# PR-ENDO: Physically Based Relightable Gaussian Splatting for Endoscopy

Joanna Kaleta[1,2,*], Weronika Smolak-Dyżewska[3,*], Dawid Malarz[3], Diego Dall'Alba[4],
Przemyslaw Korzeniowski[2], Przemysław Spurek[3]

[1]Warsaw University of Technology, [2]Sano Centre for Computational Medicine, [3]Jagiellonian University, [4]University of Verona
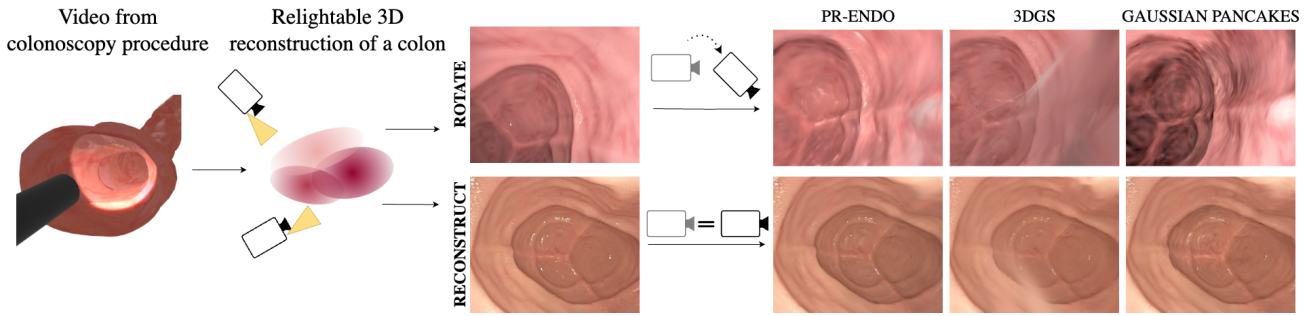j.kaleta@sanoscience.org, weronika.smolak@doctoral.uj.edu.pl

Figure 1. PR-ENDO adapts classical Gaussian Splatting for endoscopy environment. In contrast to 3DGS, we use the physical light model, which uses the connection between the camera and the light source. Moreover, we can reduce artifacts during novel view generation, which is a problem due to a limited camera trajectory during training. Finally, our model creates fewer artifacts than other state-of-the-art methods and responds to novel light conditions when we change the viewpoint or even modify geometry.

## Abstract

*Endoscopic procedures are crucial for colorectal cancer diagnosis, and three-dimensional reconstruction of the environment for real-time novel-view synthesis can significantly enhance diagnosis. We present PR-ENDO, a framework that leverages 3D Gaussian Splatting within a physically based, relightable model tailored for the complex acquisition conditions in endoscopy, such as restricted camera rotations and strong view-dependent illumination. By exploiting the connection between the camera and light source, our approach introduces a relighting model to capture the intricate interactions between light and tissue using physically based rendering and MLP. Existing methods often produce artifacts and inconsistencies under these conditions, which PR-ENDO overcomes by incorporating a specialized diffuse MLP that utilizes light angles and normal vectors, achieving stable reconstructions even with limited training camera rotations. We benchmarked our framework using a publicly available dataset and a newly introduced dataset with wider camera rotations. Our methods demonstrated superior image quality compared to baseline approaches.*

## 1. Introduction

Endoscopic images are crucial for the timely diagnosis of a wide range of pathological conditions, but their quality often suffers due to the challenging acquisition conditions [32].

Operating a flexible endoscope is complex due to the complex mapping between its tip and the non-ergonomic control interface, necessitating extensive training to master. Despite rigorous training, controlling the endoscope remains challenging and prone to human error, potentially causing patient discomfort and serious complications such as tissue perforation [32]. Furthermore, inadequate control of the endoscope can result in incomplete colon inspections, as pathological areas may fall outside the field of view of the endoscopic video frames [7, 27, 38].

These limitations have motivated advancements aimed at enhancing colonoscopy techniques [15]. Early approaches focused on automatically detecting suspicious pathological areas in endoscopic images [3]. Recent progress in Simultaneous Localization and Mapping (SLAM) within deformable environments has enabled tracking the endoscope

---

*These authors contributed equally to this work

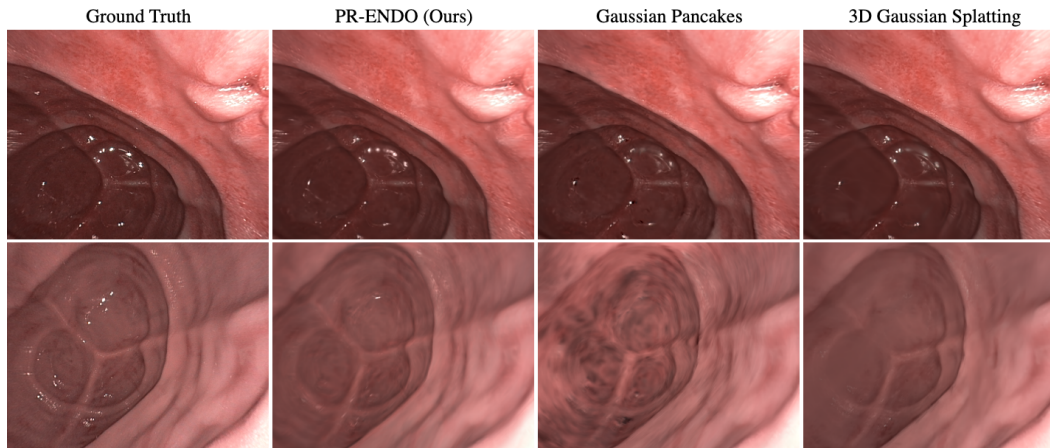| Ground Truth | PR-ENDO (Ours) | Gaussian Pancakes | 3D Gaussian Splatting |

Figure 2. **Qualitative comparison** of PR-ENDO and other methods such as Gaussian Pancakes and 3D Gaussian Splatting. Our method produces less artifacts and reconstructs geometry more accurately than other works. Please, zoom in for details.

tip's pose, thereby reconstructing a three-dimensional map of the colon [26–28].

Despite these developments, standard SLAM methods [13, 14, 31] struggle with achieving dense reconstructions, limiting the quality of three-dimensional renderings. High-quality, interactive 3D reconstructions of the colon can provide significant benefits, such as enabling endoscopists to review previously inspected areas in detail and synthesize new views for optimal diagnostic evaluation. Furthermore, accurate colon reconstructions could create realistic, interactive environments for training doctors, with models that can respond to natural body movements and simulate realistic relighting effects. In such settings, an animatable and responsive 3D colon model could be invaluable for both diagnostic and educational purposes, allowing future practitioners to practice in environments derived from real patient data.

Recent developments in Neural Radiance Fields (NeRF) [30] and Neural Implicit Surface (NeuS) [39] have been applied to colonoscopy for novel view synthesis [4, 33, 35]. However, these approaches are hindered by long training times, slow and non-interactive volumetric rendering with the absence of an explicit geometric model, limiting their clinical applicability. These challenges have been addressed by introducing 3D Gaussian Splatting (3DGS) [24], which models the environment using a set of Gaussians. While 3DGS offers fast training and real-time rendering, it relies on the quality of the explicit geometric model of the environment, typically reconstructed using SLAM techniques in endoscopic settings [27, 38]. The work in [7] presents a 3DGS approach that incorporates geometric regularization with a dense SLAM technique but does not account for the illumination model or the ability to apply movements to the reconstruction.

In this work, we address these limitations by propos-

ing PR-ENDO[1] - an approach that leverages the information it reconstructs (i.e., camera poses, depth maps, and point clouds of the environment) to introduce a physically based model capable of managing colon re-illumination. Our model is specifically designed to handle the complex conditions of endoscopic acquisitions, characterized by a camera with a limited range of motion, primarily moving forward with minimal changes in orientation. Additionally, we consider that the flexible endoscope integrates the camera and light source in a narrow space, effectively creating a single point of view for each specific illumination pattern.

Our physically based illumination framework produces a physically motivated 3D representation, enabling the generation of new views with larger camera rotation angles, improved generalization, and flexibility in adapting the 3DGS reconstruction to simulate anatomical or endoscopic movements. As noted in [12], implicit representation of complex effects like subsurface scattering is the most efficient approach. Therefore, our method extends the basic coefficient-based model by incorporating additional components to capture lighting effects that simpler models may overlook.

