipman. 2020. Implicit Geometric Regularization for Learning Shapes. In *ICML*.  
 Haoyu Guo, Sida Peng, Haotong Lin, Qianqian Wang, Guofeng Zhang, Hujun Bao, and Xiaowei Zhou. 2022b. Neural 3D Scene Reconstruction with the Manhattan-world Assumption. In *CVPR*.  
 Yu Guo, Cameron Smith, Miloš Hašan, Kalyan Sunkavalli, and Shuang Zhao. 2020. MaterialGAN: reflectance capture using a generative SVBRDF model. *ToG* 39, 6 (2020), 1–13.  
 Yuan-Chen Guo, Di Kang, Linchao Bao, Yu He, and Song-Hai Zhang. 2022a. NeRFRen: Neural radiance fields with reflections. In *CVPR*.  
 Kai Han, Kwan-Yee K Wong, Dirk Schnieders, and Miaomiao Liu. 2016. Mirror surface reconstruction under an uncalibrated camera. In *CVPR*.  
 Richard Hartley and Andrew Zisserman. 2003. *Multiple view geometry in computer vision*. Cambridge university press.  
 Jon Hasselgren, Nikolai Hofmann, and Jacob Munkberg. 2022. Shape, Light & Material Decomposition from Images using Monte Carlo Rendering and Denoising. In *NeurIPS*.  
 Asmaa Hosni, Christoph Rhemann, Michael Bleyer, Carsten Rother, and Margrit Gelautz. 2012. Fast cost-volume filtering for visual correspondence and beyond. *TPAMI* 35, 2 (2012), 504–511.  
 Eldar Insafutdinov, Dylan Campbell, João F Henriques, and Andrea Vedaldi. 2022. SNeS: Learning Probably Symmetric Neural Surfaces from Incomplete Data. In *ECCV*. 367–383.  
 Rasmus Jensen, Anders Dahl, George Vogiatzis, Engin Tola, and Henrik Aanæs. 2014. Large scale multi-view stereopsis evaluation. In *CVPR*.  
 Achuta Kadambi, Vage Taamazyan, Boxin Shi, and Ramesh Raskar. 2015. Polarized 3D: High-quality depth sensing with polarization cues. In *ICCV*.  
 James T Kajiya. 1986. The rendering equation. In *SIGGRAPH*.

