

# UC-NeRF: Uncertainty-aware Conditional Neural Radiance Fields from Endoscopic Sparse Views

Jiaxin Guo, Jiangliu Wang, Ruofeng Wei, Di Kang,  
Qi Dou, *Member, IEEE*, and Yun-hui Liu, *Fellow, IEEE*

**Abstract**—Visualizing surgical scenes is crucial for revealing internal anatomical structures during minimally invasive procedures. Novel View Synthesis is a vital technique that offers geometry and appearance reconstruction, enhancing understanding, planning, and decision-making in surgical scenes. Despite the impressive achievements of Neural Radiance Field (NeRF), its direct application to surgical scenes produces unsatisfying results due to two challenges: endoscopic sparse views and significant photometric inconsistencies. In this paper, we propose uncertainty-aware conditional NeRF for novel view synthesis to tackle the severe shape-radiance ambiguity from sparse surgical views. The core of UC-NeRF is to incorporate the multi-view uncertainty estimation to condition the neural radiance field for modeling the severe photometric inconsistencies adaptively. Specifically, our UC-NeRF first builds a consistency learner in the form of multi-view stereo network, to establish the geometric correspondence from sparse views and generate uncertainty estimation and feature priors. In neural rendering, we design a base-adaptive NeRF network to exploit the uncertainty estimation for explicitly handling the photometric inconsistencies. Furthermore, an uncertainty-guided geometry distillation is employed to enhance geometry learning. Experiments on the SCARED and Hamlyn datasets demonstrate our superior performance in rendering appearance and geometry, consistently outperforming the current state-of-the-art approaches. Our code will be released at <https://github.com/wrld/UC-NeRF>. **Index Terms**—Novel view synthesis, surgical 3D reconstruction, neural radiance fields

## I. INTRODUCTION

Minimally invasive surgery (MIS) has achieved a significant advancement in modern surgical practices. It reduces surgical

This work is supported in part by Shenzhen Portion of Shenzhen-Hong Kong Science and Technology Innovation Cooperation Zone under HZQB-KCZYB-20200089, in part by the Research Grants Council of Hong Kong under Grant T42-409/18-R, Grant 14218322, and Grant 14207320, in part by the Hong Kong Centre for Logistics Robotics, in part by the Multi-Scale Medical Robotics Centre, InnoHK, and in part by the VC Fund 4930745 of the CUHK T Stone Robotics Institute. (Corresponding author: Yun-Hui Liu)

Jiaxin Guo, Jiangliu Wang are with CUHK T Stone Robotics Institute, The Chinese University of Hong Kong, Hong Kong, China. (Emails: jxguo@mae.cuhk.edu.hk, jlwang@cuhk.edu.hk)

Ruofeng Wei, Qi Dou are with the Department of Computer Science and Engineering, The Chinese University of Hong Kong, China. (Emails: ruofenwei2-c@my.cityu.edu.hk, qidou@cuhk.edu.hk)

Di Kang is with Tencent AI Lab, Shen Zhen, China. (Email: di.kang@outlook.com)

Yun-hui Liu is with CUHK T Stone Robotics Institute, The Chinese University of Hong Kong, and with Hong Kong Center for Logistics Robotics, Hong Kong, China (Email: yhliu@mae.cuhk.edu.hk)