We experimentally evaluate the proposed framework using a set of metrics that assess the quality and speed of rendering new views, comparing it with the most relevant publicly available methods in the state of the art. The evaluation is based on two distinct datasets: one public and the other introduced in this work to provide a larger and more challenging set of camera rotations.

In summary, the contributions of this work are as follows:

- A relightning model based on light and tissue properties

---
[1]Code and dataset are available at: https://github.com/SanoScience/PR-ENDO.

in the specific context of endoscopy, capable of separating the effects of illumination on the tissue model.

- The introduction of diffuseMLP, which synthesizes new views while minimizing artifacts and visual inconsistencies.
- Demonstrating that the proposed model can handle wider camera angles and modify the scene model to replicate anatomical movements.

## 2. Related work

**Novel View Synthesis** Synthesizing new views is an essential step in modern 3D reconstruction pipelines. Many supervised approaches have been proposed, exploiting different efficient representations, in particular meshes [11, 18], voxels [19, 21] and hybrid methods [8, 20]. NeRF-based approaches demonstrate that implicit radiance fields can efficiently learn scene representations and synthesize high-quality new views [2, 30, 41]. NeRF learns a continuous function that implicitly represents the 3D scene trained from 2D images and matched camera poses. NeRF is an inefficient approach, requiring long training times and unable to achieve real-time rendering of new views.

These limitations have been recently overcome by 3DGS approaches, which propose an efficient representation of the scene via a set of 3D Gaussians [24]. 3DGS is superior in training and rendering efficiency and provides a geometric representation of the scene while NeRF relies on a volumetric representation. In our work, we, therefore, adopt a 3DGS approach, since its characteristics make it ideal for synthesizing new views in real-time.

**SLAM & Novel View Synthesis in Endoscopy** NeRF and 3DGS approaches require knowledge of the camera poses and a three-dimensional reconstruction of the scene. This information is not commonly available in the endoscopic context, and must be reconstructed with specific approaches, usually exploiting a SLAM approach [13, 14, 26, 31]. Unfortunately, deformations of the environment, poorly defined textures, and constrained and limited camera movements are very complex conditions for the use of standard SLAM approaches [28].

In this context, approaches designed explicitly for the endoscopic context have been proposed. In particular, RNNSLAM proposes the extension of the standard SLAM approach with two recurrent neural networks for pose estimation and depth maps to improve the quality of the environment reconstruction and also the tracking of camera movements [27]. Similarly, ColonNeRF proposes a specifically developed pipeline able to refine the poses and the environment reconstruction to support the subsequent NeRF modeling [35]. EndoGSLAM has recently proposed an efficient approach that integrates a SLAM approach with a 3DGS reconstruction [38]. In this work, we adopt EndoGSLAM to obtain the camera pose and the 3DGS recon-

struction of the environment, which we extend with respect to the simplified model adopted in the original work, in particular to model the effect of view-dependent illumination. However, our method is universal for other types of initialization and can be used with different SLAM modules.

The reconstruction approaches described, both NeRF and 3DGS, suffer in the endoscopic context from the limited camera movements and from the fact that the light source moves together with the camera generating a highly view-dependent illumination, which makes it difficult to synthesize new high-quality views, i.e. free from inconsistencies and artifacts. To improve the quality of the synthesis, Gaussian Pancakes extends the RNNSLAM approach to extract a 3DGS reconstruction on which geometric and depth regularization terms are introduced in order to improve the geometric quality of the reconstruction and consequently that of the images synthesized in correspondence with the new views [7]. Gaussian Pancakes does not consider the illumination model, which is instead done by REIM-NeRF, which trains a light source location-conditioned NeRF model [33]. This approach suffers from the inherent inefficiencies of the NeRF model and considers a lighting model that neglects the environmental components, which are instead considered in the proposed PR-ENDO approach.

**3DGS Relighting** There has been a lot of interest recently in tackling the problem of accurately modeling illumination within 3DGS. In its original formulation 3DGS uses a single spherical harmonic function to model both texture and color information, reducing the generalization capabilities to new views. Various methods have been proposed to indirectly include light effects [29] or decouple texture information from illumination, exploiting forward [16, 22, 23, 25, 36] or deferred shading [10, 40] to achieve effective Gaussian relighting.

GS-IR uses baking-based volumes to store occlusion and indirect illumination [25]. Relightable 3D Gaussian achieves this by employing point-based ray tracing to bake occlusion and parameterizing indirect lighting with additional spherical harmonics for each Gaussian [16]. GI-GS utilizes deferred shading to enable efficient path tracing-based occlusion and indirect illumination for rendering and relighting [10], a feature not supported by previous methods. GaussianShader introduces a shading function to 3DGS and a normal estimation framework capable of handling reflective surfaces [22], enhancing the approach proposed in GIR [36]. Together with DeferredGS [40], GaussianShader extends 3DGS to the inverse rendering task by associating BRDFs with Gaussians. GS3 proposes a triple splatting process that does not require robust normal estimation and accurately models BRDF properties using a mixture of angular Gaussians [5].

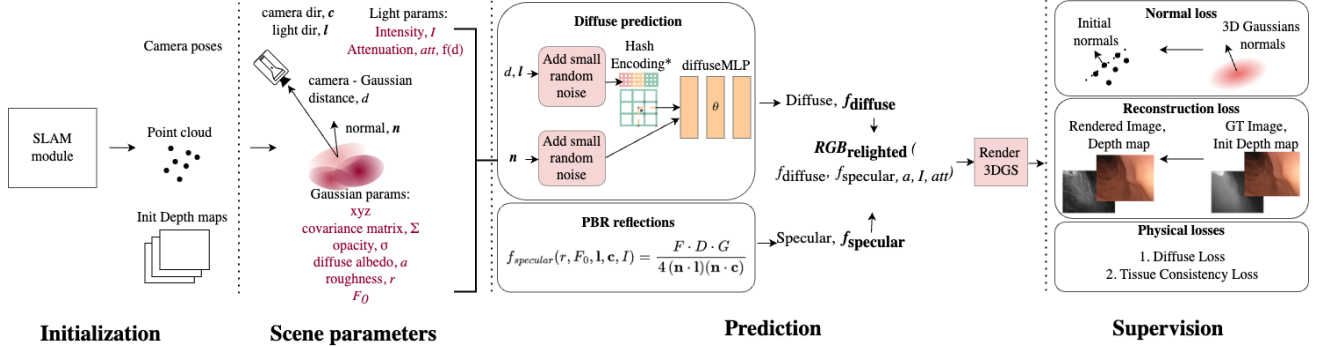PR-ENDO distinguishes itself from previous methods

Figure 3. **Overview of PR-ENDO** pipeline, which includes: (1) **Initialization**: SLAM module generates camera poses, point cloud and optionally depth maps. (2) **Scene Parameters**: Each Gaussian $i$ is defined by position, covariance $\Sigma_i$, opacity, albedo $a_i$, roughness $r_i$, and reflectivity $F0_i$. Light parameters include direction $l$, intensity $I$ and attenuation function $att(d_i)$. (3) **Prediction**: $diffuseMLP$, predicts the diffuse component, while the PBR model calculates specular reflection. The final color $\text{RGB}_{\text{relighted}}$ combines diffuse and specular terms. (4) **Supervision**: Normal, reconstruction, and physical losses are applied during training.

by employing an approach that effectively generalizes the scene, even under the specific and challenging conditions of endoscopic settings, offering compelling medical applications for such reconstructions.

## 3. Method

This section describes our model. First, we present the preliminaries on vanilla Gaussian Splatting algorithm [24] and physically based rendering. Then, we present PR-ENDO, which combines Gaussian Splatting with a relighting model, dedicated for endoscopy data.

**Preliminaries on 3D Gaussian Splatting**  3D Gaussian Splatting (3DGS) is a technique that represents scenes using explicit 3D Gaussian distributions as rendering primitives, which facilitates efficient optimization and GPU-accelerated rendering [24]. Each $i$-th Gaussian is defined by a mean position $xyz_i$ and a covariance matrix $\Sigma_i$, representing the spatial distribution in 3D. Additionally, each Gaussian is assigned an opacity $\sigma_i$ and a view-dependent color $c_i$. The covariance matrix $\Sigma_i$ is derived from a scaling matrix $S_i$ and a rotation matrix $R_i$ through the relationship $\Sigma_i = R_i S_i S_i^T R_i^T$.

