# NeILF: Neural Incident Light Field for Physically-based Material Estimation

Yao Yao, Jingyang Zhang, Jingbo Liu, Yihang Qu, Tian Fang, David McKinnon, Yanghai Tsin, and Long Quan

Apple          HKUST

**Abstract.** We present a differentiable rendering framework for material and lighting estimation from multi-view images and a reconstructed geometry. In the framework, we represent scene lightings as the Neural Incident Light Field (NeILF) and material properties as the surface BRDF modelled by multi-layer perceptrons. Compared with recent approaches that approximate scene lightings as the 2D environment map, NeILF is a fully 5D light field that is capable of modelling illuminations of any static scenes. In addition, occlusions and indirect lights can be handled naturally by the NeILF representation without requiring multiple bounces of ray tracing, making it possible to estimate material properties even for scenes with complex lightings and geometries. We also propose a smoothness regularization and a Lambertian assumption to reduce the material-lighting ambiguity during the optimization. Our method strictly follows the physically-based rendering equation, and jointly optimizes material and lighting through the differentiable rendering process. We have intensively evaluated the proposed method on our in-house synthetic dataset, the DTU MVS dataset, and real-world BlendedMVS scenes. Our method is able to outperform previous methods by a significant margin in terms of novel view rendering quality, setting a new state-of-the-art for image-based material and lighting estimation.

**Keywords:** differentiable rendering, physically-based rendering, BRDF estimation, incident light field

## 1  Introduction

Material estimation from a set of sparse images is a challenging task in both computer vision and computer graphics. The problem is usually approached by inverse rendering, where the spatially-varying bidirectional reflectance distribution functions (SV-BRDFs) and lightings of the scene are jointly optimized by minimizing the rendering loss. However, the problem is hard to solve due to the complex form of the BRDF and the high-dimensional nature of scene illuminations. To mitigate the problem, previous methods usually apply simplified material and lighting models. For example, non-spatially varying BRDF [49] is applied for certain types of objects; approximated illuminations, such as co-located flash lights [4, 5, 28, 33] and environment maps [6, 7, 27, 49, 51], are

applied to reduce the complexity of the scene lighting. In most scenarios, special capturing devices or environments are required to assist the estimation, limiting these methods to real-world applications. As the result, a practical material estimator is still missing.

On the other hand, recent progress on neural representation has shown promising results for lighting modelling. NeRF [25] jointly optimizes a neural radiance field and a density field, which has demonstrated great success for novel view synthesis. The surface light field is applied to model the outgoing light from the surface, which has been widely applied to neural surface reconstructions [45, 48]. Other methods further decompose observed lights into neural material properties and environmental lightings. However, similar to classical methods, they either use simplified lighting representations [4, 6, 7, 27], or apply approximated occlusion and indirect light handling [36, 51]. Until now, lighting modelling is still an open problem in image-based material estimation.

In this work, we address this long-standing problem by representing scene lightings as the neural incident light field. Without losing generality, the proposed NeILF is capable of modelling lighting conditions of any static scenes. Also, occlusions and indirect lights could be naturally handled in the proposed framework without the need for tracing multiple bounces of rays. For material properties, we consider a simplified Disney BRDF model [8] consisting of base color, roughness and metallic. Implementation-wise, we use multi-layer perceptrons (MLPs) to represent both the incident light field and the BRDF. The NeILF network takes a 5D vector of location and incident direction as inputs, and returns as output a RGB value of the incident light; the material network takes a 3D location as input, and outputs a 5D vector of surface BRDF properties. Meanwhile, to reduce the ambiguity between the material and the scene lighting, we propose two regularization terms, namely the bilateral smoothness and the Lambertian assumption, to constrain the optimization of roughness and metallic. Finally, we analyze similarities between NeILF for material estimation and NERF for novel view synthesis [25], providing readers an intuitive explanation of the difficulty and solvability of the problem.

We demonstrate in several datasets that our method significantly outperforms previous state-of-the-art in terms of novel view rendering accuracy. Our method is able to recover the surface BRDF even for scenes with complex lightings and geometries, which cannot be handled by previous environment map based methods. To summarize, main contributions of the paper include:

- Representing scene lightings using the neural incident light field, where occlusions and indirect lights of the scene can be naturally handled.
- A differentiable framework for joint material and lighting estimation, which significantly outperforms previous state-of-the-art in different datasets.
- A bilateral smoothness and a Lambertian assumption to constrain the roughness and the metallic, reducing the material-lighting ambiguity during the network optimization.

## 2    Related Works

### 2.1    The Rendering Equation

The rendering equation [15] computes the emitted radiance from a surface point $\mathbf{x}$ along a viewing direction $\boldsymbol{\omega_o}$:

$$L_o(\boldsymbol{\omega_o}, \mathbf{x}) = \int_\Omega f(\boldsymbol{\omega_o}, \boldsymbol{\omega_i}, \mathbf{x}) L_i(\boldsymbol{\omega_i}, \mathbf{x})(\boldsymbol{\omega_i} \cdot \mathbf{n}) d\boldsymbol{\omega_i}, \tag{1}$$

where $\mathbf{n}$ is the normal of the surface, $L_i$ is the incoming light from direction $\boldsymbol{\omega_i}$, and $f$ is the BRDF function to describe the the reflectance property, which is usually decomposed into a diffuse term and a specular term $f = f_d + f_s$. The integration is performed over all incident direction $\boldsymbol{\omega_i}$ on the hemisphere $\Omega$ where $\boldsymbol{\omega_i} \cdot \mathbf{n} > 0$.

The goal of material estimation is to recover continues functions of the scene lighting $L_i$ and the BRDF property $f$ in the above equation. Due to the complex form of the scene lighting and the material property, it is crucial to select suitable representations for $L_i$ and $f$. In this paper, we propose to use a neural incident light field to model $L_i$ (Sec. 3.1), and apply a simplified Disney BRDF [8] model to approach the BRDF $f$ (Sec. 3.2). Below we give a brief review on the physically-based material estimation from multi-view images.

### 2.2    Differentiable Rendering

Unlike classical approaches that recover 3D scene parameters in a forward reconstruction manner, differentiable rendering [2, 16] inverses the rendering process in graphics, and optimize all parameters by minimizing the difference between rendered and input images.

Recently, the technique has been combined with neural representations and has shown promising results for image-based 3D problems. NeRF [25] and followup works [21, 24, 50] decouple a 3D scene into a density field and a radiance field. Other methods also apply implicit functions to model different geometry [22, 23, 48] and appearance representations [29, 45]. In another line of works, the received radiance is further decomposed into BRDF properties and input light sources [6, 36, 49, 51]. Our method follows this practice and applies a neural BRDF model to approach the material property of the surface.

