# GaNI: Global and Near Field Illumination Aware Neural Inverse Rendering

Jiaye Wu[1], Saeed Hadadan[1], Geng Lin[1], Matthias Zwicker[1], David Jacobs[1], and Roni Sengupta[2]

[1] University of Maryland, College Park, MD 20740, USA
{jiayewu, saeedhd, geng, zwicker, dwj}@umd.edu
[2] University of North Carolina at Chapel Hill, NC, 27599
ronisen@cs.unc.edu

**Abstract.** In this paper, we present GaNI, a Global and Near-field Illumination-aware neural inverse rendering technique that can reconstruct geometry, albedo, and roughness parameters from images of a scene captured with co-located light and camera. Existing inverse rendering techniques with co-located light-camera focus on single objects only, without modeling global illumination and near-field lighting more prominent in scenes with multiple objects. We introduce a system that solves this problem in two stages; we first reconstruct the geometry powered by neural volumetric rendering NeuS, followed by inverse neural radiosity NeRad that uses the previously predicted geometry to estimate albedo and roughness. However, such a naive combination fails and we propose multiple technical contributions that enable this two-stage approach. We observe that NeuS fails to handle near-field illumination and strong specular reflections from the flashlight in a scene. We propose to implicitly model the effects of near-field illumination and introduce a surface angle loss function to handle specular reflections. Similarly, we observe that NeRad assumes constant illumination throughout the capture and cannot handle moving flashlights during capture. We propose a light position-aware radiance cache network and additional smoothness priors on roughness to reconstruct reflectance. Experimental evaluation on synthetic and real data shows that our method outperforms the existing co-located light-camera-based inverse rendering techniques. Our approach produces significantly better reflectance and slightly better geometry than capture strategies that do not require a dark room.

## 1 Introduction

Decomposing an indoor scene into geometry, material properties, and lighting, called inverse rendering [1,2,4,15,37], is a long-standing problem in computer vision. Inverse rendering has many applications in VR/AR, computational photography, and robotics perception, e.g. relighting, material editing, etc. [4,12,37,44]. However, existing methods are often tailored for individual objects [4,14,37,38] or human faces [11,26,29]. For scenes consisting of multiple objects, recent methods have shown impressive quality in reconstructing geometry [5,31], but fail to
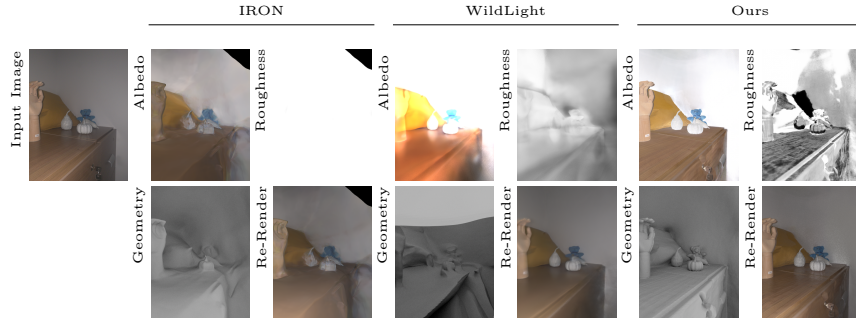
**Fig. 1:** We perform inverse rendering of a scene from multiple images captured with co-located light and camera. Our method, GaNI, produces better geometry, albedo, roughness and re-rendering in unseen views than state-of-the-art approaches, IRON [37] and WildLight [4], that also uses co-located light-camera.

reconstruct material reflectance i.e the Bi-directional Reflectance Distribution Function (BRDF), [12, 28, 35, 43, 45], since it is is highly under-constrained.

To make the inverse rendering problem well-constrained, researchers proposed capturing images using a co-located flashlight and camera [1, 4, 14, 18, 24, 37]. Such a setup is readily accessible since one can easily capture images by turning on the flashlight of one's mobile camera; light and camera location can be easily calibrated by using Structure-from-motion which imposes additional constraints on the inverse rendering problem. Recent approaches, like IRON [37] and WildLight [4], have shown impressive performance in material prediction of a single object captured using a co-located light and camera in a dark room and under ambient lighting respectively. However, these methods fail to generalize to scenes consisting of multiple objects, as noted in our experimental evaluations 1. We believe this is because the complexity of scenes with multiple objects surpasses that of individual objects, particularly with a co-located light-camera setup, which induces a pronounced near-field effect, where different regions in the scene receive different intensities of light. Moreover, interactions between objects and surfaces result in strong inter-reflections or global illumination.

In this paper, we aim to perform inverse rendering of a scene consisting of multiple objects captured using a co-located light and camera. We present GaNI, a system that can reconstruct both high-quality geometry and reflectance, by implicitly modeling both global and near-field illumination, hence overcoming the shortcomings of existing approaches [4,37]. We focus on examples that consist of multiple objects, e.g. coffee table, shoe rack, table, etc. We believe this is a first step towards developing a system that can perform accurate inverse rendering of a 360-degree room-scale scene, which is significantly more challenging.

Our proposed method develops a neural scene representation by optimizing the parameters for each scene independently. We propose a two-stage approach where we first solve for geometry followed by reflectance, instead of jointly optimizing both which is challenging in the presence of complex illumination effects. We observe that a naive combination of first reconstructing geometry with

neural volume rendering, NeuS [31], followed by reflectance estimation with InvNeRad [8], an efficient technique to model multi-bounce global illumination with known geometry, fails due to specularity, global illumination, and changing light source. We propose several technical contributions that enable geometry and reflectance estimation using a two-stage approach involving NeuS [31] and InvNeRad [8] backbones.

In the first stage, we modify NeuS [31] to handle near-field lighting and strong specular reflections which are more pronounced due to co-located light and camera. We implicitly model the near-field effect by conditioning NeuS with a parametrization that encodes the distance between the scene point from the camera. This approach enables the neural radiance representation to effectively learn and model spatially varying incident illumination fields, even in the presence of global illumination. Consequently, we achieve a robust reconstruction of scene geometry, presenting a significant improvement over NeuS [31]. To further improve reconstruction at the regions of strong specularity from the flashlight, we adaptively set weights for each point on the scene, reducing contributions from those where the light is nearly perpendicular or parallel to the surface.

In the second stage, we use surface rendering to solve the rendering equation and extract principled BRDF [3]. Here we observe that Inverse Neural Radiosity (InvNeRad) [8], which models 2nd-order light bounces using a learnable radiance cache, expects the scene illumination to remain static across the capture, which is not true for images captured with a co-located light and camera. We solve this by training a light position-aware radiance cache network that can handle moving light sources. While this leads to accurate albedo prediction, we note that the predicted roughness is often noisy for regions with faulty geometry, which we address using a total-variation smoothness prior over roughness.

We perform a detailed quantitative and qualitative analysis of our approach on 3 real and 4 synthetic scenes. First, we show that while co-located light and camera in a dark room imposes additional constraints on the capture, it produces significantly better reflectance than alternate capture possibilities, i.e. natural illumination capture [36] and co-located flashlight under ambient lighting capture (WildLight) [4]. Next, we show that in comparison to the prior co-located light and camera-based inverse rendering algorithm, IRON [37], developed for a single object, our method produces significantly better geometry and reflectance.

In summary, our contributions are ● We propose GaNI, a neural inverse rendering algorithm that reconstructs the geometry and reflectance of a scene consisting of multiple objects captured with a co-located light and camera. ● We build upon the neural volume rendering technique NeuS [31] and the neural global illumination modeling technique InvNeRad [8] by proposing key technical contributions such as: (i) implicitly model near-field lighting (ii) surface angle loss to downweight regions with strong specular reflections (iii) learning radiance cache with varying light source to model global illumination (iv) total-variation smoothness prior on roughness. ● Our evaluations show that the proposed approach reconstructs significantly better geometry and reflectance compared to the state-of-the-art object-centric inverse rendering technique IRON [37]; and
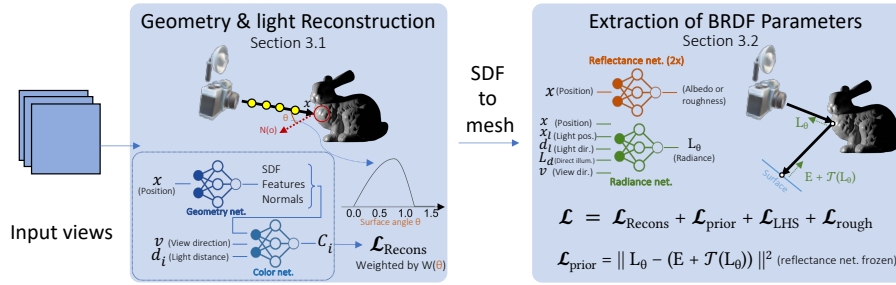
Fig. 2: **Overview of our pipeline**. Our system consists of two stages. In the first stage, we reconstruct geometry with volume rendering under near-field and global illumination. In the second stage, we extract accurate material properties with surface rendering while accounting for multi-bounce global illumination while solving the rendering equation by minimizing the radiometric prior. Our output is principled BRDF [3] (albedo and roughness) represented by two separate neural networks and a neural signed distance field.

predicts better reflectance and similar geometry compared to alternate capture setup, such as natural illumination [36] and WildLight [4].

Our code and data will be available upon acceptance.

## 2    Related Works

**Role of capture setup.** Researchers have explored different capture configurations with specific advantages that can prove beneficial depending on the intended applications. For example, some works [12,28] attempt to solve inverse rendering from unconstrained images captured in-the-wild but produce low-quality results. On the other hand, researchers have used Light Stage [7,26,39], a spherical gantry consisting of well-calibrated lights and cameras, to reconstruct very high-quality geometry and reflectance but can rarely be used in-the-wild.