- Brian Karis and Epic Games. 2013. Real shading in unreal engine 4. *Proc. Physically Based Shading Theory Practice* 4, 3 (2013), 1.
- Hiroharu Kato, Yoshitaka Ushiku, and Tatsuya Harada. 2018. Neural 3d mesh renderer. In *CVPR*.
- Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014).
- Georgios Kopanas, Thomas Leimkühler, Gilles Rainer, Clément Jambon, and George Drettakis. 2022. Neural point catacaustics for novel-view synthesis of reflections. *ToG* 41, 6 (2022), 1–15.
- Zhengfei Kuang, Kyle Olszewski, Menglei Chai, Zeng Huang, Panos Achlioptas, and Sergey Tulyakov. 2022. NeROIC: Neural Rendering of Objects from Online Image Collections. In *SIGGRAPH*.
- Hai Li, Xingrui Yang, Hongjia Zhai, Yuqian Liu, Hujun Bao, and Guofeng Zhang. 2022. Vox-Surf: Voxel-based implicit surface representation. *IEEE Transactions on Visualization and Computer Graphics* (2022).
- Junxuan Li and Hongdong Li. 2022a. Neural Reflectance for Shape Recovery with Shadow Handling. In *CVPR*.
- Junxuan Li and Hongdong Li. 2022b. Self-calibrating photometric stereo by neural inverse rendering. In *ECCV*.
- Zhengqin Li, Mohammad Shafiei, Ravi Ramamoorthi, Kalyan Sunkavalli, and Manmohan Chandraker. 2020. Inverse rendering for complex indoor scenes: Shape, spatially-varying lighting and svbrdf from a single image. In *CVPR*.
- Zhengqin Li, Zexiang Xu, Ravi Ramamoorthi, Kalyan Sunkavalli, and Manmohan Chandraker. 2018. Learning to reconstruct shape and spatially-varying reflectance from a single image. *ToG* 37, 6 (2018), 1–11.
- Shichen Liu, Tianye Li, Weikai Chen, and Hao Li. 2019. Soft rasterizer: A differentiable renderer for image-based 3d reasoning. In *CVPR*.
- Yang Liu, Alexandros Neophytou, Sunando Sengupta, and Eric Sommerlade. 2021. Relighting images in the wild with a self-supervised siamese auto-encoder. In *CVPR*.
- Xiaoxiao Long, Cheng Lin, Peng Wang, Taku Komura, and Wenping Wang. 2022. Sparseneus: Fast generalizable neural surface reconstruction from sparse views. In *ECCV*.
- Linjie Lyu, Ayush Tewari, Thomas Leimkühler, Marc Habermann, and Christian Theobalt. 2022. Neural Radiance Transfer Fields for Relightable Novel-view Synthesis with Global Illumination. In *ECCV*.
- B Mildenhall, PP Srinivasan, M Tancik, JT Barron, R Ramamoorthi, and R Ng. 2020. NeRF: Representing scenes as neural radiance fields for view synthesis. In *ECCV*.
- Jacob Munkberg, Jon Hasselgren, Tianchang Shen, Jun Gao, Wenzheng Chen, Alex Evans, Thomas Müller, and Sanja Fidler. 2022. Extracting Triangular 3D Models, Materials, and Lighting From Images. In *CVPR*.
- Giljoo Nam, Joo Ho Lee, Diego Gutierrez, and Min H Kim. 2018. Practical svBRDF acquisition of 3d objects with unstructured flash photography. *ToG* 37, 6 (2018), 1–12.
- Thomas Nestmeyer, Jean-François Lalonde, Iain Matthews, and Andreas Lehrmann. 2020. Learning physics-guided face relighting under directional light. In *CVPR*.
- Fred E Nicodemus. 1965. Directional reflectance and emissivity of an opaque surface. *Applied optics* 4, 7 (1965), 767–775.
- Michael Niemeyer, Lars Mescheder, Michael Oechsle, and Andreas Geiger. 2020. Differentiable volumetric rendering: Learning implicit 3d representations without 3d supervision. In *CVPR*.
- Merlin Nimier-David, Delio Vicini, Tizian Zeltner, and Wenzel Jakob. 2019. Mitsuba 2: A retargetable forward and inverse renderer. *ToG* 38, 6 (2019), 1–17.
- Michael Oechsle, Songyou Peng, and Andreas Geiger. 2021. Unisurf: Unifying neural implicit surfaces and radiance fields for multi-view reconstruction. In *ICCV*.
- Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. 2019. DeepSDF: Learning continuous signed distance functions for shape representation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 165–174.
- Julien Philip, Michaël Gharbi, Tinghui Zhou, Alexei A Efros, and George Drettakis. 2019. Multi-view relighting using a geometry-aware network. *ACM Transactions on Graphics (TOG)* 38, 4 (2019), 78–1.
- Julien Philip, Sébastien Mordenthaler, Michaël Gharbi, and George Drettakis. 2021. Free-viewpoint indoor neural relighting from multi-view stereo. *ACM Transactions on Graphics (TOG)* 40, 5 (2021), 1–18.
- Poly Heaven. 2022. Poly Heaven. <https://polyhaven.com/>.
- Stefan Rahmann and Nikos Canterakis. 2001. Reconstruction of specular surfaces using polarization imaging. In *CVPR*.
- Christian Richardt, Douglas Orr, Ian Davies, Antonio Criminisi, and Neil A Dodgson. 2010. Real-time spatiotemporal stereo matching using the dual-cross-bilateral grid. In *ECCV*.
- Simon Rodriguez, Siddhant Prakash, Peter Hedman, and George Drettakis. 2020. Image-based rendering of cars using semantic labels and approximate reflection flow. *Proceedings of the ACM on Computer Graphics and Interactive Techniques* 3 (2020).
- Stefan Roth and Michael J Black. 2006. Specular flow and the recovery of surface structure. In *CVPR*.
- Viktor Rudnev, Mohamed Elgharib, William Smith, Lingjie Liu, Vladislav Golyanik, and Christian Theobalt. 2022. NeRF for outdoor scene relighting. In *ECCV*.
- Daniel Scharstein and Richard Szeliski. 2002. A taxonomy and evaluation of dense two-frame stereo correspondence algorithms. *IJCV* 47, 1 (2002), 7–42.
- Carolin Schmitt, Simon Donne, Gernot Riegler, Vladlen Koltun, and Andreas Geiger. 2020. On joint estimation of pose, geometry and svBRDF from a handheld scanner. In *CVPR*.
- Johannes Lutz Schönberger and Jan-Michael Frahm. 2016. Structure-from-Motion Revisited. In *CVPR*.
- Johannes Lutz Schönberger, Enliang Zheng, Marc Pollefeys, and Jan-Michael Frahm. 2016. Pixelwise View Selection for Unstructured Multi-View Stereo. In *ECCV*.
- Tianchang Shen, Jun Gao, Kangxue Yin, Ming-Yu Liu, and Sanja Fidler. 2021. Deep Marching Tetrahedra: a hybrid representation for high-resolution 3d shape synthesis. In *NeurIPS*.
- Yichang Shih, Sylvain Paris, Frédo Durand, and William T Freeman. 2013. Data-driven hallucination of different times of day from a single outdoor photo. *ACM Transactions on Graphics (TOG)* 32, 6 (2013), 1–11.
- Sudipta N Sinha, Johannes Kopf, Michael Goesele, Daniel Scharstein, and Richard Szeliski. 2012. Image-based rendering for scenes with reflections. *ToG* 31, 4 (2012), 1–10.
- Vincent Sitzmann, Michael Zollhöfer, and Gordon Wetzstein. 2019. Scene representation networks: Continuous 3d-structure-aware neural scene representations. *Advances in Neural Information Processing Systems* 32 (2019).
- Christoph Strecha, Rik Fransens, and Luc Van Gool. 2006. Combined depth and outlier estimation in multi-view stereo. In *CVPR*.
- Jiaming Sun, Xi Chen, Qianqian Wang, Zhengqi Li, Hadar Averbuch-Elor, XiaoWei Zhou, and Noah Snavely. 2022. Neural 3D reconstruction in the wild. In *SIGGRAPH*.
- Ayush Tewari, Ohad Fried, Justus Thies, Vincent Sitzmann, Stephen Lombardi, Kalyan Sunkavalli, Ricardo Martin-Brualla, Tomas Simon, Jason Saragih, Matthias Nießner, et al. 2020. State of the art on neural rendering. In *Computer Graphics Forum*, Vol. 39. Wiley Online Library, 701–727.
- Ayush Tewari, Justus Thies, Ben Mildenhall, Pratul Srinivasan, Edgar Tretschk, W Yifan, Christoph Lassner, Vincent Sitzmann, Ricardo Martin-Brualla, Stephen Lombardi, et al. 2022. Advances in neural rendering. In *Computer Graphics Forum*, Vol. 41. Wiley Online Library, 703–735.
- Siu-Kei Tin, Jinwei Ye, Mahdi Nezamabadi, and Can Chen. 2016. 3d reconstruction of mirror-type objects using efficient ray coding. In *ICCP*. IEEE, 1–11.
- Kushagra Tiwary, Askhat Dave, Nikhil Behari, Tzofi Klinghoffer, Ashok Veeraraghavan, and Ramesh Raskar. 2022. ORCa: Glossy Objects as Radiance Field Cameras. *arXiv preprint arXiv:2212.04531* (2022).
- Kenneth E Torrance and Ephraim M Sparrow. 1967. Theory for off-specular reflection from roughened surfaces. *Josa* 57, 9 (1967), 1105–1114.
- Dor Verbin, Peter Hedman, Ben Mildenhall, Todd Zickler, Jonathan T Barron, and Pratul P Srinivasan. 2022. Ref-NeRF: Structured view-dependent appearance for neural radiance fields. In *CVPR*.
- Fangjinhua Wang, Silvano Galliani, Christoph Vogel, Pablo Speciale, and Marc Pollefeys. 2021a. PatchMatchNet: Learned multi-view patchmatch stereo. In *CVPR*.
- Jiepeng Wang, Peng Wang, Xiaoxiao Long, Christian Theobalt, Taku Komura, Lingjie Liu, and Wenping Wang. 2022c. Neuris: Neural reconstruction of indoor scenes using normal priors. In *ECCV*.
- Peng Wang, Lingjie Liu, Yuan Liu, Christian Theobalt, Taku Komura, and Wenping Wang. 2021b. NeuS: Learning neural implicit surfaces by volume rendering for multi-view reconstruction. In *NeurIPS*.
- Yiming Wang, Qin Han, Marc Habermann, Kostas Daniilidis, Christian Theobalt, and Lingjie Liu. 2022a. NeuS2: Fast Learning of Neural Implicit Surfaces for Multi-view Reconstruction. *arXiv preprint arXiv:2212.05231* (2022).
- Yiqun Wang, Ivan Skorokhodov, and Peter Wonka. 2022b. HF-NeuS: Improved Surface Reconstruction Using High-Frequency Details. In *NeurIPS*.
- Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. 2004. Image quality assessment: from error visibility to structural similarity. *Transactions on Image Processing* 13, 4 (2004), 600–612.
- Thomas Whelan, Michael Goesele, Steven J Lovegrove, Julian Straub, Simon Green, Richard Szeliski, Steven Butterfield, Shobhit Verma, Richard A Newcombe, M Goesele, et al. 2018. Reconstructing scenes with mirror and glass surfaces. *ToG* 37, 4 (2018), 102–1.
- Felix Wimbauer, Shangzhe Wu, and Christian Rupprecht. 2022. De-rendering 3D Objects in the Wild. In *CVPR*.
- Shihao Wu, Hui Huang, Tiziano Portenier, Matan Sela, Daniel Cohen-Or, Ron Kimmel, and Matthias Zwicker. 2018. Specular-to-diffuse translation for multi-view reconstruction. In *ECCV*.
- Tong Wu, Jiaqi Wang, Xingang Pan, Xudong Xu, Christian Theobalt, Ziwei Liu, and Dahua Lin. 2022. Voxurf: Voxel-based Efficient and Accurate Neural Surface Reconstruction. *arXiv preprint arXiv:2208.12697* (2022).
- Jianfeng Yan, Zizhuang Wei, Hongwei Yi, Mingyu Ding, Runze Zhang, Yisong Chen, Guoping Wang, and Yu-Wing Tai. 2020. Dense hybrid recurrent multi-view stereo net with dynamic consistency checking. In *ECCV*.

