

---

# NFL-BA: Near-Field Light Bundle Adjustment for SLAM in Dynamic Lighting

---

Andrea Dunn Beltran<sup>\*1</sup>, Daniel Rho<sup>\*1</sup>, Marc Niethammer<sup>2</sup>, Roni Sengupta<sup>1</sup>

<sup>1</sup> University of North Carolina at Chapel Hill    <sup>2</sup> University of California San Diego

<sup>\*</sup> Equal contribution

{asdunnbe, dn103c1, ronisen}@cs.unc.edu, mniethammer@ucsd.edu



Figure 1: **NFL-BA enhances tracking and mapping in neural rendering-based SLAM** (e.g., MonoGS [27]) by explicitly modeling dynamic near-field lighting, with applications in endoscopy.

## Abstract

Simultaneous Localization and Mapping (SLAM) systems typically assume static, distant illumination; however, many real-world scenarios, such as endoscopy, subterranean robotics, and search & rescue in collapsed environments, require agents to operate with a co-located light and camera in the absence of external lighting. In such cases, dynamic near-field lighting introduces strong, view-dependent shading that significantly degrades SLAM performance. We introduce Near-Field Lighting Bundle Adjustment Loss (NFL-BA) which explicitly models near-field lighting as a part of Bundle Adjustment loss and enables better performance for scenes captured with dynamic lighting. NFL-BA can be integrated into neural rendering-based SLAM systems with implicit or explicit scene representations. Our evaluations mainly focus on endoscopy procedure where SLAM can enable autonomous navigation, guidance to unsurveyed regions, blindspot detections, and 3D visualizations, which can significantly improve patient outcomes and endoscopy experience for both physicians and patients. Replacing Photometric Bundle Adjustment loss of SLAM systems with NFL-BA leads to significant improvement in camera tracking, 37% for MonoGS and 14% for EndoGS, and leads to state-of-the-art camera tracking and mapping performance on the C3VD colonoscopy dataset. Further evaluation on indoor scenes captured with phone camera with flashlight turned on, also demonstrate significant improvement in SLAM performance due to NFL-BA.

## 1 Introduction

Simultaneous Localization and Mapping (SLAM) enables autonomous agents to build a spatial map of an unknown environment while estimating their own poses within it, with wide-ranging applications in robotics, computer vision, autonomous vehicles, and scientific imaging. Most SLAM systems [38, 34, 7, 58, 60, 56, 49, 20, 27, 16] assume an autonomous agent navigating an environment with distant, static illumination, e.g., a self-driving car in the streets, and they optimize a Photometric Bundle Adjustment loss where they minimize an error between the captured image and the re-rendered image using estimated 3D scene and camera poses.

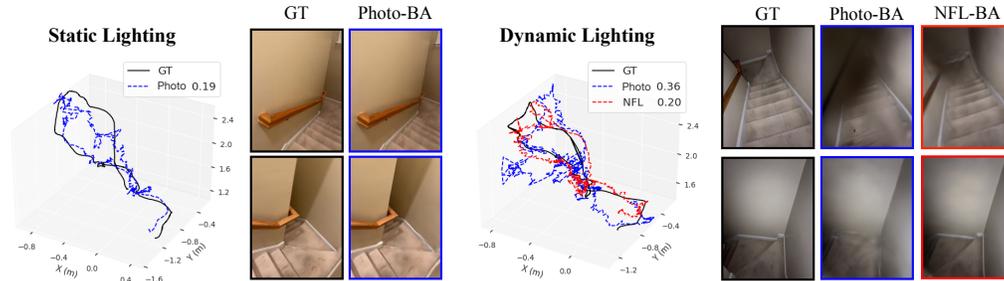


Figure 2: **MonoGS performance under (1) distant static lighting and (2) dynamic near-field lighting from a co-located flashlight.** Standard photometric BA performs well under static lighting but fails under dynamic lighting, degrading both trajectory and map quality. NFL-BA restores performance under dynamic lighting, matching the quality of the static-light setup.

However, many scientific and safety-critical applications demand that autonomous agents operate in environments devoid of external illumination, relying instead on self-mounted light sources. For example, in endoscopy procedures, a slender flexible tube with a co-located light and camera is used to inspect internal organs such as the airway and the colon [9, 45, 33, 23, 51, 17]. Accurate trajectory estimation is crucial for reliably guiding instruments to areas of interest, mapping anomalies, and avoiding tissue damage during navigation. In subterranean search-and-rescue or collapsed-building inspection, robots rely on onboard lamps to explore unstable voids; slight errors in pose estimation can accumulate into large drift, leading to misaligned maps, missed victims, or costly back-tracking.

Despite the prevalence of these use cases, current SLAM systems perform poorly under such conditions (see Fig. 2). This performance drop is primarily due to the effects of *dynamic near-field lighting*, where the only illumination is co-located with the camera and moves with it. Dynamic near-field lighting causes different points of the surface to receive different intensities of light at each time step, depending on the distance and orientation of the point to the camera, introducing strong, view-dependent shading. These lighting artifacts significantly impair both feature-based and direct (photometric) tracking, resulting in substantial failures in mapping accuracy and pose estimation.

To alleviate these issues, we propose a new Bundle Adjustment loss that accounts for dynamic near-field lighting. Our key intuition is that the shading effect of the captured image can provide valuable information about the relative distance and orientation between the surface and the camera. With this, we formulate a Near-Field Lighting Bundle Adjustment loss, NFL-BA, where we optimize the surface geometry and the camera parameters such that the rendered image has shading variations that match the relative distance and orientation between the surface and the camera. Our NFL-BA loss can be applied to any neural rendering-based SLAM algorithm, i.e., with neural implicit and explicit 3D Gaussian scene representation.

In this paper, we specifically focus on demonstrating how NFL-BA can improve the performance of existing SLAM systems for 3D reconstruction and localization from endoscopy videos. SLAM can enable autonomous navigation through internal organs and guide physicians to unsurveyed regions to improve physicians’ situational awareness by providing 3D visualizations, and can help measure organ shapes. We evaluated NFL-BA with two state-of-the-art 3DGS-based SLAM systems, general-purpose MonoGS [27] and endoscopy-specific EndoGSLAM [45], and one neural implicit SLAM, NICE-SLAM [58], by replacing their Photometric Bundle Adjustment loss with NFL-BA loss. We observe that the NFL-BA loss improves the performance of all SLAM algorithms on average when using ground-truth or estimated depth maps on the C3VD colonoscopy dataset. For example, NFL-BA significantly improves MonoGS by reducing camera tracking error by 37% (3.48 to 2.18 mm) and camera mapping error by 38% (1.59 to 0.99 mm) when initialized by PPSNet depth[35].

Additionally, we also demonstrate the effectiveness of NFL-BA on indoor rooms captured with a moving co-located light and camera without any external light source, mimicking agent navigation during search & rescue and covert military operations. By replacing incorporating our NFL-BA loss, we see an average improvement of  $\sim 35\%$  in pose estimation across all scenes.

## 2 Related Works

**Dense SLAM and Bundle Adjustment.** Early SLAM pipelines focused on sparse feature matching for pose estimation and mapping [30, 6, 42, 11]. With advancements in neural scene representations

several proposed SLAM frameworks [58, 2] generate dense, pixel-level that yield more detailed and robust reconstructions. More recently, 3D Gaussian surface methods have demonstrated real-time rendering with high-fidelity mapping [20, 27, 49, 16, 10, 53].

These dense SLAM approaches all rely on a core Bundle Adjustment step. Bundle Adjustment (BA) alternatively optimizes camera parameters and surface geometry by minimizing errors across multiple frames. Traditional geometric BA aligns detected 2D feature points to their 3D counterparts by minimizing reprojection error, assuming static lighting and Lambertian surfaces [14]. Although effective in controlled environments, it struggles in complex or low-texture scenes. Photometric BA (Photo-BA) [1] incorporates pixel intensities into the optimization process, minimizing photometric re-projection errors and proving advantageous in environments where feature matching fails [11]. However, Photo-BA does not exploit the correspondence cues provided by dynamic or near-field lighting where image intensities vary across frames.

**Near-field Lighting models.** Near-field lighting has been leveraged for 3D reconstruction tasks like monocular depth and surface normal estimation [35, 57] and Photometric Stereo [21]. Some of these approaches [35, 21] use a near-field lighting representation as input to a CNN along with captured images for predicting surface normal and geometry. In the context of Endoscopy, LightDepth [37] and PPSNet [35] demonstrated the effectiveness of near-field lighting to enhance depth estimation. LightNeus [3] exploited the inverse-square law for light decay to improve endoscopic surface reconstruction, however with known camera parameters and pre-operative 3D CT scan.

*It has never, however, been used for Simultaneous Localization & Mapping (SLAM) problems, let alone in combination with neural rendering methods.* To this end, we propose a Bundle Adjustment Loss with Near-Field Lighting (NFL-BA), considering the most commonly available single co-located camera & light in the endoscope or other autonomous agents.

**Dynamic Lighting in SLAM.** Visual SLAM performance often degrades under illumination changes such as exposure shifts, specularities, and varying color temperature. Early photometric calibration methods jointly optimize camera intrinsics, exposure, and scene depths to normalize brightness variations in real time [11] while probabilistic SLAMs with unscented filtering further stabilizes pose estimates under uncertain lighting conditions [26]. More recently, learning-based matchers [22, 48] adapt descriptors to cope with complex lighting variations. None of these methods, however, explicitly model near-field lighting geometry to handle this co-located light setting.

**SLAM in endoscopy.** Early works [40, 12] demonstrated the feasibility of applying SLAM in such environments by addressing dynamic lighting and tissue deformation. Researchers have often used a mixture of supervised learning on synthetic and self-supervised learning on real endoscopy datasets for tailoring SLAM frameworks to endoscopy with complex camera motion [25, 55, 46] and developed novel endoscopy SLAM frameworks [36, 29, 18]. However these techniques often struggle with challenging sequences from both synthetic and clinical data. Recently, neural rendering-based methods [39, 23, 45, 51, 13, 15] have proved especially effective in generating high-quality details and modeling textureless regions with a large number of Gaussians. In this work, we adopt neural rendering approaches and explicitly model the near-field lighting effects, alleviating dynamic lighting challenges and improving performance.