Recently researchers have focused on simpler, 'at-home' configurations, that can be easily replicated in the wild but also provide sufficient constraints to generate better reconstructions than unconstrained approaches [12, 28]. These 'at-home' capture setups can be under unconstraint natural illumination [9, 17, 23, 32, 38, 40, 41], a co-located light and camera [1, 4, 18, 24, 27, 37] or a moving camera and a flashlight [15, 42]. Additionally, WildLight [4] has explored inverse rendering with co-located light and camera in the presence of ambient illumination. Recent researches have shown that the choice of illumination condition can heavily impact the quality of BRDF recovery. A natural illumination capture setup is most popular due to its simplicity, as it only requires multi-view images. However, such methods often cannot extract accurate reflectance, especially in multi-object scenes with complicated geometry due to fundamental ambiguity between geometry, lighting and reflectance. On the other hand, a co-located setup often gives high quality inverse rendering results, but requires a darkroom. WildLight combines the best of the two worlds, but requires roughly equal ambient illumination energy and flashlight energy, which is difficult to

satisfy in multi-object scenes. Our work is focused on co-located light and camera setup, which allows accurate inverse rendering results. We show that the state-of-the-art inverse rendering technique with co-located light and camera fails to model near-field and global illumination effects that are more visible in multi-object scenes than objects. Our approach is the first to enable co-located light-camera-based inverse rendering in scenes by implicitly modeling near-field and global illumination effects.

**Geometry Reconstruction.** Recently, NeuS [31] and VolSDF [33] have gained popularity since they achieve high-accuracy geometry by adopting a volumetric approach to optimize a neural implicit surface representation as a Signed Distance Field (SDF). We adopt NeuS [31] for reconstructing the geometry. However, NeuS cannot handle near-field illumination and strong specularity. Therefore, we propose to model near-field illumination and use surface angle loss to handle this scenario.

**Material Estimation Under Co-located Light and Camera.** Most previous inverse rendering algorithms, such as IRON [37], WildLight [4], for Co-located Light and Camera do not attempt to handle global illumination. While global illumination is often negligible for objects, such inter-reflection is much more prominent for multi-object scenes. Some previous methods [36,41] targeted at natural illumination estimate global illumination by leveraging a radiance field fitted to image observations. However such approaches are problematic for co-located light and camera setup as the radiance of only a subset of the scene can be observed.

To properly define radiance for the entire scene requires solving the rendering equation. Such a solution typically requires solving path integrals, i.e. integration over the path space and building paths that connect the camera to the light source. The most naive approach that serves this purpose is *differentiable path tracing*, which is known to have an intractable memory footprint in complex scenes and linear time complexity in the path length. To alleviate the time and memory requirements, a recent work, InvNeRad [8], adopted radiance caching techniques and only computes the primary ray intersection, and queries the cache for the contribution of the rest of the path, thus accounting for global illumination. In particular, InvNeRad [8] uses a neural network to represent the radiance cache and solve the rendering equation in self-supervised fashion. It shows significant improvements in terms of memory and time compared to previous work for inverse rendering.

However, InvNeRad assumes static illumination, and cannot be directly applied to co-located light and camera setup. Therefore, we propose a light source conditioned neural radiance cache to avoid creating hundreds of radiance caches and allow sharing computed radiance across different light positions. On real data, we also found small errors in geometry often lead to errors in roughness. Therefore, we propose a roughness regularization to reduce the error.

## 3   Method

**Capture Setup.** Similar to [37], we capture images with a smartphone with a flashlight turned on in a dark room. Such a setup provides multi-illumination images, which helps resolve ambiguity between illumination and material properties that occur in a natural illumination setup, allowing our method to obtain significantly better material estimation than capture under natural illumination. While this means our method can only be used at night time, it produces significantly better BRDF without requiring any sophisticated hardware equipment or calibration.

**Overview.** We observe that existing one-stage colocated inverse rendering systems, such as WildLight [4], which represents radiance as principled BRDF (Bidirectional Reflectance Distribution Function) parameters under a co-located flashlight, are often not robust against global illumination effects, such as inter-reflection between objects, and fails for real geometry reconstruction. Therefore, our reconstruction process consists of two optimization stages, as depicted in Figure 2. The first stage (Section A3) performs volume rendering by improving NeuS [31] to reconstruct the geometry with implicit near-field illumination and handle strong specular reflection with a surface angle loss. The second stage (Section A4) builds on InvNeRAD [8] and uses the predicted geometry to solve the rendering equation at all moving light positions and extracts principled BRDF parameters [3] with a smoothness prior on roughness.

### 3.1   Stage 1: Geometry Recovery with Volume Rendering

We build our volume rendering procedure on top of NeuS [31], and we similarly represent the scene by two neural networks.

$\quad$ **Geometry Network** $\mathrm{S}_{\Theta_\mathrm{S}}(\mathbf{x}) \to \{s, \mathbf{f}\}$, which maps a 3D position $\mathbf{x}$ to its signed distance to the closest surface, and a feature vector $\mathbf{f}$. The network represents the geometry of the scene.

$\quad$ **Color Network** $\mathrm{C}_{\Theta_\mathrm{C}}(\mathbf{x}, \mathbf{n}, \mathbf{v}, \mathbf{f}) \to \mathbf{c}$, where $\mathbf{x}$ is the 3D position of the point, $\mathbf{n}$ is the normal at position $\mathbf{x}$, $\mathbf{v}$ is the view direction of the camera and $\mathbf{f}$ is the feature vector output by the geometry network. The network represents the radiance of the scene.

$\quad$ The images are rendered with volume rendering as defined by NeuS [31]. Using their unbiased and occlusion-aware weight $\mathrm{w}(\cdot)$, with $\mathbf{o}$ as camera position, $\mathbf{v}$ as view direction, t as distance along camera ray, and $\mathrm{p}(t)$ as 3D positions along camera ray, we have the following equation.

$$\mathrm{C}(\mathbf{o}, \mathbf{v}) = \int_0^{+\infty} w(t)\, \mathrm{C}_{\Theta_\mathrm{C}}(\mathrm{p}(t), \mathbf{v}) dt \tag{1}$$

$\quad$ For a larger scene, the distance from co-located light to different parts of the scene will be different. However, as shown previously, the color network of NeuS does not model such differences. As shown in table 1 and Figure 4, previous inverse rendering methods such as IRON [37], which directly use NeuS [31] often fails in our multi-object dataset.

We model the co-located flashlight as a point light source. To handle the change of radiance, we additionally parameterize the color network with the flashlight position $\mathbf{x_i}$:

$$\mathrm{C}_{\Theta_{\mathrm{C}}}(\mathbf{x}, \mathbf{n}, \mathbf{v}, \mathbf{f}, \mathbf{x_i}) \rightarrow \mathbf{c} \tag{2}$$

Since the camera is co-located with the flashlight, $\mathbf{x_i} = \mathbf{x} + t\mathbf{v}$, where $t$ is the distance between the flashlight and queried point. We can uniquely identify a $(\mathbf{x_i}, \mathbf{x}, \mathbf{v})$ combination by just $(t, \mathbf{x}, \mathbf{v})$. Additionally, since the flashlight is a point light source, the irradiance follows the inverse quadratic rule $\frac{1}{t^2}$ as distance increases from the flashlight. As such, we propose the following parametrization:

$$\mathrm{C}_{\Theta_{\mathrm{C}}}(\mathbf{x}, \mathbf{n}, \mathbf{v}, \mathbf{f}, \frac{1}{t^2}) \rightarrow \mathbf{c} \tag{3}$$

By encouraging but not forcing the light conditioned radiance network to model near field light effects from a point light source, we can robustly and accurately reconstruct geometry even in the presence of global illumination.

Similar to previous methods, we encode view direction with spherical harmonics [22] and every other input with positional encoding [21].

We supervise the pixels with ground-truth image pixel values via reconstruction loss. However, as our color network models the illumination change, we supervise the network in linear RGB space to bypass the non-linearity introduced by gamma function. Together, we use linearized log loss introduced by RawNeRF [20] to better learn radiance in dark regions. Denote $y_i$ as the i-th pixel of groundtruth.

$$\tilde{L}_{\mathrm{recons}}(\hat{y}, y) = \sum_i \left( \frac{\mathrm{C}(\mathbf{o}, \mathbf{v}) - y_i}{sg(\mathrm{C}(\mathbf{o}, \mathbf{v})) + \epsilon} \right)^2 \tag{4}$$

**Surface Angle Weighting** We observe that strong reflections from the flashlight heavily deteriorate geometry prediction. While several works, such as neural-pbir [30] or Ref-NeuS, have explored loss weighting schemes to improve geometry reconstruction for specular objects, we found these methods often tend to degrade the geometry for concave scenes.

Our weighting scheme is based on the fact that the light position is known, and co-located with the camera, so we can calculate the angle the light ray forms with the surface normal. Our intuition is to downweight very small angles (perpendicular direction) that produce strong specular reflections and very large angles (grazing direction) that cause the diffuse component to become minimal, while specular inter-reflection become even stronger due to fresnel effects.

Since we do not have ground truth geometry, we use the neural sdf during training as proxy geometry. We render the surface normal similar to color:

$$\mathrm{N}(\mathbf{o}) = \int_0^{+\infty} w(t) \nabla \mathrm{S}_{\Theta_{\mathrm{S}}}(\mathbf{x}) dt \tag{5}$$

Given angle $\theta = \arccos(N \cdot L)$ between normal $N$ and light ray orientation $L$, we propose to weight the reconstruction loss $L_{\mathrm{recons}}$ of each pixels by:

$$W_{a,b}(\theta) = \begin{cases} \max(\cos(a(\theta - \frac{\pi}{4})), 0) & x \leq \frac{\pi}{4} \\ \max(\cos(b(\theta - \frac{\pi}{4})), 0) & x > \frac{\pi}{4} \end{cases} \qquad (6)$$

Our choice of function is motivated by the following requirements: decrease slowly in the middle near $\frac{\pi}{4}$, where neither specularity nor fresnel effects are prominent, and decrease rapidly on the two extremes where those effects quickly dominate observation; the function needs to be asymmetrical, as retroreflective specular highlights are only observable from a very narrow range of angles w.r.t surface normal, but specular inter-reflection is prominent over a significantly larger range of angles. We use $a = 2$, $b = 3$ for all our experiments.