- Jiayu Yang, Wei Mao, Jose M Alvarez, and Miaomiao Liu. 2020. Cost volume pyramid based depth inference for multi-view stereo. In *CVPR*.
- Wenqi Yang, Guanying Chen, Chaofeng Chen, Zhenfang Chen, and Kwan-Yee K Wong. 2022a. Ps-nerf: Neural inverse rendering for multi-view photometric stereo. In *ECCV*.
- Wenqi Yang, Guanying Chen, Chaofeng Chen, Zhenfang Chen, and Kwan-Yee K Wong. 2022b. S3-NeRF: Neural Reflectance Field from Shading and Shadow under a Single Viewpoint. In *NeurIPS*.
- Yao Yao, Zixin Luo, Shiwei Li, Tian Fang, and Long Quan. 2018. MVSNet: Depth inference for unstructured multi-view stereo. In *ECCV*.
- Yao Yao, Jingyang Zhang, Jingbo Liu, Yihang Qu, Tian Fang, David McKinnon, Yanghai Tsin, and Long Quan. 2022. NeLF: Neural incident light field for physically-based material estimation. In *ECCV*.
- Lior Yariv, Jiatao Gu, Yoni Kasten, and Yaron Lipman. 2021. Volume rendering of neural implicit surfaces. In *NeurIPS*.
- Lior Yariv, Yoni Kasten, Dror Moran, Meirav Galun, Matan Atzmon, Basri Ronen, and Yaron Lipman. 2020. Multiview neural surface reconstruction by disentangling geometry and appearance. In *NeurIPS*.
- Weicai Ye, Shuo Chen, Chong Bao, Hujun Bao, Marc Pollefeys, Zhaopeng Cui, and Guofeng Zhang. 2022. Intrinsicnerf: Learning intrinsic neural radiance fields for editable novel view synthesis. *arXiv preprint arXiv:2210.00647* (2022).
- Ye Yu, Abhimitra Meka, Mohamed Elgharib, Hans-Peter Seidel, Christian Theobalt, and William AP Smith. 2020. Self-supervised outdoor scene relighting. In *ECCV*.
- Ye Yu and William AP Smith. 2019. Inverserendernet: Learning single image inverse rendering. In *CVPR*.
- Jason Zhang, Gengshan Yang, Shubham Tulsiani, and Deva Ramanan. 2021c. NeRS: Neural reflectance surfaces for sparse-view 3d reconstruction in the wild. In *NeurIPS*.
- Kai Zhang, Fujun Luan, Zhengqi Li, and Noah Snavely. 2022a. IRON: Inverse Rendering by Optimizing Neural SDFs and Materials from Photometric Images. In *CVPR*.
- Kai Zhang, Fujun Luan, Qianqian Wang, Kavita Bala, and Noah Snavely. 2021a. PhySG: Inverse rendering with spherical gaussians for physics-based material editing and relighting. In *CVPR*.
- Kai Zhang, Gernot Riegler, Noah Snavely, and Vladlen Koltun. 2020. Nerf++: Analyzing and improving neural radiance fields. *arXiv preprint arXiv:2010.07492* (2020).
- Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. 2018. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*.
- Xiuming Zhang, Pratul P Srinivasan, Boyang Deng, Paul Debevec, William T Freeman, and Jonathan T Barron. 2021b. NeRFactor: Neural factorization of shape and reflectance under an unknown illumination. In *SIGGRAPH*.
- Yuanqing Zhang, Jiaming Sun, Xingyi He, Huan Fu, Rongfei Jia, and Xiaowei Zhou. 2022b. Modeling Indirect Illumination for Inverse Rendering. In *CVPR*.
- Boming Zhao, Bangbang Yang, Zhenyang Li, Zuoyue Li, Guofeng Zhang, Jiashu Zhao, Dawei Yin, Zhaopeng Cui, and Hujun Bao. 2022b. Factorized and controllable neural re-rendering of outdoor scene for photo extrapolation. In *ACM MM*.
- Fuqiang Zhao, Yuheng Jiang, Kaixin Yao, Jiakai Zhang, Liao Wang, Haizhao Dai, Yuhui Zhong, Yingliang Zhang, Minye Wu, Lan Xu, et al. 2022a. Human performance modeling and rendering via neural animated mesh. In *SIGGRAPH Asia*.
- Quan Zheng, Gurprit Singh, and Hans-Peter Seidel. 2021. Neural Relightable Participating Media Rendering. *NeurIPS* (2021).