### 2.3    Material and Lighting Estimation

Due to the difficulty of joint material and lighting estimation, previous methods usually apply additional sensors or controlled lightings to facilitate the optimization process. For example, additional sensors [3, 12, 31], co-located flash lights [4, 5, 28, 33], or turn-table settings [10, 41] are applied to capture image or scene lightings. Moreover, simplified material and lighting models, e.g., the non-spatially varying BRDF [49] or approximated illuminations of environment
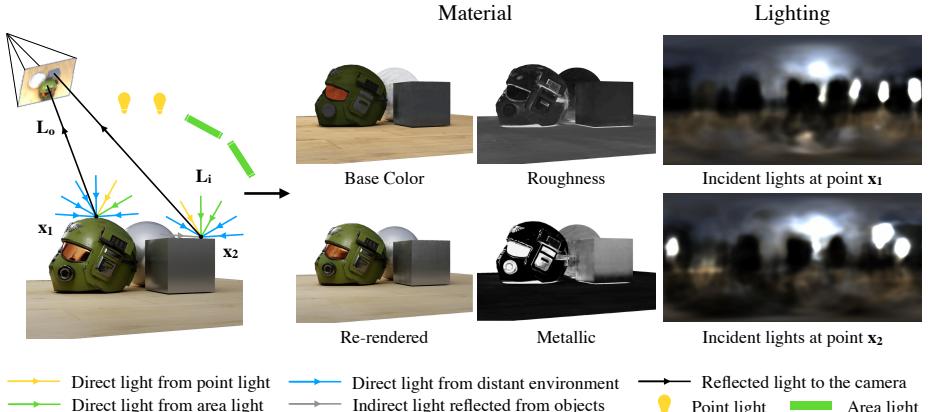
Fig. 1: Illustrations of the proposed method and our material and lighting estimation results. NeILF is capable of modelling the joint illumination of direct/indirect lights from different sources. Estimated incident lights at point $\mathbf{x_1}$ and point $\mathbf{x_2}$ well explain the mixed lighting of the scene, including an environment map with high-radiance sun light, two near-range point lights, and two near-range area lights.

maps [6, 49, 51], are applied to reduce the complexity of the problem. NeRV [36] introduces the visibility field to model indirect lights, however, requires environmental lightings to be known in advance. Nevertheless, such mitigations will inevitably limit these methods to real-world applications. In contrast, our method applies a unified incident light field to represent different light sources in the scene, and is capable of jointly estimating material and lighting under any lighting conditions.

## 3    Method

### 3.1    Neural Incident Light Field

One of the keys to invert the rendering equation is to model the incoming light $L_i$ in a correct way. Ideally, $L_i$ should take into account 1) *direct lights* from light sources in the scene, 2) *occlusions* that block the surface point from receiving direct lights, and 3) *indirect lights* that are reflected from other surface points. However, each of the three components is hard to model. Previous methods [6, 20, 49, 51] usually approximate direct lights as an environment map and hardly handle indirect lights as they require multi-bounce raytracing.

In contrast, we formulate incoming lights in the scene directly as a neural incident light field, where an MLP takes a point location $\mathbf{x}$ and an incident direction $\boldsymbol{\omega}$ as inputs, and returns an incident light radiance $L$ as output:

$$\mathbb{L} : \{\mathbf{x}, \boldsymbol{\omega}\} \to \mathbf{L}. \tag{2}$$

Without losing generality, the proposed NeILF representation is capable of modelling the joint illumination effect of direct/indirect lights and occlusions of *any static scenes*. An illustration is shown in Fig. 1. Compared with the commonly used environment map, NeILF is able to handle the spatially-varying illumination effect, making it possible to estimate material for scenes with complex geometries and lightings.

## 3.2 Simplified Disney BRDF

In this section, we describe the BRDF representation used in the proposed framework. We apply a simplified Disney principled BRDF model, where the BRDF of a surface point $\mathbf{x}$ is parameterized by a base color $\mathbf{b}(\mathbf{x}) \in [0,1]^3$, a roughness $r(\mathbf{x}) \in [0,1]$ and a metallic $m(\mathbf{x}) \in [0,1]$, which is a subset of the full Disney model [8]. Similar to the neural incident lighting field, BRDF parameters are also stored using multi-layer perceptrons:

$$\mathbb{B} : \mathbf{x} \to \{\mathbf{b}, r, m\}, \tag{3}$$

where the MLP takes a 3D surface point $\mathbf{x}$ as input, and returns the 5-channel BRDF parameters as output. Note that other representations, e.g. UV atlas or per-vertex BRDF parameters, can also be applied. Here we choose the neural representation because it has been proven to be effective for modelling continuous functions in 3D space [25, 49], and its derivative can be easily and analytically derived for our regularization computation (Sec. 3.3).

**The rendering equation**     Given the BRDF parameterization, we now describe the concrete formulation of $f$ in Equation 1. In the following equations, we omit notations of surface point $\mathbf{x}$ and normal $\mathbf{n}$ as the geometry of the scene is assumed to be given. The diffuse term can be calculated as $f_d = \frac{1-m}{\pi} \cdot \mathbf{b}$, and the specular term as:

$$f_s(\boldsymbol{\omega}_o, \boldsymbol{\omega}_i) = \frac{D(\mathbf{h}; r) \cdot F(\boldsymbol{\omega}_o, \mathbf{h}; \mathbf{b}, m) \cdot G(\boldsymbol{\omega}_i, \boldsymbol{\omega}_o, \mathbf{h}; r)}{(\mathbf{n} \cdot \boldsymbol{\omega}_i) \cdot (\mathbf{n} \cdot \boldsymbol{\omega}_o)}, \tag{4}$$

where $\mathbf{h}$ is the half vector between the incident direction $\boldsymbol{\omega}_i$ and the viewing direction $\boldsymbol{\omega}_o$. The first term $D$ is the normal distribution function of the microfacets in the surface. It is related to the roughness $r$ and we use Spherical Gaussians to model this function as in previous methods [38, 49]. The second Fresnel term $F$ models the portion of light that can be reflected from the surface, which is determined by the surface metallic $m$ and the base color $\mathbf{b}$. The final geometry term $G$ handles the shadow and occlusion of the microfacets, which is parameterized on the roughness $r$ and is approximated using the GGX distribution [37]. Details of $D$, $F$ and $G$ are provided in the supplementary material.

## 3.3 Material-Lighting Ambiguity and Regularizations

While the Disney BRDF and incident light field are capable of representing materials and lightings of different scenes, jointly optimizing both would inevitably

lead to ambiguous solutions among them. One degenerate case could be that we can force a pure reflective BRDF to all surface points, and then only optimize the incident lights to adjust the input image. Theoretically, we can still find a perfect solution that fits the given BRDF and input images: for each 3D point, whenever there is a visible camera, we set its mirror symmetric incident light equal to the viewing out-going light, and set other incident lights equal to zero. It is also reported in previous work [32] that even human observers cannot distinguish the two confounded components from only image observations.

In the proposed framework, we can still manage to recover reasonable material and lighting results as MLPs can implicitly enforce a spatial smoothness constraint [50] on the two components. However, for robust material and lighting estimation, additional regularizations are desired. In this paper, we propose two regularizations for roughness $r$ and metallic $m$:

**Bilateral Smoothness**    We encourage $r$ and $m$ not to change rapidly in space, and the gradient of the input image $\mathbf{I}$ can be used as a hint to guide the smoothing process. Thus, we define the bilateral smoothness cost of $r$ and $m$ as:

$$l_{smooth} = \frac{1}{|S_I|} \sum_{\mathbf{p} \in S_I} (\|\nabla_{\mathbf{x}} r(\mathbf{x_p})\| + \|\nabla_{\mathbf{x}} m(\mathbf{x_p})\|) e^{-\|\nabla_{\mathbf{p}} \mathbf{I}(\mathbf{p})\|}, \tag{5}$$

where $S_I$ is the set of all sampled pixels and $\mathbf{x_p}$ is the corresponding 3D point of the sampled pixel $\mathbf{p}$. The image gradient $\nabla_{\mathbf{p}} \mathbf{I}(\mathbf{p})$ can be pre-calculated from the input image, and the roughness gradient $\nabla_{\mathbf{x}} r(\mathbf{x_p})$ and metallic gradient $\nabla_{\mathbf{x}} m(\mathbf{x_p})$ can be derived analytically by back-propagating the neural network.