### 3 Background

In this section, we review the general framework of neural rendering-based SLAM. We represent the camera at time  $t$  by its extrinsics  $P_t = [R_t, T_t] \in \mathbf{SE}(3)$  and known intrinsics  $K$ , yielding the projection  $\pi_t = KP_t$ . We assume the camera intrinsic  $K$  to be the same for all frames and known or calibrated ahead of time. Pixels are denoted  $p$  and 3D camera-space points by  $x$ .

In neural rendering, scene parameters  $\Theta$ , whether in the form of neural networks or primitives, encode visual and geometric information, such as colors  $c_i$  and occupancy  $\alpha_i$ . Given  $\Theta$  and  $P_t$ , we can get the color  $\hat{C}(\cdot)$  and the depth  $\hat{D}(\cdot)$  of a pixel  $p$  from a frame at time  $t$  as follows [28, 20]:

$$\hat{C}(p) = \sum_{i \in \mathcal{N}} c_i \alpha_i \prod_{j=1}^{i-1} (1 - \alpha_j), \quad \hat{D}(p) = \sum_{i \in \mathcal{N}} z_i \alpha_i \prod_{j=1}^{i-1} (1 - \alpha_j) \quad (1)$$

where  $\mathcal{N}$  denotes the group of samples for a pixel  $p$ , with  $\alpha_i$  representing the occupancy of the  $i$ -th sample, and  $z_i$  denotes its distance from the camera center.

To optimize  $P_t$  and  $\Theta$ , dense SLAM methods typically use rendering loss  $\mathcal{L}_{ren}$ , reducing the rendering errors between the rendered and captured images [58, 27, 49] and, if estimated or ground

truth depth maps are available, an additional depth loss  $\mathcal{L}_{geo}$  can be added [43]. Typically, these losses take the form of  $L^p$  norm as follows with variations with  $M_t$  as a pixel-wise mask:

$$\mathcal{L}_{ren} = \|M_t \odot (\hat{C} - C)\|_p, \quad \mathcal{L}_{geo} = \|M_t \odot (\hat{D} - D)\|_p \quad (2)$$

Bundle adjustment optimizes both  $P_t$  and  $\Theta$  using the following combined loss:

$$\text{Photo-BA:} \quad \min \sum_{t \in \mathcal{W}} \lambda_{ren} \mathcal{L}_{ren}(\hat{C}, C; M_t) + \lambda_{geo} \mathcal{L}_{geo}(\hat{D}, D; M_t) \quad (3)$$

where  $\mathcal{W}$  denotes the set of frames used for the bundle adjustment and the hyperparameters  $\lambda_{ren}$  and  $\lambda_{geo}$  are the loss weights. Additionally, the objective function can include any other regularization terms, such as artifact suppressing [27] or opacity regularization [59].

During the Mapping stage, both  $\Theta$  and  $P_t$  are optimized over a set of keyframes. The exact algorithm for keyframe selection, keyframe update and optimization strategies for tracking and mapping phase vary between different SLAM approaches and their specific objectives.

**Implicit Neural Representations.** Neural field-based SLAM methods [41, 58, 44, 60, 38] uses a set of neural networks  $F(x, d; \Theta) \rightarrow (c_i, \sigma_i)$ , optimized to estimate the color  $c_i$  and the volume density  $\sigma_i$  for an input 3D coordinate  $x$  and the view direction  $d$ . The the occupancy can be calculated from the volume density  $\sigma_i$  and the distance between adjacent samples  $\delta_i$  as  $\alpha_i = 1 - \exp(-\sigma_i \delta_i)$ .

**3D Gaussian Splatting.** For 3D Gaussian Splatting [20] SLAM methods, the scene is represented by a set of Gaussians with mean  $\mu^i$ , covariance  $\Sigma^i$  in world space, color  $c_i$ , and opacity  $\alpha^i$ . The shape parameters and occupancy  $\alpha^i$  of the *splatted* 2D Gaussians are computed as follows:

$$\bar{\mu}_t^i = \pi_t \mu^i, \quad \bar{\Sigma}_t^i = J_t R_t \Sigma^i R_t^T J_t^T, \quad \alpha_i = \alpha^i \exp\left(-\frac{1}{2}(p - \bar{\mu}_t^i)^\top (\bar{\Sigma}_t^i)^{-1} (p - \bar{\mu}_t^i)\right) \quad (4)$$

where  $J_t$  is the Jacobian of the projection  $\pi_t$ ,  $p$  denotes a pixel coordinate, and  $\bar{\mu}_t^i, \bar{\Sigma}_t^i$  are the splatted mean and covariance of Gaussian  $\mathcal{G}^i$  in pixel space.

## 4 Near-Field Light Bundle Adjustment

We introduce a novel Near-Field Lighting based Bundle Adjustment loss, NFL-BA, that integrates near-field lighting with neural-rendering 3D scene representations to improve performance of existing SLAM systems on images captured with dynamic lighting co-located with the camera. Our proposed NFL-BA can replace commonly used Photometric Bundle adjustment loss, defined in Eq. 3, within neural-rendering based SLAM framework. Photo-BA typically optimizes scene appearance parameter as RGB color, which is sufficient when the illumination on each scene point remains the constant throughout the capture. However, for scenes with a dynamic light co-located with a moving camera, the illumination received at each point varies per frame as the camera and the light moves through the scene. In this setting, the illumination received at each point depends on the relative distance and orientation between the point and the camera, as conceptualized in Fig. 3. Thus continuing to model scene appearance as simple RGB color is inaccurate for dynamic near-field lighting as it doesn't separate effects of illumination due to camera movement from the intrinsic view-independent color of the scene, i.e. albedo.

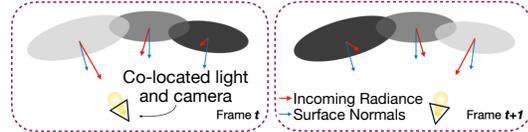


Figure 3: **Illustration of our key idea.** As the co-located light and camera, moves through the scene, different 3D Gaussians on the surface receive different intensities of light (red arrow), dependent on the relative distance and orientation between the 3D Gaussian and the camera.

Our goal is to explicitly model surface appearance as albedo and separate near-field lighting effects from it. To accurately model dynamic lighting we then represent near-field illumination effects with camera pose and scene geometry. In sec. 4.1 we describe our image formation model using neural rendering framework that will decompose the surface appearance into albedo and incoming lighting, which will be further represented as a function of scene geometry and camera pose. Then in sec. 4.2, we will use this image formation to create the Near-Field Bundle Adjustment loss and show how it can be easily integrated into neural rendering based SLAM framework.

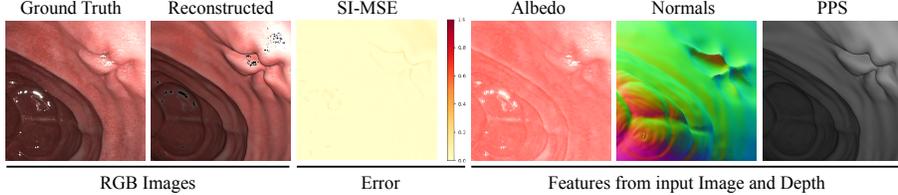


Figure 4: **Image Formation Validation.** We show that C3VD images captured with a real endoscope conform to our co-located light-camera and zero attenuation  $\beta$  image formation model, as indicated by very low per-pixel scale-invariant MSE between the original image and the reconstructed image with masked-out specular regions.

#### 4.1 Image Formation with Near-Field Lighting

We consider an image-formation model under near-field lighting for a single image following previous works [19, 35]. Each pixel  $p$  and the corresponding three-dimensional point  $x_p$  in the camera space receives different light intensities and directions, characterized by the light source to surface direction  $L^d(\cdot)$  and attenuation term  $L^a(\cdot)$ , as follows:

$$L^d(x_p) = \frac{x_p - x_L}{\|x_p - x_L\|}, \quad L^a(x_p) = \frac{(L^d(x_p)^\top f)^\beta}{\|x_p - x_L\|^2}, \quad (5)$$

where  $x_L$  is the location of the light source,  $f$  is the forward (optical axis) vector.  $\beta$  is an angular attenuation coefficient, and will be discussed in sec. 4.2.

Assuming a diffuse reflectance model, which has proven effective for depth estimation in endoscopic scenes [35], we can approximate the rendered image at each pixel  $\hat{C}(\cdot)$  as:

$$PPS(x_p) = L^a(x_p) \cdot (L^d(x_p)^\top n(x_p)), \quad \hat{C}(p) = \rho(x_p) PPS(x_p), \quad (6)$$

where  $\rho(\cdot)$  and  $n(\cdot)$  are albedo and normal at position  $x_p$  of pixel  $p$  respectively.  $PPS(\cdot)$  is a per-pixel shading term. Note that existing approaches that uses this near-field light image formation model [19, 35] uses pixel-based representation to predict depth map or surface geometry from images captured from a single viewpoint only. In this paper, we extend the Near-Field Image Formation model beyond single-view pixel-based representation to multi-view 3D representation.

Our key insight is that the standard volumetric rendering equation can be modified to incorporate the near-field lighting model described in eq. 6, while keeping the overall SLAM pipeline intact. In our framework, we reinterpret the direct color ( $c_i$  in eq. 3) as the product of the albedo  $\rho(\cdot)$  and the shading term  $PPS(\cdot)$ , which models dynamic near-field lighting. Note that both albedo and shading is defined directly on the 3D neural representations, i.e. neural radiance field or 3D gaussians, and not in pixel-space. This leads to the modified rendering equation under near-field lighting:

$$\hat{C}_{pps}(p) = \sum_{i \in \mathcal{N}} \rho(x_i) PPS(x_i) \alpha_i \prod_{j=1}^{i-1} (1 - \alpha_j) \quad (7)$$

Note that eq. 6 represents a special case of eq. 7 where a single sample is considered and the occupancy  $\alpha_i$  equals one. Our image formation model assumes diffuse reflectance and no angular attenuation, to reduce the complexity of the modeling. While it is easy to extend the image formation model to handle specular reflectance and angular attenuation of lighting, this leads to additional parameters that needs to be optimized during the Bundle Adjustment.