**Structure-from-Motion (SfM) pointcloud supervision on real data.** Similar to other neural reconstruction techniques, we require camera poses to be known ahead of time before reconstruction, which is commonly estimated from structure-from-motion (SfM). Similar to prior works such as Geo-NeuS [5], Neilf++ [36], instead of discarding the point cloud estimated in SfM, our system can optionally use the point cloud to supervise geometry by requiring the neural signed distance field to be zero at point cloud locations $X_{\text{sfm}}$. We use SfM point cloud supervision on real data.

### 3.2  Stage 2: BRDF Estimation with Surface Rendering

In the second stage, we assume a surface based rendering model, and by properly modeling global illumination effects, we recover material properties represented by popular principled BRDF [3], which allows intuitive human editing for downstream applications. Specifically, we optimize for spatially varying albedo and roughness. We represent them using separate neural networks, $\phi_a(x)$ and $\phi_r(x)$, respectively, where $x$ is any 3D position.

With a co-located light and camera capture setup, illumination changes in every image, and we observe a subset of the radiance field. Hence, we cannot compute global illumination directly from image observations similar to some previous inverse rendering methods [41]. To recover unobserved radiance, we leverage an efficient global illumination inverse rendering technique InvNeRad [8] that solves the rendering equation directly on geometry defined by the signed distance field.

**Background: Inverse Neural Radiosity.** InvNeRad [8] is an inverse rendering method that efficiently accounts for global illumination without the need to explicitly simulate multiple scattering events through path tracing. It uses a radiance cache $L_\theta$ represented as a neural network with parameter set $\theta$, and after bouncing the ray only once, it queries the cache to collect the contribution of the rest of the path. More formally, the incident radiance at pixel $k$, denoted as $I_k$, is determined by the measurement equation,

$$I_k = \int_{\mathcal{A}} \int_{\mathcal{H}^2} W_k(x, \omega).(E + \mathcal{T}(L_\theta)(x, \omega)) dx d\omega^\perp \qquad (7)$$

where $W_k(x, \omega)$ models the response of a sensor pixel to incident radiance over its area $\mathcal{A}$ and the hemisphere of directions $\mathcal{H}^2$, $E$ is the emitted radiance distribution $E(x, \omega_o)$, $\mathcal{T}$ is the transport operator.

To optimize the scene parameters $\phi$ from the ground truth images, we denote the image rendered using $\phi$ as $I(\phi)$, then we can define a reconstruction loss as

$$\mathcal{L}_{\text{recons}}(I(\phi)) = \|I(\phi) - I^{\text{GT}}\|. \tag{8}$$

To train the radiance cache, InvNeRad [8] introduces a radiometric prior

$$\mathcal{L}_{\text{prior}}(\theta) = \|L_\theta(x, \omega_o) - (E(x, \omega_o) + \mathcal{T}(L_\theta)(x, \omega_o))\|. \tag{9}$$

Following InvNeRad [8], we additionally constrain the neural radiance cache with groundtruth radiance from images by directly rendering the neural radiance cache into an image $I_k^{LHS}$. Then we define an additional loss term to match the image to groundtruth.

$$\mathcal{L}_{\text{LHS}}(\theta) = \left\| I^{LHS}(\theta) - I^{\text{GT}} \right\|^2. \tag{10}$$

**Co-located Light Inverse Neural Radiosity.** InvNeRad assumes illumination is constant across the input views. With our co-located light capture setup, the lighting changes with each view. A naive approach is to have a radiance cache for every flashlight $l$, and formulate the radiometric prior as:

$$\mathcal{L}_{\text{prior}}(\theta) = \|L_{l,\theta}(x, \omega_o) - (E_l(x, \omega_o) + \mathcal{T}(L_{l,\theta})(x, \omega_o))\| \tag{11}$$

where for each captured image with flashlight $l$, $L_{l,\theta}$ requires a separate network. This is computationally prohibitive as we could easily have close to 1000 different light positions, thus 1000 different networks. Instead, we propose to use a single radiance cache network to represent radiance at all flashlight positions, which will take additional parameters $x_l$ and $d_l$, the position of the flashlight and orientation, as input. We also noticed that under a co-located light scenario, direct illumination often has strong discontinuities that change with light positions, making it hard for the network to model. However, it can be computed easily. Therefore, we additionally condition the neural radiance network on direct illumination $L_{\text{direct}}$. Denote visibility of surface from flashlight as $\mathbb{1}_{\text{vis}}$, flashlight radiant intensity as $E_{\text{flash}}$, and distance between flashlight and surface as $d$.

$$L_{\text{direct}} = \mathbb{1}_{\text{vis}} \frac{E_{\text{flash}}}{d^2} \tag{12}$$

**Roughness Regularization.** Unlike the original InvNeRad setup, we do not have access to groundtruth geometry. We found small imperfections will often lead to significant errors in roughness. Therefore, we additionally incorporate a total variation loss for roughness, and we set $\lambda = 0.02$ in all our experiments.

$$L_{rough} = \lambda|\nabla\phi_r| \tag{13}$$

Finally, the optimization task is

$$\phi^*, \theta^* = \arg\min_{\phi,\theta} \mathcal{L}_{\text{recons}}(I(\phi)) + \mathcal{L}_{\text{prior}}(\theta) + \mathcal{L}_{\text{rough}}(\phi) + \mathcal{L}_{\text{LHS}}(\theta) \tag{14}$$

**Implementation Details** Similar to InvNeRad, we implement our second stage on top of Mitsuba 3 [10], and use hashgrid encoding to encode the 3D position inputs to the reflectance fields $\phi_a$ and $\phi_r$. We also use the hashgrid encoding for flashlight positions, and spherical harmonics for flashlight orientation to encourage smoothness between contiguous flashlight positions but still allow high frequency changes.

To further speedup global illumination computation, we extract the neural sdf at $2048^3$ resolution and store as a mesh. We found the resolution sufficient to represent the details of the scenes we tested.

**Table 1: Quantitative comparison of geometry on synthetic data.** Our method produces significantly better geometry than IRON [37] and slightly better geometry than WildLight [4] with images captured with co-located lighting and camera in a darkroom.

|  | Chamfer ($\times 100$) | | | Depth Map L1 | | |
|---|---|---|---|---|---|---|
|  | Bedroom | Shelf | Counter | Bedroom | Shelf | Counter |
| IRON | 1.012 | 0.194 | 1.672 | 0.124 | 0.076 | 0.174 |
| WildLight | **0.000** | <u>0.013</u> | <u>0.076</u> | **0.004** | <u>0.026</u> | <u>0.024</u> |
| Ours | <u>0.032</u> | **0.011** | **0.023** | <u>0.012</u> | **0.014** | **0.015** |

## 4  Evaluation

We perform qualitative evaluation on real data and quantitative and qualitative evaluation on synthetic data.
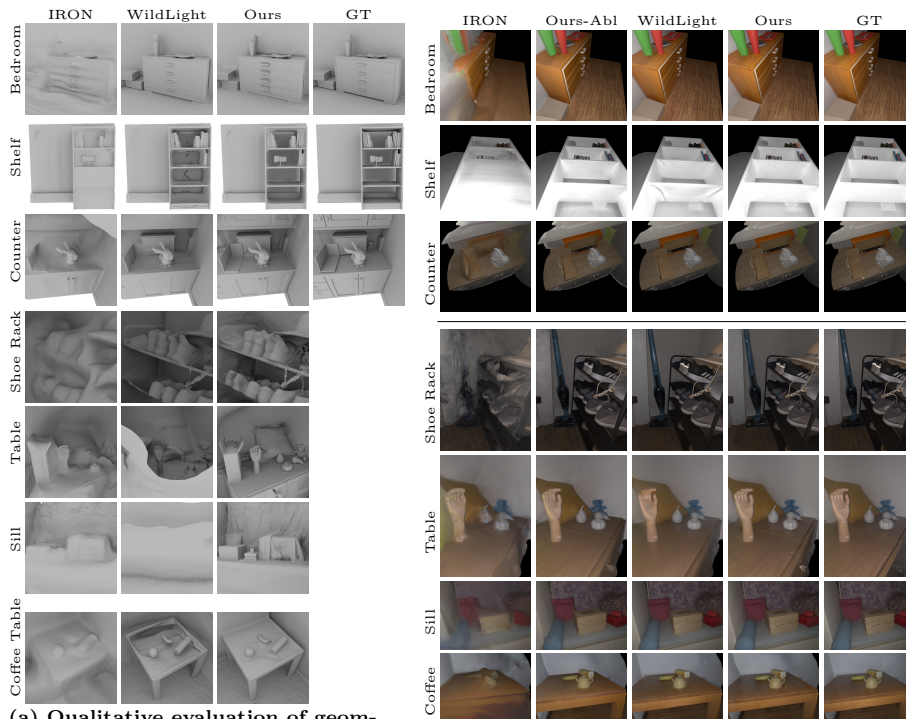
**Synthetic data.** We created three scenes, **bedroom**, **shelf**, and **kitchen counter** based on openly available Blender and Mitsuba scenes. We render the dataset with a moving co-located point light and camera as HDR images with the Mitsuba 3 renderer. Each synthetic scene contains 1000 training images and 50 validation images with randomly generated camera poses.

**Real data.** We captured four real scenes using co-located light and camera, **window sill**, **table**, **shoe rack**, and **coffee table** that contain 1006, 704, 564, and 650 images respectively. We reserve 5% of the images of each scene as the validation split.