**Lambertian Assumption**    We also assume that all surfaces tend to be Lambertian if no view-dependent lighting is observed, which leads to high roughness and low metallic, and we define the Lambertian cost as:

$$l_{lambertian} = \frac{1}{|S_I|} \sum_{\mathbf{p} \in S_I} (|r(\mathbf{x_p}) - 1| + |m(\mathbf{p})|). \tag{6}$$

The proposed two regularizations will be minimized during network training. It is noteworthy that the two losses may not necessarily improve quantitative results as they are heuristically defined for robust material and lighting estimation. We show in a later ablation study that the bilateral smoothness will lead to visually much more pleasing results for real-world reconstructions.

### 3.4   Loss

Similar to other differentiable rendering framework, we compute the L1 loss between the rendered image and the input image:

$$l_{image} = \frac{1}{|S_I|} \sum_{\mathbf{p} \in S_I} \|\mathbf{I}(\mathbf{p}) - L_o(\mathbf{x_p}, \boldsymbol{\omega}_o)\|_1. \tag{7}$$

The final loss of the proposed system is a weighted sum of the image loss and the two regularization losses: $l = l_{image} + w_s l_{smooth} + w_l l_{lambertian}$, where the two weights are empirically set to $w_s = 10^{-4}$ and $w_l = 10^{-3}$ in all our experiments.

## 4    Implementations

### 4.1    Sphere Sampling

To compute $L_o$ using a finite number of incident lights, we need to discretize Equation 1 as: $L_o(\boldsymbol{\omega_o}, \mathbf{x}) = \sum_{i \in S_L} f(\boldsymbol{\omega_o}, \boldsymbol{\omega_i}, \mathbf{x}) L_i(\boldsymbol{\omega_i}, \mathbf{x})(\boldsymbol{\omega_i} \cdot \mathbf{n}) \cdot A(\boldsymbol{\omega_i})$, where $S_L$ is the set of incident lights sampled for point $\mathbf{x}$ and $A(\boldsymbol{\omega_i})$ is the solid angle that corresponds to the incident light. In computer graphics, randomized Monto-Carlo Samplings are usually applied in ray-tracing, and the solid angle $A(\boldsymbol{\omega_i})$ is approximated by the probability distribution $P(\boldsymbol{\omega_i})$ of ray samples.

However, in differentiable rendering, it is critical to accurately compute the solid angle $A(\boldsymbol{\omega_i})$ for each light sample as we need to correctly pass loss gradients to network parameters. We found that using random sampling and approximating $A(\boldsymbol{\omega_i})$ as the probability distribution $P(\boldsymbol{\omega_i})$ will lead to erroneous BRDF results. Thus, we apply a fixed Fibonacci sampling over the half sphere to get all samples. In this case, $A(\boldsymbol{\omega_i}) = \frac{2\pi}{|S_L|}$ and the rendering equation becomes:

$$L_o(\boldsymbol{\omega_o}, \mathbf{x}) = \frac{2\pi}{|S_L|} \sum_{i \in S_L} f(\boldsymbol{\omega_o}, \boldsymbol{\omega_i}, \mathbf{x}) L_i(\boldsymbol{\omega_i}, \mathbf{x})(\boldsymbol{\omega_i} \cdot \mathbf{n}). \tag{8}$$

### 4.2    Learned HDR-LDR Mapping

For real-world datasets with low dynamic range (LDR) images, we need to convert the high dynamic range (HDR) output from our renderer to LDR before computing the image loss. As such transformation is unavailable in previous MVS datasets, we apply a learned HDR-LDR mapping to mimic the conversion in our network. Note that linear transformations, including exposure and white balance, can be embedded into the incident light. Thus, we only explicitly model the gamma correction with a learnable parameter:

$$L_o^{LDR} = (L_o^{HDR})^\gamma. \tag{9}$$

### 4.3    Training Details

We use an 8-layer Siren [35] with feature size of 512 and a skip connection in the middle to represent the BRDF MLP. Also, the positional encoding [24] is applied to further strengthen the network. The NeILF MLP shares the same implementation as BRDF, except that 1) the feature size is downsized to 128 to reduce the VRAM usage and 2) the last layer activation function is changed from *tanh* to *exp* in order to guarantee non-negative and unbounded light intensities.

In the experiment, we use $|S_L| = 128$ incident lights to compute the output radiance during training, and use $|S_L| = 256$ incident lights to evaluate the rendered image during testing. For each training iteration, we randomly sample 16000 pixels from all images, and the network is optimized for a total of 15000 iterations. The Adam optimizer [19] with an initial learning rate of $10^{-3}$ is applied in our network, and the learning rate is scaled down by $\sqrt{0.1}$ at 5000 and 10000 iterations. The training process takes around 1.5 hours to finish on a Tesla V100 GPU and the VRAM consumption is around 30 GB.

## 5   Experiments

### 5.1   Baseline Methods

We compare our method with the following baselines:

**PhySG**$^*$    Firstly, we consider the recent PhySG [49] for material estimation. The original PhySG jointly optimizes the non-spatially varying BRDF, the environment map, and the geometry of the object. To fairly compare with the method, we fix the given geometry and optimize only the uniform BRDF and the environment map.

**SG-Env**    This baseline is another variant of PhySG [49]. Compared with PhySG$^*$, SG-Env applies a SV-BRDF model and a slightly different rendering formulation (the same $f_d$, $D$ and $G$ as ours). We use this baseline to directly compare NeILF with the SG environment map representation.

**Pix-Env**    Compared with SG-Env, this baseline uses a 2D image of resolution $32 \times 16$ to represent the environment map, which can be viewed as a variant of NeRFactor [51] but without visibility handling. We sample all pixels in the environment image to compute the rendering equation, where pixel values are optimized directly during training.

**Ne-Env**    Lastly, we compare our method with the neural environment map representation. This baseline shares the same implementation of the proposed NeILF, except that the positional input in the incident light field is omitted such that the incident light is only related to the incoming direction: $\{\boldsymbol{\omega}\} \to \mathbf{L}$.

### 5.2   Benchmark on Synthetic Scenes

To quantitatively evaluate our method under different lighting conditions, we generate a set of synthetic data and compare our method with the above baselines.

**Data Preparation**    The synthetic dataset contains three objects and their combinations: a single rough metallic sphere, a single rough metallic cube, and a helmet with spatially variant materials. The objects are placed on a plane to model the real-world object capture. We also create six lighting conditions to lit the objects, including three environment maps and three mixed lightings:

- *Env–city*: an environment map of a city;
- *Env–studio*: an environment map of a studio;
- *Env–castel*: an environment map of a castel;
- *Mix–city*: Env–city plus two point lights and two area lights;
- *Mix–studio*: Env–studio plus two point lights and two area lights;
- *Mix–castel*: Env–castel plus two point lights and two area lights.