**Angular attenuation.** Following previous works [19, 35], we simplify the near-field light image formation model by setting the attenuation coefficient  $\beta$  in eq. 5 to zero. This effectively ignores the directional fall-off component, reducing the light attenuation term to a simple inverse-square fall-off  $L^a(x_p) = 1/\|x_p - x_L\|^2$ . This simplification is justified because the angular attenuation in settings like endoscopy is often negligible compared to the inverse square law attenuation, and estimating  $\beta$  accurately can be challenging due to variations in endoscope designs. In future work, we plan determine the optimal value of  $\beta$  for different systems and incorporate the light direction vector  $r_t^e$  for more accurate modeling which can further improve camera rotation during Bundle Adjustment.

**Empirical validation of our image-formation model on colonoscopy image.** Fig. 4 provides an example colonoscopy image from the C3VD dataset [4], showing the accuracy of the near-light field model (Eq. 6) with  $\beta$  of 0. Albedo was estimated by converting each RGB image to HSV color

space, setting the value channel to 1 across all pixels, then converting the modified image back to RGB space. This standardizes pixel intensity variations, approximating a reflectance map where illumination effects are minimized, but does not strictly represent ground truth albedo. As shown, the image formulation model is sufficient to represent endoscopic scenes with low reconstruction errors.

## 4.2 Near-Field Light Bundle Adjustment Loss

Next, we will re-define the Photometric Bundle Adjustment loss of eq. 3 using the near-field lighting based image formation model defined in eq. 7 expressed as follows:

$$\text{NFL-BA: } \min_{t \in \mathcal{W}} \sum \lambda_{ren} \mathcal{L}_{ren}(\hat{C}_{pps}, C; M_t) + \lambda_{geo} \mathcal{L}_{geo}(\hat{D}, D; M_t) \quad (8)$$

where  $\hat{C}_{pps}$  denotes the rendered image with near-field lighting-incorporated volumetric rendering equation (Eq. 7). This reformulation seamlessly integrates near-field lighting cues into the neural rendering framework without altering the rest of the SLAM framework. Since our formulation is confined solely to the rendering process, and thus to the bundle adjustment, we do not modify or replace any other SLAM components for fair comparison. This design choice enables easier integration with existing neural rendering-based SLAM methods.

**Choice of image space in optimization.** Note that many settings, and especially endoscopy, frames are stored in standard sRGB color space, whereas our near-field shading term  $PPS(\cdot)$  is computed in a linear space. To ensure consistency, we apply an inverse gamma correction of  $\gamma = 2.2$  to the sRGB images before computing  $PPS$ , or equivalently, gamma-correct the linear PPS output by  $\gamma = 1/2.2$  when rendering back to sRGB. This step aligns the lighting model with the true photometric intensities and prevents bias from the nonlinear sRGB transfer function.

**Normal calculation during Bundle Adjustment.** To calculate the normals  $n(\cdot)$  from neural fields, we utilize the direction of the gradient of the occupancy with respect to the spatial coordinates as follows [5]:  $n(x_i) = -\nabla\sigma(x_i)/\|\nabla\sigma(x_i)\|$ . For Gaussian Splatting, we use the shortest axis of each Gaussian as its normal, following [52, 8, 47]. In both cases, we ensure the computed normal is oriented towards the camera by enforcing  $n(x)^\top L^d(x)$  to be positive. Otherwise, we flip the normals by multiplying them by -1 for stability.

## 5 Evaluation

Our proposed method is a plug-in approach that can be applied to any existing neural-rendering-based SLAM framework. We first test our method on endoscopy videos using one neural implicit SLAM, NICE-SLAM [58], as well as two existing 3DGS-SLAM frameworks: the general-purpose MonoGS [27] and the endoscopy-specific EndoGSLAM [45]. In each case, we replace the standard Photometric Bundle Adjustment loss (3) with our proposed equation NFL-BA loss (8). Additionally, we also test MonoGS [27] on self-captured indoor scenes with a co-located light and camera.

Table 1: **Quantitative Evaluation on the C3VD [4] dataset with oracle depth map.** Replacing Photometric BA with NFL-BA significantly improves tracking quality of two state-of-the-art 3D Gaussian SLAMs, MonoGS [27] and EndoGS [59], and one neural implicit SLAM, NICE-SLAM [58].

Method	BA	Tracking		Mapping
		ATE <sub>t</sub> (mm)↓	ATE <sub>r</sub> (°)↓	Chamfer (mm)↓
NICE-SLAM, <i>CVPR'22</i>	Photo	4.16	2.68	1.95
	NFL	<b>2.88</b>	2.81	<b>1.70</b>
EndoGSLAM, <i>MICCAI'24</i>	Photo	1.93	1.81	0.85
	NFL	2.04	<b>1.13</b>	0.97
MonoGS, <i>CVPR'24</i>	Photo	2.90	1.11	1.16
	NFL	<b>1.60</b>	1.49	<b>0.79</b>

### 5.1 Evaluation Setting

**Datasets.** We evaluate our method on three datasets that reflect different challenges in handling near-field dynamic lighting: (1) a phantom endoscopy dataset, (2) a clinical endoscopy dataset, and (3) a dataset of indoor scenes captured with phone camera with flashlight turned on.

**C3VD.** The C3VD dataset [4] (CC BY-NC-SA 4.0) was created using a phantom colon with synthetic materials to simulate realistic tissue geometry. **Colon10K.** To test generalization in real-world clinical endoscopy settings, we evaluate on Colon10K [24], a large-scale video dataset without depth or pose supervision. Videos are sampled from actual The endoscopy video was captured by a surgeon who performs different endoscopy procedures on the phantom colon with a real endoscope capturing RGB images coupled with corresponding depth maps. We focus on 8 sequences ranging from 70 to 800 frames from different regions of the colon, for more details, please see supplementary. We evaluate using both ground truth and predicted depths.

Table 2: **Quantitative evaluation on the C3VD [4] dataset** using depth maps estimated by SOTA techniques, PPSNet [35] and DA-Hybrid [50, 32]. Replacing Photometric BA with NFL-BA significantly improves tracking for both MonoGS [27] and EndoGSLAM [45], and mapping and rendering quality for MonoGS [27]. Note that SOTA performance for each of the tracking, mapping, and rendering metrics is observed when NFL-BA is used.

Method	Depth	BA	Tracking		Mapping	Rendering
			$ATE_t$ (mm)↓	$ATE_r^\circ$ ↓	Chamfer (mm) ↓	LPIPS ↓
EndoGSLAM [45] <i>MICCAI'24</i>	PPS-Net	Photo	3.03	1.73	1.23	0.39
	PPS-Net	NFL	<b>2.62</b>	<b>1.24</b>	1.25	<b>0.39</b>
	DA-Hybrid	Photo	6.67	2.26	2.12	0.43
	DA-Hybrid	NFL	<b>3.91</b>	<b>1.58</b>	2.39	<b>0.42</b>
MonoGS [27] <i>CVPR'24</i>	PPS-Net	Photo	3.48	1.70	1.59	0.56
	PPS-Net	NFL	<b>2.18</b>	<b>1.65</b>	<b>0.99</b>	<b>0.53</b>
	DA-Hybrid	Photo	4.63	1.69	1.34	0.52
	DA-Hybrid	NFL	<b>2.35</b>	<b>1.14</b>	<b>1.13</b>	0.52

procedures and are typically around 300-600 frames. This setting is significantly more challenging than the phantom setting, since frames may contain motion blur, specular highlights, and fluid occlusions. Sequences are uniformly sampled and fisheye corrected.

*Self-Captured Indoor Scenes.* To study the role of near-field lighting in a controlled non-clinical setting, we capture a dataset of four indoor scenes (*Guitars, Porch, Pool, and Stairs*) using an iPhone 15 Pro. Scenes include objects with varied geometry and reflectance (diffuse, specular), imaged under dynamic motion. Scenes are captured using a co-located point light source mounted to camera. Ground-truth camera trajectories were recorded via motion capture, but no reference point clouds are available; hence, we report only trajectory error (ATE<sub>t</sub>) and perceptual quality (LPIPS), omitting Chamfer distance. Details and additional visualizations are included in the supplementary.

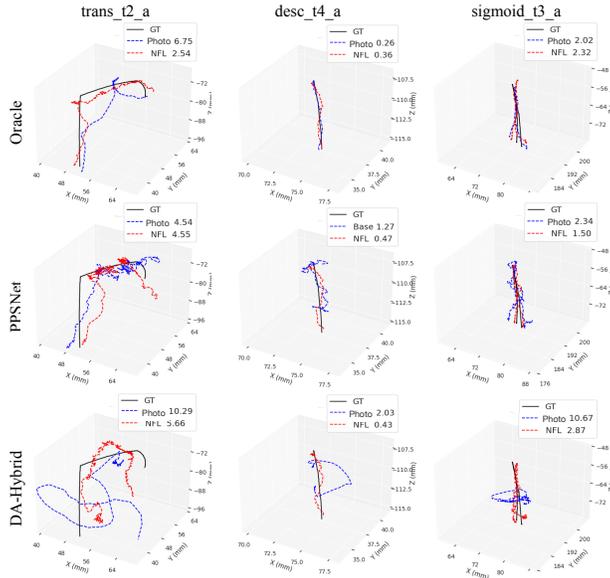


Figure 5: **Camera tracking improvement over EndoGSLAM [45].** Replacing the Photo-BA loss (in blue) with NFL-BA loss (in red) significantly improves camera tracking for different depth initialization. Average tracking error ATE<sub>t</sub> for each sequence is reported in the inset. (zoom for details)

**Metrics.** For evaluation, we basically followed other neural rendering SLAM algorithms [27, 45]. For tracking performance, we measure the root mean square error of the Absolute Trajectory Error (ATE) for both translation and rotation across all frames. Translation error ATE<sub>t</sub> is in millimeters (mm) for the endoscopy scenes and meters (m) for the in-door scenes. And rotation error ATE<sub>r</sub> is in degrees. To assess the mapping quality, we use the Chamfer distance from ground truth point clouds to the nearest points in the estimated point clouds [46], for more details please see supplementary. In addition, we evaluate rendering quality using the Learned Perceptual Image Patch Similarity (LPIPS) [54]. We note that for many endoscopic SLAM applications, tracking and mapping accuracies are more important than photorealism of the rendered images, unlike many indoor or outdoor scenes.