The rendering process involves projecting 3D Gaussians onto a 2D image plane by calculating both the 2D positions and the transformed covariance matrix $\Sigma'$ in screen space. The covariance matrix in camera coordinates is computed as:

$$\Sigma'_i = JW\Sigma_i W^T J^T, \tag{1}$$

where $W$ is the viewing transformation matrix, and $J$ is the Jacobian matrix of the affine approximation of $W$.

Rendering is performed through a differentiable point-based alpha blending technique, where the color of each

pixel $C$ is accumulated along the ray by blending ordered 2D Gaussians from front to back:

$$C = \sum_{i \in N} T_i \alpha_i c_i, \quad T_i = \prod_{j=1}^{i-1}(1 - \alpha_j). \tag{2}$$

Here, $\alpha$ is computed by combining the opacity $\sigma$ with contributions from the 2D covariance matrix $\Sigma'$ and the pixel location in image space. The attributes of each Gaussian, including color $c_i$, are encoded with spherical harmonics (SH) to capture view-dependent effects [34].
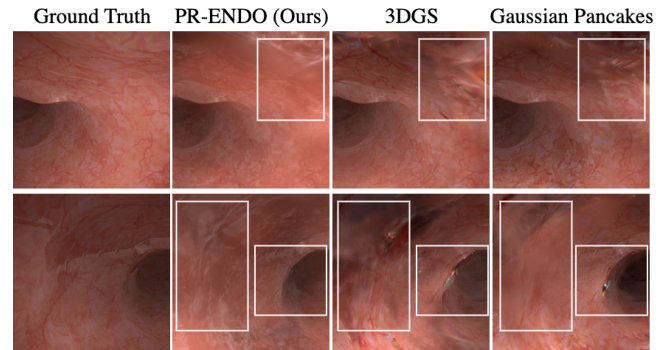


Figure 4. **Qualitative comparison** of PR-ENDO and other methods such as Gaussian Pancakes and 3D Gaussian Splatting. A comparison was made with our dataset ColonRotate, which has ground truth renders for various camera rotation angles.

**Preliminaries on Physically-Based Rendering**
Physically-Based Rendering (PBR) [1, 9] simulates light-material interactions based on physical principles, enabling realistic image synthesis. Material properties in PBR govern how surfaces interact with light, typically modeled through diffuse and specular reflections.

The **Bidirectional Reflectance Distribution Function (BRDF)** is defined as the sum of the diffuse and specular reflection terms:

$$\text{BRDF} = k_d f_{\text{diffuse}} + f_{\text{specular}} \tag{3}$$

where $k_d$ is the diffuse reflection coefficient, and $f_{\text{diffuse}}$ is the diffuse reflection term, $f_{\text{specular}}$ is the specular reflection term. The latter term includes the specular reflection coefficient $k_s$, which is related to $k_d$ by $k_d = 1 - k_s$.

*Diffuse Component* The diffuse term represents the scattering of light uniformly in all directions on a rough surface. This is typically modeled using the *Lambertian reflectance model*, which defines the diffuse reflection as:

$$f_{\text{diffuse}} = a \cdot (\mathbf{n} \cdot \mathbf{l}) \tag{4}$$

where: - $\mathbf{n}$ is the surface normal, - $\mathbf{l}$ is the light direction, - $a$ is the diffuse albedo (reflectivity of the material in diffuse reflections).

*Specular Component* Specular reflection models mirror-like reflections influenced by view angle and surface roughness. In PBR, the Cook-Torrance model computes the specular term as:

$$f_{\text{specular}} = \frac{D\,G\,F}{4\,(\mathbf{l} \cdot \mathbf{n})(\mathbf{c} \cdot \mathbf{n})} \tag{5}$$

where: $\mathbf{c}$ is the viewing (camera) direction, $D$ is the microfacet distribution function, describing how microfacets are oriented on the surface, $G$ is the geometric attenuation term, accounting for shadowing and masking between microfacets, $F$ is the Fresnel term, representing the angle-dependent reflectivity of the surface.

*Fresnel Effect* The Fresnel effect describes how reflectivity changes with the viewing angle, especially pronounced at grazing angles. It is computed as:

$$F = F_0 + (1 - F_0)(1 - (\mathbf{l} \cdot \mathbf{H}))^5 \tag{6}$$

where $F_0$ is the base reflectivity at normal incidence, $\mathbf{H}$ is the half-vector between the light direction $\mathbf{l}$ and view direction $\mathbf{c}$. $F$ relates to the specular reflection coefficient $k_s$ as follows $k_s = F$.

*Microfacet Distribution* The microfacet distribution function $D$ models surface roughness, commonly using the Trowbridge-Reitz GGX distribution:

$$D = \frac{\alpha^2}{\pi\,((\mathbf{n} \cdot \mathbf{H})^2(\alpha^2 - 1) + 1)^2} \tag{7}$$

where $\alpha = \text{roughness}^2$, representing the surface roughness. The GGX distribution captures both sharp and rough reflections, making it versatile for a range of materials. Here, $\mathbf{H}$ is the half-vector, and $\mathbf{n} \cdot \mathbf{H}$ represents the angle between the surface normal and the half-vector.

*Geometric Attenuation* The *geometric attenuation term* $G$ accounts for shadowing and masking effects on microfacets, which reduces the amount of light reflected due to self-shadowing on rough surfaces. The *Schlick-Beckmann approximation* simplifies this calculation:

$$G = \frac{\mathbf{n} \cdot \mathbf{l}}{(\mathbf{n} \cdot \mathbf{l})(1 - k) + k} \cdot \frac{\mathbf{n} \cdot \mathbf{c}}{(\mathbf{n} \cdot \mathbf{c})(1 - k) + k} \tag{8}$$

where $k = \frac{\alpha}{2}$. This approximation efficiently models the effects of surface roughness on reflection attenuation, especially when combined with the GGX distribution.

**PR-ENDO** Our method, PR-ENDO, introduces a 3DGS physically-based relighting model tailored for endoscopic scenes, capturing both diffuse and specular interactions between light and surfaces.

We show the overview of our method in Fig. 3. PR-ENDO separates the scene representation into a light model and set of relightable Gaussians with additional parameters: base color $a$, roughness $r$ and base reflectivity $F0$. We use a dedicated light model, which assumes that the light source and camera are co-located.

Endoscopy videos are characterized by limited training views due to specific and challenging training camera trajectory. Therefore, generalization to novel viewpoints (such as previously unseen rotations) leads to artifacts. To tackle this problem, instead of simple coefficient-based diffuse relightning model, we introduce diffuseMLP $\theta$. The MLP is fed with Gaussian normal vector $\mathbf{n_i}$ and light position parameters $(\mathbf{l_i}, d_i)$ relative to the $i$-th Gaussian. We add small
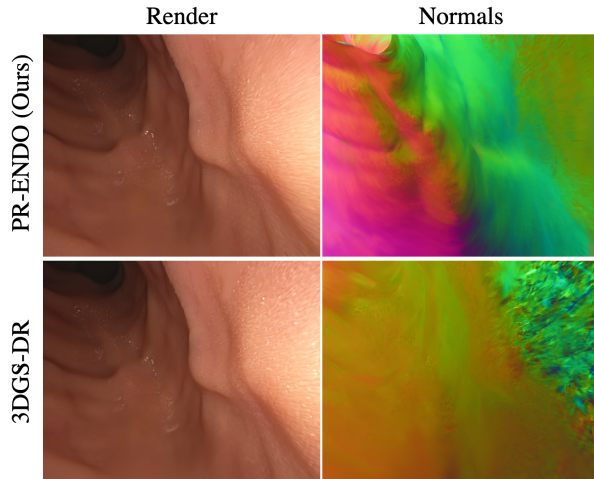


Figure 5. **Normal comparison between PR-ENDO and 3DGS-DR [40].** Although both methods produce plausible renderings, our method has a superior normals reconstruction.