Each scene contains 96 images, where 87 images are used for training and the left out 9 images are used for evaluation. The image trajectory forms three loops
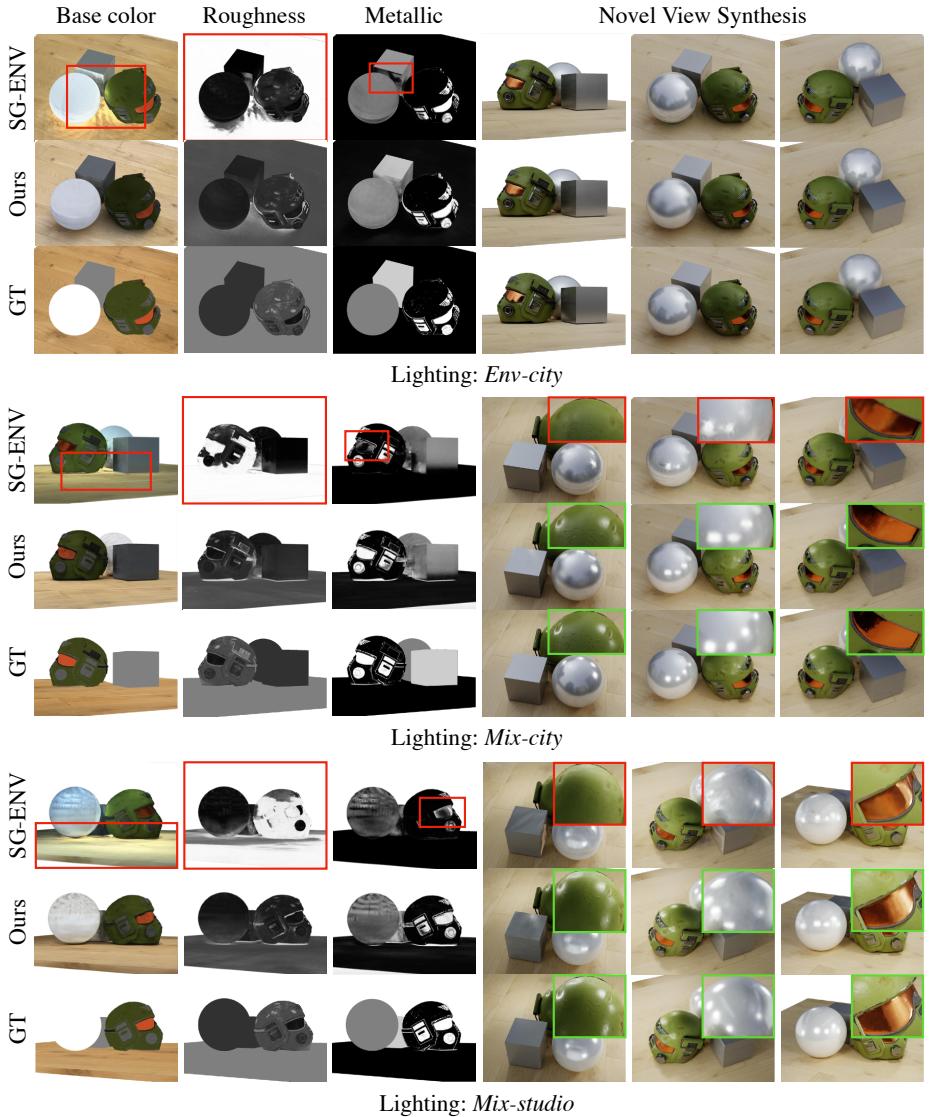
Fig. 2: Comparative results on BRDF estimation and novel view synthesis on the synthetic dataset. From left to right are images of *base color*, *roughness*, *metallic* and synthesized testing views. Our method is able to generate high-quality BRDF and novel view synthesis results under different lighting conditions. In contrast, the environment map based SG-ENV [49] produces noisy BRDF outputs especially in occluded regions. And also, high lights are wrongly recovered in novel view renderings if mixed lightings occur.

Table 1: Quantitative results on Synthetic scenes. We compare the proposed NeILF with four baseline methods described in Sec. 5.1 using PSNR scores. Our method generates consistently the best novel view rendering for all scenes. Also, our method produces significantly better BRDF results than other methods if multiple objects and mixed lightings are given.

| Scene Geometry | | Single-helmet | | | Combined-objects | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| Scene Lighting | | Env-city | Env-studio | Env-castel | Env-city | Env-studio | Env-castel | Mix-city | Mix-studio | Mix-castel |
| Base Color | PhySG*[49] | 13.43 | 14.87 | 13.95 | 15.01 | 16.96 | 16.13 | 14.16 | 12.43 | 14.29 |
| | SG-ENV[49] | **20.61** | **18.45** | 15.99 | **22.38** | **20.74** | **22.21** | 16.92 | 13.16 | 16.69 |
| | Pix-ENV[51] | 14.20 | 13.37 | 12.18 | 13.18 | 15.09 | 7.94 | 12.16 | 11.44 | 11.79 |
| | Ne-ENV | 13.43 | 12.65 | 12.45 | 11.68 | 11.98 | 7.66 | 11.87 | 10.90 | 9.54 |
| | Ours | 16.36 | 16.36 | **18.28** | 15.59 | 15.48 | 12.95 | **17.39** | **16.88** | **17.37** |
| Metallic | PhySG*[49] | 7.57 | 7.48 | 7.83 | 8.72 | 7.97 | 8.35 | 8.67 | 8.95 | 8.76 |
| | SG-ENV[49] | **21.19** | **21.31** | **21.79** | 17.01 | 16.40 | **16.39** | 15.44 | 14.25 | 14.49 |
| | Pix-ENV[51] | 14.28 | 18.03 | 18.44 | 16.15 | 15.61 | 11.42 | 15.27 | 14.24 | 14.84 |
| | Ne-ENV | 9.31 | 17.40 | 6.15 | 15.86 | 16.10 | 11.42 | 15.43 | 15.35 | 15.49 |
| | Ours | 17.79 | 18.52 | 16.82 | **18.22** | **19.11** | 10.29 | **18.42** | **18.43** | **17.34** |
| Roughness | PhySG*[49] | 6.91 | 11.88 | 6.75 | 6.62 | 11.29 | 6.22 | 6.27 | 6.83 | 6.14 |
| | SG-ENV[49] | 14.77 | 15.77 | 9.64 | 9.61 | 17.64 | 9.74 | 8.77 | 12.58 | 9.14 |
| | Pix-ENV[51] | 13.88 | 14.36 | 15.92 | 15.55 | 13.95 | 13.06 | 16.11 | 14.12 | 13.41 |
| | Ne-ENV | 11.95 | 14.84 | 9.26 | 15.56 | 14.48 | 12.94 | 16.20 | 14.43 | 14.14 |
| | Ours | **16.13** | **16.19** | **17.16** | **17.48** | **18.30** | **13.40** | **17.05** | **16.27** | **16.44** |
| Rendering | PhySG*[49] | 24.59 | 24.77 | 26.52 | 24.82 | 25.65 | 27.24 | 24.38 | 24.04 | 25.81 |
| | SG-ENV[49] | 29.73 | 29.86 | 32.13 | 31.01 | 29.46 | 32.34 | 27.20 | 25.88 | 27.70 |
| | Pix-ENV[51] | 29.28 | 29.03 | 31.09 | 30.81 | 29.06 | 32.30 | 27.88 | 25.85 | 27.97 |
| | Ne-ENV | 28.60 | 29.56 | 29.76 | 30.75 | 29.07 | 32.05 | 28.07 | 26.01 | 28.33 |
| | Ours | **31.57** | **30.84** | **34.43** | **33.77** | **31.07** | **35.28** | **31.11** | **28.59** | **32.11** |

around the object at altitude 0, 22.5 and 45 and the image resolution is set to $1600 \times 1200$. We use Blender [9] to render the HDR images by ray tracing. Position maps and normal maps at all viewpoints are rendered to serve as the geometry input for the system. Meanwhile, per-view ground truth base color, metallic and roughness maps are provided for quantitative evaluation.

**Results**    We use the PSNR score as our evaluation metric. Quantitative comparisons on 1) base color, 2) metallic, 3) roughness and 4) novel view synthesis are shown in Table 1. Our method consistently outperforms other methods with a large margin in terms of the novel view rendering quality. For material estimation, we found that if single objects and environment map light sources are given, SG-Env [49] is able to generate comparable results with ours. However, if multiple objects or mixed light sources are given, its quality will drop significantly. This is because the environment map representation cannot model mixed light sources of point and area lights. Also, indirect lights and occlusions within multiple objects are not handled by SG-Env. In contrast, our NeILF representation can robustly deal with mixed lightings and complex scene geometries. Qualitative results are shown in Fig. 2.