**Table 2: Qualitative evaluation of reflectance on synthetic data (Stage 2).** We present MSE ($\times 10$) of re-rendering in unseen views, albedo, and roughness on the validation set of synthetic data. Best is shown in **bold** and second places <u>underlined</u>.

|  | Validation View | | | | Albedo | | | | Roughness | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
|  | Bedroom | Shelf | Counter | *Mean* | Bedroom | Shelf | Counter | *Mean* | Bedroom | Shelf | Counter | *Mean* |
| IRON | 0.122 | 0.212 | 0.029 | 0.121 | 0.565 | <u>0.278</u> | <u>0.271</u> | <u>0.371</u> | 1.003 | 2.530 | 2.645 | 2.059 |
| WildLight | **0.006** | <u>0.040</u> | <u>0.005</u> | <u>0.017</u> | 0.101 | 1.158 | 0.603 | 0.621 | <u>0.636</u> | <u>1.340</u> | <u>1.933</u> | <u>1.303</u> |
| Ours | <u>0.010</u> | **0.023** | **0.005** | **0.013** | **0.034** | **0.051** | **0.036** | **0.041** | **0.380** | **1.276** | **1.118** | **0.924** |
| Ours-Abl | 0.016 | 0.046 | 0.014 | 0.025 | <u>0.079</u> | 1.404 | 0.665 | 0.716 | 0.677 | 2.656 | 2.662 | 1.998 |

**Comparison.** We compare our approach with that of WildLight [4], which does not handle global illumination, and IRON [37], which additionally does

(a) **Qualitative evaluation of geometry.** Our method produces significantly better geometry than IRON [37] on synthetic and real data. Compared to Wild-Light [4] we produce slightly better geometry on synthetic data and significantly better on real data.

(b) **Qualitative comparison of re-rendering.** We present re-rendering in validation views for the synthetic scenes (row 1-3) and real scenes (row 4-7). Our method produces better re-rendering w.r.t. IRON [37] due to our ability to better model near-field and global illumination.

**Fig. 3:** Qualitative comparison of geometry and re-rendering.

not model near-field illumination. We test all compared methods in a dark-room. While WildLight is designed for capture under ambient ilumination, we found such design has trouble converging correctly for multi-object scenes under ambient illumination. We will show the results in supplementary. The original implementation of IRON uses Mitsuba roughplastic BRDF. Since all of our synthetic scenes use principled BRDF, we modified IRON to use principled BRDF for a fair comparison. We also compare with Nelif++ [36] which performs inverse rendering under natural illumination without requiring a dark room.

We reconstruct better geometry and reflectance than IRON [37] and Wild-Light [4] (see Fig. 3a, 4 and Tab. 1, 2). Compared to other capture setups that do not require dark room, Nelif++ [36] we improve the geometry slightly but significantly improve BRDF prediction.

### 4.1  Comparison with Inverse Rendering Techniques that use Co-located Light & Camera

**Evaluation of Geometry.** We compare the geometry reconstructed by our method with state of art inverse rendering methods using a co-located capture
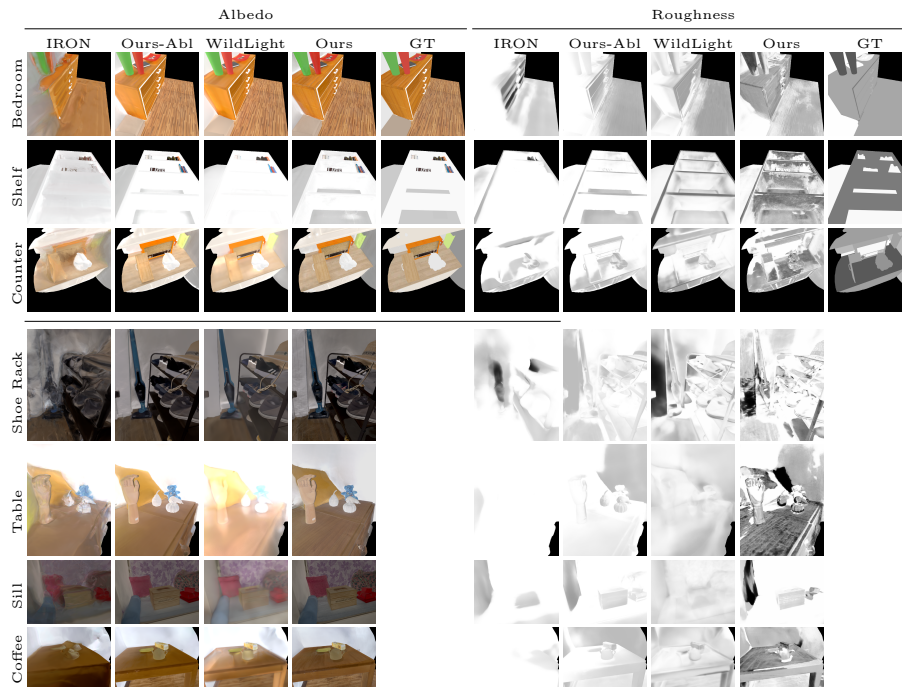
Fig. 4: **Qualitative comparison of reflectance estimation.** We present estimated albedo, and roughness in validation views for the synthetic scenes (row 1-3) and real scenes (row 4-7). Our method produces significantly better albedo, roughness and re-rendering w.r.t. IRON [37] due to our ability to better model near-field and global illumination. **(Please zoom in for better visualization)**

setup. We compare to IRON [37] and WildLight [4] quantitatively on synthetic data in Tab. 1 and qualitatively on both synthetic and real data in Fig. 3a.

As the mesh models of our synthetic scene contain hidden parts not visible in any camera views, we cannot naively compute the distance between reconstructed geometry and ground-truth meshes. Therefore, we rendered a depth map of ground-truth mesh and reconstructed mesh on validation views, and compute the L1 distance between the two depth maps as our depth map score. Additionally, we can convert the depth map of ground-truth mesh and of reconstructed mesh as two-point clouds. Then we compute the Chamfer distance from the ground-truth point cloud to the point cloud of the reconstructed mesh.

IRON [37] does not attempt to model near-field illumination and performs poorly in multi-object scenes. WildLight [4] assumes a point light source model which does not model global illumination. Therefore, while WildLight performed well on an open scene without much inter-reflection, such as 'bedroom', the performance is worse on scenes with strong concavity such as 'shelf', 'kitchen counter' and complex geometry, e.g. the sign in 'shelf' scene, the teddy bear in the back of the kitchen 'counter', or the left-most shoe on the second row of the 'shoe rack'. Additionally, WildLight often catastrophically fails on many

real scenes with prominent inter-reflection, such as the 'table', 'window sill', or 'coffee table'.

In comparison, our method is more robust in the presence of global illumination. However, since we do not strictly enforce a physically based BRDF model but implicitly regularize the radiance with a neural network, we perform slightly worse in textureless areas such as the wall.

**Quantitative Evaluation of Reflectance on Synthetic data** We compare the reflectance estimation of our method with that of IRON [37] quantitatively on synthetic data in 2 and qualitatively on both synthetic and real data in Figure 4 on held-out validation set. Along with albedo and roughness estimation, we also re-render the scene in validation views, unseen during training, to visualize the overall effect of inverse rendering. We compute Mean Squared Error (MSE) for re-rendering in validation views and prediction of albedo and roughness. During re-rendering in validation views, we clip the range of both ground-truth and prediction to $[0, 1]$ to prevent high-intensity pixels from dominating the error.

For quantitative evaluation, we optimize the point flashlight energy



**Fig. 5:** Comparison of our method with Neilf++ [36], a state-of-the-art inverse rendering algorithm for natural illumination. We captured the same scene under both natural and co-located illuminations with similar number of images and camera poses. We found that our method significantly outperforms Neilf++, epsecially in albedo.

during training for all methods. Therefore, we rescale the predicted albedo of each method by the flashlight energy. For qualitative evaluation, we mask out pixels that do not have the corresponding geometry reconstructed.

Overall, our method significantly outperforms IRON [37] on average for albedo estimation (MSE 0.066 vs 0.552), roughness estimation (MSE: 0.846 vs 1.792), and re-rendering in validation views (MSE: 0.014 vs 0.129). Similarly on real data, we significantly outperform IRON [37]. Re-rendering images of IRON [37] contain many artifacts due to poorly reconstructed geometry, e.g. Cabinet gets fused with in wall in Bedroom; two layers of 'shelf' are reconstructed as one solid block. On the other hand, our method has very minimal errors in re-rendering of novel views, and hence predicted geometry and reflectance.

## 4.2 Comparison with Natural Illumination and Hybrid Illumination

Both our approach and IRON expects dark room for capture, which means it can be only performed at nighttime. Recent methods like Neilf++ [36] can perform inverse rendering in the presence of natural illumination. that our method outperforms Neilf++ [36], which can capture scenes under natural illumination.
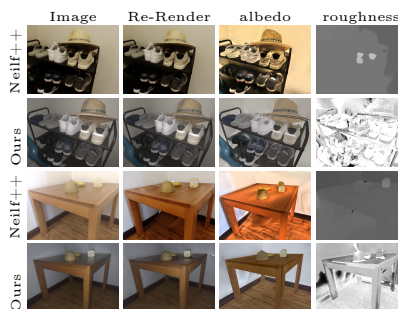
We captured the same scene with similar number of images and poses under both natural and co-located illumination. We captured 750 images under both illuminations for the first example shoe rack, and 649 and 686 for co-located and natural illumination respectively for the second example coffee table. However, in Fig. 5 we show prominent errors exist on the right side of the white shoe under the hat, as incorrect estimation of local illumination results in very high intensity albedo. Shading also gets baked into albedo of the hat. Similarly shading are baked into the albedo of the table for the second example. The difference is more obvious for reflectance estimation, and hence generating re-rendering views, highlighting the importance of using co-located capture for high-quality reflectance estimation.