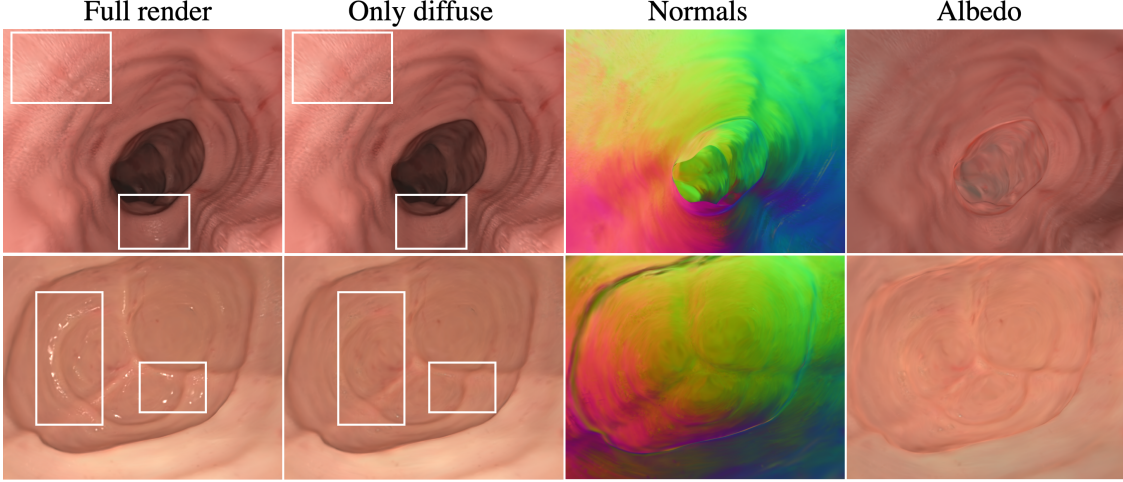| Full render | Only diffuse | Normals | Albedo |

Figure 6. **Decomposition.** Thanks to our physically motivated model, we are able to decompose the render into diffuse, normals and albedo views.

random noise to the inputs to improve robustness in different lighting and viewpoint variations, encouraging consistent diffuse predictions under similar conditions. The output returns diffuse part of the $i$-th Gaussian color:

$$f_{\text{diffuse},i} = \theta(\mathbf{l_i}, d_i, \mathbf{n_i}) \tag{9}$$

The result of diffuseMLP not only allows for better generalization to novel rotations but can also incorporate additional lightning effects not included in the simple coefficient-based diffuse model.

The specular component $f_{specular,i}$ for the $i$-th Gaussian is computed using the PBR model as in Eq. (5), based on $r_i$ and $F0_i$. In our method we clamp $F_0$ values to 0.03, which represent realistic values of the colon refractive index [17].

The final relighted color, $RGB_{\text{relight},i}$ is obtained by combining the diffuse and specular components, scaled by the base color $a_i$, light intensity $I$ and a distance-based attenuation factor $att(d_i)$:

$$RGB_{\text{relight},i} = I \cdot att(d_i) \cdot$$
$$(f_{\text{diffuse},i} \cdot (1 - F_i) \cdot a_i + f_{\text{specular},i}) \tag{10}$$

In our data-driven PBR model, all material, as well as light properties, are optimized via gradient descent.

**Optimization** Our model's total loss function $\mathcal{L}$ combines image reconstruction, depth, geometric, diffuse, and tissue consistency losses, designed to optimize visual fidelity, depth alignment, and physical realism.

*Image Loss*: An $L_1$ loss, $\mathcal{L}_{\text{Image}}$, measures the pixel-wise difference between rendered and ground truth images.
*Structural Similarity Loss*: A structural similarity index (SSIM) loss, $\mathcal{L}_{\text{D-SSIM}}$, compares the rendered and ground truth images, balancing pixel accuracy with perceived similarity.

*Depth Loss*: To address the minimal depth diversity in endoscopic images, we apply a depth loss based on the Huber loss function following [7]. For depth maps $\{D_i\}_{i=1}^{M}$, where $M$ is the number of viewpoints, the depth loss $\mathcal{L}_D$ is defined as:

$$\mathcal{L}_D(i) = \begin{cases} 0.5\Delta D_i^2, & \text{if } |\Delta D_i| < \delta \\ \delta\left(|\Delta D_i| - 0.5\delta\right), & \text{otherwise} \end{cases} \tag{11}$$

where $\Delta D_i = |D_i - \hat{D}_i|$ denotes the difference between true and rendered depth.

*Geometric Regularization*: To enforce structural consistency, a geometric loss $\mathcal{L}_G$ [7] is applied, aligning Gaussians with the surface's principal curvature. The cosine similarity-based geometric loss is:

$$\mathcal{L}_G = 1 - \frac{|A \cdot B|}{\|A\|\|B\|} \tag{12}$$

where $A$ and $B$ are the normal vectors of nearest neighbors on the surface.

The overall incorporated loss function $\mathcal{L}_{\text{inc}}$ is a weighted sum:

$$\mathcal{L}_{\text{inc}} = (1 - \lambda_1)\mathcal{L}_{\text{Image}} + \lambda_1\mathcal{L}_{\text{D-SSIM}} + \lambda_2\mathcal{L}_D + \lambda_3\mathcal{L}_G. \tag{13}$$

Additionally, we designed physically-based losses.
*Diffuse Loss*: To ensure consistency between the diffuse lighting predicted by the diffuseMLP $f_{\text{diffuse,MLP},i}$ and the coefficients $f_{\text{diffuse,coeff},i}$, we apply a mean squared error loss summed over all ($N$) Gaussians:

$$\mathcal{L}_{\text{Diffuse}} = \sum_{i=1}^{N} \left(f_{\text{diffuse,MLP},i} - f_{\text{diffuse,coeff},i}\right)^2 \tag{14}$$

Table 1. Performance Comparison on C3VD and RotateColon Datasets. Our method achieves superior performance over existing SOTA in both endoscopy and relighting tasks. *While Gaussian Shader produces reasonable renderings in terms of metrics, it significantly fails to accurately reconstruct correct scene geometry.

| Model | C3VD | | | RotateColon | | |
|---|---|---|---|---|---|---|
| | PSNR ↑ | SSIM ↑ | LPIPS ↓ | PSNR ↑ | SSIM ↑ | LPIPS ↓ |
| 3DGS [24] | $33.90 \pm 1.97$ | $0.89 \pm 0.03$ | $0.28 \pm 0.07$ | $20.29 \pm 6.80$ | $0.82 \pm 0.11$ | $\mathbf{0.25 \pm 0.11}$ |
| EndoGSLAM [38] | $22.16 \pm 2.66$ | $0.77 \pm 0.08$ | $\mathbf{0.22 \pm 0.05}$ | - | - | - |
| GaussianPancakes [7] | $33.12 \pm 2.89$ | $0.89 \pm 0.03$ | $0.30 \pm 0.05$ | $20.10 \pm 4.92$ | $0.88 \pm 0.05$ | $0.27 \pm 0.07$ |
| GaussianShader* [22] | $29.82 \pm 2.44$ | $0.86 \pm 0.03$ | $0.40 \pm 0.04$ | $21.25 \pm 3.79$ | $0.87 \pm 0.05$ | $0.38 \pm 0.11$ |
| 3DGS-DR [40] | $33.77 \pm 1.83$ | $0.89 \pm 0.09$ | $0.31 \pm 0.04$ | $21.49 \pm 5.50$ | $0.89 \pm 0.06$ | $0.28 \pm 0.08$ |
| PR-ENDO (ours) | $\mathbf{34.00 \pm 2.16}$ | $\mathbf{0.90 \pm 0.03}$ | $0.29 \pm 0.05$ | $\mathbf{21.90 \pm 5.32}$ | $\mathbf{0.87 \pm 0.05}$ | $0.28 \pm 0.06$ |

*Tissue Consistency Loss*: To enforce consistency in tissue characteristics, we constrain the albedo $a$, roughness $r$ and base reflectance $F0$ to remain close to their respective means with each component measured as an MSE and summed over all ($N$) Gaussians:

$$\mathcal{L}_{\text{Tissue}} = \sum_{i=1}^{N}(a_i - a_{\text{mean}})^2$$
$$+ \sum_{i=1}^{N}(r_i - r_{\text{mean}})^2 + \sum_{i=1}^{N}(F0_i - F0_{\text{mean}})^2. \quad (15)$$

Each part of the equation can be thought as a separate loss for each component respectively ($\mathcal{L}_a$, $\mathcal{L}_r$, $\mathcal{L}_{F0}$).

The final loss function $\mathcal{L}$ is the sum of all loss terms, balanced with appropriate weights:

$$\mathcal{L} = \mathcal{L}_{\text{inc}} + \mathcal{L}_{\text{Diffuse}} + \mathcal{L}_{\text{Tissue}} \quad (16)$$

## 4. Experiments

**Datasets** We evaluate our proposed method on the Colonoscopy 3D Video Dataset (C3VD) [6], preprocessed by EndoGSLAM [38] with their camera pose estimations and initial point cloud. Additionally, this dataset provides ground-truth RGB images and depth maps. The dataset images are pre-undistorted and have a resolution of 675×540.