## A APPENDIX

### A.1 BRDF model

We adopt the Cook-Torrance BRDF [Cook and Torrance 1982]. The basic reflection ratio  $F_0 = (m * \mathbf{a} + (1 - m) * 0.04)$  where  $\mathbf{a}$  is the albedo and  $m$  is the metalness. Then, the Fresnel term is

$$F = F_0 + (1 - F_0)(1 - (\mathbf{h} \cdot \omega_o))^5, \quad (20)$$

where  $\mathbf{h}$  is the half-way vector and  $\omega_o$  is the viewing direction. The Geometry function is based on the Schlick-GGX Geometry function:

$$G(\mathbf{n}, \omega_o, \omega_i, k) = G_{\text{sub}}(\mathbf{n}, \omega_o, k)G_{\text{sub}}(\mathbf{n}, \omega_i, k), \quad (21)$$

$$G_{\text{sub}}(\mathbf{n}, \mathbf{v}, k) = \frac{\mathbf{n} \cdot \mathbf{v}}{(\mathbf{n} \cdot \mathbf{v})(1 - k) + k}, \quad (22)$$

where  $k = \rho^4/2$  and  $\rho$  is the roughness. The normal distribution for Stage II is the Trowbridge-Reitz GGX distribution

$$N(\mathbf{n}, \mathbf{h}, \alpha) = \frac{\alpha^2}{\pi((\mathbf{n} \cdot \mathbf{h})(\alpha^2 - 1) + 1)^2}, \quad (23)$$

where  $\alpha = \rho^2$ .

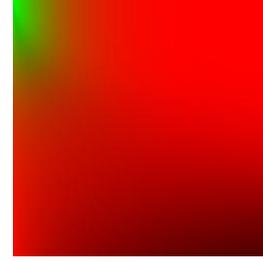


Fig. 21. The prefiltered  $F_1$  and  $F_2$  for the integral of the BRDF.

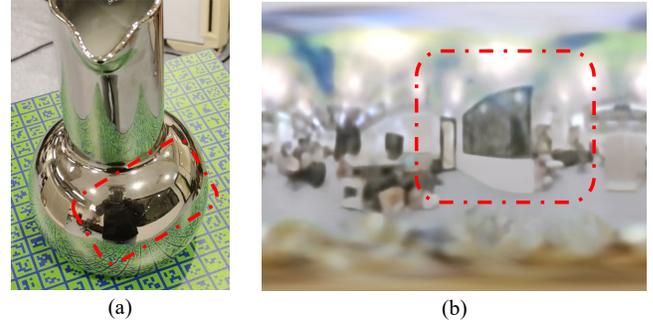


Fig. 22. **Visualization of direct environment lights.** (a) An Image captured in an indoor environment. (b) Estimated direct lights. Even in an indoor environment, modeling the direct environment lights with only directions produces a satisfactory approximation for surface reconstruction.

The prefiltered BRDF  $F_1$  and  $F_2$  used in the split-sum are stored in an image as shown in Fig. 21, where  $F_1$  is the red color,  $F_2$  is the green color, x-y axis represents the roughness  $\rho$  and the  $\mathbf{n} \cdot \omega_o$ , respectively. Given  $\rho$  and  $\mathbf{n} \cdot \omega_o$ , we interpolate on the image of Fig. 21 to get  $F_1$  and  $F_2$ .

### A.2 Discussion on the direct light representation

In Eq. 9, the direct light is represented by  $g_{\text{direct}}(\omega_i)$ , which only takes a direction  $\omega_i$  as input. This direct light implicitly contains an assumption that direct lights all come from the light sources located at infinity. Such an assumption is accurate enough for surface reconstruction even in a challenging indoor environment. We provide an example in Fig. 22 to show the environment lights estimated by our method in Stage I.

However, when using such direct light representation in the BRDF estimation of Stage II, we find the approximation is not accurate enough. The main reason for this is that the BRDF estimation is very sensitive to the locations of strong light sources at a finite distance. As shown in Fig. 23 (a), the estimated strong light sources in red circles are all enlarged. The reason is illustrated by Fig. 23 (b). There are three rays (A, B, C) corresponding to three surface points (A, B, C). Since both Ray A and Ray C are pointing to the light source, the network will increase the light intensity on the reflective directions of Ray A and C, which causes the strong light source to enlarge. Meanwhile, Ray B is not pointing to the light source but has the same reflective direction as Ray A. Increasing the light intensity

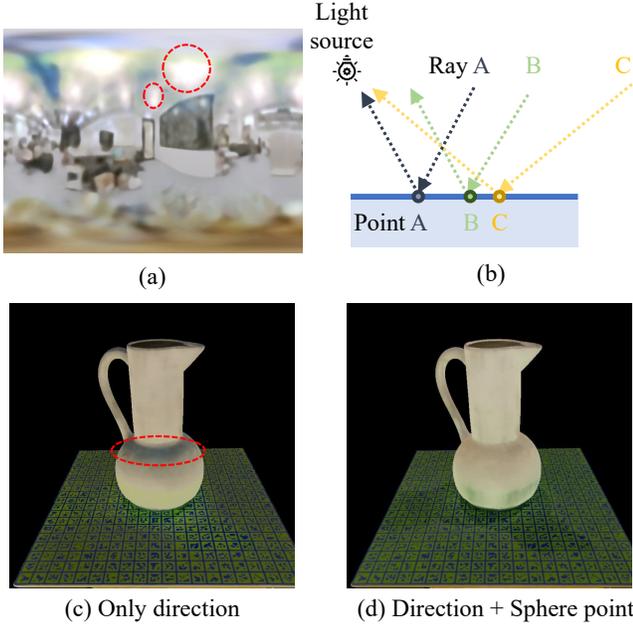


Fig. 23. (a) The reconstructed direct environment lights using only directions as input. Note that the strong light sources in the red circles are enlarged. (b) A diagram containing three rays and a light source to show the sensitiveness to strong light sources. (c) Estimated albedo using only directions as the direct light representation. Red circles show the inaccurate albedo estimation. (d) Estimated albedo using both directions and sphere intersection points as the direct light representation.