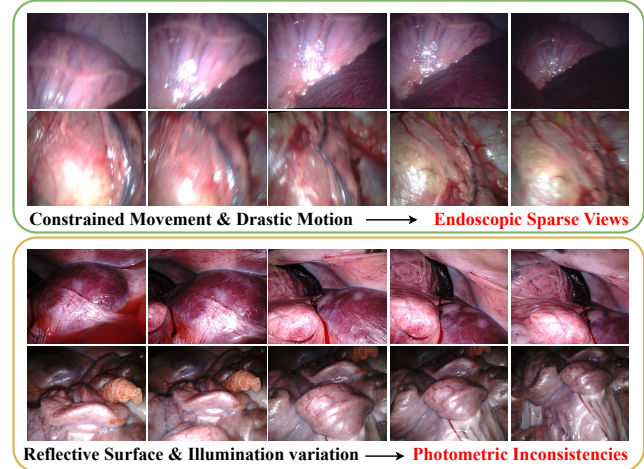
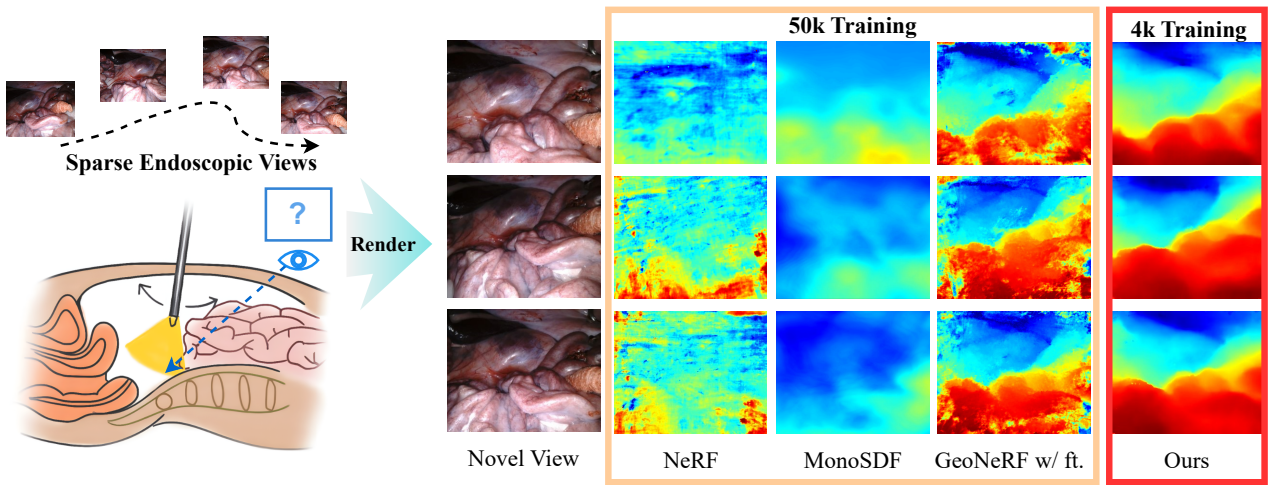


Fig. 1. Employing NeRF [1] in surgical scene encounters two main challenges, i.e. endoscopic sparse views and photometric inconsistency.

trauma, lessens post-operative discomfort, and shortens recovery time [2]. The endoscope allows surgeons to inspect internal structures, facilitating precise navigation through complex anatomical landscapes [3]. However, the 3D perception of the endoscope is impeded by the limited range of viewpoint changes inherent to endoscopic procedures, as well as the limited view area and two-dimensional imaging. These limitations make it challenging for surgeons to perceive depth and fully assess the surrounding conditions, which can hinder their understanding of the internal anatomy and potentially impact the effectiveness of surgical interventions. Additionally, traditional visualization systems in surgery, such as Ultrasound, Magnetic Resonance Imaging (MRI), or Computed Tomography (CT), increase the cost and complexity of medical imaging [6]. While 3D reconstruction methods allow surface geometry learning, they fall short of predictive ability and flexibility in exploring and visualizing surgical scenes from different perspectives with high-fidelity details.

Utilizing endoscopic multi-view images as input, the novel view synthesis technique generates photo-realistic free-view images of surgical scenes and intricate abdominal structures. This technique provides advantages for various applications, including virtual reality interactions, intra-operative surgical navigation, and autonomous robotic surgery [7]–[9]. In this area, the Neural Radiance Fields (NeRF) approach [1] has shown remarkable success. It synthesizes novel viewpoints from dense sets of input images using implicit volumetric representations, paving the way for an era of visually immersive



**Fig. 2. Training NeRF on sparse surgical views is challenging.** NeRF [1] fails to produce desirable views given sparse surgical scenes as inputs. State-of-the-art few-shot NeRF methods MonoSDF [4] and GeoNeRF [5] show degeneration in geometry rendering results. In contrast, our approach presents consistent improvement and achieves faster convergence in 4k compared to the 50k optimization of other baselines.

and interactive surgical practices.