## 5.3    Test on Real-world Scenes

We then test our method on two real-world datasets, namely DTU [14] and BlendedMVS [43] datasets. DTU dataset is captured in a lab setting with a fixed

Table 2: Quantitative results on DTU [14] and BlendedMVS [43] Datasets. The table shows PSNR scores of novel view renderings of test images. Our method consistently outperforms the other methods in terms of the rendering quality.

| | DTU [14] | | | | | BlendedMVS [43] | | | | | |
| | scan-1 | scan-11 | scan-37 | scan-75 | scan-97 | bull | cam | dog | gold | statue | stone |
|---|---|---|---|---|---|---|---|---|---|---|---|
| PhySG* [49] | 20.40 | 20.78 | 20.30 | 16.03 | 19.86 | 21.64 | 18.11 | 20.70 | 19.06 | 19.74 | 21.22 |
| SG-ENV [49] | 22.18 | 21.56 | 21.71 | 18.06 | 21.09 | 22.51 | 20.14 | 22.06 | 19.44 | 20.79 | 22.31 |
| Pix-ENV [51] | 23.61 | 23.61 | 22.68 | 19.54 | 21.52 | 21.58 | 19.98 | 21.36 | 19.28 | 20.46 | 22.89 |
| Ne-ENV | 23.77 | 23.79 | 22.87 | 19.52 | 21.51 | 22.17 | 20.17 | 21.73 | 19.66 | 20.55 | 23.08 |
| NeILF (Ours) | **24.79** | **24.33** | **24.44** | **23.46** | **23.96** | **24.93** | **22.10** | **22.36** | **20.80** | **21.51** | **24.22** |



Fig. 3: Qualitative results on DTU [14] and BlendedMVS [43] datasets. Our method successfully removes high lights in the base color and produces visually plausible results of roughness and metallic. Note that brand names in BlendedMVS-camera have been masked out from images.

lighting and camera trajectory, while BlendedMVS contains a variety of indoor and outdoor scenes captured by different users. As the two datasets provide only LDR images, the learned HDR-LDR mapping described in Sec. 4.2 is applied for the loss computation. For each scene, we select 5 images for testing and use remaining images for network training.

It is noteworthy that unlike in the synthetic dataset, here we use multi-view stereo methods to generate the geometry input rather than directly using the ground truth. For DTU datatset, we use Vis-MVSNet [47] to generate the dense 3D point cloud and SPSR [18] to recover the mesh surface. For BlendedMVS dataset, we use original images and the provided reference mesh geometry as

| Input image | Base color | Roughness | Metallic | Base color | Roughness | Metallic |

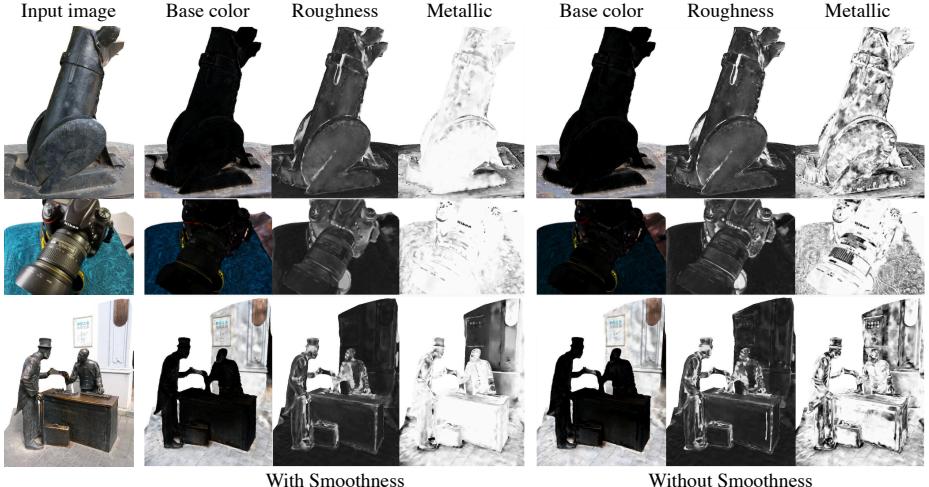With Smoothness                         Without Smoothness

Fig. 4: Qualitative comparisons on with and without the proposed bilateral smoothness regularization.

our inputs. By doing so, our method can be viewed as an extension to nowadays 3D reconstruction pipelines.

Quantitative results are shown in Table 2 and qualitative results compared with SG-ENV are shown in Fig. 3. The proposed method produces both the best rendering PSNR and the most visually pleasing BRDF in all selected scenes. We believe the proposed method can be integrated into traditional 3D reconstruction pipelines [11, 13, 17, 34, 42] for relightable mesh model reconstruction.

## 5.4   Ablation Study

In this section, we analyze several design choices of the proposed framework. The ablation studies are conducted on the synthetic dataset, and we report the average scores over all scenes to compare different settings.

**Ray Sample Number**    We first study the influence of the ray sample number for material estimation quality. The ray sample number is decreased from $S_L = 128$ to $S_L = 64$ and $S_L = 32$. As shown in Table 3, higher sampling number will lead to better reconstruction results. In our default setting, we choose $S_L = 128$ to better balance the quality and the VRAM/runtime consumption.

**Random Sample**    Next, we compare the fixed Fibonacci sample described in Sec. 4.1 with the random uniform sample commonly used in computer graphics. It is shown in Table 3 that the random sample would lead to worse results, showing that it is crucial to precisely discretize the rendering equation in the differentiable rendering.

**Regularizations**    Lastly, we study the effectiveness of the two regularizations proposed in Sec. 3.3. We find that the bilateral smoothness is essential for material estimation of real-world scenes, where the roughness and metallic will be significantly improved if the smoothness is applied (Fig. 4).

On the other hand, we also find that the two heuristics have limited influence to quantitative results of the synthetic dataset. We believe this is because the vanilla NeILF already produces high-quality estimations for synthetic scenes. In our default setting, we keep the two regularizations for all scenes but we encourage users to selectively apply the two terms depending on different characteristics of input scenes.

| | Base. | Meta. | Roug. | Rend. |
|---|---|---|---|---|
| S = 128 | **16.30** | **17.22** | 16.49 | **32.09** |
| S = 64 | 15.20 | 16.47 | 18.88 | 31.40 |
| S = 32 | 13.40 | 16.27 | **19.49** | 30.57 |
| Rand. Samp. | 12.45 | 15.11 | 18.10 | 29.73 |

Table 3: Ablation studies. Average scores among all synthetic scenes are reported.

## 6    Discussions

### 6.1    Comparison with NeRF Optimization

In this section, we compare the proposed NeILF framework with the neural radiance field [24] optimization. We show that the two frameworks share similarities in many aspects, and thus provide readers an intuitive explanation why the proposed NeILF can successfully disentangle the complex material and lighting in the joint optimization.

**Lighting Representations**    NeRF [24] represents the scene appearance as the neural radiance field. While the radiance field is physically different with the incident light field, its complexity is completely the same as ours: both NeILF and NeRF take a 3D position $\mathbf{x}$ along with a direction $\boldsymbol{\omega}$ as inputs, and returns a RGB value as output.

**Spatially-varying Properties**    NeILF aims to recover surface materials as BRDF properties, while NeRF jointly optimizes the scene geometry as a density field. Both our BRDF and NeRF's density MLPs take only a 3D position $\mathbf{x}$ as input, and return different spatial properties as outputs. Implementation-wise, the only difference is that our BRDF is consist of a 5D parameter vector, while the density value is a 1D scalar. Nevertheless, the two spatially-varying properties are very similar and their complexities are comparable.