Fig. 24. Results with or without Sphere Intersection Encoding (SIE).

on the reflective direction of Ray A will also make Ray B brighter. Therefore, the model will tend to decrease the albedo of Point B. This causes the phenomenon shown in Fig. 23 (c) that the region in the red circle has a darker albedo than other regions. To resolve this problem, we use the direct light representation  $g_{\text{direct}}(\mathbf{q}(\mathbf{p}, \omega_i), \omega_i)$  in Stage II, where  $\mathbf{q}(\mathbf{p}, \omega_i)$  is the intersection point on the bounding sphere of the ray emitting from  $\mathbf{p}$  along the direction  $\omega_i$ . We call this sphere intersection encoding. With sphere intersection encoding, we are able to more accurately estimate the albedo as shown in Fig. 23 (d).

An optional choice is to add the sphere intersection encoding in Stage I. We conduct experiments on the “Bunny” and “Bear” using this sphere intersection encoding. The results are shown in Fig. 24. We find that adding such sphere intersection encoding degenerates the performance. The main reason is that adding a sphere

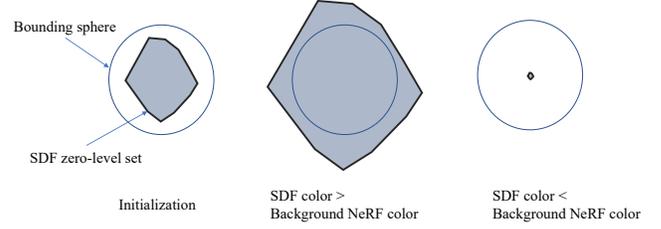


Fig. 25. Imbalanced convergence speed between the outer background NeRF color and the SDF shading color will make the reconstruction fail. (left) The initialized zero-level set of neural SDF locates inside the bounding sphere. Note that only the sample points inside the bounding sphere will be used in the neural SDF to compute the opacity density [Wang et al. 2021b]. (middle) If the SDF shading color converges faster than the background NeRF color, then the surface of the neural SDF will dilate to exceed the bounding sphere and is unable to shrink back to the bounding sphere. (right) If the SDF shading color converges slower than the background NeRF color, then the surface of the neural SDF will shrink and eventually disappear. In both cases, the training of NeuS will fail.

intersection improves the fitting ability, which makes the network indiscriminately focus on color fitting even using distorted surfaces instead of recovering the faithful geometry. Stage II is not affected by such an overfitting problem because Stage II uses Monte Carlo sampling to sample many rays to render a single pixel, which avoids overfitting in a single direction.

### A.3 Rationale of light integral approximations

The rationale of the two light integral approximations is that:

- (1) On a smooth surface with a small specular lobe, the light integral is mainly determined by the light from the reflective direction. In this case, the occlusion probability of the reflective direction will be the dominant factor in the integral. Therefore, the first approximation of using the probability of the reflective direction is a good approximation for all other directions. Meanwhile, as shown in [Verbin et al. 2022], an integrated directional encoding (IDE) with a small roughness  $\rho$  will produce a high-frequency directional encoding that is suitable to represent strong view-dependent colors on smooth surfaces. Therefore, the IDE approximation also provides a good estimation.
- (2) On a rough surface with a large specular lobe, the light integral is not only affected by the reflective direction so the first approximation seems to be too assumptive. However, since the light integral on a large specular lobe will mainly be white and change slowly with the view direction, it is not essential that we use the  $g_{\text{direct}} (s = 0)$  or the  $g_{\text{indirect}} (s = 1)$  to predict such a slowly-changing white light. Meanwhile, the IDE with a large  $\rho$  will also produce a low-frequency directional encoding for this case.

### A.4 Stabilization loss

As shown in Fig. 25, the initialization of the neural SDF follows [Atzmon and Lipman 2020]. We follow NeuS [Wang et al. 2021b] to use a background NeRF to render the background of the image, which

is actually a NeRF++ [Zhang et al. 2020] and should be referred to NeuS [Wang et al. 2021b] for more details. However, the imbalanced convergence speed between the shading color of the neural SDF and the background color provided by the background NeRF will cause the training process to collapse. If the convergence speed of the foreground shading color is faster, then the surface tends to enlarge to exceed the bounding sphere. Otherwise, it would shrink and eventually disappear. Both cases lead to the failure of training. To avoid these, we prevent the zero-level set from expanding outside the bounding sphere or shrinking to disappear by a stabilization loss  $\ell_{\text{stable}}$ . Note that the bounding sphere is normalized to a unit sphere at the origin. We first sample some points  $S_1 = \{\mathbf{p} \in \mathbb{R}^3 \mid \|\mathbf{p}\| < 0.02\}$  and  $S_2 = \{\mathbf{p} \in \mathbb{R}^3 \mid \|\mathbf{p}\| > 0.98\}$  in the bounding sphere. Then, we penalize the SDF values of these sample points to get close to 0 in the first 1000 steps. We find that such a stabilization loss is essential for a correct convergence of NeuS [Wang et al. 2021b] and our method on these challenging reflective objects.

#### A.5 Network architectures and implementation details

**Architectures.** We illustrate the architectures of all MLPs and implementation details in the following and more detailed structures can be referred to the codes at <https://github.com/liuyuan-pal/NeRO>. Same as NeuS [Wang et al. 2021b],  $g_{\text{sdf}}$  uses a positional encoding with a frequency of 6 as inputs and contains 8 linear layers with the channel number of 256 and a skip connection on the 4th layer. The output of  $g_{\text{sdf}}$  is an SDF value and a 256-dim feature vector. The 256-dim feature vector is fed into  $g_{\text{material}}$  with 4 linear layers of 256 channels to output roughness, metalness, and albedo. Both the direct light MLP  $g_{\text{direct}}$  and indirect light MLP  $g_{\text{indirect}}$  use the spherical harmonics directional encoding up to degree 5. The occlusion prob MLP  $g_{\text{occ}}$  and the indirect light MLP  $g_{\text{indirect}}$  both use the positional encoding of frequency 8.  $g_{\text{occ}}$ ,  $g_{\text{direct}}$ , and  $g_{\text{indirect}}$  all contain 4 layers with a width of 256. The 2D NeRF built on the XoY plane of the camera system also contains 4 layers with a width of 256. In Stage II, we use a base feature extraction network with the same structure as the  $g_{\text{sdf}}$  but use a higher positional encoding of 8 to extract 256-dim feature vectors. All other MLPs ( $g_{\text{direct}}$ ,  $g_{\text{indirect}}$  and  $g_{\text{material}}$ ) in Stage II have the same structures as Stage I. We summarize the trainable components for both stages in Table 7.