**Computational costs** We trained all models on a single NVIDIA RTX A6000 GPU. The per-scene optimization takes ~1 FPS. For more information on runtime speed, please see supplementary materials.

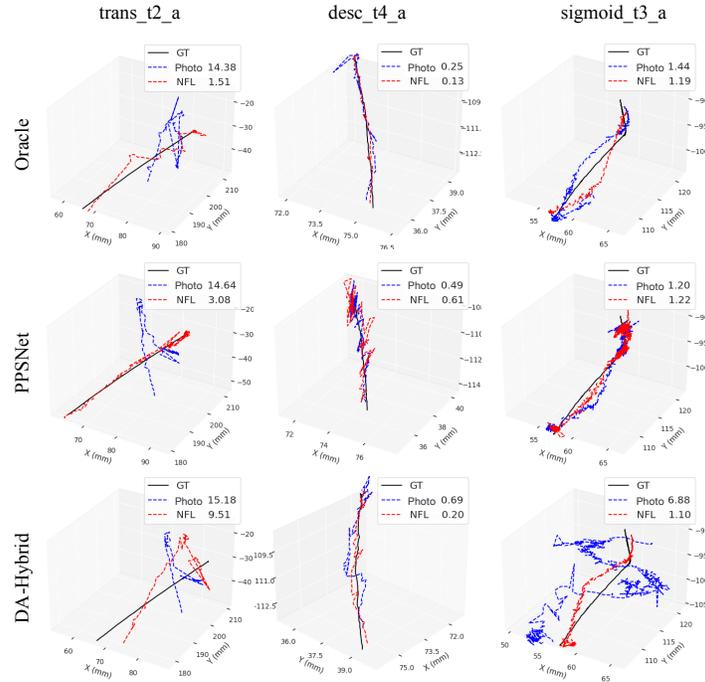


Figure 6: **Camera tracking improvement over MonoGS [27]**. Replacing the Photo BA loss (in blue) with NFL-BA loss (in red) significantly improves camera tracking for different depth initialization. Average tracking error  $ATE_t$  for each sequence is reported in the inset.

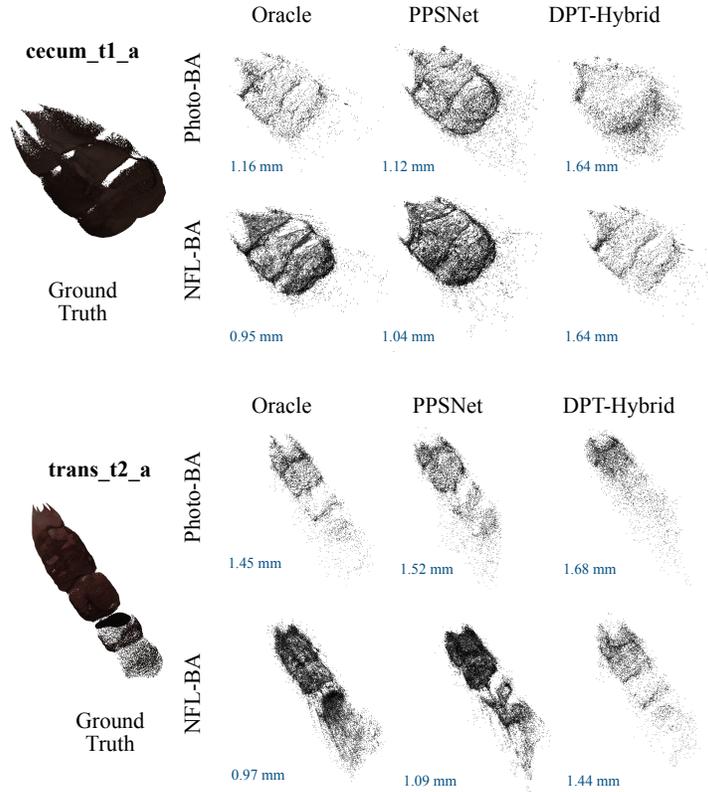


Figure 7: **Reconstructed point clouds using MonoGS [27]** show that NFL-BA improves coverage and density while reducing scatter compared to Photometric BA, as measured by Chamfer distance.

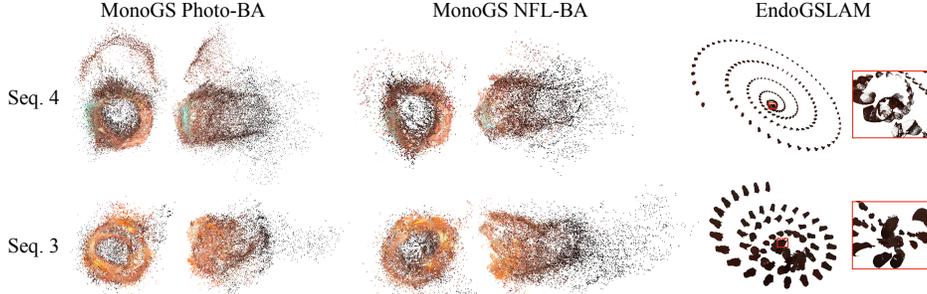


Figure 8: **Results on real endoscopy from Colon10k dataset.** On Sequences 3 and 4 with PPSNet depth, NFL-BA improves MonoGS tracking and mapping, yielding more coherent, elongated colon structures, while EndoGSLAM fails under sudden camera motion.

## 5.2 Evaluation on C3VD Endoscopy Data

Because NFL-BA is designed as a drop-in replacement for photometric bundle adjustment, we only adjusted the two associated loss weights; all other hyperparameters remain identical between the Photo-BA and NFL-BA experiments. Please see supplemental for detailed hyperparameter settings.

**SLAM with oracle depth map.** In Tab. 1 we replace Photometric Bundle Adjustment loss with NFL-BA loss for depth is initialized with ground-truth or oracle. NFL-BA significantly improves camera localization ( $ATE_t$ ) and mapping for NICE-SLAM and MonoGS, and only camera rotation ( $ATE_r$ ) for EndoGSLAM. EndoGSLAM was specifically designed for synthetic data with an oracle depth map, and we will show later that for estimated depth maps or real endoscopy videos, it performs significantly worse than MonoGS Ground-truth depths are never available during endoscopy, and the majority of endoscopes hardly have any depth sensors.

**SLAM with predicted depth map.** Under realistic conditions with estimated depths, NFL-BA’s impact is even more pronounced. In Tab. 2 we replace Photometric BA loss with NFL-BA loss for MonoGS [27] and EndoGSLAM [45] for depth maps we use PPSNet [35], a state-of-the-art monocular depth estimation algorithm for endoscopy, and fine-tuned general-purpose depth estimator, which we will call it as DA-Hybrid - DepthAnything[50] with DINOv2 encoder [32]. NFL-BA significantly improves camera localization ( $ATE_t$ ) and camera rotation ( $ATE_r$ ) for both MonoGS [27] and EndoGSLAM [45] while producing similar rendering quality. For example, camera localization for MonoGS is improved by 37% for PPSNet and 49% for DA-Hybrid depth initialization. Mapping accuracy of MonoGS also improves by 37% for PPSNet and 16% for DA-Hybrid depth maps. Overall, these results demonstrate that NFL-BA can compensate for noisy depth estimation and improves performance. Across all four metrics, for tracking, mapping, and rendering, the SOTA performance on the C3VD dataset is in fact achieved when NFL-BA loss is used in the SLAM framework.

## 5.3 Evaluation on Real Endoscopy Data

We show results on real endoscopy sequence from Colon10k sequence 3 and 4 in Fig. 8. EndoGSLAM fails to construct any real structure, with many disconnected regions along a spiral trajectory. EndoGSLAM assumes constant velocity and is not robust to the sudden motion common in endoscopy procedures, which is significantly more in real data than C3VD. This results extremely poor or failed reconstructions.

**Sequence 4.** This pull-back “down-the-barrel” sequence exposes a clear cylindrical lumen. With Photo-BA, MonoGS captures the overall shape but produces a broken segment due to trajectory drift. NFL-BA corrects this, yielding a continuous “hollow-center” reconstruction. Minor artifacts from extreme specular highlights remain (green points), as detailed in the supplement.

**Sequence 3.** In the extended traversal, both Photo-BA and NFL-BA recover the colon’s general geometry, but NFL-BA produces a longer, tighter model with less point scatter. It also better preserves interior ridges (interactive point clouds in the supplement).

## 5.4 Evaluation on Indoor Scene

To validate NFL-BA in a non-medical setting, we evaluate on four indoor scenes. Table 3 shows that replacing standard Photometric BA with NFL-BA yields substantial reductions in  $ATE_t$  across

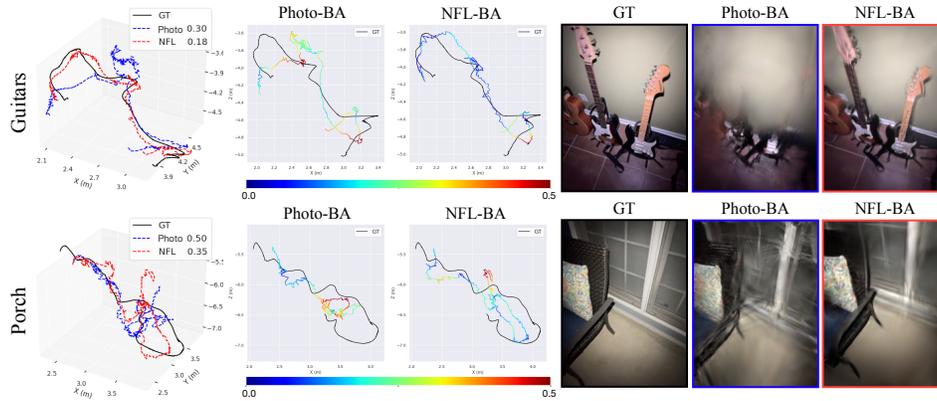


Figure 9: **Results on indoor scenes captured with co-located flashlight and phone camera.** Qualitative comparison on two self-captured indoor scenes using MonoGS with standard Photo-BA versus NFL-BA. (left) Estimated camera trajectories overlaid on ground truth. (center) Per-frame tracking error relative to ground truth. (right) Example re-rendered views, illustrating the sharper, more accurate reconstructions enabled by NFL-BA.