**Rendering Formulations**    Our method applies the physically-based rendering to compute the reflected light from a surface point, while NeRF adopts the volume rendering to get the accumulated color along a viewing ray. On the one hand, NeILF requires incident light integration over the hemisphere; on the other hand, NeRF requires alpha composition along the viewing ray. Implementation-wise, to render a pixel, NeILF needs to sample the BRDF MLP once and the

incident light MLP for multiple times, while NeRF does the same operations on the density and radiance MLPs.

**Reconstruction Ambiguities**    The geometry-appearance ambiguity is addressed in NeRF++ [50]. Similarly, we analyze the material-lighting ambiguity in our NeILF optimization. It has also been reported that with proper regularizations on density [30] or converted SDF [39, 44], the NeRF framework is able to produce high-quality geometry reconstructions. In contrast, we also show that the proposed bilateral smoothness can significantly improve the roughness and metallic quality for real-world scenes (Fig. 4).

## 6.2   Limitations and Future Works

While the proposed method has already shown promising results for material estimation, the current pipeline still contains several limitations that could be further addressed in future works.

**Running Speed**    Similar to NeRF, our method requires multiple samples of the light MLP to render one pixel, which makes the training process time consuming. Our current implementation takes around 1.5 hours to estimate the BRDF of a given scene (details in Sec. 4.3). We hope that in the future explicit Octree [46], spherical harmonics [1, 40] or neural hashing [26] could be applied to speed up the NeILF optimization.

**Geometry Optimization**    Our method assumes that the geometry of the scene should be given in advance. Although we have shown that multi-view reconstructed meshes are qualified enough for real-world DTU and BlendedMVS scenes, it would be better if we could jointly optimize the geometry during training. Possible directions include displacement/normal map estimation and recent differentiable surface refinements [45, 48, 49].

## 7   Conclusions

We have presented a differentiable rendering framework for material and lighting estimation. Compared with the environment map approximation, the proposed neural incident light field is capable of modelling the lighting condition of any static scenes, making it possible to estimate qualified material properties even for scenes with complex lightings and geometries. The proposed method strictly follows the physically-based rendering equation, and jointly optimizes material and lighting through the differentiable rendering process. We have intensively evaluated our method on our in-house synthetic dataset, the DTU MVS dataset, and the real-world BlendedMVS scenes. Our method is able to outperform previous environment map based methods by a significant margin in terms of the novel view rendering quality, setting a new state-of-the-art for image-based material and lighting estimation.

# Bibliography

[1] Alex Yu and Sara Fridovich-Keil, Tancik, M., Chen, Q., Recht, B., Kanazawa, A.: Plenoxels: Radiance fields without neural networks. arXiv:2112.05131 (2021)

[2] Azinovic, D., Li, T.M., Kaplanyan, A., Niessner, M.: Inverse path tracing for joint material and lighting estimation. In: CVPR (June 2019)

[3] Azinovic, D., Li, T.M., Kaplanyan, A., Nießner, M.: Inverse path tracing for joint material and lighting estimation. In: CVPR (2019)

[4] Bi, S., Xu, Z., Sunkavalli, K., Hašan, M., Hold-Geoffroy, Y., Kriegman, D., Ramamoorthi, R.: Deep reflectance volumes: Relightable reconstructions from multi-view photometric images. In: ECCV (2020)

[5] Bi, S., Xu, Z., Sunkavalli, K., Kriegman, D., Ramamoorthi, R.: Deep 3d capture: Geometry and reflectance from sparse multi-view images. In: CVPR (2020)

[6] Boss, M., Braun, R., Jampani, V., Barron, J.T., Liu, C., Lensch, H.: Nerd: Neural reflectance decomposition from image collections. In: ICCV (2021)

[7] Boss, M., Jampani, V., Braun, R., Liu, C., Barron, J., Lensch, H.P.: Neural-pil: Neural pre-integrated lighting for reflectance decomposition. In: NeurIPS 2021 (2021)

[8] Burley, B., Studios, W.D.A.: Physically-based shading at disney. In: ACM SIGGRAPH (2012)

[9] Community, B.O.: Blender - a 3D modelling and rendering package. Blender Foundation, Stichting Blender Foundation, Amsterdam (2018), `http://www.blender.org`

[10] Dong, Y., Chen, G., Peers, P., Zhang, J., Tong, X.: Appearance-from-motion: Recovering spatially varying surface reflectance under unknown lighting. TOG (2014)

[11] Furukawa, Y., Ponce, J.: Accurate, dense, and robust multiview stereopsis. PAMI (2009)

[12] Guo, K., Lincoln, P., Davidson, P., Busch, J., Yu, X., Whalen, M., Harvey, G., Orts-Escolano, S., Pandey, R., Dourgarian, J., et al.: The relightables: Volumetric performance capture of humans with realistic relighting. TOG (2019)

[13] Hiep, V.H., Keriven, R., Labatut, P., Pons, J.P.: Towards high-resolution large-scale multi-view stereo. In: CVPR (2009)

[14] Jensen, R., Dahl, A., Vogiatzis, G., Tola, E., Aanæs, H.: Large scale multi-view stereopsis evaluation. In: CVPR (2014)

[15] Kajiya, J.T.: The rendering equation. In: Proceedings of the 13th annual conference on Computer graphics and interactive techniques (1986)

[16] Kato, H., Beker, D., Morariu, M., Ando, T., Matsuoka, T., Kehl, W., Gaidon, A.: Differentiable rendering: A survey. arXiv preprint arXiv:2006.12057 (2020)

[17] Kazhdan, M., Bolitho, M., Hoppe, H.: Poisson surface reconstruction. In: Proceedings of the fourth Eurographics symposium on Geometry processing (2006)

[18] Kazhdan, M., Hoppe, H.: Screened poisson surface reconstruction. TOG (2013)

[19] Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. In: ICLR (2015)

[20] Li, Z., Xu, Z., Ramamoorthi, R., Sunkavalli, K., Chandraker, M.: Learning to reconstruct shape and spatially-varying reflectance from a single image. TOG (2018)

[21] Liu, L., Gu, J., Lin, K.Z., Chua, T.S., Theobalt, C.: Neural sparse voxel fields. NeurIPS (2020)

[22] Liu, S., Zhang, Y., Peng, S., Shi, B., Pollefeys, M., Cui, Z.: Dist: Rendering deep implicit signed distance function with differentiable sphere tracing. In: CVPR (2020)

[23] Liu, S., Saito, S., Chen, W., Li, H.: Learning to infer implicit surfaces without 3d supervision. arXiv preprint arXiv:1911.00767 (2019)

[24] Martin-Brualla, R., Radwan, N., Sajjadi, M.S., Barron, J.T., Dosovitskiy, A., Duckworth, D.: Nerf in the wild: Neural radiance fields for unconstrained photo collections. In: CVPR (2021)

[25] Mildenhall, B., Srinivasan, P.P., Tancik, M., Barron, J.T., Ramamoorthi, R., Ng, R.: Nerf: Representing scenes as neural radiance fields for view synthesis. In: ECCV (2020)

[26] Müller, T., Evans, A., Schied, C., Keller, A.: Instant neural graphics primitives with a multiresolution hash encoding. arXiv:2201.05989 (2022)

[27] Munkberg, J., Hasselgren, J., Shen, T., Gao, J., Chen, W., Evans, A., Mueller, T., Fidler, S.: Extracting Triangular 3D Models, Materials, and Lighting From Images. arXiv:2111.12503 (2021)

[28] Nam, G., Lee, J.H., Gutierrez, D., Kim, M.H.: Practical svbrdf acquisition of 3d objects with unstructured flash photography. TOG (2018)

[29] Niemeyer, M., Mescheder, L., Oechsle, M., Geiger, A.: Differentiable volumetric rendering: Learning implicit 3d representations without 3d supervision. In: CVPR (2020)