**Training details.** We train all MLPs in Stage II from scratch. All activation functions in MLPs are ReLUs and we follow NeuS [Wang et al. 2021b] to add weight normalization on all linear weights except the SDF MLP. The final activation for the  $g_{\text{material}}$  is Sigmoid to get values in  $[0, 1]$  while the final activation function for  $g_{\text{direct}}$  and  $g_{\text{indirect}}$  is the exponential function to get light radiance in  $[0, \infty)$ . We apply the standard gamma correction to get colors in the sRGB space before computing the rendering loss. On all objects except the bunny object, we freeze the variance used in the computation of opacity density for the first 15k steps for better convergence. On the bunny object, we find such a freezing operation will make the SDF unable to reconstruct the hole between the legs of the bunny. The weights used in the loss computation are  $\lambda_{\text{eikonal}} = 0.1$ ,  $\lambda_{\text{occ}} = 1.0$ ,  $\lambda_{\text{smooth}} = 0.05$  and  $\lambda_{\text{light}} = 0.1$  for all experiments.

Object	Bear	Bunny	Coral	Maneki	Vase
Image Num.	97	129	126	128	128

Table 6. Image numbers of each object in the Glossy-Real dataset.

	$g_{\text{sdf}}$	$g_{\text{material}}$	$g_{\text{direct}}$	$g_{\text{indirect}}$	$g_{\text{occ}}$	$g_{\text{camera}}$
Stage I	✓	✓	✓	✓	✓	✓
Stage II		✓	✓	✓		✓

Table 7. Trainable components in Stage I and Stage II.

	Ref-NeRF	NeuS	NeuS-Large	Ours
Model size	5.2M	5.3M	8.7M	7.7M
Chamfer distance↓	0.0169	0.0212	0.0200	<b>0.0038</b>

Table 8. Model sizes and average CDs on “Bell” and “Cat” of Ref-NeRF [Verbin et al. 2022], NeuS [Wang et al. 2021b] and our method. “NeuS-Large” means that we use a deeper and wider color network for NeuS to make the model larger.

#### A.6 Dataset statistics

The Glossy-Blender dataset contains 128 training images for each object, which are uniformly distributed on the upper hemisphere. In order to evaluate the NVS quality, we additionally render 8 evenly-distributed test images to compute the metrics of PSNR, SSIM, and LPIPS. Both training and test images have a resolution of 800×800. The number of images for each object in the Glossy-Real dataset is shown in Table 6. All images have a resolution of 1024×768 and are used for training because our target is shape and material reconstruction.

#### A.7 Model size

In this section, we compare the model sizes of NeRO with Ref-NeRF [Verbin et al. 2022] and NeuS [Wang et al. 2021b] in Table 8. NeRO uses the same networks as NeuS to model SDF and background. The network of NeRO is slightly larger due to the decomposition of the color function into materials and lighting. We further try to increase the width and the depth of the color network of NeuS to get an ~8M model called “NeuS-Large”. However, even with large model size, NeuS-Large is still unable to correctly reconstruct the reflective surfaces and shows similar results to the original NeuS model.

#### A.8 Relighting results in SSIM/LPIPS

We provide the evaluation of relighting quality on the Glossy-Blender dataset in terms of SSIM [Wang et al. 2004] and LPIPS [Zhang et al. 2018] in Table 10 and Table 9 respectively. On both metrics, we observe similar results to PSNR where our method outperforms the baselines by a significant margin on these reflective objects.

	NDR	NDRMC	MII	NeILF	Ours
Bell	0.913	0.903	0.882	0.916	<b>0.965</b>
Cat	0.901	0.907	0.889	0.921	<b>0.962</b>
Teapot	0.844	0.899	0.884	0.918	<b>0.977</b>
Potion	0.824	0.858	0.885	0.903	<b>0.950</b>
TBell	0.739	0.883	0.870	0.897	<b>0.968</b>
Angel	0.864	0.865	0.867	0.889	<b>0.911</b>
Horse	0.930	0.935	0.893	0.945	<b>0.954</b>
Luyu	0.854	0.859	0.850	0.877	<b>0.914</b>
Avg.	0.859	0.889	0.878	0.908	<b>0.950</b>

Table 9. SSIM $\uparrow$  of NDR [Munkberg et al. 2022], NDRMC [Hasselgren et al. 2022], MII [Zhang et al. 2022b], NeILF [Yao et al. 2022] and our method on the Glossy-Blender dataset.

	NDR	NDRMC	MII	NeILF	Ours
Bell	0.0989	0.1180	0.1477	0.1056	<b>0.0557</b>
Cat	0.1104	0.1150	0.1437	0.0737	<b>0.0523</b>
Teapot	0.1385	0.1207	0.1373	0.0980	<b>0.0283</b>
Potion	0.2062	0.2027	0.1723	0.1443	<b>0.0843</b>
TBell	0.2632	0.1703	0.2117	0.1407	<b>0.0460</b>
Angel	0.1066	0.1213	0.1183	0.0940	<b>0.0790</b>
Horse	0.0492	0.0530	0.0723	0.0457	<b>0.0403</b>
Luyu	0.1085	0.1080	0.1300	0.0960	<b>0.0723</b>
Avg.	0.1352	0.1261	0.1417	0.0996	<b>0.0573</b>

Table 10. LPIPS $\downarrow$  of NDR [Munkberg et al. 2022], NDRMC [Hasselgren et al. 2022], MII [Zhang et al. 2022b], NeILF [Yao et al. 2022] and our method on the Glossy-Blender dataset.

	NeuS	Ref-NeRF	Ours
PSNR $\uparrow$	27.80	27.86	<b>29.73</b>
SSIM $\uparrow$	0.875	0.878	<b>0.904</b>
LPIPS $\downarrow$	0.365	0.375	<b>0.324</b>

Table 11. NVS quality of NeuS [Wang et al. 2021b] and Ref-NeRF [Verbin et al. 2022] on the Glossy-Blender dataset with PSNR, SSIM and LPIPS.