### 4.3   Ablation studies

We perform ablation of individual component of our system, surface angle weighting, and our second stage to show their effectiveness.



**Ablation study of Surface Angle Weighting** In Figure 6, we conduct ablation study of our proposed weighting scheme in the first stage. The artifacts on the table highlighted in red are present in the ablated version which is caused by interreflection. Such artifacts become significantly less pronounced in the full version as we put less weight on pixels at parallel angle.
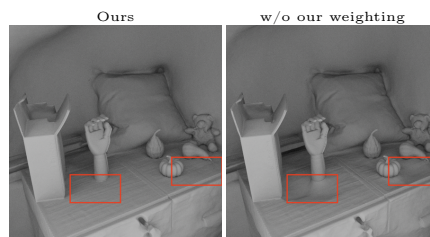
**Fig. 6:** Ablation of our proposed surface angle weighting. Red box highlights areas where specular inter-reflection causes artifacts in geometry without surface angle weighting. Such error become significantly less pronounced in our full variant.

**Ablation study of Ours Second Stage** To separate the effect of geometry and reflectance, we perform an ablation study where we replace the second stage of our system with IRON's second stage, which we denote as Ours-Abl. Ours-Abl uses the same reconstructed geometry, only differing from ours in material properties extraction. As shown in Figure 4, even with good geometry initialization, Ours-Abl often leaves strong artifacts in the albedo, such as on the side of cabinet of **kitchen counter** or the top right corner of the pillow in **table**. Without proper handling global illumination, Ours-Abl often heavily bakes global illumination effects into albedo. Moreover, Ours-Abl often has trouble recovering roughness in regions where global illumination is prominent, such as the side of cabinet, or the inside of shelf.

# 5  Conclusion

We introduce a novel system for inverse rendering for scenes using co-located light and camera. By modeling changing point light source near field illumination, we are able to obtain robust and accurate reconstruction of geometry under co-located flashlight. Then we propose light position conditioned radiance cache which extends InvNeRad to co-located light and extracts accurate material properties under global illumination. Our experiments show strong performance against state of the art methods both in geometry reconstruction and reflectance recovery. We believe our work will be an important step toward accurate inverse rendering of scenes.

**Limitations** Similar to previous co-located light and camera approaches, our approach requires capturing at nightime with lights off. Our current approach is focused on reconstruction of multi-object scenes rather than room scale scenes. One potential future direction is to extend our approach to full room scale scenes.

# References

1. Bi, S., Xu, Z., Sunkavalli, K., Kriegman, D., Ramamoorthi, R.: Deep 3d capture: Geometry and reflectance from sparse multi-view images. In: 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 5959–5968. IEEE Computer Society, Los Alamitos, CA, USA (jun 2020). https://doi.org/10.1109/CVPR42600.2020.00600, https://doi.ieeecomputersociety.org/10.1109/CVPR42600.2020.00600 1, 2, 4

2. 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: Computer Vision – ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part III. p. 294–311. Springer-Verlag, Berlin, Heidelberg (2020). https://doi.org/10.1007/978-3-030-58580-8_18, https://doi.org/10.1007/978-3-030-58580-8_18 1

3. Burley, B., Studios, W.D.A.: Physically-based shading at disney. In: Acm Siggraph. vol. 2012, pp. 1–7. vol. 2012 (2012) 3, 4, 6, 8

4. Cheng, Z., Li, J., Li, H.: Wildlight: In-the-wild inverse rendering with a flashlight. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 4305–4314 (June 2023) 1, 2, 3, 4, 5, 6, 10, 11, 12

5. Fu, Q., Xu, Q., Ong, Y.S., Tao, W.: Geo-neus: Geometry-consistent neural implicit surfaces learning for multi-view reconstruction. In: Koyejo, S., Mohamed, S., Agarwal, A., Belgrave, D., Cho, K., Oh, A. (eds.) Advances in Neural Information Processing Systems. vol. 35, pp. 3403–3416. Curran Associates, Inc. (2022), https://proceedings.neurips.cc/paper_files/paper/2022/file/16415eed5a0a121bfce79924db05d3fe-Paper-Conference.pdf 1, 8

6. Goldman, D.B.: Vignette and exposure calibration and compensation. IEEE Transactions on Pattern Analysis and Machine Intelligence **32**(12), 2276–2288 (2010). https://doi.org/10.1109/TPAMI.2010.55 2

7. Guo, K., Lincoln, P., Davidson, P., Busch, J., Yu, X., Whalen, M., Harvey, G., Orts-Escolano, S., Pandey, R., Dourgarian, J., Tang, D., Tkach, A., Kowdle, A., Cooper, E., Dou, M., Fanello, S., Fyffe, G., Rhemann, C., Taylor, J., Debevec, P., Izadi, S.: The relightables: Volumetric performance capture of humans with realistic relighting. ACM Trans. Graph. **38**(6) (nov 2019). https://doi.org/10.1145/3355089.3356571, https://doi.org/10.1145/3355089.3356571 4

8. Hadadan, S., Lin, G., Novák, J., Rousselle, F., Zwicker, M.: Inverse global illumination using a neural radiometric prior. In: ACM SIGGRAPH 2023 Conference Proceedings. SIGGRAPH '23, Association for Computing Machinery, New York, NY, USA (2023). https://doi.org/10.1145/3588432.3591553, https://doi.org/10.1145/3588432.3591553 3, 5, 6, 8, 9

9. Hasselgren, J., Hofmann, N., Munkberg, J.: Shape, Light, and Material Decomposition from Images using Monte Carlo Rendering and Denoising. arXiv:2206.03380 (2022) 4

10. Jakob, W., Speierer, S., Roussel, N., Nimier-David, M., Vicini, D., Zeltner, T., Nicolet, B., Crespo, M., Leroy, V., Zhang, Z.: Mitsuba 3 renderer (2022), https://mitsuba-renderer.org 10, 1

11. Kim, H., Zollöfer, M., Tewari, A., Thies, J., Richardt, C., Theobalt, C.: Inverse-facenet: Deep single-shot inverse face rendering from a single image. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (2018) 1

12. Li, Z., Shafiei, M., Ramamoorthi, R., Sunkavalli, K., Chandraker, M.: Inverse rendering for complex indoor scenes: Shape, spatially-varying lighting and svbrdf from

a single image. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (June 2020) 1, 2, 4

13. LibRaw LLC: Libraw. https://github.com/LibRaw/LibRaw 2
14. Lichy, D., Sengupta, S., Jacobs, D.W.: Fast light-weight near-field photometric stereo. In: 2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 12602–12611. IEEE Computer Society, Los Alamitos, CA, USA (jun 2022). https://doi.org/10.1109/CVPR52688.2022.01228, https://doi.ieeecomputersociety.org/10.1109/CVPR52688.2022.01228 1, 2
15. Lichy, D., Wu, J., Sengupta, S., Jacobs, D.W.: Shape and material capture at home. In: 2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 6119–6129. IEEE Computer Society, Los Alamitos, CA, USA (jun 2021). https://doi.org/10.1109/CVPR46437.2021.00606, https://doi.ieeecomputersociety.org/10.1109/CVPR46437.2021.00606 1, 4
16. Lindenberger, P., Sarlin, P.E., Larsson, V., Pollefeys, M.: Pixel-Perfect Structure-from-Motion with Featuremetric Refinement. In: ICCV (2021) 2
17. Liu, Y., Wang, P., Lin, C., Long, X., Wang, J., Liu, L., Komura, T., Wang, W.: Nero: Neural geometry and brdf reconstruction of reflective objects from multiview images (2023) 4
18. Luan, F., Zhao, S., Bala, K., Dong, Z.: Unified shape and svbrdf recovery using differentiable monte carlo rendering. Computer Graphics Forum 40 (2021), https://api.semanticscholar.org/CorpusID:232404668 2, 4
19. Maik Riechert: rawpy. https://github.com/letmaik/rawpy (2022) 2
20. Mildenhall, B., Hedman, P., Martin-Brualla, R., Srinivasan, P.P., Barron, J.T.: NeRF in the dark: High dynamic range view synthesis from noisy raw images. CVPR (2022) 7
21. Mildenhall, B., Srinivasan, P.P., Tancik, M., Barron, J.T., Ramamoorthi, R., Ng, R.: Nerf: Representing scenes as neural radiance fields for view synthesis. Commun. ACM 65(1), 99–106 (dec 2021). https://doi.org/10.1145/3503250, https://doi.org/10.1145/3503250 7
22. Müller, T., Evans, A., Schied, C., Keller, A.: Instant neural graphics primitives with a multiresolution hash encoding. ACM Trans. Graph. 41(4) (jul 2022). https://doi.org/10.1145/3528223.3530127, https://doi.org/10.1145/3528223.3530127 7
23. Munkberg, J., Hasselgren, J., Shen, T., Gao, J., Chen, W., Evans, A., Müller, T., Fidler, S.: Extracting Triangular 3D Models, Materials, and Lighting From Images. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 8280–8290 (June 2022) 4
24. Nam, G., Lee, J.H., Gutierrez, D., Kim, M.H.: Practical svbrdf acquisition of 3d objects with unstructured flash photography. ACM Trans. Graph. 37(6) (dec 2018). https://doi.org/10.1145/3272127.3275017, https://doi.org/10.1145/3272127.3275017 2, 4
25. NVIDIA: Nvidia optix ray tracing engine, https://developer.nvidia.com/rtx/ray-tracing/optix 1
26. Pandey, R., Orts-Escolano, S., LeGendre, C., Haene, C., Bouaziz, S., Rhemann, C., Debevec, P., Fanello, S.: Total relighting: Learning to relight portraits for background replacement. vol. 40 (August 2021). https://doi.org/10.1145/3450626.3459872 1, 4
27. Schmitt, C., Donne, S., Riegler, G., Koltun, V., Geiger, A.: On joint estimation of pose, geometry and svbrdf from a handheld scanner. In: Proceedings IEEE Conf. on Computer Vision and Pattern Recognition (CVPR) (2020) 4