While NeRF has shown remarkable effectiveness in handling natural images, its application in surgical scenes often results in subpar rendering, due to two challenges: endoscopic sparse views and photometric inconsistency, as illustrated in Fig. 1. Endoscopic sparse views are primarily due to constrained camera movement and drastic motion. The movement of the endoscope is inherently restricted by the narrow internal spaces within the body and the limited flexibility of the endoscope. These constraints confine the camera to a limited number of views, making it difficult to observe the same spot from different viewpoints. Moreover, when surgeons maneuver the endoscope, drastic movements under low camera frame rates can cause intermittent visibility and motion blur, hindering the acquisition of clear and consecutive frames. Both two factors reduce the overlap between captured views, leading to the sparsity of useful visual data. This is particularly challenging for NeRF, which relies on dense views for accurate 3D geometry inference and novel view synthesis. Besides, surgical scenes frequently exhibit considerable photometric inconsistencies, due to non-Lambertian reflections and fluctuating illumination. Such inconsistencies complicate the accurate capture of scene geometry and appearance, leading to potential artifacts and incorrect density distributions in NeRF synthesized views due to its fundamental design around minimizing RGB error. These factors contribute to a pressing question: *How can we enhance NeRF's ability to handle the severe shape-radiance ambiguity problem caused by the endoscopic sparse views and photometric inconsistency?*

To address this problem, we take the inspiration from few-shot NeRF to solve the problems induced by endoscopic sparse view. Some methods are proposed for few-shot NeRF to enhance the rendering performance and reduce the shape-radiance ambiguity, which could be roughly categorized into two classes: NeRF with pre-training [5], [10], [11] and NeRF with geometry guidance [4], [12]–[15]. The former class requires large datasets to pre-train the generalizable NeRF and fine-tune it on similar target scenes. The latter class

employs geometric information (depths, normals, pointclouds) to supervise and regularize the neural rendered depth or control the ray sampling range to get rid of outliers. While these methods achieve improvements compared to NeRF, they ignore the extent of photometric inconsistencies and directly inject the constraints equally in spatial, leading to unsatisfying performance as presented in Fig. 2. To tackle the severe shape-radiance ambiguity in surgical scene, we aim to explicitly detect the photometric inconsistency from the multi-view inputs and enable the neural radiance fields to adaptively model the regions with the uncertainty estimation.

In this paper, we propose a new network for surgical novel view synthesis, to empower NeRF with great robustness and efficiency to tackle the challenging surgical scene with inconsistencies and sparse views, i.e. UC-NeRF. The key novelty of our UC-NeRF lies in incorporating the multi-view uncertainty information with geometry and appearance priors to condition the neural radiance field for improved accuracy and robustness to sparse views. Our network has three essential designs to exploit the uncertainty information: i) A **consistency learner** to build the geometry correspondence and learn the uncertainty estimation across sparse multi-views. ii) An **uncertainty-aware dual-branch NeRF** designed with base-adaptive architecture utilizing the uncertainty information to handle photometric inconsistencies and solve shape-radiance ambiguity. iii) The **distillation from monocular geometry priors** to optimize the neural rendering with the uncertainty guidance to improve the accuracy in rendered depth. We validate our proposed approach on the SCARED [16] and Hamlyn datasets [17]–[19]. The extensive experiments demonstrate the state-of-the-art performance of UC-NeRF in novel view synthesis from sparse endoscopic views, with consistent improvement in effectiveness and efficiency.

In summary, our contributions are three folds: 1) To the best of our knowledge, we are the first to present a NeRF-based method addressing the challenging problem of novel view synthesis from endoscopic sparse images. 2) We devise an uncertainty-aware dual-branch NeRF to exploit the learned

uncertainty information, to recover view-dependent appearance while reducing the shape radiance ambiguity, with the distillation from the monocular geometry priors to enhance the accuracy and robustness. 3) The experiment results demonstrate that our method outperforms the previous state-of-the-art baselines, showing the superior efficiency and robustness to endoscopic sparse views.

## II. RELATED WORKS

**Novel View Synthesis.** Novel view synthesis focuses on generating new images or views of a scene from viewpoints not captured in the original imagery. Traditionally, novel view synthesis is based on geometric methods like image-based rendering [20]–[22] and light field rendering [23]–[26], which require precise camera calibration and struggle with complex scenes. The integration of deep learning allows a significant advancement for more realistic view generation [27]–[29]. A major breakthrough in this area was the development of NeRF [1], which used a fully connected deep network to model volumetric scene functions, greatly enhancing the quality of synthesized views. Following NeRF, various extensions have emerged by improving speed, quality, and generalization [4], [10], [11], [13], [14].