Table 3: Quantitative results on four self-captured indoor scenes under dynamic lighting, comparing MonoGS with standard Photo-BA versus NFL-BA. For each scene, the best of each metric is bold.

BA	Guitars		Porch		Pool		Stairs	
	ATE <sub>t</sub> (m)↓	LPIPS ↓						
Photo	0.30	0.39	0.50	<b>0.49</b>	0.41	0.46	0.36	0.40
NFL	<b>0.18</b>	<b>0.37</b>	<b>0.35</b>	0.50	<b>0.30</b>	<b>0.44</b>	<b>0.20</b>	<b>0.31</b>

all scenes: from 0.30m to 0.18m (40%) in *Guitar*, 0.50m to 0.35m (30%) in *Outdoor*, 0.41m to 0.30m (27%) in *Pool*, and 0.36m to 0.20m (44%) in *Stair*. On average, NFL-BA reduces tracking error by  $\sim 35\%$ , demonstrating that near-field shading cues greatly enhance pose estimation even in richly textured, well-lit indoor environments. While LPIPS remains largely comparable, with slight improvements in *Guitars* and *Stairs* and minor variations in *Porch* and *Pool*, the primary benefit of NFL-BA is clear in trajectory accuracy (see Fig. 9).

## 6 Conclusions

In this paper, we presented a novel bundle adjustment loss that explicitly models dynamic near-field lighting by incorporating light intensity fall-off based on the relative distance and orientation between the surface and the co-located light and camera. This formulation is especially effective for endoscopic scenes, where traditional geometric or photometric bundle adjustment losses struggle under dynamic near-field lighting conditions on textureless surfaces. We demonstrated the general applicability of our approach by integrating it into three different neural rendering-based SLAM methods, improving performance on a challenging endoscopy dataset and indoor scenes captured with a phone camera with a flashlight turned on.

**Limitations.** While our new formulation for SLAM effectively represents scenes with co-located and dynamic lighting environments, it is currently limited in handling specular reflections, sub-surface scattering, and inter-reflections. Incorporating a more complex image formulation is beyond the scope of the current work, and addressing these remains a promising direction for future research.

## 7 Acknowledgments

This work is supported by a National Institute of Health (NIH) project #R21EB035832 "Next-gen 3D Modeling of Endoscopy Videos" and #R21EB037440 "Gen-AI Airway Simulator for 3D Endoscopy". We also thank Stephen M. Pizer, Ron Alterovitz, and Dr. Sarah McGill for helpful discussions during the project.

## References

- [1] Hatem Alismail, Brett Browning, and Simon Lucey. Photometric bundle adjustment for vision-based slam, 2016.
- [2] Dejan Azinović, Ricardo Martin-Brualla, Dan B Goldman, Matthias Nießner, and Justus Thies. Neural rgb-d surface reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6290–6301, June 2022.
- [3] Víctor M Batlle, José MM Montiel, Pascal Fua, and Juan D Tardós. LightNeuS: Neural surface reconstruction in endoscopy using illumination decline. In *Medical Image Computing and Computer Assisted Intervention (MICCAI)*, 2023.
- [4] 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.
- [5] Mark Boss, Raphael Braun, Varun Jampani, Jonathan T. Barron, Ce Liu, and Hendrik P.A. Lensch. NerD: Neural reflectance decomposition from image collections. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 12684–12694, October 2021.
- [6] C. Campos, R. Elvira, J. J. G. Rodríguez, J. M. M. Montiel, and J. D. Tardós. Orb-slam3: An accurate open-source library for visual, visual-inertial, and multi-map slam. *IEEE Transactions on Robotics*, 37(6):1874–1890, 2021.
- [7] Devendra Singh Chaplot, Ruslan Salakhutdinov, Abhinav Gupta, and Saurabh Gupta. Neural topological slam for visual navigation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 12875–12884, 2020.
- [8] Hanlin Chen, Fangyin Wei, Chen Li, Tianxin Huang, Yunsong Wang, and Gim Hee Lee. Vcr-gaus: View consistent depth-normal regularizer for gaussian surface reconstruction. *arXiv preprint arXiv:2406.05774*, 2024.
- [9] B. Cui, H. Zhang, X. Li, W. Zhou, and T. Cheng. Endodac: Efficient adapting foundation model for self-supervised depth estimation from any endoscopic camera. In *Proceedings of the Medical Image Computing and Computer-Assisted Intervention Conference (MICCAI)*, 2024.
- [10] Tianchen Deng, Yaohui Chen, Leyan Zhang, Jianfei Yang, Shenghai Yuan, Jiuming Liu, Danwei Wang, Hesheng Wang, and Weidong Chen. Compact 3d gaussian splatting for dense visual slam, 2024.
- [11] J. Engel, V. Koltun, and D. Cremers. Direct sparse odometry. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 40(3):611–625, 2018.
- [12] Oscar G. Grasa, Javier Civera, and J. M. M. Montiel. Ekf monocular slam with relocalization for laparoscopic sequences. In *2011 IEEE International Conference on Robotics and Automation*, pages 4816–4821, 2011.
- [13] Jiaxin Guo, Jiangliu Wang, Di Kang, Wenzhen Dong, Wenting Wang, and Yun-hui Liu. Free-SurGS: SfM-Free 3D Gaussian Splatting for Surgical Scene Reconstruction. In *proceedings of Medical Image Computing and Computer Assisted Intervention – MICCAI 2024*, volume LNCS 15007. Springer Nature Switzerland, October 2024.
- [14] R. Hartley and A. Zisserman. *Multiple View Geometry in Computer Vision*. Cambridge University Press, 2003.
- [15] Michel Hayoz, Christopher Hahne, Thomas Kurmann, Max Allan, Guido Beldi, Daniel Candinas, Pablo Márquez-Neila, and Raphael Sznitman. Online 3D reconstruction and dense tracking in endoscopic videos. In *proceedings of Medical Image Computing and Computer Assisted Intervention – MICCAI 2024*, volume LNCS 15006. Springer Nature Switzerland, October 2024.
- [16] Huajian Huang, Longwei Li, Cheng Hui, and Sai-Kit Yeung. Photo-slam: Real-time simultaneous localization and photorealistic mapping for monocular, stereo, and rgb-d cameras. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024.
- [17] Yiming Huang, Beilei Cui, Long Bai, Ziqi Guo, Mengya Xu, Mobarakol Islam, and Hongliang Ren. Endo-4DGS: Endoscopic Monocular Scene Reconstruction with 4D Gaussian Splatting. In *proceedings of Medical Image Computing and Computer Assisted Intervention – MICCAI 2024*, volume LNCS 15006. Springer Nature Switzerland, October 2024.

- [18] Raúl Iranzo, Víctor M Batlle, Juan D Tardós, and José MM Montiel. Endometric: Near-light metric scale monocular slam. *arXiv preprint arXiv:2410.15065*, 2024.
- [19] Y. Iwahori, H. Sugie, and N. Ishii. Reconstructing shape from shading images under point light source illumination. In *[1990] Proceedings. 10th International Conference on Pattern Recognition*, volume i, pages 83–87 vol.1, 1990.
- [20] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkuehler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Trans. Graph.*, 42(4), July 2023.
- [21] Daniel Lichy, Soumyadip Sengupta, and David W. Jacobs. Fast light-weight near-field photometric stereo, 2022.
- [22] Philipp Lindenberger, Paul-Edouard Sarlin, and Marc Pollefeys. LightGlue: Local Feature Matching at Light Speed. In *ICCV*, 2023.
- [23] Yifan Liu, Chenxin Li, Chen Yang, and Yixuan Yuan. Endogaussian: Gaussian splatting for deformable surgical scene reconstruction. *arXiv preprint arXiv:2401.12561*, 2024.
- [24] Ruibin Ma, Sarah K. McGill, Rui Wang, Julian Rosenman, Jan-Michael Frahm, Yubo Zhang, and Stephen Pizer. Colon10k: A benchmark for place recognition in colonoscopy. In *2021 IEEE 18th International Symposium on Biomedical Imaging (ISBI)*, pages 1279–1283, 2021.
- [25] Ruibin Ma, Rui Wang, Stephen Pizer, Julian Rosenman, Sarah K McGill, and Jan-Michael Frahm. Real-time 3d reconstruction of colonoscopic surfaces for determining missing regions. In *Medical Image Computing and Computer Assisted Intervention–MICCAI 2019: 22nd International Conference, Shenzhen, China, October 13–17, 2019, Proceedings, Part V 22*, pages 573–582. Springer, 2019.
- [26] R. Martínez-Cantín and J. D. Tardós. Unscented kalman filter for visual slam with photometric calibration. In *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 1234–1241, 2016.
- [27] Hidenobu Matsuki, Riku Murai, Paul H.J. Kelly, and Andrew J. Davison. Gaussian splatting slam. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 18039–18048, June 2024.
- [28] 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. *Commun. ACM*, 65(1):99–106, December 2021.
- [29] Javier Morlana, Juan D Tardós, and José MM Montiel. Topological slam in colonoscopies leveraging deep features and topological priors. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 733–743. Springer, 2024.
- [30] R. Mur-Artal, J. M. M. Montiel, and J. D. Tardos. Orb-slam: A versatile and accurate monocular slam system. In *IEEE Transactions on Robotics*, volume 31, pages 1147–1163, 2015.
- [31] Andriy Myronenko and Xubo Song. Point set registration: Coherent point drift. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 32(12):2262–2275, 2010.
- [32] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy V. Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel HAZIZA, Francisco Massa, Alaaeldin El-Nouby, Mido Assran, Nicolas Ballas, Wojciech Galuba, Russell Howes, Po-Yao Huang, Shang-Wen Li, Ishan Misra, Michael Rabbat, Vasu Sharma, Gabriel Synnaeve, Hu Xu, Herve Jegou, Julien Mairal, Patrick Labatut, Armand Joulin, and Piotr Bojanowski. DINOv2: Learning robust visual features without supervision. *Transactions on Machine Learning Research*, 2024. Featured Certification.
- [33] Kutsev Bengisu Ozyoruk, Guliz Irem Gokceler, Taylor L. Bobrow, Gulfize Coskun, Kagan Incetan, Yasin Almalioglu, Faisal Mahmood, Eva Curto, Luis Perdigoto, Marina Oliveira, Hasan Sahin, Helder Araujo, Henrique Alexandrino, Nicholas J. Durr, Hunter B. Gilbert, and Mehmet Turan. Endoslam dataset and an unsupervised monocular visual odometry and depth estimation approach for endoscopic videos. *Medical Image Analysis*, 71:102058, 2021.
- [34] Albert Palomer, Pere Ridao, and David Ribas. Inspection of an underwater structure using point-cloud slam with an auv and a laser scanner. *Journal of field robotics*, 36(8):1333–1344, 2019.
- [35] Akshay Paruchuri, Samuel Ehrenstein, Shuxian Wang, Inbar Fried, Stephen M Pizer, Marc Niethammer, and Roni Sengupta. Leveraging near-field lighting for monocular depth estimation from endoscopy videos. In *Computer Vision – ECCV 2024*, Cham, 2024. Springer Nature Switzerland.