[30] Oechsle, M., Peng, S., Geiger, A.: Unisurf: Unifying neural implicit surfaces and radiance fields for multi-view reconstruction. In: ICCV (2021)

[31] Park, J.J., Holynski, A., Seitz, S.M.: Seeing the world in a bag of chips. In: CVPR (2020)

[32] Pont, S.C., Te Pas, S.F.: Material—illumination ambiguities and the perception of solid objects. Perception (2006)

[33] Schmitt, C., Donne, S., Riegler, G., Koltun, V., Geiger, A.: On joint estimation of pose, geometry and svbrdf from a handheld scanner. In: CVPR (2020)

[34] Schönberger, J.L., Zheng, E., Frahm, J.M., Pollefeys, M.: Pixelwise view selection for unstructured multi-view stereo. In: ECCV (2016)

[35] Sitzmann, V., Martel, J., Bergman, A., Lindell, D., Wetzstein, G.: Implicit neural representations with periodic activation functions. NeurIPS (2020)

[36] Srinivasan, P.P., Deng, B., Zhang, X., Tancik, M., Mildenhall, B., Barron, J.T.: Nerv: Neural reflectance and visibility fields for relighting and view synthesis. In: CVPR (2021)

[37] Walter, B., Marschner, S., Li, H., Torrance, K.: Microfacet models for refraction through rough surfaces. In: EGSR (2007)

[38] Wang, J., Ren, P., Gong, M., Snyder, J., Guo, B.: All-frequency rendering of dynamic, spatially-varying reflectance. In: ACM SIGGRAPH Asia (2009)

[39] Wang, P., Liu, L., Liu, Y., Theobalt, C., Komura, T., Wang, W.: Neus: Learning neural implicit surfaces by volume rendering for multi-view reconstruction (2021)

[40] Wizadwongsa, S., Phongthawee, P., Yenphraphai, J., Suwajanakorn, S.: Nex: Real-time view synthesis with neural basis expansion. In: CVPR (2021)

[41] Xia, R., Dong, Y., Peers, P., Tong, X.: Recovering shape and spatially-varying surface reflectance under unknown illumination. TOG (2016)

[42] Yao, Y., Luo, Z., Li, S., Fang, T., Quan, L.: Mvsnet: Depth inference for unstructured multi-view stereo. In: ECCV (2018)

[43] Yao, Y., Luo, Z., Li, S., Zhang, J., Ren, Y., Zhou, L., Fang, T., Quan, L.: Blendedmvs: A large-scale dataset for generalized multi-view stereo networks. In: CVPR (2020)

[44] Yariv, L., Gu, J., Kasten, Y., Lipman, Y.: Volume rendering of neural implicit surfaces. arXiv preprint arXiv:2106.12052 (2021)

[45] Yariv, L., Kasten, Y., Moran, D., Galun, M., Atzmon, M., Ronen, B., Lipman, Y.: Multiview neural surface reconstruction by disentangling geometry and appearance. NeurIPS (2020)

[46] Yu, A., Li, R., Tancik, M., Li, H., Ng, R., Kanazawa, A.: PlenOctrees for real-time rendering of neural radiance fields. In: ICCV (2021)

[47] Zhang, J., Yao, Y., Li, S., Luo, Z., Fang, T.: Visibility-aware multi-view stereo network. BMVC (2020)

[48] Zhang, J., Yao, Y., Quan, L.: Learning signed distance field for multi-view surface reconstruction. In: ICCV (2021)

[49] Zhang, K., Luan, F., Wang, Q., Bala, K., Snavely, N.: Physg: Inverse rendering with spherical gaussians for physics-based material editing and relighting. In: CVPR (2021)

[50] Zhang, K., Riegler, G., Snavely, N., Koltun, V.: Nerf++: Analyzing and improving neural radiance fields. arXiv preprint arXiv:2010.07492 (2020)

[51] Zhang, X., Srinivasan, P.P., Deng, B., Debevec, P., Freeman, W.T., Barron, J.T.: NeRFactor: Neural Factorization of Shape and Reflectance Under an Unknown Illumination. TOG (2021)

[52] Zhou, K., Synder, J., Guo, B., Shum, H.Y.: Iso-charts: stretch-driven mesh parameterization using spectral analysis. In: Proceedings of the 2004 Eurographics/ACM SIGGRAPH symposium on Geometry processing (2004)

# Supplementary Material for NeILF

In this supplementary material, we show implementation details and additional results of the proposed method. For relighting results, please refer to our **supplementary video**.

## 1   Implementation Details

### 1.1   Network Architecture

The detailed architecture of the proposed method is shown in Fig. 1. The overall design of our network is simple yet effective. We believe the architecture can be easily re-implemented or extended by other researchers.
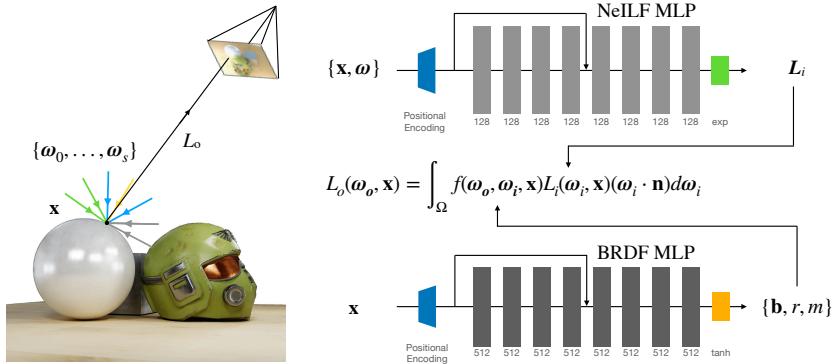


Fig. 1: Network architecture of the proposed NeILF.

### 1.2   Detailed BRDF Functions

In this section, we describe detailed implementations of our normal distribution term D, Fresnel term F, and geometry term G. The normal distribution function D is approximated by the Spherical Gaussian function:

$$D(\mathbf{h}; r) = S(\mathbf{h}, \frac{1}{\pi r^2}, \mathbf{n}, \frac{2}{r^2}) = \frac{1}{\pi r^2} e^{\frac{2}{r^2}(\mathbf{h}\cdot\mathbf{n}-1)}.$$

The Fresnel term is given as:

$$F(\boldsymbol{\omega}_o, \mathbf{h}; \mathbf{b}, m) = F_0 + (1 - F_0)(1 - (\boldsymbol{\omega}_o \cdot \mathbf{h})^5),$$
$$\text{where } F_0 = 0.04(1 - m) + \mathbf{b}m.$$

Finally, the geometry term is approximated by the GGX function [37]:

$$G(\boldsymbol{\omega}_i, \boldsymbol{\omega}_o, \mathbf{n}; r) = G_{GGX}(\boldsymbol{\omega}_i \cdot \mathbf{n})G_{GGX}(\boldsymbol{\omega}_o \cdot \mathbf{n}),$$
$$\text{where } G_{GGX}(z) = \frac{2z}{z + \sqrt{r^2 + (1 - r^2)z^2}}.$$

## 2    Synthetic Dataset

### 2.1    Dataset Setup

In this section, we illustrate the camera and lighting setup of our in-house synthetic dataset. The camera trajectory is shown in Fig. 2, which contains 96 camera positions with the outside-look-in trajectory. The mixed lighting set up is also shown in Fig. 2, which consists of two near-range point lights, two near-range area lights and one background environment map from {*Env-city, Env-studio, Env-castel*} (see main paper Sec. 5.2 for details).



Fig. 2: The camera and the near-range light setting of our synthetic dataset. We add two point lights and two area lights in the mixed lighting to lit the scene.