**NeRF from Surgical Scenes.** While NeRF has shown its potential in novel view synthesis, current approaches focus on 3D reconstruction using neural implicit fields from stereo endoscope videos. EndoNeRF [30] explores neural rendering for deformable tissue reconstruction from stereo endoscope inputs and devised a mask-guided ray-casting strategy to address the tool occlusion challenge. Sun et al. [31] propose a depth estimation network and a reconstruction network utilizing neural radiance fields for dynamic reconstruction. EndoSurf [32] proposes to model the deformation, geometry, and appearance separately to represent the deforming surfaces. Unlike these approaches that focus on single-view stereo video reconstruction, to the best of our knowledge, we are the first to present a sparse-view NeRF approach to tackle the challenge of novel view synthesis from endoscopic sparse views.

**Few-shot Neural Rendering.** Recently, the few-shot NeRF techniques that utilize sparse views for novel view synthesis have opened up new possibilities for improving the robustness and efficiency of NeRF, which can be classified into two classes: NeRF with pre-training and NeRF with geometry guidance. NeRF with pre-training in the first category pretrain a generalizable NeRF across multiple large datasets before fine-tuning on specific targets. PixelNeRF [10] leverages CNN features from input images to predict a continuous neural scene representation. IBRNet [33] introduces a network architecture to estimate radiance and volume density with appearance information from multiple source views. MVSNeRF [11] takes usage of cost volumes to reason the prior features to enhance the neural radiance field reconstruction. GeoNeRF [5] incorporates a transformer-based attention mechanism with volume rendering to manage complex occlusion conditions. However, these methods are prone to degeneration during significant data domain shifts, such as surgical scenes that differ vastly from general ones. Moreover, the need for extensive and

diverse pre-training datasets results in substantial time and computational costs. In this paper, we advocate the conditional NeRF by exploiting features and uncertainty estimation from sparse source images and realize generalizable training on the multiple surgical scenes, outperforming previous methods even without fine-tuning. NeRF with geometry guidance focuses on guiding NeRF training with geometric information. For example, NerfingMVS [13] takes MVS depth from COLMAP [34] as a prior and employs a monocular depth network to guide NeRF optimization. RegNeRF [14] regularizes geometry and appearance from unobserved viewpoints, refining ray sampling space over time. DS-NeRF [15] uses sparse 3D points from Structure-from-Motion (SfM) to monitor and adjust NeRF ray termination, enabling faster training with sparse views. MonoSDF [4] integrates multiple monocular geometric priors into neural implicit surface reconstruction. Roessle et al. [12] upsample SfM sparse points into dense depth maps to guide NeRF optimization. SparseNeRF [35] employs a local depth ranking regularization and a spatial continuity regularization to distill the depth priors. ConsistentNeRF [36] learns the 3D consistency by leveraging depth information to regularize the multi-view and single-view among pixels. Despite these advances, they ignore the extent of photometric inconsistencies and inject the constraints equally in spatial. Therefore, the injected hard geometry constraints easily cause over-fitting in sparse training data, leading to degeneration in test novel views. Our proposed UC-NeRF introduces an uncertainty-aware conditional NeRF to address these challenges, by leveraging the uncertainty information to model the 3D surgical scene and guide the optimization.