We expect our method to work exceptionally well with SLAM approaches specifically developed for colon geometry reconstruction such as RNNSLAM [27], utilized in [7]. The model for [27] was not open sourced at the time of submission. All of the comparisons on C3VD made in this paper are done using the same initialization from EndoGSLAM.

In addition to C3VD, we introduce a second dataset, RotateColon, developed specifically for evaluating novel view synthesis under extended rotations using our in-house simulator. Unlike traditional evaluation that samples every n-th frame from the camera trajectory, this dataset includes intense rotations not observed during the training sequence,

enabling a more rigorous assessment of generalization. We do not employ a SLAM-based approach for this dataset and use ground truth camera poses. The resolution of RotateColon images used for training is 640x640.
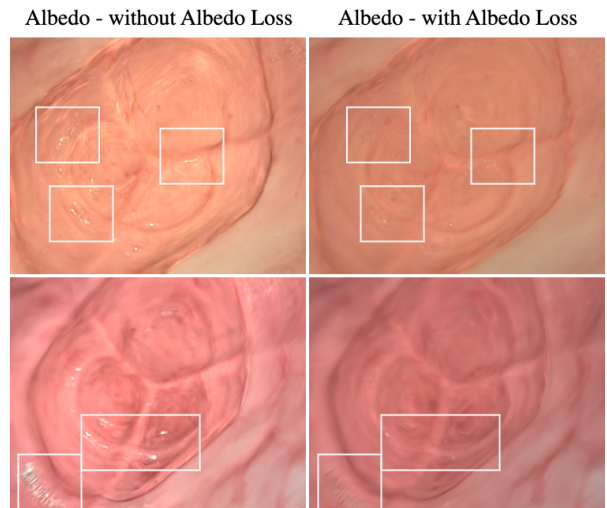


Figure 7. **Ablation study on albedo loss.** Albedo loss prevents reflections and lighting effects from becoming part of the base color.

We measure rendering performance using widely approved set of metrics: peak signal-to-noise ratio (PSNR), structural similarity index (SSIM), and learned perceptual image patch similarity (LPIPS).

**Training details** We train PR-ENDO on C3VD for 40000 iterations using Adam optimizer for Gaussian's attributes and diffuseMLP. Our hyperparameters are as follows: diffuseMLP lr 0.002, $L_G$ weight 0.1, $L_D$ weight 0.1, $L_a$ weight 1e-5, $F_0$ and $L_r$ weights 1e-6. $L_{Diffuse}$ weight starts at 1e-6 and decreases linearly during training. Learning rates for roughness and $F_0$ are at 0.003. For RotateColon dataset, we set $L_a$ weight at 1e-4.

We run other baseline methods using our EndoGSLAM

initialization. We use configurations suggested in original repositories, with exception for GaussianPancakes. Given that GaussianPancakes relies heavily on the initialization method, we extended training to 30k iterations and densification to 7k iterations for a fair comparison, while keeping other parameters at default. This 30k threshold was set experimentally, as results were observed to stabilize by this point.

**Reconstruction** We compare reconstruction capabilities of PR-ENDO and baselines quantitatively in Tab. 1. On both datasets, the proposed method outperforms the baseline methods, including the vanilla 3DGS [24], EndoGSLAM [38], GaussianPankake [7] (developed specifically for endoscopy applications), and two representative methods from 3DGS relighting works: GaussianShader [22] and 3DGS-DR [40].

Fig. 2 showcases our superior reconstruction on C3VD dataset. Using the same initialization, we achieve fewer artifacts on the surface and more accurate reflections. Fig. 4 shows how our method behaves during challenging rotation angles and is more efficient at generalization.

Table 2. Study on hash-grid (HG) influence in the considered datasets.

| Dataset | HG | PSNR $\uparrow$ | SSIM $\uparrow$ | LPIPS $\downarrow$ |
|---|---|---|---|---|
| C3VD | ✗ | $33.45 \pm 2.40$ | $0.94 \pm 0.02$ | $0.32 \pm 0.05$ |
| | ✓ | $34.00 \pm 2.16$ | $0.90 \pm 0.03$ | $0.29 \pm 0.05$ |
| Rotate Colon | ✗ | $21.90 \pm 5.32$ | $0.87 \pm 0.05$ | $0.28 \pm 0.06$ |
| | ✓ | $20.45 \pm 4.87$ | $0.89 \pm 0.04$ | $0.27 \pm 0.06$ |

**Decomposition and relightning** Our method enables the decomposition of diffuse, albedo, and normals as seen in Fig. 6. Compared with other state-of-the-art methods in decomposition, our model achieves more plausible normals (Fig. 5). Additionally, thanks to the fact that PR-ENDO responds to light, it enables an efficient way to change light properties (such as position) in novel views. In Apppendix we demonstrate that we can separate light from camera position to confirm the effectiveness of our method.

**Novel view and anatomy movement** Thanks to PR-ENDO superior abilities of generalization (see Fig. 4, it is suitable for employing simulations and adding simple colon movement animations. By simply transforming the colon reconstructed by our method into GaMeS representation [37], we evaluate the quality of renderings and reduced artifacts during body movement simulations. See supplemented videos for details.

Table 3. Ablation study on C3VD set. We run ablation study for a base model without HashGrid

| PR-ENDO without | PSNR $\uparrow$ | SSIM $\uparrow$ | LPIPS $\downarrow$ |
|---|---|---|---|
| - | $33.45 \pm 2.40$ | $0.94 \pm 0.02$ | $0.32 \pm 0.05$ |
| diffuseMLP | $32.48 \pm 2.77$ | $0.89 \pm 0.03$ | $0.31 \pm 0.04$ |
| $L_{\text{Tissue}}$ | Failed | Failed | Failed |
| $L_{\text{Diffuse}}$ | $32.98 \pm 2.57$ | $0.89 \pm 0.03$ | $0.31 \pm 0.04$ |

**HashGrid (HG) component** Examination of our architecture shows (see Tab. 2) that inclusion of HG boosts model's ability to correctly reconstruct. However, it hinders PR-ENDO's ability to generalize during challenging rotations. Due to that, we decide that inclusion of HG in our architecture is optional and should be tweaked depending on user's preferences.

**Ablation** We test effectiveness our designed components and losses quantitatively in Tab. 3 on the whole C3VD dataset. During this study, we disable HashGrid for the basic version of our method. PR-ENDO without $L_{\text{Tissue}}$ suffers from the vanishing gradient problem, making it difficult to train effectively. During our ablation study, we also show that albedo loss is crucial for proper separation between components and prevents the albedo from taking over the representation of the reflections and lightning effects (see Fig. 7).

**Limitations** While our model outperforms existing methods, it has some limitations. It still struggles with extreme rotations and when the light source is positioned significantly far from the camera. Training time is around 2 times longer than for vanilla 3DGS. For more textured tissues a more complex constraint for $\mathcal{L}_a$ may be required.

## 5. Conclusion

PR-ENDO is a model specifically designed for reconstructing endoscopic videos. It leverages the unique property of endoscopic videos, which typically feature a single light source whose position closely aligns with the camera location. This allows us to model reflections in endoscopic videos accurately through physical principles. Such separation enhances PR-ENDO's ability to generalize and produce realistic novel views, even at challenging angles. Qualitative and quantitative comparisons demonstrate that PR-ENDO achieves superior results compared to previous state-of-the-art methods.