### A.9 Novel-view synthesis quality

To show the quality of novel view synthesis (NVS), we additionally render 8 novel-view images on the Glossy-Blender dataset and report the NVS quality on these images in Table 11 in terms of PSNR, SSIM, and LPIPS.

### A.10 Results on less- or non-reflective objects

NeRO is also able to reconstruct less or non-reflective objects, which we demonstrate from three aspects. First, we show the reconstruction results on three less reflective objects in Fig 26, where we capture  $\sim 100$  images with a resolution of  $1024 \times 768$  on each object and recover the camera poses by COLMAP. Second, we show that NeRO can simultaneously reconstruct reflective and non-reflective

	scan24	scan37	scan110	scan114	scan118	scan122
NeuS	1.00	1.37	1.20	0.35	0.49	0.54
Ours	1.10	1.13	1.14	0.39	0.52	0.57

Table 12. CDs $\downarrow$  on the DTU dataset.

Description	Angel	Bell	Cat	Teapot	Avg.
Only indirect lights	0.0037	0.0043	0.0210	0.0064	0.0089
Ref-NeRF+NeuS	0.0038	0.0062	0.0219	0.0064	0.0096
Ours	0.0034	0.0032	0.0044	0.0037	0.0037

Table 13. Comparison with the direct combination of Ref-NeRF and NeuS in terms of CD.

objects in Fig. 27, where we capture  $\sim 200$  images for each set of objects. Finally, we further evaluate NeRO on the DTU dataset. Since NeRO assumes a static light environment and the DTU dataset contains images with inconsistent light environments and shadows, we manually remove images with inconsistent light environments for training. The quantitative and qualitative results are shown in Table 12 and Fig. 28.

### A.11 Direct combination of Ref-NeRF with NeuS

To improve the reflective color fitting of NeuS [Wang et al. 2021b], an alternative solution is to combine the color function of Ref-NeRF [Verbin et al. 2022] with the neural SDF of NeuS. This combination leads to a model similar to Model 2 of Table 4 in the geometry ablation study, both of which predict specular lights from a PE and an IDE. The differences are that Ref-NeRF does not consider the integral of material  $M_{\text{material}}$  in Eq. 7 and that Ref-NeRF directly predicts the diffuse color from an MLP while Model 2 relies the IDE on the normal direction (Eq. 8 and Eq. 11) to compute diffuse color. We report the CDs of this model in Table 13, which are comparable to Model 2 and worse than NeRO. The qualitative results of this combination model are shown in Fig. 29.

### A.12 Copyrights

The Glossy-Blender dataset is created from the models listed in Table 14. On some objects, we modified their appearances and geometry to make the dataset. All the HDR images used in relighting or rendering are downloaded from <https://polyhaven.com> under the CC0 license.



Fig. 26. Reconstruction results on less reflective objects.



Fig. 27. Reconstruction results on both reflective and non-reflective objects.



Fig. 28. Qualitative reconstruction results on the DTU dataset.



Fig. 29. Reconstructed surfaces of directly combining Ref-NeRF [Verbin et al. 2022] with NeuS [Wang et al. 2021b].

Model	Creator	License	Link
Bell	jQueary	CC BY 4.0	here
Cat	Suushimi	CC BY-NC 4.0	here
Teapot	Martin Newell	N/A	here
Luyu	romullus	CC BY-SA 4.0	here
TBell	gla_bot	CC BY 4.0	here
Horse	halimi13744	CC BY 4.0	here
Potion	Blender3D	CC BY 4.0	here
Angel	SebastianSosnowski	CC BY 4.0	here

Table 14. Copyrights of all models used in the Glossy-Blender dataset.