Fig. 3: Two selected points for incident light visualization. Estimated incident lights for $\mathbf{x}_1$ and $\mathbf{x}_2$ are visualized in Fig. 4

### 2.2    Lighting Estimation

Due to the space limit, we did not analyze the lighting estimation result in the main paper. In this section, we demonstrate the powerful capability of the proposed NeILF for lighting modelling. As shown in Fig. 3, we select two surface points on top of the helmet and the cube to show the estimated incident lights at these two points. The lighting estimation results under six lighting conditions are visualized in Fig. 4. Our lighting estimations enjoy the following properties:

- For environment map based lightings (*Env-city*, *Env-studio* and *Env-studio*), our incident light estimations successfully recovery corresponding environment maps at both points.
- For mixed lightings (*Mix-city*, *Mix-studio* and *Mix-studio*), our results well explain all light sources, including background environment maps, two near-range point lights, and two near-range area lights.
- For mixed lightings, we generate consistent point and area light estimations (high lights in images) across different background environments.
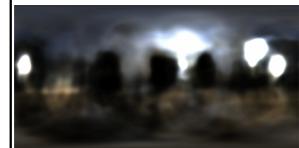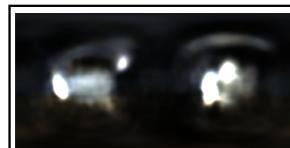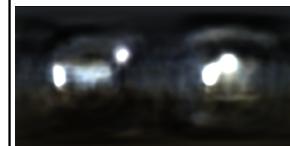
Fig. 4: Visualizations of our lighting estimations. For environment map based lightings (*Env-city*, *Env-studio* and *Env-studio*), our incident light estimations correctly recovery GT environment maps at both points (see Fig. 3). For mixed lightings with point and area lights (*Mix-city*, *Mix-studio* and *Mix-studio*), our results well explain all light sources. Note that the three mixed lightings share the same near-range light settings as shown in Fig. 2, and we are able to generate consistent point and area light estimations (high lights in images) across different background environments.

# 3     DTU and BlendedMVS Datasets

We show BRDF estimation of all DTU [14] scenes in Fig. 5 and all Blended-MVS [43] scenes in Fig. 6.
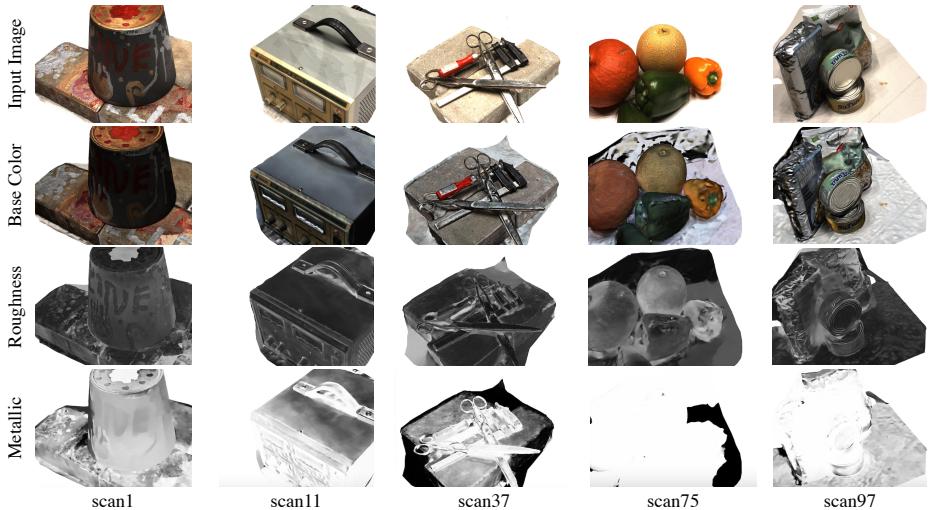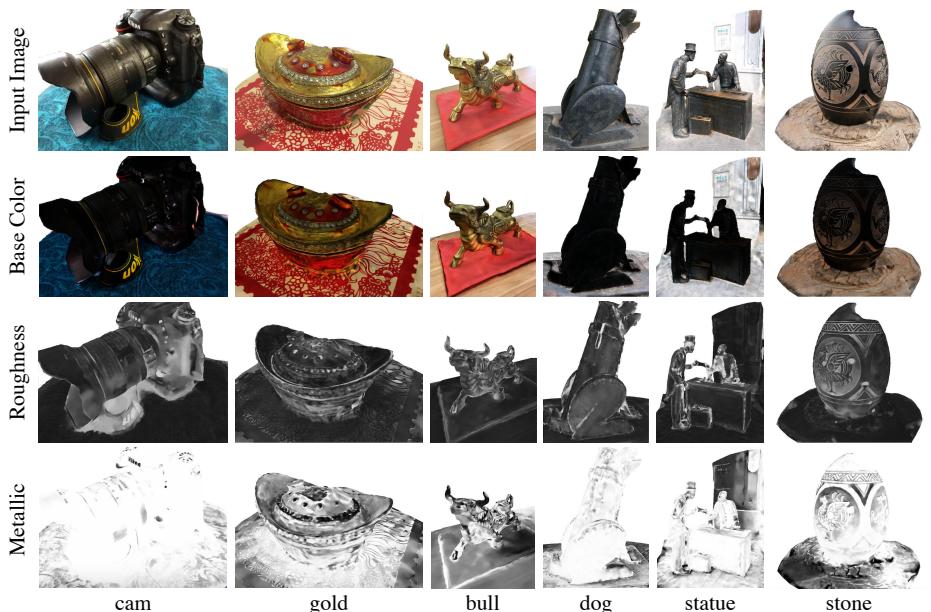


Fig. 5: Results of DTU [14] scenes.



Fig. 6: Results of BlendedMVS [43] scenes.

# 4    Relighting

In this section, we show the relighting result of our BRDF estimation. Given a mesh geometry input, we apply the the Iso-charts[52] algorithm to paramiterize the mesh surface as a 2D UV map. Then, for each pixel in the UV map, we pass the corresponding surface point to the BRDF MLP to obtain its BRDF values. The resulting BRDF texture maps (Fig. 7) can be directly used in rendering pipelines for relighting. Please refer to our **supplementary video** for convincing results.
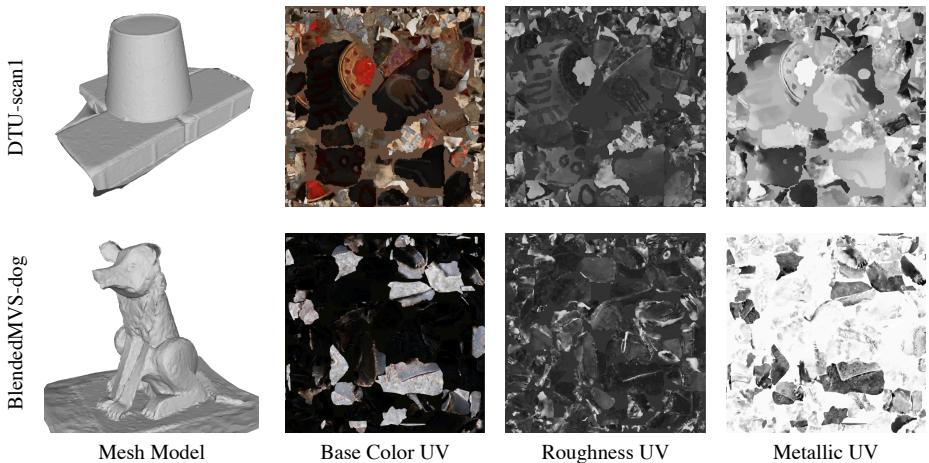


Fig. 7: Exported BRDF UV maps.