## 6. Acknowledgments

## References

[1] 09 - reflection models. In *Physically Based Rendering*, pages 409–459. Morgan Kaufmann, Burlington, 2004. 4

[2] Jonathan T Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and Pratul P Srinivasan. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 5855–5864, 2021. 3

[3] Ishita Barua, Daniela Guerrero Vinsard, Henriette C Jodal, Magnus Løberg, Mette Kalager, Øyvind Holme, Masashi Misawa, Michael Bretthauer, and Yuichi Mori. Artificial intelligence for polyp detection during colonoscopy: a systematic review and meta-analysis. *Endoscopy*, 53(03):277–284, 2021. 1

[4] Víctor M Batlle, José MM Montiel, Pascal Fua, and Juan D Tardós. Lightneus: Neural surface reconstruction in endoscopy using illumination decline. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 502–512. Springer, 2023. 2

[5] Zoubin Bi, Yixin Zeng, Chong Zeng, Fan Pei, Xiang Feng, Kun Zhou, and Hongzhi Wu. Gsˆ3: Efficient relighting with triple gaussian splatting. *arXiv preprint arXiv:2410.11419*, 2024. 3

[6] Taylor L Bobrow, Mayank Golhar, Rohan Vijayan, Venkata S Akshintala, Juan R Garcia, and Nicholas J Durr. Colonoscopy 3d video dataset with paired depth from 2d-3d registration. *Medical Image Analysis*, page 102956, 2023. 7

[7] Sierra Bonilla, Shuai Zhang, Dimitrios Psychogyios, Danail Stoyanov, Francisco Vasconcelos, and Sophia Bano. Gaussian pancakes: Geometrically-regularized 3d gaussian splatting for realistic endoscopic reconstruction. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 274–283. Springer, 2024. 1, 2, 3, 6, 7, 8

[8] Michael Broxton, John Flynn, Ryan Overbeck, Daniel Erickson, Peter Hedman, Matthew Duvall, Jason Dourgarian, Jay Busch, Matt Whalen, and Paul Debevec. Immersive light field video with a layered mesh representation. *ACM Transactions on Graphics (TOG)*, 39(4):86–1, 2020. 3

[9] Brent Burley and Walt Disney Animation Studios. Physically-based shading at disney. In *Acm Siggraph*, pages 1–7. vol. 2012, 2012. 4

[10] Hongze Chen, Zehong Lin, and Jun Zhang. Gigs: Global illumination decomposition on gaussian splatting for inverse rendering. *arXiv preprint arXiv:2410.02619*, 2024. 3

[11] Alvaro Collet, Ming Chuang, Pat Sweeney, Don Gillett, Dennis Evseev, David Calabrese, Hugues Hoppe, Adam Kirk, and Steve Sullivan. High-quality streamable free-viewpoint video. *ACM Transactions on Graphics (ToG)*, 34(4):1–13, 2015. 3

[12] Jan-Niklas Dihlmann, Arjun Majumdar, Andreas Engelhardt, Raphael Braun, and Hendrik P.A. Lensch. Subsurface scattering for gaussian splatting, 2024. 2

[13] Jakob Engel, Thomas Schöps, and Daniel Cremers. Lsd-slam: Large-scale direct monocular slam. In *European conference on computer vision*, pages 834–849. Springer, 2014. 2, 3

[14] Jakob Engel, Vladlen Koltun, and Daniel Cremers. Direct sparse odometry. *IEEE transactions on pattern analysis and machine intelligence*, 40(3):611–625, 2017. 2, 3

[15] Zuoming Fu, Ziyi Jin, Chongan Zhang, Zhongyu He, Zhenzhou Zha, Chunyong Hu, Tianyuan Gan, Qinglai Yan, Peng Wang, and Xuesong Ye. The future of endoscopic navigation: a review of advanced endoscopic vision technology. *IEEE Access*, 9:41144–41167, 2021. 1

[16] Jian Gao, Chun Gu, Youtian Lin, Hao Zhu, Xun Cao, Li Zhang, and Yao Yao. Relightable 3d gaussian: Real-time point cloud relighting with brdf decomposition and ray tracing. *arXiv preprint arXiv:2311.16043*, 2023. 3

[17] Panagiotis Giannios, Spyridon Koutsoumpos, Konstantinos G. Toutouzas, Maria Matiatou, George C. Zografos, and Konstantinos Moutzouris. Complex refractive index of normal and malignant human colorectal tissue in the visible and near-infrared. *Journal of Biophotonics*, 10(2):303–310, 2017. 6, 12

[18] Kaiwen Guo, Feng Xu, Yangang Wang, Yebin Liu, and Qionghai Dai. Robust non-rigid motion tracking and surface reconstruction using l0 regularization. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3083–3091, 2015. 3

[19] Kaiwen Guo, Peter Lincoln, Philip Davidson, Jay Busch, Xueming Yu, Matt Whalen, Geoff Harvey, Sergio Orts-Escolano, Rohit Pandey, Jason Dourgarian, et al. The relightables: Volumetric performance capture of humans with realistic relighting. *ACM Transactions on Graphics (ToG)*, 38(6):1–19, 2019. 3

[20] Yuan-Chen Guo, Yan-Pei Cao, Chen Wang, Yu He, Ying Shan, and Song-Hai Zhang. Vmesh: Hybrid volume-mesh representation for efficient view synthesis. In *SIGGRAPH Asia 2023 Conference Papers*, pages 1–11, 2023. 3

[21] Tao Hu, Tao Yu, Zerong Zheng, He Zhang, Yebin Liu, and Matthias Zwicker. Hvtr: Hybrid volumetric-textural rendering for human avatars. In *2022 International Conference on 3D Vision (3DV)*, pages 197–208. IEEE, 2022. 3

[22] Yingwenqi Jiang, Jiadong Tu, Yuan Liu, Xifeng Gao, Xiaoxiao Long, Wenping Wang, and Yuexin Ma. Gaussianshader: 3d gaussian splatting with shading functions for reflective surfaces. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5322–5332, 2024. 3, 7, 8

[23] Joanna Kaleta, Kacper Kania, Tomasz Trzcinski, and Marek Kowalski. Lumigauss: High-fidelity outdoor relighting with 2d gaussian splatting, 2024. 3

[24] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Trans. Graph.*, 42(4):139–1, 2023. 2, 3, 4, 7, 8

[25] Zhihao Liang, Qi Zhang, Ying Feng, Ying Shan, and Kui Jia. Gs-ir: 3d gaussian splatting for inverse rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21644–21653, 2024. 3

[26] Xingtong Liu, Zhaoshuo Li, Masaru Ishii, Gregory D Hager, Russell H Taylor, and Mathias Unberath. Sage: slam with appearance and geometry prior for endoscopy. In *2022 International conference on robotics and automation (ICRA)*, pages 5587–5593. IEEE, 2022. 2, 3

[27] Ruibin Ma, Rui Wang, Yubo Zhang, Stephen Pizer, Sarah K McGill, Julian Rosenman, and Jan-Michael Frahm. Rnnslam: Reconstructing the 3d colon to visualize missing regions during a colonoscopy. *Medical image analysis*, 72:102100, 2021. 1, 2, 3, 7

[28] Nader Mahmoud, Alexandre Hostettler, Toby Collins, Luc Soler, Christophe Doignon, and Jose Maria Martinez Montiel. Slam based quasi dense reconstruction for minimally invasive surgery scenes. *arXiv preprint arXiv:1705.09107*, 2017. 2, 3

[29] Dawid Malarz, Weronika Smolak, Jacek Tabor, Sławomir Tadeja, and Przemysław Spurek. Gaussian splatting with nerf-based color and opacity, 2024. 3

[30] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM*, 65(1): 99–106, 2021. 2, 3

[31] Raul Mur-Artal, Jose Maria Martinez Montiel, and Juan D Tardos. Orb-slam: a versatile and accurate monocular slam system. *IEEE transactions on robotics*, 31(5):1147–1163, 2015. 2, 3

[32] Ameya Pore, Martina Finocchiaro, Diego Dall'Alba, Albert Hernansanz, Gastone Ciuti, Alberto Arezzo, Arianna Menciassi, Alicia Casals, and Paolo Fiorini. Colonoscopy navigation using end-to-end deep visuomotor control: A user study. In *2022 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 9582–9588. IEEE, 2022. 1

[33] Dimitrios Psychogyios, Francisco Vasconcelos, and Danail Stoyanov. Realistic endoscopic illumination modeling for nerf-based data generation. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 535–544. Springer, 2023. 2, 3

[34] Ravi Ramamoorthi and Pat Hanrahan. An efficient representation for irradiance environment maps. In *Proceedings of the 28th annual conference on Computer graphics and interactive techniques*, pages 497–500, 2001. 4

[35] Yufei Shi, Beijia Lu, Jia-Wei Liu, Ming Li, and Mike Zheng Shou. Colonnerf: Neural radiance fields for high-fidelity long-sequence colonoscopy reconstruction. *arXiv preprint arXiv:2312.02015*, 2023. 2, 3