28. Sengupta, S., Gu, J., Kim, K., Liu, G., Jacobs, D.W., Kautz, J.: Neural inverse rendering of an indoor scene from a single image. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 8598–8607 (2019) 2, 4

29. Sengupta, S., Kanazawa, A., Castillo, C.D., Jacobs, D.W.: Sfsnet: Learning shape, reflectance and illuminance of facesin the wild'. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 6296–6305 (2018) 1

30. Sun, C., Cai, G., Li, Z., Yan, K., Zhang, C., Marshall, C., Huang, J., Zhao, S., Dong, Z.: Neural-pbir reconstruction of shape, material, and illumination. In: 2023 IEEE/CVF International Conference on Computer Vision (ICCV). pp. 18000–18010. IEEE Computer Society, Los Alamitos, CA, USA (oct 2023). https://doi.org/10.1109/ICCV51070.2023.01654, https://doi.ieeecomputersociety.org/10.1109/ICCV51070.2023.01654 7

31. Wang, P., Liu, L., Liu, Y., Theobalt, C., Komura, T., Wang, W.: Neus: Learning neural implicit surfaces by volume rendering for multi-view reconstruction. In: Ranzato, M., Beygelzimer, A., Dauphin, Y., Liang, P., Vaughan, J.W. (eds.) Advances in Neural Information Processing Systems. vol. 34, pp. 27171–27183. Curran Associates, Inc. (2021), https://proceedings.neurips.cc/paper_files/paper/2021/file/e41e164f7485ec4a28741a2d0ea41c74-Paper.pdf 1, 3, 5, 6, 2

32. Yao, Y., Zhang, J., Liu, J., Qu, Y., Fang, T., McKinnon, D., Tsin, Y., Quan, L.: Neilf: Neural incident light field for physically-based material estimation. In: European Conference on Computer Vision (ECCV) (2022) 4

33. Yariv, L., Gu, J., Kasten, Y., Lipman, Y.: Volume rendering of neural implicit surfaces. In: Thirty-Fifth Conference on Neural Information Processing Systems (2021) 5

34. Yariv, L., Kasten, Y., Moran, D., Galun, M., Atzmon, M., Ronen, B., Lipman, Y.: Multiview neural surface reconstruction by disentangling geometry and appearance. Advances in Neural Information Processing Systems **33** (2020) 2

35. Yu, Y., Smith, W.A.: Inverserendernet: Learning single image inverse rendering. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (2019) 2

36. Zhang, J., Yao, Y., Li, S., Liu, J., Fang, T., McKinnon, D., Tsin, Y., Quan, L.: Neilf++: Inter-reflectable light fields for geometry and material estimation. In: International Conference on Computer Vision (ICCV) (2023) 3, 4, 5, 8, 11, 13

37. Zhang, K., Luan, F., Li, Z., Snavely, N.: Iron: Inverse rendering by optimizing neural sdfs and materials from photometric images. In: 2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 5555–5564 (2022). https://doi.org/10.1109/CVPR52688.2022.00548 1, 2, 3, 4, 5, 6, 10, 11, 12, 13

38. Zhang, K., Luan, F., Wang, Q., Bala, K., Snavely, N.: Physg: Inverse rendering with spherical gaussians for physics-based material editing and relighting. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 5453–5462 (June 2021) 1, 4

39. Zhang, X., Fanello, S., Tsai, Y.T., Sun, T., Xue, T., Pandey, R., Orts-Escolano, S., Davidson, P., Rhemann, C., Debevec, P., Barron, J.T., Ramamoorthi, R., Freeman, W.T.: Neural light transport for relighting and view synthesis. ACM Trans. Graph. **40**(1) (jan 2021). https://doi.org/10.1145/3446328, https://doi.org/10.1145/3446328 4

40. 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. ACM Trans. Graph. **40**(6) (dec 2021). https://doi.org/10.1145/3478513.3480496, https://doi.org/10.1145/3478513.3480496 4

41. Zhang, Y., Sun, J., He, X., Fu, H., Jia, R., Zhou, X.: Modeling indirect illumination for inverse rendering (2022). `https://doi.org/10.48550/ARXIV.2204.06837`, `https://arxiv.org/abs/2204.06837` 4, 5, 8

42. Zhao, D., Lichy, D., Perrin, P.N., Frahm, J.M., Sengupta, S.: Mvpsnet: Fast generalizable multi-view photometric stereo. arXiv preprint arXiv:2305.11167 (2023) 4

43. Zhu, J., Huo, Y., Ye, Q., Luan, F., Li, J., Xi, D., Wang, L., Tang, R., Hua, W., Bao, H., Wang, R.: I$^2$-sdf: Intrinsic indoor scene reconstruction and editing via raytracing in neural sdfs. In: CVPR (2023) 2

44. Zhu, J., Luan, F., Huo, Y., Lin, Z., Zhong, Z., Xi, D., Wang, R., Bao, H., Zheng, J., Tang, R.: Learning-based inverse rendering of complex indoor scenes with differentiable monte carlo raytracing. In: SIGGRAPH Asia 2022 Conference Papers. SA '22, Association for Computing Machinery, New York, NY, USA (2022). `https://doi.org/10.1145/3550469.3555407`, `https://doi.org/10.1145/3550469.3555407` 1

45. Zhu, R., Li, Z., Matai, J., Porikli, F., Chandraker, M.: Irisformer: Dense vision transformers for single-image inverse rendering in indoor scenes. In: 2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 2812–2821. IEEE Computer Society, Los Alamitos, CA, USA (jun 2022). `https://doi.org/10.1109/CVPR52688.2022.00284`, `https://doi.ieeecomputersociety.org/10.1109/CVPR52688.2022.00284` 2

## A1 Summary

Here we summarize our supplementary.

- We render all our scenes under natural lighting conditions with our recovered geometry and material parameters using path tracing renderer mitsuba [10], and denoise with its built-in NVIDIA OptiX denoiser [25]. The videos feature ambient lighting with moving point light sources and cameras to demonstrate the practical applications of our method. They are included as separate mp4 files.
- In Figure 7, we show qualitative results of WildLight on WildLight style capture setup (Co-Located Light and Camera under ambient natural illumination). We found WildLight often fail to converge on our multi-object scenes.
- In section A2, we provide additional details regarding how we capture and process our real data.
- In section A3, we provide additional details for our stage 1 architecture and training.
- In section A4, we provide additional details for our stage 2 training.
- In Figure 8, Figure 9, Figure 10, we provide additional visualization of qualitative comparisons on synthetic data. In Figure 11, Figure 12, Figure 13, Figure 14, we provide additional visualization of qualitative comparisons on real data. Additional details of the comparison are discussed in section A5.
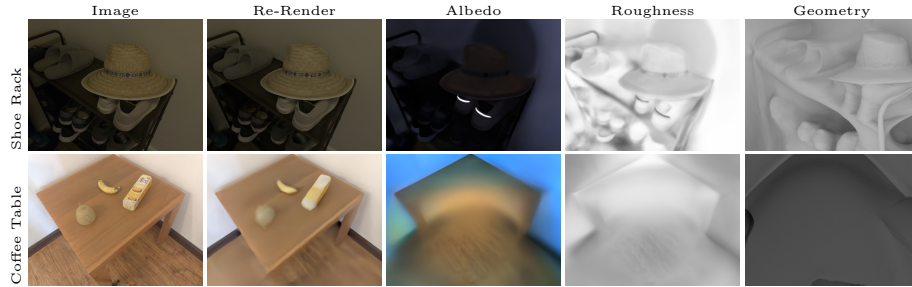
**Fig. 7:** Qualitative results of WildLight on Co-located Light and Camera under ambient natural illumination.

## A2    Real Data Capture Setup and Post-processing

We capture all of our real data using an iPhone XS Max and an iPhone 11 Pro. We capture all the image with ProCamera app on iOS as raw dng file. During capture, we keep manual and fixed white balance, focus and exposure. For co-located capture, we also keep flashlight constantly on through the capture session. We process the raw files with RawPy [19], which is a python interface around libraw [13]. We perform structure-from-motion reconstruction using pixel perfect sfm [16]. We apply camera undistortion parameters estimated by pixel perfect sfm to our captured images. We found that both iPhone XS Max and iPhone 11 Pro experience significant vignetting. To calibrate for vignetting, we use a piece of white paper on a sunny day under direct sunlight as the calibration target. We model the vignetting as 6-th degree even order polynomial [6], and apply vignetting correction accordingly. We store all final processed images as 16-bit unsigned png images with linear response curve without any gamma curve applied, which are used for all following experiments.

## A3    Additional Details of Stage 1 Architecture and Training

In section 3.1 of the main paper, we described some changes to original NeuS architecture, including introducing light position conditioned radiance field. Here we described additional details of implementation.

    With our synthetic dataset, the images can contain "background pixels" where the primary ray from the camera does not intersect with any scene geometry during rendering. Since we are not using environment map for our scenes, the values of these pixels are undefined, and we use a per image binary mask to ignore these pixels during training. Consequently, we do not put any supervision in the background region. To prevent the network from producing arbitrary values for the background, we adopt mask loss commonly used by prior works [31, 34]. NeuS [31] defines unbiased weights $w_{k,i}$ along the k-th camera ray based on the underlying signed distance field. Denote $\hat{O}_k = \sum_{i=1}^{n} w_i$ as the the sum of weights

along the k-th camera ray, $M_k = \{0, 1\}$ as the value of the binary mask on the k-th pixel, and BCE as the binary cross entropy loss, we have the following equation.

$$L_{\text{mask}} = \text{BCE}(M_k, \hat{O}_k) \tag{15}$$

Such mask loss is only used for synthetic data, and not used for real data.