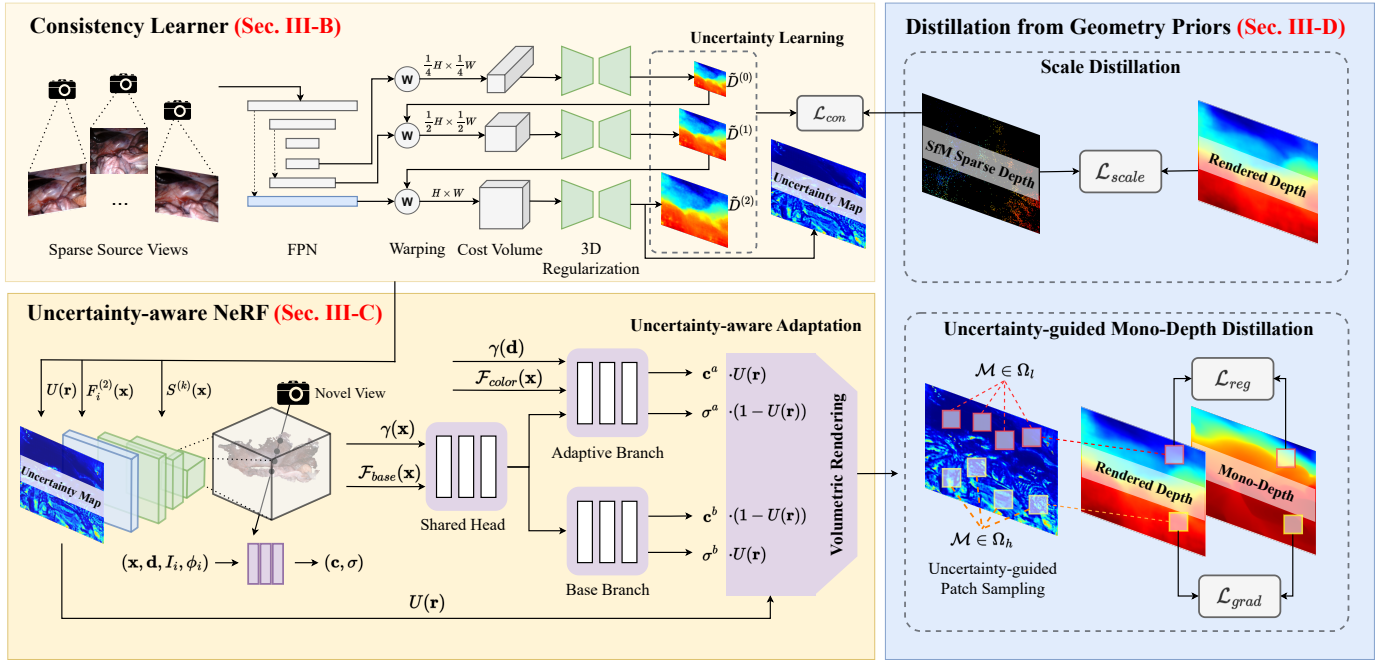
**Uncertainty Estimation in NeRF.** Uncertainty estimation has been adopted in diverse areas of computer vision to improve the interpretability and reduce the risk of the model. NeRF-W [37] leverages the uncertainty to tackle the transient object problem. S-NeRF [38] learns a probability distribution to quantify the uncertainty associated with the scene information. CF-NeRF [39] introduces latent variable modeling and conditional normalized flow to incorporate uncertainty quantification into NeRF. ActiveNeRF [40] incorporates the uncertainty estimation into a NeRF model by modeling the radiance values as a Gaussian distribution. In this paper, we aim to address the challenge of novel view synthesis from endoscopic sparse views. Unlike previous works, we incorporate the multi-view uncertainty from the consistency learner with dual branch NeRF, to explicitly handle the photometric inconsistencies from endoscopic sparse views.

## III. UNCERTAINTY-AWARE CONDITIONAL NERF

Given a set of sparse input images  $\{I_1, \dots, I_N\}$  and their corresponding camera parameters  $\{\phi_1, \dots, \phi_N\}$  as input, our goal is to reconstruct a radiance field that can faithfully capture the view-independent effects and the underlying true geometry, so that we can volume-render an image from any novel viewpoint and thus facilitate the diagnosis. Mathematically, this process is denoted as follows:

$$(\mathbf{c}, \sigma) = \text{UC-NeRF}(\gamma(\mathbf{x}), \gamma(\mathbf{d}); I_i, \phi_i), \quad (1)$$





**Fig. 3. Overview of our Uncertainty-aware Conditional NeRF (UC-NeRF).** We first build a consistency learner upon the multi-view stereo network, to capture the view-consistent constraints to generate the uncertainty map. Then, the uncertainty-aware NeRF takes image features (from FPN and 3D regularization module) and the uncertainty map as input to predict the radiance field, resulting in reduced shape-radiance ambiguity and improved rendering accuracy. Finally, we introduce the distillation from geometry priors for further optimizing the neural rendering results.

where  $\gamma$  is the encoding function to map position  $\mathbf{x}$  and view direction  $\mathbf{d}$  to a higher dimensional space. Through conditional inputs, our method enables generalizable training across multiple surgical scenes, thereby promoting training efficiency and robustness.