[36] Yahao Shi, Yanmin Wu, Chenming Wu, Xing Liu, Chen Zhao, Haocheng Feng, Jingtuo Liu, Liangjun Zhang, Jian Zhang, Bin Zhou, et al. Gir: 3d gaussian inverse rendering for relightable scene factorization. *arXiv preprint arXiv:2312.05133*, 2023. 3

[37] Joanna Waczyńska, Piotr Borycki, Sławomir Tadeja, Jacek Tabor, and Przemysław Spurek. Games: Mesh-based adapting and modification of gaussian splatting. *arXiv preprint arXiv:2402.01459*, 2024. 8, 12

[38] Kailing Wang, Chen Yang, Yuehao Wang, Sikuang Li, Yan Wang, Qi Dou, Xiaokang Yang, and Wei Shen. Endogslam: Real-time dense reconstruction and tracking in endoscopic surgeries using gaussian splatting. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 219–229. Springer, 2024. 1, 2, 3, 7, 8

[39] Peng Wang, Lingjie Liu, Yuan Liu, Christian Theobalt, Taku Komura, and Wenping Wang. Neus: Learning neural implicit surfaces by volume rendering for multi-view reconstruction. *arXiv preprint arXiv:2106.10689*, 2021. 2

[40] Tong Wu, Jia-Mu Sun, Yu-Kun Lai, Yuewen Ma, Leif Kobbelt, and Lin Gao. Deferredgs: Decoupled and editable gaussian splatting with deferred shading. *arXiv preprint arXiv:2404.09412*, 2024. 3, 5, 7, 8, 14, 15

[41] Kai Zhang, Gernot Riegler, Noah Snavely, and Vladlen Koltun. Nerf++: Analyzing and improving neural radiance fields. *arXiv preprint arXiv:2010.07492*, 2020. 3

# Appendix

## A. Detailed Endoscopic Setup

The exact positioning of the light source relative to the lens is illustrated in Fig. 8. Our setup supports seamless optimization of the spotlight angle for the light source. While this feature was not utilized in our experiments due to the characteristics of our datasets, it can be jointly optimized alongside other light parameters. Additionally, experiments involving adjustments to the spotlight angle during the inference stage are shown in Appendix C.
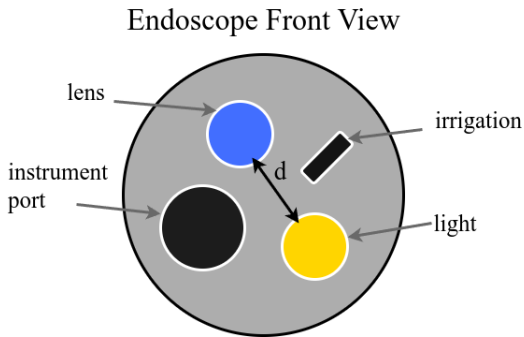


Figure 8. **Detailed colonoscope architecture.** The light and camera sources are colocated, with a small distance between them, denoted as $d$. We optimize $d$ alongside other light parameters.

## B. Physical Properties of Optimized Gaussians

Table 4 demonstrates that the optimized roughness values are comparable to those estimated in [17] when $F_0$ is clamped at 0.03. This highlights the physical plausibility of the optimized parameters in our model.

Table 4. Clamping the $F_0$ value to 0.03 reflects realistic refractive properties of the colon, as reported in [17]. This approach ensures that the optimized values align with expected physical properties.

| Sequence | Roughness | F0 |
| --- | --- | --- |
| Expected physical value | $\sim$0.25 | $\sim$0.021 |
| cecum_t1_b | 0.2722 | 0.0300 |
| cecum_t2_b | 0.2585 | 0.0300 |
| cecum_t3_a | 0.2474 | 0.0300 |
| sigmoid_t1_a | 0.2800 | 0.0300 |
| sigmoid_t2_a | 0.2675 | 0.0300 |
| sigmoid_t3_a | 0.2600 | 0.0300 |
| trans_t1_b | 0.2741 | 0.0300 |
| trans_t2_c | 0.2630 | 0.0300 |
| trans_t4_a | 0.2562 | 0.0300 |
| trans_t4_b | 0.2549 | 0.0300 |

## C. Additional Experiments on Relighting

We demonstrate the manipulation of the light source and physical parameters. For this experiment, we decouple the light angle from the camera to showcase the effectiveness of our method. By adjusting the spotlight angle, we control the light's coverage area. Additionally, modifying the roughness of the Gaussians results in changes to specularity: smoother surfaces exhibit greater reflectivity, while rougher surfaces reduce specularity. We strongly recommend reviewing the supplementary videos, which effectively demonstrate the relighting capabilities of our method.

## D. Additional Experiments on Simulating Body Movements

Capturing detailed changes during body movements in static images is challenging; hence, we provide videos for visualization. We highly recommend reviewing these videos, as they effectively demonstrate how our method enables realistic body movement simulations with minimal artifacts. To simulate body movements, we utilize GaMeS [37] reparameterization combined with parametric or cage-based physical simulation.

## E. Additional Results for Geometry

In Fig. 10, Fig. 11, Fig. 12, we present an overview of the generated normals for each C3VD scene, alongside a comparison to normals produced by 3DGS-DR. The results demonstrate that our normals are more accurate and closely aligned with the ground truth geometry.
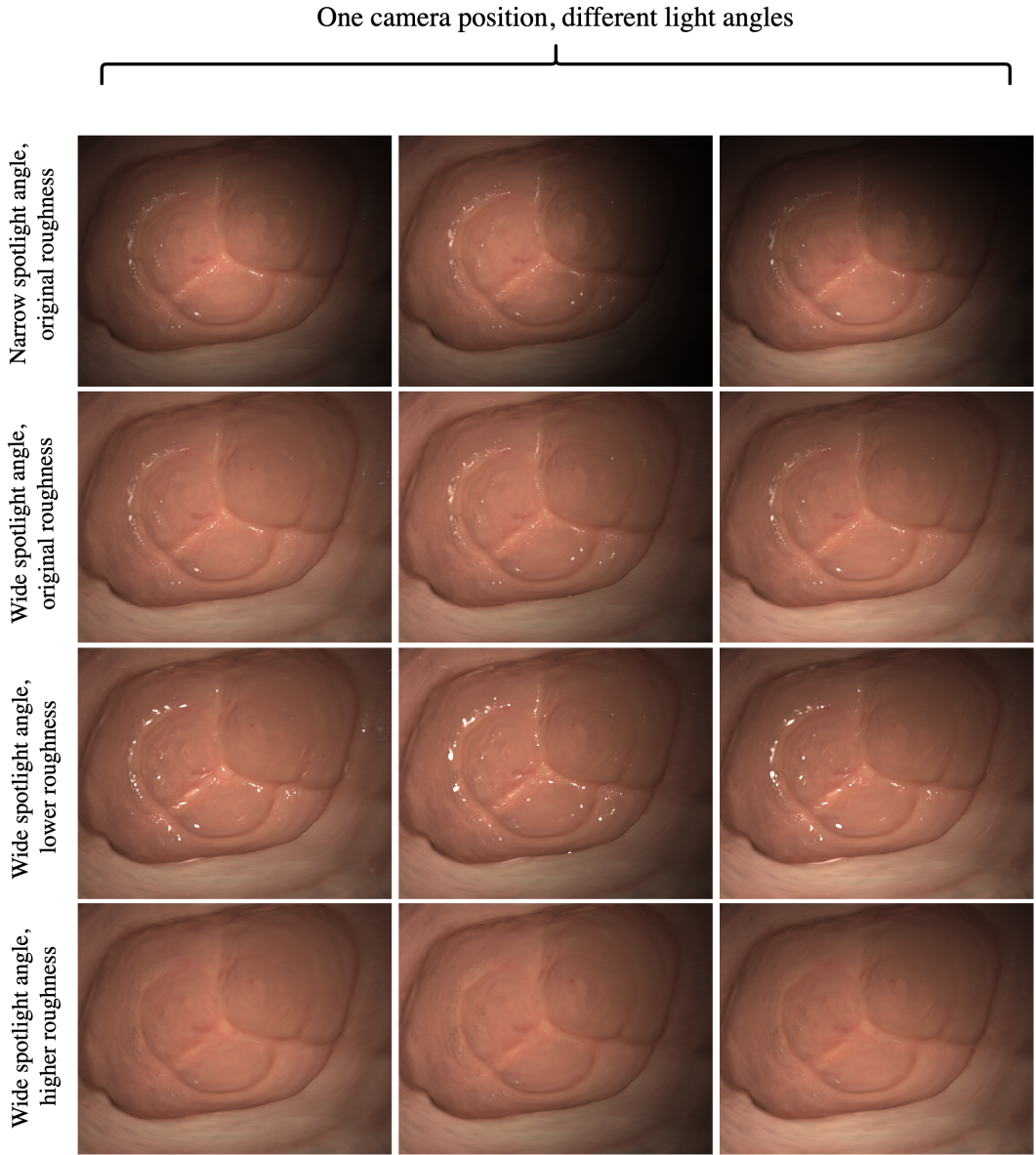
Figure 9. **Additional relighting results for Scene 1.** Light angles and Gaussian roughness are manipulated to explore their effect on specularity and surface appearance.