- [36] Juan J Gómez Rodríguez, José MM Montiel, and Juan D Tardós. Nr-slam: Non-rigid monocular slam. *IEEE Transactions on Robotics*, 2024.
- [37] Javier Rodríguez-Puigvert\*, Víctor M. Batlle\*, José María M. Montiel, Rubén Martínez-Cantín, Pascal Fua, Juan D. Tardós, and Javier Civera. LightDepth: single-view depth self-supervision from illumination decline. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2023.
- [38] Erik Sandström, Yue Li, Luc Van Gool, and Martin R Oswald. Point-slam: Dense neural point cloud-based slam. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 18433–18444, 2023.
- [39] Yufei Shi, Beijia Lu, Jia-Wei Liu, Ming Li, and Mike Zheng Shou. Colonerf: Neural radiance fields for high-fidelity long-sequence colonoscopy reconstruction. *arXiv preprint arXiv:2312.02015*, 2023.
- [40] D. Stoyanov, G. Mylonas, F. Deligianni, A. Darzi, and G.-Z. Yang. Soft-tissue motion tracking and structure estimation for robotic assisted mis procedures. In *Medical Image Computing and Computer-Assisted Intervention (MICCAI)*, pages 139–146, 2005.
- [41] E. Sucar, S. Liu, J. Ortiz, and A. J. Davison. imap: Implicit mapping and positioning in real-time. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 6229–6238, 2021.
- [42] Zachary Teed and Jia Deng. Droid-slam: Deep visual slam for monocular, stereo, and rgb-d cameras. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, volume 34, pages 16558–16569. Curran Associates, Inc., 2021.
- [43] Fabio Tosi, Youmin Zhang, Ziren Gong, Erik Sandström, Stefano Mattocchia, Martin R Oswald, and Matteo Poggi. How nerfs and 3d gaussian splatting are reshaping slam: a survey. *arXiv preprint arXiv:2402.13255*, 4, 2024.
- [44] Hengyi Wang, Jingwen Wang, and Lourdes Agapito. Co-slam: Joint coordinate and sparse parametric encodings for neural real-time slam. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 13293–13302, June 2023.
- [45] 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 *proceedings of Medical Image Computing and Computer Assisted Intervention – MICCAI 2024*, volume LNCS 15006. Springer Nature Switzerland, October 2024.
- [46] Shuxian Wang, Yubo Zhang, Sarah K. McGill, Julian G. Rosenman, Jan-Michael Frahm, Soumyadip Sengupta, and Stephen M. Pizer. A surface-normal based neural framework for colonoscopy reconstruction. In Alejandro Frangi, Marleen de Bruijne, Demian Wassermann, and Nassir Navab, editors, *Information Processing in Medical Imaging*, pages 797–809, Cham, 2023. Springer Nature Switzerland.
- [47] Qianyi Wu, Jianmin Zheng, and Jianfei Cai. Surface reconstruction from 3d gaussian splatting via local structural hints. In *Computer Vision – ECCV 2024*, Cham, 2024. Springer Nature Switzerland.
- [48] Kuan Xu, Yuefan Hao, Shenghai Yuan, Chen Wang, and Lihua Xie. AirSLAM: An efficient and illumination-robust point-line visual slam system. *IEEE Transactions on Robotics (TRO)*, 2024.
- [49] Chi Yan, Delin Qu, Dan Xu, Bin Zhao, Zhigang Wang, Dong Wang, and Xuelong Li. Gs-slam: Dense visual slam with 3d gaussian splatting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19595–19604, 2024.
- [50] Lihe Yang, Bingyi Kang, Zilong Huang, Xiaogang Xu, Jiashi Feng, and Hengshuang Zhao. Depth anything: Unleashing the power of large-scale unlabeled data. In *CVPR*, 2024.
- [51] Shuoqun Yang, Qian Li, Daiyun Shen, Bingchen Gong, Qi Dou, and Yueming Jin. Deform3DGS: Flexible Deformation for Fast Surgical Scene Reconstruction with Gaussian Splatting . In *proceedings of Medical Image Computing and Computer Assisted Intervention – MICCAI 2024*, volume LNCS 15006. Springer Nature Switzerland, October 2024.
- [52] Keyang Ye, Qiming Hou, and Kun Zhou. 3d gaussian splatting with deferred reflection. In *ACM SIGGRAPH 2024 Conference Papers*, SIGGRAPH '24, New York, NY, USA, 2024. Association for Computing Machinery.
- [53] Vladimir Yugay, Yue Li, Theo Gevers, and Martin R. Oswald. Gaussian-slam: Photo-realistic dense slam with gaussian splatting, 2024.

- [54] Richard Zhang, Phillip Isola, Alexei A. Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
- [55] Yubo Zhang, Jan-Michael Frahm, Samuel Ehrenstein, Sarah K McGill, Julian G Rosenman, Shuxian Wang, and Stephen M Pizer. Colde: a depth estimation framework for colonoscopy reconstruction. *arXiv preprint arXiv:2111.10371*, 2021.
- [56] S. Zhi, J. Lai, A. Kundu, M. Bloesch, A. Davison, and A. Zisserman. Nerf-slam: Real-time dense monocular slam with neural radiance fields. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2022.
- [57] L. Zhou, R. Klette, and K. Scheibe. A multi-image shape-from-shading framework for near-lighting perspective endoscopes. *International Journal of Computer Vision*, 82:1–24, 2009.
- [58] A. Zhu, Z. Zhang, H. Su, L. Li, G. Wang, and X. Li. Nice-slam: Neural implicit scalable encoding for slam. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022.
- [59] Lingting Zhu, Zhao Wang, Jiahao Cui, Zhenchao Jin, Guying Lin, and Lequan Yu. Endogs: Deformable endoscopic tissues reconstruction with gaussian splatting. In M. Emre Celebi, Mauricio Reyes, Zhen Chen, and Xiaoxiao Li, editors, *Medical Image Computing and Computer Assisted Intervention – MICCAI 2024 Workshops*, pages 135–145, Cham, 2025. Springer Nature Switzerland.
- [60] Zihan Zhu, Songyou Peng, Viktor Larsson, Zhaopeng Cui, Martin R Oswald, Andreas Geiger, and Marc Pollefeys. Nicer-slam: Neural implicit scene encoding for rgb slam. In *International Conference on 3D Vision (3DV)*, March 2024.

---

# Supplementary Materials

---

## Overview of Appendices

Our appendices contain the following additional details:

- Appendix A: Dataset Processing and Collection
- Appendix B: Implementation Details and Computational Costs
- Appendix C: Point Clouds and Metrics
- Appendix D: Long Sequence Validation
- Appendix E: Evaluation of Center-Crop Baseline
- Appendix F: Lighting model and Angular Attenuation
- Appendix G: Results on Endoscopy Datasets

## A Dataset Processing and Collection

**C3VD.** We test our method on a subset of 8 videos with at least one video from each section of the colon, with varying camera motion, and anomalies. We choose these 8 videos from the test split of PPSNet to avoid any bias when the SLAM is initialized with PPSNet predicted depth map. The names of the sequences are as follows: cecum\_t1\_a, cecum\_t2\_a, cecum\_t3\_a, sigmoid\_t3\_a, desc\_t4\_a\_p2, trans\_t2\_a, trans\_t3\_a, and trans\_t4\_a.

All images were cropped to remove any artifacts resulting from fish-eye correction then downscale and crop the images. Specifically, we resize each image to have a height of 384 pixels while maintaining the aspect ratio, then crop the central region to obtain a 384×384 pixel image.

**Colon10k.** Since Colon10K provides no ground-truth depths, we compute per-frame estimates with PPSNet. Each image is center-cropped and uniformly resized to 384 × 384 px to match our SLAM input requirements.

**Indoor Self Captures.** We recorded four indoor scenes using an iPhone 15 with LiDAR, capturing synchronized RGB (1440 × 1920 px) and depth (256 × 192 px) streams. Raw RGB frames are downsampled to 256 × 192 px to align with the depth map resolution. Depth maps are stored as 16-bit values up to 10 m. We logged camera poses via Apple’s ARKit framework, code for our custom capture app and preprocessing scripts will be released alongside the dataset.

## B Implementation Details and Computational Costs

As mentioned in the main paper, we ran all models on a single NVIDIA RTX A6000 GPU. The per-scene optimization takes approximately 1 FPS. We found that NFL-BA only reduces the fps runtime by a small amount on the C3VD dataset.

**NICE-SLAM.** Since the scene is encoded using neural networks, we extract normals from the occupancy grid, as described in Sec. 4.2 to calculate the shading term. NICE-SLAM requires a well-defined bounding box which we obtained from the ground truth point clouds (see appendix C). We also used the default loss weights of NICE-SLAM, setting  $\lambda_{ren}$  to 0.5 during tracking and 0.2 during mapping, and  $\lambda_{geo}$  to 1 in both phases.