As in Fig. 3, UC-NeRF contains two major components: 1) a consistency learner that builds the geometry correspondences across multi-view inputs and generates an uncertainty map (Sec. III-B); and 2) a conditional NeRF that enables uncertainty-aware neural radiance fields reconstruction (Sec. III-C). To maximize generalization capability and reduce errors caused by shape-radiance ambiguity, we propose to distill geometric priors from the estimated sparse SfM points and a monocular depth estimator into our UC-NeRF (Sec. III-D), resulting in further improved robustness and accuracy in the rendered depth.

### A. Preliminaries

We first briefly introduce some NeRF basics [1] which are used in this paper. NeRF takes as input a set of posed images and represents a scene as a continuous volumetric function parameterized by MLPs. Given a 3D point  $\mathbf{x} \in \mathbb{R}^3$  and a viewing direction  $\mathbf{d} \in \mathbb{R}^2$ , NeRF learns to map from  $(\mathbf{x}, \mathbf{d})$  to the volume density  $\sigma$  and the emitted color  $\mathbf{c} = (r, g, b)$ :

$$(\mathbf{c}, \sigma) = \text{MLP}(\gamma(\mathbf{x}), \gamma(\mathbf{d})). \quad (2)$$

The color of an image pixel is calculated with the volume rendering [1]. Specifically, a ray  $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$  is determined by the camera origin  $\mathbf{o}$  and the pixel location  $\mathbf{p}$ , where  $\mathbf{d}$  is the unit direction vector passing from the camera origin  $\mathbf{o}$  to the pixel  $\mathbf{p}$ ,  $t$  is the distance of a sampling point to the origin

on this ray. The volume rendering equation [1] to obtain a pixel's color is defined as follows:

$$\hat{\mathbf{C}}(\mathbf{r}) = \sum_{i=1}^N T_i (1 - \exp(-\sigma_i \delta_i)) \mathbf{c}_i, \quad (3)$$

$$\text{where } T_i = \exp(-\sum_{j=1}^{i-1} \sigma_j \delta_j), \delta_i = t_{i+1} - t_i,$$

where  $N$  is the number of sample points along each ray,  $\sigma_i$  is the density value of point  $\mathbf{x}_i$ ,  $\delta_i$  is the distance between two consecutive sample points along the ray,  $T_i$  is the accumulated transparency from the camera origin. Following [1], [12], the depth of an image pixel can be similarly calculated by integrating every sample point's distance to the camera origin:

$$\hat{D}(\mathbf{r}) = \sum_{i=1}^N T_i (1 - \exp(-\sigma_i \delta_i)) t_i. \quad (4)$$

With the above differential volume rendering process, NeRF optimizes the radiance fields by minimizing the reconstruction error between the rendered color and the ground truth color:

$$\mathcal{L}_{rgb} = \|\hat{\mathbf{C}}(\mathbf{r}) - \mathbf{C}(\mathbf{r})\|_2^2. \quad (5)$$

However, NeRF reconstruction degrades drastically, especially the geometry part, due to the sparse endoscopic views and the photometric inconsistencies across images caused by varying lighting conditions during the surgical operation.

### B. Consistency Learner

In UC-NeRF, a consistency learner is utilized to exploit robust geometric information from sparse SfM depth and learn the consistency. Specifically, it adopts the CasMVSNet [41]



as the backbone to extract intermediate image features and construct cascade cost volumes to predict a dense depth map and an uncertainty map. The 2D image features and the 3D neural volumes, which contain geometry and appearance information of the target surgical scenes, are used as the condition input of the later NeRF network (detailed in Sec. III-C.1). The uncertainty map is used to adaptively re-weight the predictions before radiance integral (detailed in Sec. III-C.2). With the help of this consistency learner, our UC-NeRF can use better conditional information as input and estimate more accurate novel view images and depth maps.

1) *Cascade Neural Volumes*: Given  $N$  sparse source views  $\{I_i\}_{i=1}^N$  with resolution size  $H \times W$  as input, we utilize a Feature Pyramid Network (FPN) [42] to extract image features in different spatial resolutions across three stages.

$$F_i^{(k)} = \text{FPN}(I_i), \quad k = \{0, 1, 2\}, \quad (6)$$

where  $F_i^{(k)} \in \mathbb{R}^{\frac{H}{2^{2-k}} \times \frac{W}{2^{2-k}} \times 2^{-k}Z}$  represents the 2D feature maps extracted at the stage  $k$  from the  $i_{th}$  input view,  $Z$  is the feature dimension of stage 0.