**Sigmoid sequences**

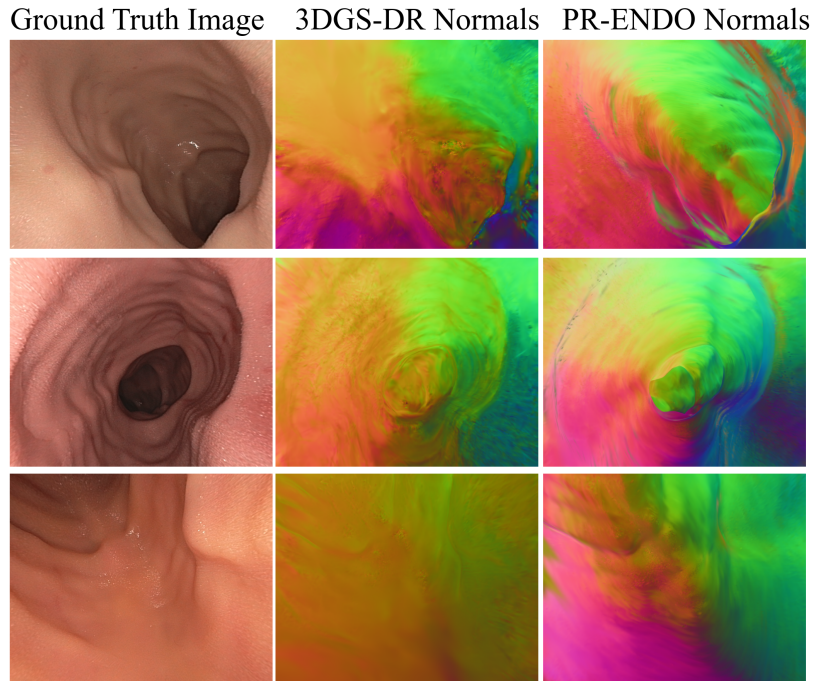| Ground Truth Image | 3DGS-DR Normals | PR-ENDO Normals |



Figure 10. **Normal comparison across sigmoid C3VD scenes.** A single random capture is presented for each C3VD sequence. Our method generates higher-quality normals aligned with the ground truth geometry compared to the state-of-the-art 3DGS-DR [40].

**Cecum sequences**

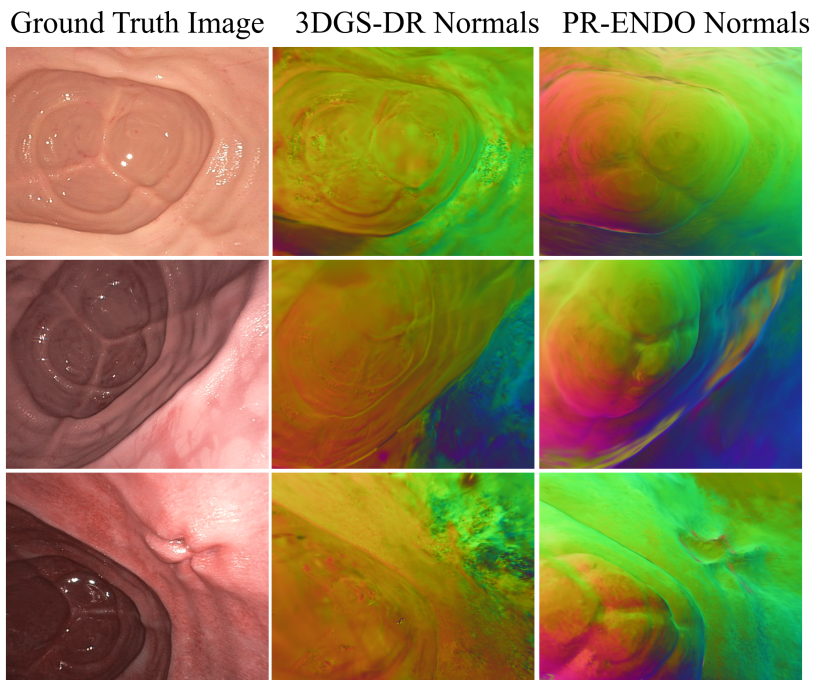| Ground Truth Image | 3DGS-DR Normals | PR-ENDO Normals |



Figure 11. **Normal comparison across cecum C3VD scenes.** A single random capture is presented for each C3VD sequence. Our method generates higher-quality normals aligned with the ground truth geometry compared to the state-of-the-art 3DGS-DR [40].

**Trans sequences**

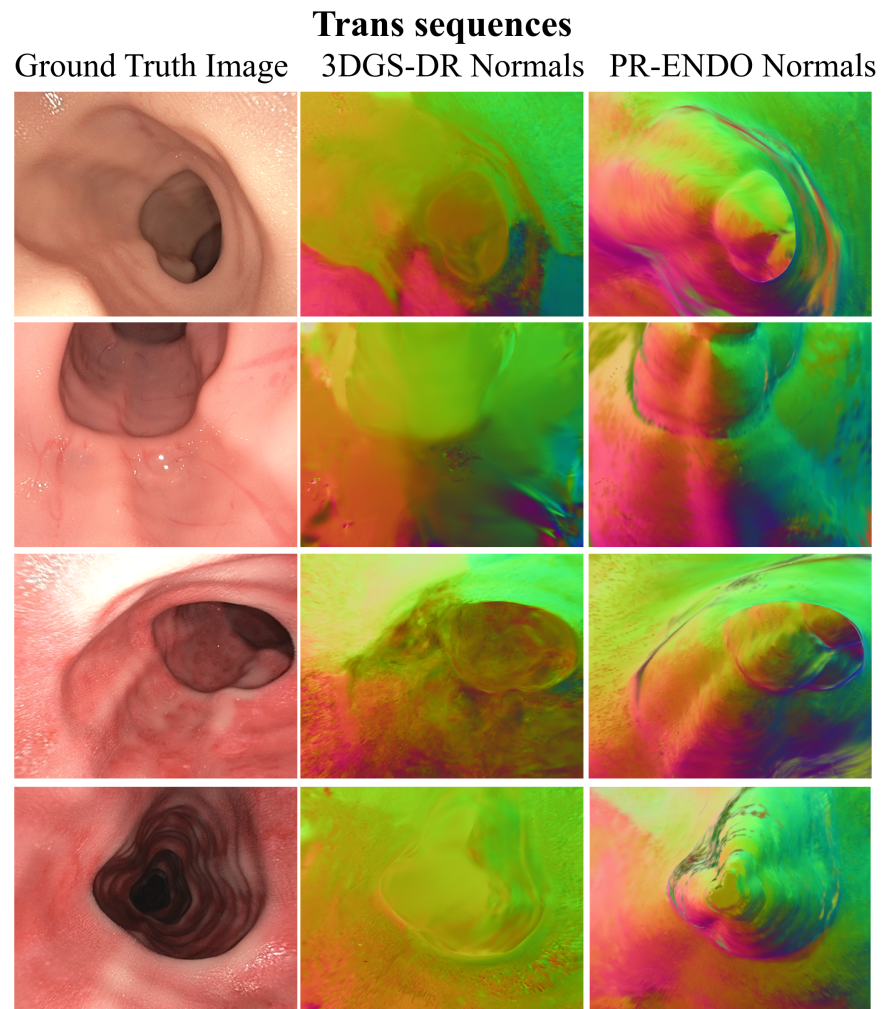Ground Truth Image | 3DGS-DR Normals | PR-ENDO Normals

Figure 12. **Normal comparison across trans C3VD scenes.** A single random capture is presented for each C3VD sequence. Our method generates higher-quality normals aligned with the ground truth geometry compared to the state-of-the-art 3DGS-DR [40].