Stage 1 is trained with batch size of 512 rays, learning rate of $5 \times 10^{-4}$ for 500K steps on synthetic data, and 1M steps on real data.

## A4 Additional Details of Stage 2 Architecture and Training

Here we provide additional details regarding the second stage of our system. We use Principled BRDF [3] in stage 2. Similar to InvNeRad [8], we only optimize for albedo and roughness, while keeping other parameters fixed. We set all fixed parameters to zeros, except for "specular" (which sometimes called "specular albedo"), which we set to 0.6 for synthetic scenes and 0.5 for real scenes. For fairness of comparison, we also fix "specular" to the same values for IRON [37].

We train our stage 2 for 48000 iterations at learning rate of $5 \times 10^{-4}$ with batch size of 16384 for synthetic data and 15625 for real data.

## A5 Additional Visualization of Qualitative Results

In Figure 8, Figure 9, Figure 10, we provide additional visualization of qualitative results on synthetic data. In Figure 11, Figure 12, Figure 13, Figure 14, we provide additional visualization of qualitative results on real data.
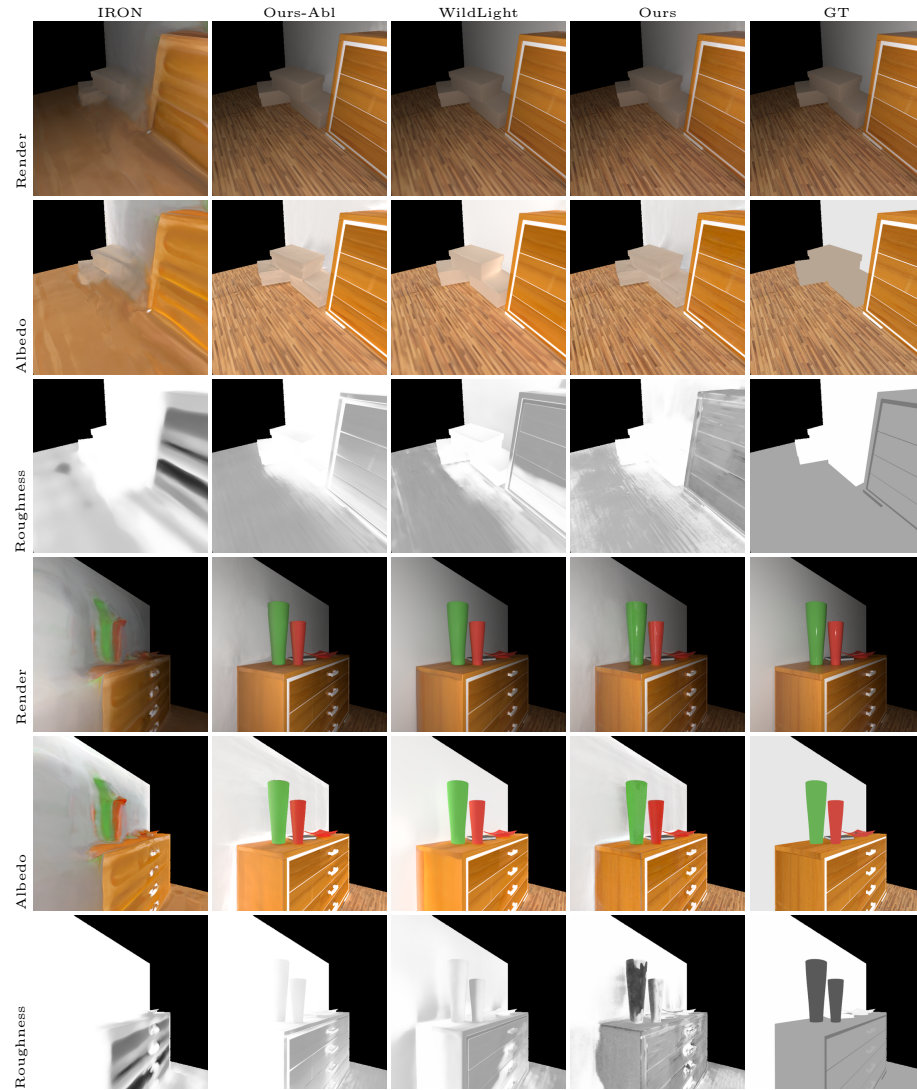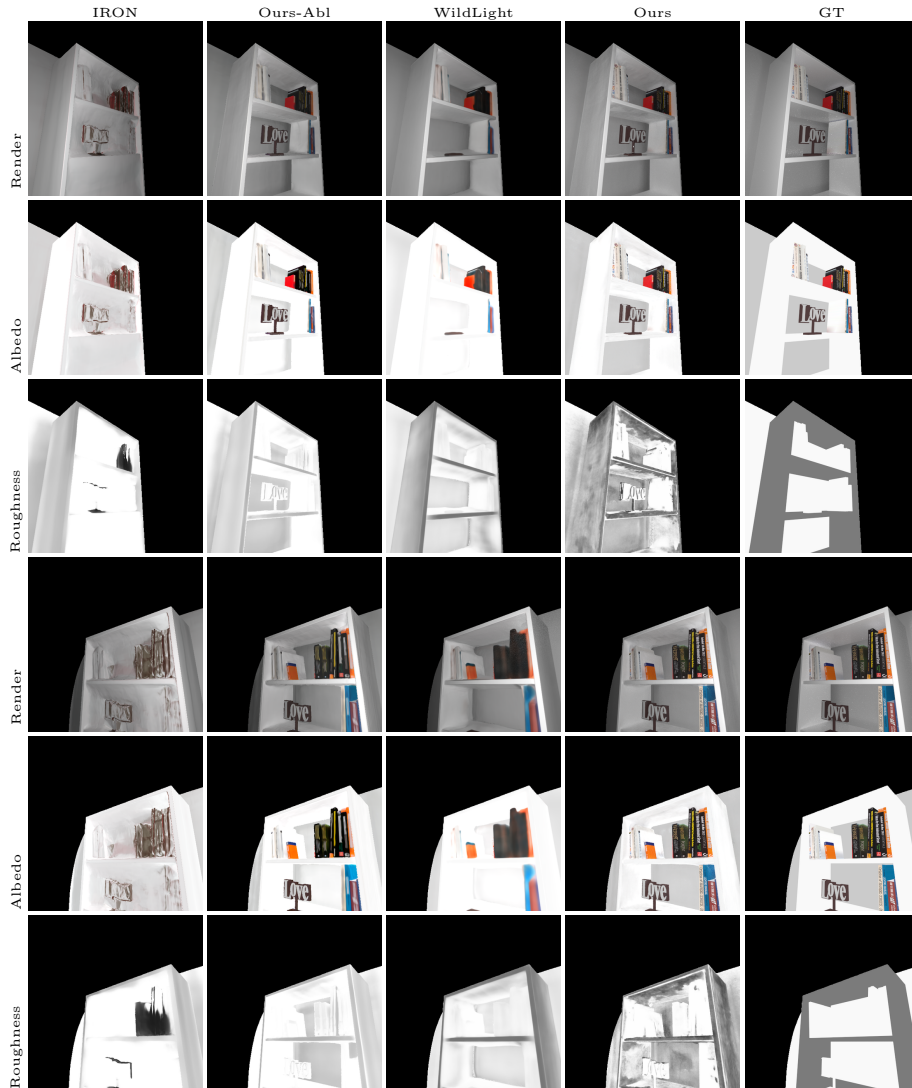
**Fig. 8:** Qualitative comparison on synthetic scene bedroom.

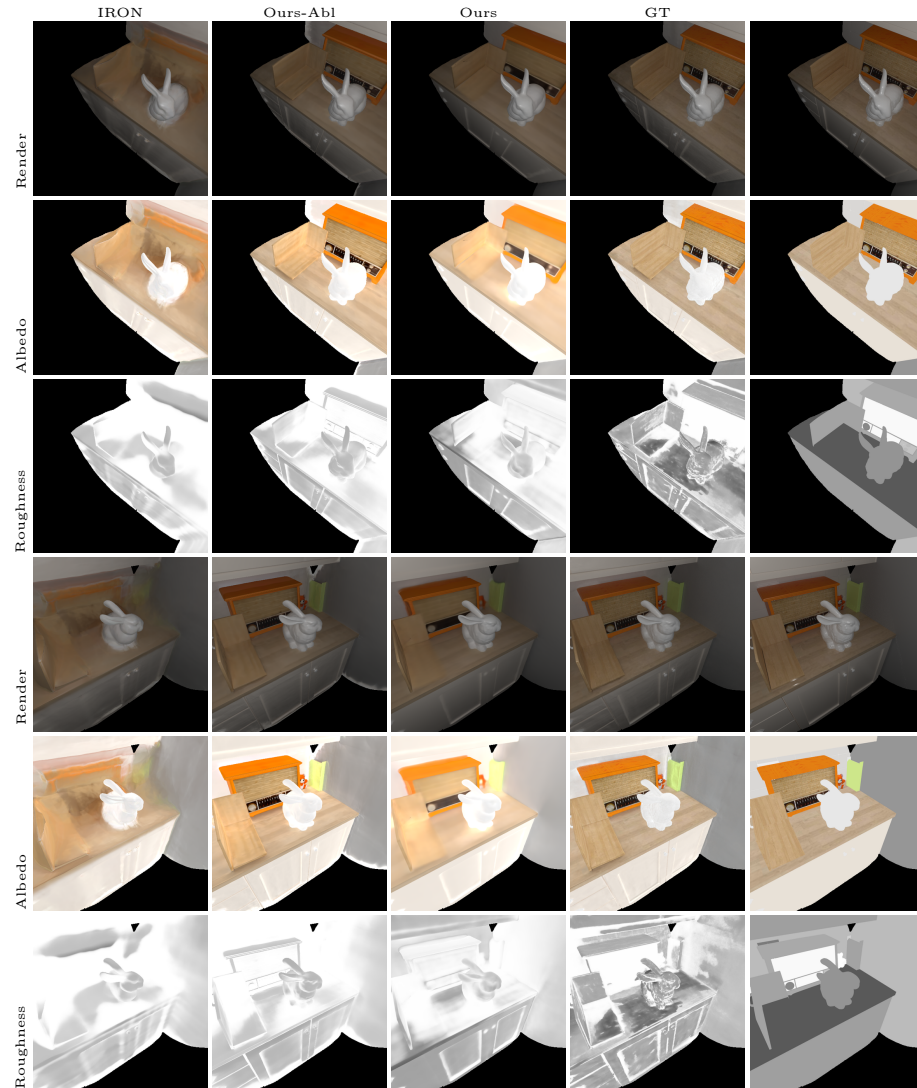**Fig. 9:** Qualitative comparison on synthetic scene coffee table.