Next, we warp the 2D feature maps from different source views to the plane sweeping volume feature  $V_i^{(k)}$  on the frustum of the target view following [41], [43]. Given the camera parameters  $\{\phi_i\}_{i=1}^N$  for the input source images, and  $\phi_0$  for the target view, we can apply homography warping to warp the 2D feature maps  $F_i^{(k)}$  from source views into hypothetical planes of the target view, forming 3D features  $V_i^{(k)}$ . Note that the target view is the novel view to be rendered, which is the same as the reference camera in MVS methods [43]. Following CasMVSNet [41], the hypothesis depth planes range from  $d_{max}^{(k)}$  to  $d_{min}^{(k)}$  with  $Y^{(k)}$  discrete depth values evenly spaced in-between. The depth planes are configured in a coarse-to-fine manner from stage 0 which uses the largest depth range and the most planes, to stage 2.

The cost volume is calculated with a variance-based metric from 3D feature  $V_i^{(k)}$  [41], [43]. Note that we only use features from the neighboring source views during feature extraction and the following depth estimation, since we do not have access to the image of the novel view to be synthesized. We further regularize the cost volume with a 3D-CNN to generate a 3D neural volume  $S^{(k)}$  and a probability volume  $P^{(k)}$ :

$$S^{(k)}, P^{(k)} = \text{3D-CNN}(\text{Var}(V_i^{(k)})), \quad (7)$$

where  $S^{(k)} \in \mathbb{R}^{Y^{(k)} \times \frac{H}{2^{2-k}} \times \frac{W}{2^{2-k}} \times 2^{-k}Z}$  denotes the 3D neural volume storing the geometry information of the target view, and is used as one conditional input to build the neural radiance fields.  $P^{(k)} \in \mathbb{R}^{Y^{(k)} \times \frac{H}{2^{2-k}} \times \frac{W}{2^{2-k}}}$  represents the probability volume specifying every spatial location's depth probability among all possible depth planes. After computing  $P^{(k)}$ , we also perform a softmax operation on the depth plane dimension of  $P^{(k)}$  to ensure that the probability of the depth hypothesis remains in the range  $[0, 1]$ . The final depth estimation  $\tilde{D}^{(k)}$  of the target view is obtained as the expectation of different depth hypothesis [43]:

$$\tilde{D}^{(k)} = \sum_{d=d_{min}^{(k)}}^{d_{max}^{(k)}} d \times P^{(k)}(d). \quad (8)$$

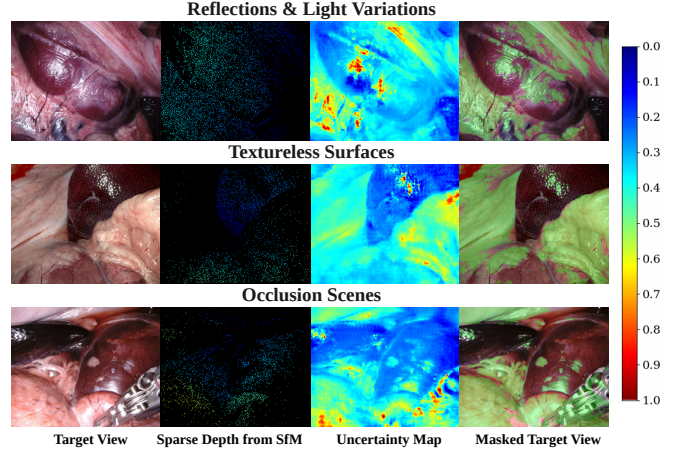


Fig. 4. Visualization of the Sparse SfM depth and the estimated uncertainty map. With the guidance from SfM, the uncertainty map measures the extent of photometric inconsistency. The masked target view indicates the region with uncertainty larger than the mean value.

We further enhance the consistency learner with view-consistent points from SfM. Specifically, we utilize the sparse depth  $D_{sfm}$  projected from SfM points in every view to supervise the predicted depth from the consistency learner. Taking into account existing noise, we use the SfM reprojection error  $\omega$  as the weight in the sparse depth for the regression:

$$L_{con} = \sum_{k=0}^2 \alpha^{(k)} \exp(-(\omega/\bar{\omega})^2) \|\tilde{D}^{(k)} - D_{sfm}\|_1, \quad (9)$$

where  $\exp(-(\omega/\bar{\omega})^2)$  is the weight for depth loss at stage  $k$ , which injects larger weight on the depth loss for the point with a smaller reprojection error. Given that SfM leverages SIFT points and bundle adjustment optimization, it excels in accurately capturing geometric correspondences. Consequently, the consistency learner can exploit the reliable geometric correspondences captured by SfM, i.e. light-invariant and visible points among the source view inputs.