**EndoGSLAM.** The main difference is the weight map  $M_t$  in the bundle adjustment loss (Eq. 3) to exclude over-exposed pixels that can arise in endoscopy-specific lighting conditions. Furthermore, given that the shading term  $PPS(\cdot)$  is sensitive to depth scales, we rescaled the depth maps so that their maximum values are approximately 5. Notably, scaling the depth maps did not improve baseline performance (when using PPS depth maps, the average  $ATE_t$  went from 3.03 to 3.38). We used the

Table 4: Runtime: We evaluate the frames per second (**FPS**) for all methods on the C3VD dataset.

Method	Depth	Photo-BA	NFL-BA
NICE-SLAM	Oracle	$\ll 1$	$\ll 1$
	PPSNet	$\ll 1$	$\ll 1$
EndoGSLAM	Oracle	1.79	1.35
	PPSNet	1.53	1.22
	DPT-Hybrid	1.20	0.90
MonoGS	Oracle	1.38	1.09
	PPSNet	1.06	0.93
	DPT-Hybrid	0.99	0.83

default loss weights of EndoGSLAM,  $\lambda_{ren}$ ;  $\lambda_{geo}$ , set to 0.5 and 1 during tracking and 1 and 1 during mapping, respectively.

**MonoGS.** To integrate our method, we treat the Gaussian color features as albedo features and multiply them with the shading term before rasterization, and then we use the rendered output colors for bundle adjustment. For all input depths, we set  $\lambda_{ren}$  and  $\lambda_{geo}$  to 0.8 and 0.5, respectively.

## C Point Clouds and Metrics

**Chamfer Distances.** Since ground truth point clouds are unavailable for C3VD, we generate them by unprojecting 2D images into 3D space with the correct camera configuration and the oracle depth maps, provided in the C3VD dataset. For neural fields-based SLAM, we use the vertices of the output meshes as the estimated point clouds, while for Gaussian Splatting-based SLAMs, we use the Gaussian positions. For point cloud alignment, we use Coherent Point Drift [31] and the Chamfer distances are also in millimeters.

**Coloring Point Clouds.** We use extracted point clouds from 3D Gaussian positions for visualization. For colors, we directly use Gaussian color features.

## D Long Sequence Validation

To evaluate NFL-BA’s performance on longer sequence, we test our method on a longer screening video for C3VD, specifically **c0\_full\_t2\_v2**, which spans over 4,000 frames with ground-truth poses, allowing us to measure cumulative drift. In the main paper, we only evaluate on registered videos with ground truth depth, not screening videos. We show how NFL-BA performs relative to Photo-BA using the MonoGS backbone in Table 5.

NFL-BA matches Photometric BA on short sequences and increasingly outperforms it as trajectory length grows. As a plug-and-play bundle-adjustment loss, NFL-BA enhances long-term robustness under dynamic near-field lighting.

Table 5: ATE<sub>T</sub> (cm) on c0\_full\_t2\_v2 at varying lengths. NFL-BA matches or beats Photo-BA on short sequences and shows long-term stability.

BA	500 frames	1,000 frames	2,000 frames	4,000 frames
Photo	2.599	1.937	8.459	14.790
NFL	2.422	2.254	5.742	11.368

## E Evaluation of Center-Crop Baseline

We compared MonoGS Photo-BA optimized on only the central 75% and 50% of each frame against full-frame MonoGS + NFL-BA on two sequences, measuring translational ATE (ATE<sub>T</sub>), Chamfer distance (CD), and reconstructed point count in Table 6.

Although center-cropping improves camera tracking performance (ATE\_T) of MonoGS and gets close to NFL-BA, it results in significantly worse reconstruction in terms of quality (CD) and density (number of points). This is because for camera localization, not all pixels are essential, and focusing only on central pixels can eliminate near-field lighting effects. In contrast, for reconstruction, all pixels matter. This highlights the need for a principled mechanism for handling dynamic near-field lighting, as proposed by NFL-BA.

Table 6: Center-cropping improves trajectory error but reduces map density, while full-frame NFL-BA maintains both accuracy and dense reconstructions.

Method	trans_t3_a			desc_t4_p2		
	ATE <sub>T</sub>	CD	Points	ATE <sub>T</sub>	CD	Points
MonoGS + NFL-BA	0.26	0.60	28,196	0.13	0.67	8,879
MonoGS (full frame)	0.31	0.76	7,095	0.25	0.76	6,398
MonoGS (75% center)	0.35	1.06	5,761	0.13	0.97	4,329
MonoGS (50% center)	0.40	2.82	1,865	0.15	1.56	2,100

## F Lighting Model and Angular Attenuation

We clarify that NFL-BA models only diffuse reflectance and direct illumination, ignoring specular and subsurface effects. Explicitly modeling these introduces non-differentiable and highly non-convex terms, remaining an open challenge.

To mitigate specular highlights, we apply an intensity mask discarding pixels above 0.9 grayscale intensity. This handles most artifacts but cannot correct reflective surfaces beyond the mask. While effective on matte datasets (C3VD, in-the-wild), degradation occurs on Colon10K with more specular highlights. Future work will address specular reflections under dynamic lighting.

Regarding angular attenuation,  $\beta = 0$  is justified for tightly collimated endoscopic LEDs, we must validate this for non endoscopy scenes. Using NFL-BA on the MonoGS backbone, we validate various  $\beta$  values on two indoor sequences as seen in Table 7

$\beta = 0$  provides competitive mean accuracy and low variance. Non-zero values may yield small gains but introduce instability. Given minimal benefit versus tuning cost, we retain  $\beta = 0$ , and plan to explore learning  $\beta$  as a per-scene parameter in future work.

Table 7: Ablation of the angular-attenuation coefficient  $\beta$  on two indoor sequences, reporting translational ATE and LPIPS.

Scene	Metric	$\beta=0.00$	$\beta=0.25$	$\beta=0.50$	$\beta=0.75$	$\beta=1.00$
Guitars	ATE_T	0.17±0.01	0.17±0.003	0.16±0.006	0.17±0.006	0.17±0.011
	LPIPS	0.36±0.02	0.37±0.006	0.40±0.030	0.36±0.017	0.36±0.020
Porch	ATE_T	0.33±0.16	0.25±0.013	0.27±0.045	0.36±0.158	0.22±0.012
	LPIPS	0.50±0.03	0.53±0.008	0.48±0.005	0.51±0.035	0.49±0.010

## G Results on Endoscopy Datasets

For all experiments for each of the slam systems, we report the median of three runs for all tables and figures. We have included per-sequence metrics for the median run below. Additionally, we include point clouds for EndoGSLAM in Figure 10, more can be found on our project page.

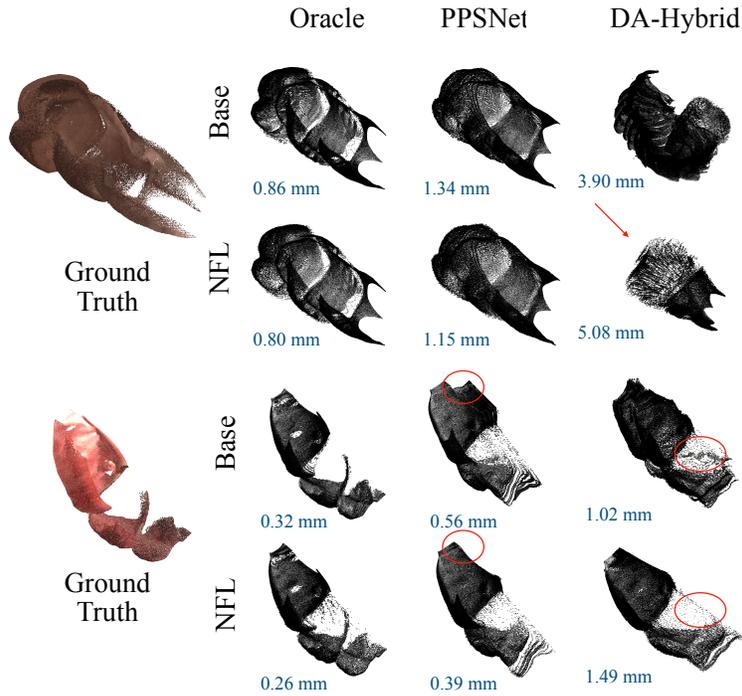


Figure 10: EndoGSLAM Point cloud results for 2 sequences: cecum\_t1\_a\_under\_review (top) and desc\_t4\_a\_p2\_under\_review (bottom).

Table 8: Results for sequence **cecum\_t1\_a\_under\_review**

Method	Depth	BA	$ATE_t$ (mm)↓	$ATE_r$ (°)↓	Chamfer (mm)↓
NICE-SLAM	Oracle	Photo	2.65	2.82	0.08
		NFL	1.39	2.81	0.15
	PPS-Net	Photo	10.14	2.71	0.51
		NFL	4.14	2.75	0.19
EndoGSLAM	Oracle	Photo	1.39	0.24	0.12
		NFL	0.91	0.31	0.16
	PPS-Net	Photo	2.79	0.60	0.30
		NFL	2.93	0.70	0.35
	DA-Hybrid	Photo	2.64	0.08	0.04
		NFL	8.44	2.42	1.21
MonoGS	Oracle	Photo	1.16	0.39	1.56
		NFL	1.22	0.35	0.95
	PPS-Net	Photo	2.30	0.65	1.19
		NFL	2.45	0.76	1.04
	DA-Hybrid	Photo	5.44	0.30	1.64
		NFL	1.09	0.40	1.65

Table 9: Results for sequence **cecum\_t2\_a\_under\_review**

Method	Depth	BA	$ATE_t$ (mm)↓	$ATE_r$ (°)↓	Chamfer (mm)↓
NICE-SLAM	Oracle	Photo	8.13	2.19	0.49
		NFL	1.15	2.82	0.11
	PPS-Net	Photo	1.11	2.13	0.11
		NFL	7.98	2.82	0.50
EndoGSLAM	Oracle	Photo	2.32	1.06	0.53
		NFL	6.75	2.79	1.40
	PPS-Net	Photo	4.55	1.18	0.59
		NFL	4.55	2.82	1.41
	DA-Hybrid	Photo	5.66	0.89	0.45
		NFL	9.47	2.25	1.13
MonoGS	Oracle	Photo	4.04	0.34	1.47
		NFL	6.84	0.48	1.24
	PPS-Net	Photo	6.72	2.71	1.74
		NFL	6.52	2.73	1.70
	DA-Hybrid	Photo	5.20	0.66	2.18
		NFL	4.26	0.32	1.35