**Fig. 10:** Qualitative comparison on synthetic scene shelf.

**Fig. 11:** Qualitative comparison on real scene shoe rack.
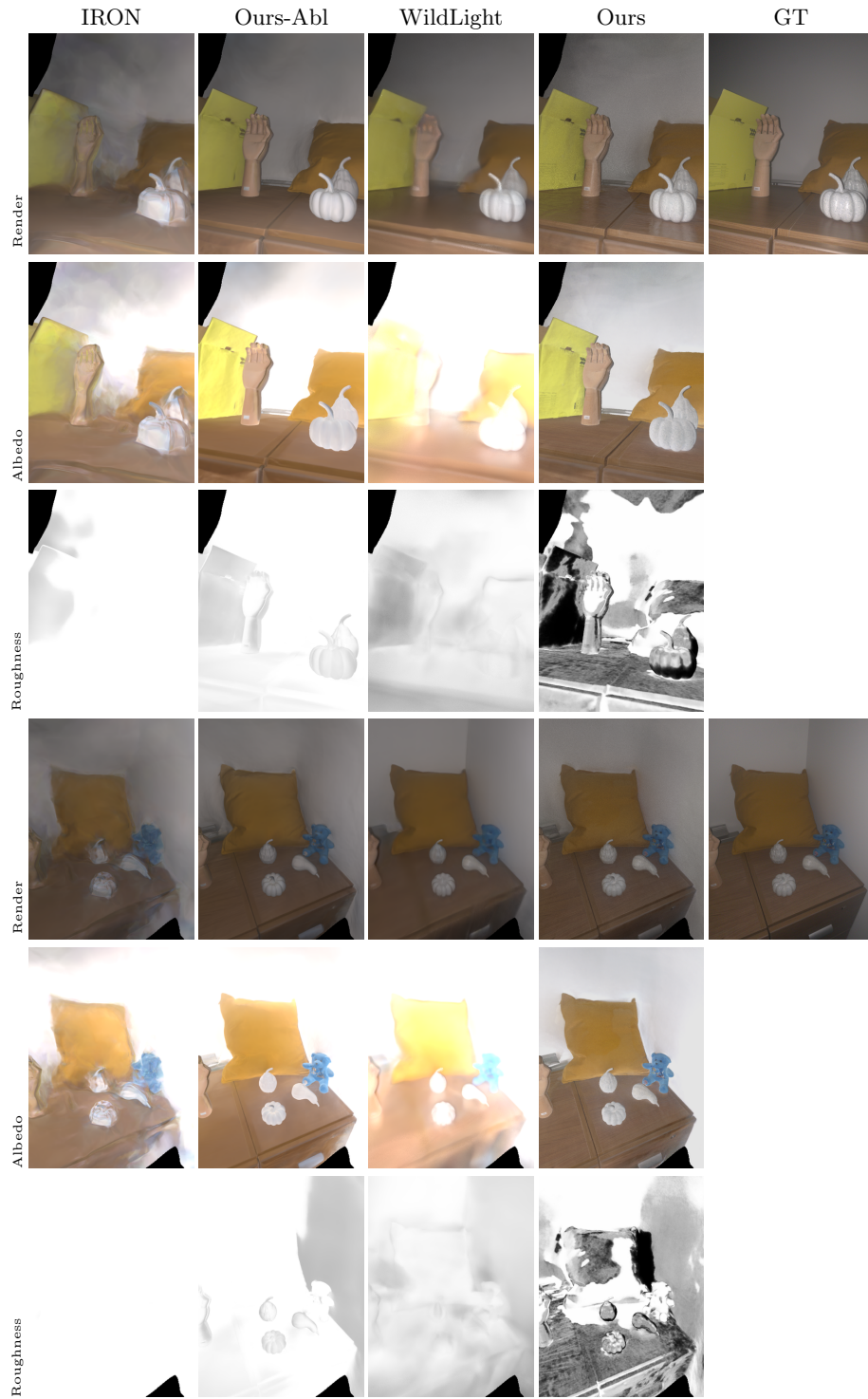
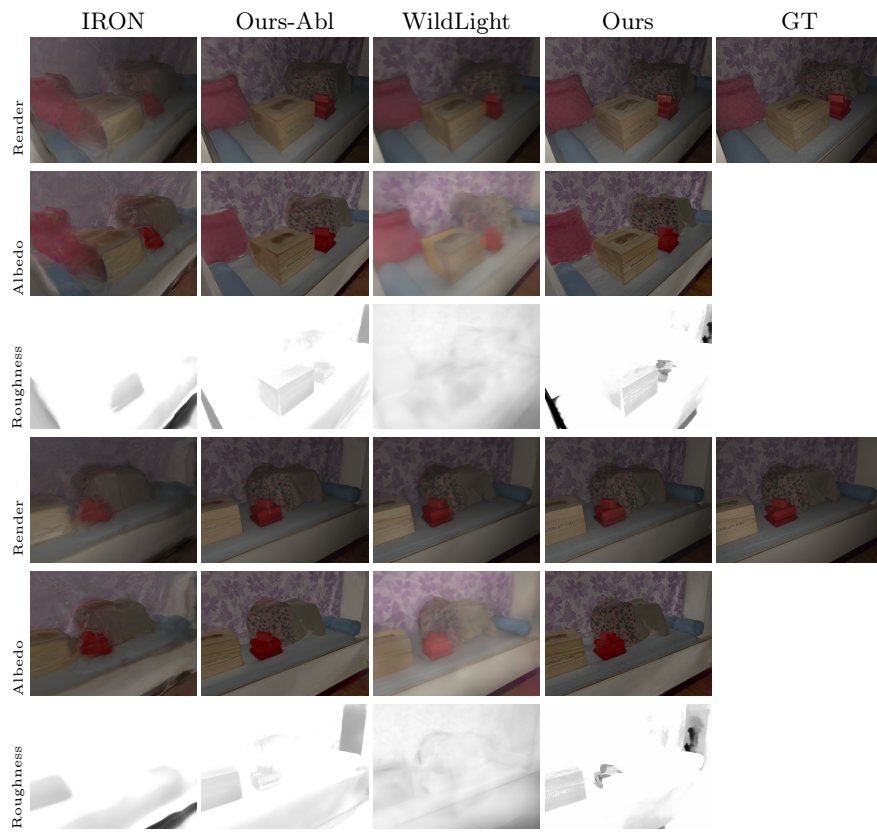**Fig. 12:** Qualitative comparison on real scene table.
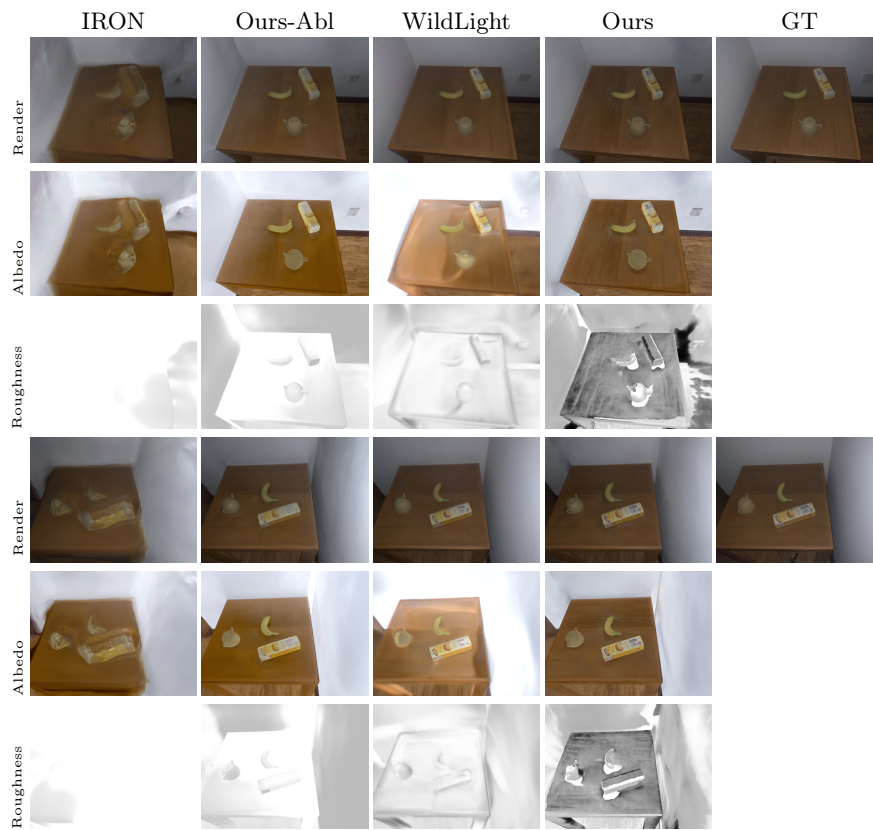
**Fig. 13:** Qualitative comparison on real scene window sill.

**Fig. 14:** Qualitative comparison on real scene coffee table.