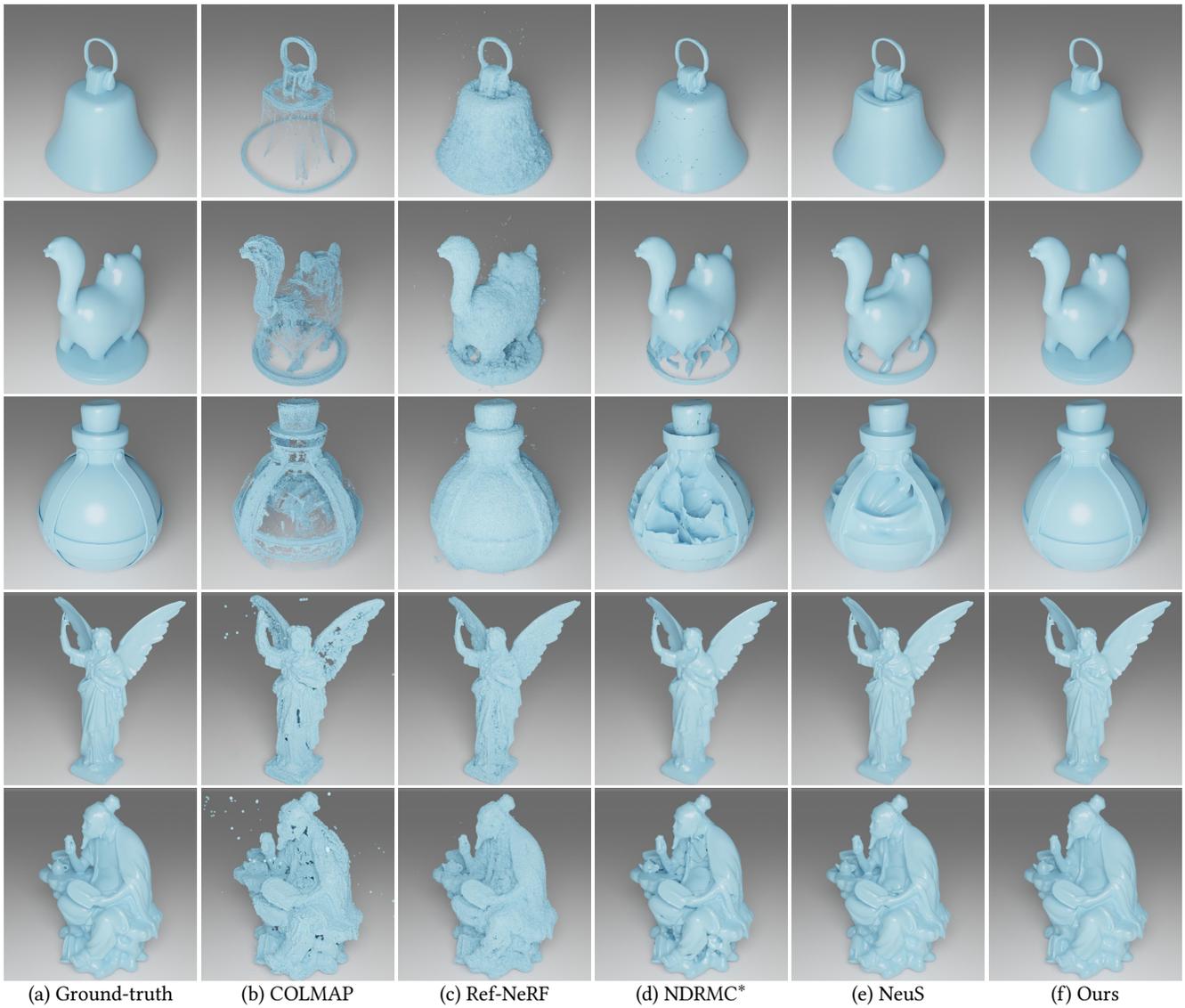


Fig. 30. **Ground-truth and reconstructed surfaces of the Glossy-Blender dataset.** We compare our results with COLMAP [Schönberger et al. 2016], Ref-NeRF [Verbin et al. 2022], NDRMC [Hasselgren et al. 2022], and NeuS [Wang et al. 2021b]. \*NDRMC [Hasselgren et al. 2022] is trained with ground-truth object masks while the other methods do not use object masks. The supplementary video contains more qualitative results.



Fig. 31. **Relighting objects in the Glossy-Blender dataset.** We compare our method with NDR [Munkberg et al. 2022], NDRMC [Hasselgren et al. 2022], MII [Zhang et al. 2022b] and NeLF [Yao et al. 2022]. Note that all relighted images are normalized to match the average colors of the ground-truth images. The supplementary video contains more qualitative results.



Fig. 32. **Relighting objects from the Glossy-Real dataset.** We provide a visual comparison with NDR [Munkberg et al. 2022], NDRMC [Hasselgren et al. 2022], MII [Zhang et al. 2022b], and NeLF [Yao et al. 2022]. We provide the input image with the nearest viewpoint and the relighting HDR map as a reference. The supplementary video contains more qualitative results.

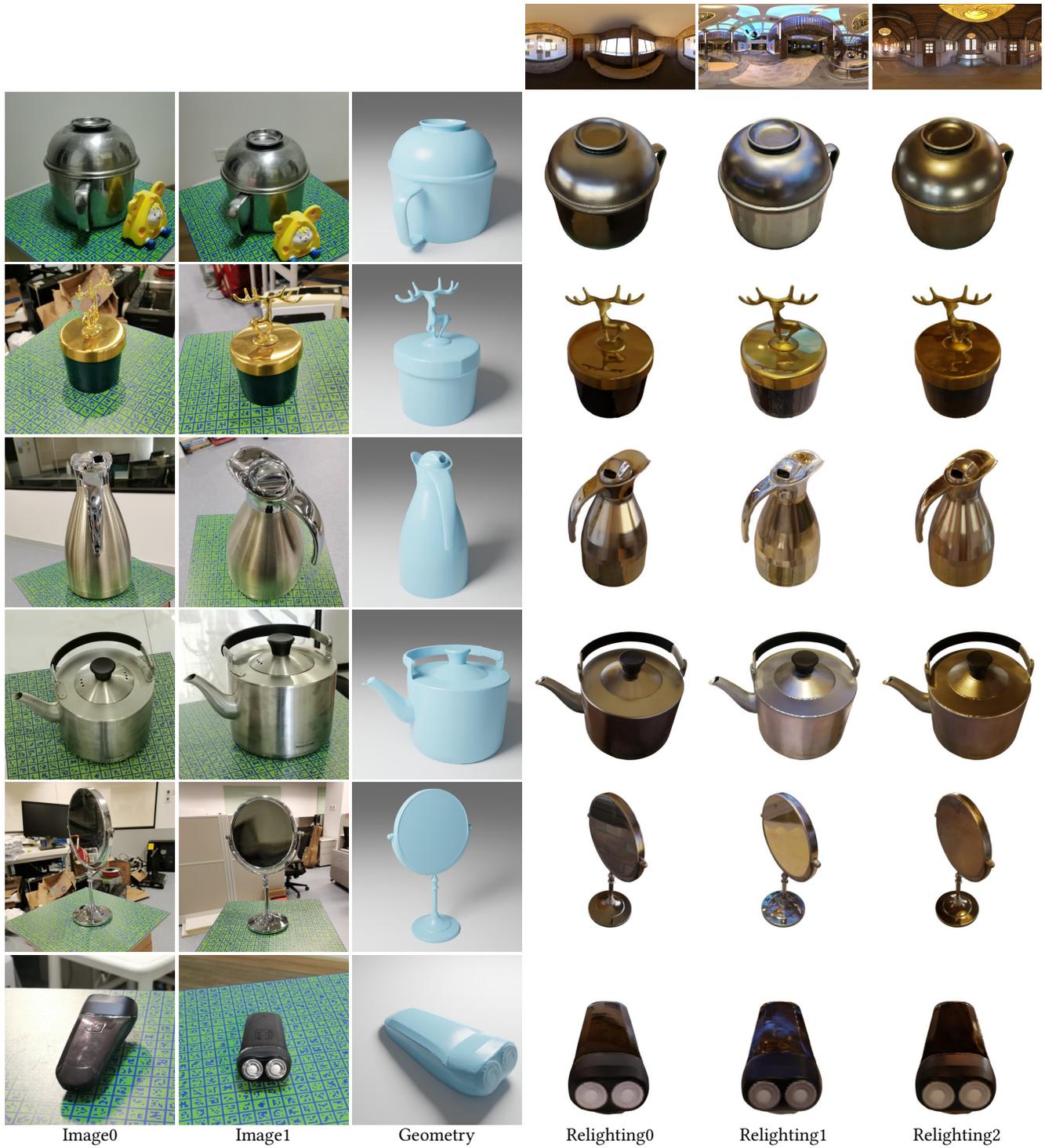


Fig. 33. **Reconstruction results on other real objects.** The supplementary video contains more qualitative results.