Table 10: Results for sequence **cecum\_t3\_a\_under\_review**

Method	Depth	BA	$ATE_t$ (mm)↓	$ATE_r$ (°)↓	Chamfer (mm)↓
NICE-SLAM	Oracle	Photo	3.46	2.82	0.11
		NFL	3.42	2.82	0.07
	PPS-Net	Photo	2.80	2.64	0.17
		NFL	3.83	2.59	0.14
EndoGSLAM	Oracle	Photo	0.75	0.15	0.08
		NFL	0.74	0.27	0.14
	PPS-Net	Photo	0.79	0.16	0.08
		NFL	2.18	0.22	0.11
	DA-Hybrid	Photo	1.45	0.62	0.31
		NFL	1.52	0.18	0.09
MonoGS	Oracle	Photo	0.36	0.16	1.23
		NFL	0.87	0.31	0.62
	PPS-Net	Photo	1.07	0.17	1.12
		NFL	1.24	0.14	0.81
	DA-Hybrid	Photo	1.06	0.33	0.95
		NFL	1.16	0.30	0.93

Table 11: Results for sequence **desc\_t4\_a\_p2\_under\_review**

Method	Depth	BA	$ATE_t$ (mm)↓	$ATE_r$ (°)↓	Chamfer (mm)↓
NICE-SLAM	Oracle	Photo	15.90	2.65	0.59
		NFL	8.88	2.75	0.50
	PPS-Net	Photo	0.87	2.80	0.10
		NFL	0.68	2.80	0.10
EndoGSLAM	Oracle	Photo	0.36	0.56	0.28
		NFL	0.26	0.94	0.47
	PPS-Net	Photo	0.48	0.58	0.29
		NFL	1.00	0.83	0.42
	DA-Hybrid	Photo	0.45	1.16	0.58
		NFL	1.80	2.81	1.40
MonoGS	Oracle	Photo	0.25	1.05	0.76
		NFL	0.13	0.98	0.67
	PPS-Net	Photo	0.49	1.78	0.77
		NFL	0.61	1.67	0.70
	DA-Hybrid	Photo	0.69	2.63	0.80
		NFL	0.20	0.97	0.77

Table 12: Results for sequence **sigmoid\_t3\_a\_under\_review**

Method	Depth	BA	$ATE_t$ (mm)↓	$ATE_r$ (°)↓	Chamfer (mm)↓
NICE-SLAM	Oracle	Photo	1.04	2.47	0.05
		NFL	0.84	2.83	0.05
	PPS-Net	Photo	9.47	2.79	0.28
		NFL	6.60	2.55	0.35
EndoGSLAM	Oracle	Photo	4.81	1.14	0.57
		NFL	4.11	2.69	1.35
	PPS-Net	Photo	3.13	1.40	0.70
		NFL	9.82	2.56	1.29
	DA-Hybrid	Photo	8.51	2.36	1.18
		NFL	8.15	2.80	1.41
MonoGS	Oracle	Photo	1.44	0.46	1.17
		NFL	1.19	2.33	0.67
	PPS-Net	Photo	1.20	1.53	4.63
		NFL	1.22	2.22	0.89
	DA-Hybrid	Photo	6.88	2.18	1.42
		NFL	1.10	0.67	1.33

Table 13: Results for sequence **trans\_t2\_a\_under\_review**

Method	Depth	BA	$ATE_t$ (mm)↓	$ATE_r$ (°)↓	Chamfer (mm)↓
NICE-SLAM	Oracle	Photo	0.77	2.83	0.05
		NFL	0.94	2.80	0.05
	PPS-Net	Photo	9.18	2.82	0.26
		NFL	9.85	2.81	0.31
EndoGSLAM	Oracle	Photo	4.21	1.80	0.90
		NFL	0.49	2.82	1.41
	PPS-Net	Photo	10.05	1.66	0.84
		NFL	0.90	1.46	0.73
	DA-Hybrid	Photo	9.14	2.77	1.39
		NFL	14.59	2.18	1.10
MonoGS	Oracle	Photo	14.38	1.42	1.49
		NFL	1.51	2.72	0.97
	PPS-Net	Photo	14.64	1.91	1.52
		NFL	3.08	1.02	1.09
	DA-Hybrid	Photo	15.18	2.12	1.68
		NFL	9.51	1.41	1.45

Table 14: Results for sequence **trans\_t3\_a\_under\_review**

Method	Depth	BA	$ATE_t$ (mm)↓	$ATE_r$ (°)↓	Chamfer (mm)↓
NICE-SLAM	Oracle	Photo	0.97	2.81	0.15
		NFL	6.26	2.79	0.68
	PPS-Net	Photo	3.78	2.80	0.22
		NFL	1.12	2.65	0.13
EndoGSLAM	Oracle	Photo	0.17	2.70	1.35
		NFL	0.20	2.70	1.35
	PPS-Net	Photo	0.30	2.80	1.40
		NFL	0.50	2.81	1.41
	DA-Hybrid	Photo	0.46	2.81	1.41
		NFL	0.33	2.74	1.37
MonoGS	Oracle	Photo	0.31	2.67	0.76
		NFL	0.26	2.66	0.60
	PPS-Net	Photo	0.33	2.79	0.89
		NFL	0.61	2.83	0.82
	DA-Hybrid	Photo	0.29	2.81	0.91
		NFL	0.38	2.71	0.71

Table 15: Results for sequence **trans\_t4\_a\_under\_review**

Method	Depth	BA	$ATE_t$ (mm)↓	$ATE_r$ (°)↓	Chamfer (mm)↓
NICE-SLAM	Oracle	Photo	0.39	2.81	0.05
		NFL	0.23	2.83	0.05
	PPS-Net	Photo	7.26	2.77	0.14
		NFL	7.88	2.73	0.23
EndoGSLAM	Oracle	Photo	2.64	1.47	0.74
		NFL	1.99	2.02	1.01
	PPS-Net	Photo	2.02	1.73	0.86
		NFL	2.58	1.78	0.89
	DA-Hybrid	Photo	2.99	2.08	1.04
		NFL	10.67	2.72	1.37
MonoGS	Oracle	Photo	1.23	2.41	0.83
		NFL	0.76	2.06	0.62
	PPS-Net	Photo	1.12	2.10	0.89
		NFL	1.70	1.80	0.86
	DA-Hybrid	Photo	2.30	2.46	1.10
		NFL	1.14	2.36	0.85

## NeurIPS Paper Checklist

The checklist is designed to encourage best practices for responsible machine learning research, addressing issues of reproducibility, transparency, research ethics, and societal impact. Do not remove the checklist: **The papers not including the checklist will be desk rejected.** The checklist should follow the references and follow the (optional) supplemental material. The checklist does NOT count towards the page limit.

Please read the checklist guidelines carefully for information on how to answer these questions. For each question in the checklist:

- You should answer [Yes], [No], or [NA].
- [NA] means either that the question is Not Applicable for that particular paper or the relevant information is Not Available.
- Please provide a short (1–2 sentence) justification right after your answer (even for NA).

**The checklist answers are an integral part of your paper submission.** They are visible to the reviewers, area chairs, senior area chairs, and ethics reviewers. You will be asked to also include it (after eventual revisions) with the final version of your paper, and its final version will be published with the paper.

The reviewers of your paper will be asked to use the checklist as one of the factors in their evaluation. While "[Yes]" is generally preferable to "[No]", it is perfectly acceptable to answer "[No]" provided a proper justification is given (e.g., "error bars are not reported because it would be too computationally expensive" or "we were unable to find the license for the dataset we used"). In general, answering "[No]" or "[NA]" is not grounds for rejection. While the questions are phrased in a binary way, we acknowledge that the true answer is often more nuanced, so please just use your best judgment and write a justification to elaborate. All supporting evidence can appear either in the main paper or the supplemental material, provided in appendix. If you answer [Yes] to a question, in the justification please point to the section(s) where related material for the question can be found.

### 1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: Yes, they are.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

### 2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: In Sec. 6, we discuss limitations.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.

- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

### 3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [NA]

Justification: Our paper does not present formal theorems or proofs; it focuses on an engineering loss formulation validated empirically. We include the equations necessary to formulate and justify our loss function but no explicit proof is needed.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

### 4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: We will put that information in the supplementary due to the page limits.

Guidelines: We provide full algorithmic details (loss terms, hyperparameters, depth preprocessing), dataset splits (C3VD sequences, indoor scenes) in the supplemental material.

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.

- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
  - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
  - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
  - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
  - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

## 5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: We are willing to make captured data and code publicly available once accepted.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

## 6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: They are described in Sec. 5.1.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

## 7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [No]

Justification: Due to page limits, we will provide those in the supplementary.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

## 8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: It is included in Sec. 5.1.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

## 9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

Answer: [Yes]

Justification: Our work conforms with the NeurIPS Code of ethics.

Guidelines: Our work involves algorithmic development and non-sensitive imaging data; no ethical issues beyond standard academic practice. All endoscopy data are publicly available and properly referenced.

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

## 10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: They are described in Sec. 1.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

## 11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: There are no high-risk pretrained models or sensitive datasets requiring special access controls.

Guidelines:

- The answer NA means that the paper poses no such risks.

- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

## 12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: We cited all the datasets and included license if possible.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets) has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

## 13. New assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [No]

Justification: We intend to release our code and captured data upon acceptance, but structured documentation (e.g., README, data schema) is not yet provided; this will be included with the release.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

## 14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: No human participants or crowdsourced annotations were involved.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

**15. Institutional review board (IRB) approvals or equivalent for research with human subjects**

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: Our paper does not involve crowdsourcing or research with human subjects. We only used publicly available data that already passed IRB approvals.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

**16. Declaration of LLM usage**

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigor, or originality of the research, declaration is not required.

Answer: [NA]

Justification: Large language models were not used for any part of the methodology

Guidelines:

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (<https://neurips.cc/Conferences/2025/LLM>) for what should or should not be described.