2) *Uncertainty Learning*: To model the extent of the photometric inconsistencies, an uncertainty map of the target view is generated according to the probability distribution along the depth hypothesis. Specifically, for every pixel  $(u, v)$ , its uncertainty is calculated according to four probability values from the probabilistic volume  $P^{(2)}$  that are closest to its depth estimation ( $\tilde{D}^{(2)}$  in Eq. (8)) following [41], [43]:

$$U(u, v) = 1 - \sum_{i=j-1}^{j+2} P^{(2)}(u, v, i) \Big|_{j=\text{Index}(\tilde{D}^{(2)}(u, v))}, \quad (10)$$

where  $\text{Index}(\cdot)$  indicates the index of the hypothesis depth plane. If a point exhibits consistency across multiple views, it is prone to have a unimodal distribution in its depth estimation probability. This results in the estimated probability  $P^{(2)}$ , being close to the peak of this unimodal distribution, indicating low uncertainty. In contrast, multimodal distribution signifies higher uncertainty, as there is a lack of consensus among the different views on the correct depth value. As the uncertainty map shown in Fig. 4, regions with high uncertainty tend to have minimal texture, large occlusion, and severe reflective surfaces. This suggests that our consistency learner

is able to produce a sensible measure of uncertainty based on the guidance of SfM sparse points.

### C. Uncertainty-aware dual-branch NeRF

We propose a dual-branch NeRF utilizing the uncertainty information to explicitly handle photometric inconsistencies caused by moving light and non-Lambertian surfaces. Specifically, our dual-branch NeRF contains a base branch  $\text{MLP}_{\theta_b}$  spatially weighted by the confidence score (i.e. 1 - uncertainty score), which aims at modeling view-consistent appearance and geometry in the surgical scene, and an adaptive branch  $\text{MLP}_{\theta_a}$  weighted by the uncertainty score, which aims at modeling view-dependent effects and details, i.e. illumination variations, non-Lambertian texture. The reason is that regions with higher uncertainty values from the depth estimator usually lack consistency across neighboring multi-views, which are supposed to be modeled by the adaptive branch.

**1) Base-Adaptive Rendering Network:** The NeRF is conditioned on the feature priors from the consistency learner to generalize in different scenes with better efficiency and rendering performance (see Fig. 7). Concretely, we directly build the cost volume upon the target (i.e. novel) view frustum to collect more “original” appearance and geometry features from the target view rather than highly processed abstract features from other views. This is different from MVSNeRF [11], which builds cost volume upon the neighboring input view and then warps the 3D points to this different view for sampling condition features. This modification of the cost volume reconstruction process results in better geometry accuracy and improved image rendering performance (See Tab. I).

Specifically, a sample point is warped to input views to obtain appearance prior (i.e.  $\{I_i(\mathbf{x}_{0 \rightarrow i})\}_{i=1}^N$ ) and geometry prior  $S^{(k)}\{\mathbf{x}\}_{k=0}^2$ . Note that  $S^{(k)}$  is a 3D volume space in NDC coordinate frame and we use trilinear interpolation. Mathematically, the condition input  $\mathcal{F}_{base}(\mathbf{x})$  and the conditional NeRF are defined as:

$$\mathcal{F}_{base}(\mathbf{x}) = \text{Concat}(\{S^{(k)}(\mathbf{x})\}_{k=0}^2, \{I_i(\mathbf{x}_{0 \rightarrow i})\}_{i=1}^N), \quad (11)$$

$$h = \text{MLP}_{\theta_s}(\gamma(\mathbf{x}), \mathcal{F}_{base}(\mathbf{x})), \quad (12)$$

where  $h$  is the shared latent feature used by both base and adaptive branches.

$$\mathbf{c}^b, \sigma^b = \text{MLP}_{\theta_b}(h), \quad (13)$$

$$\mathbf{c}^a, \sigma^a = \text{MLP}_{\theta_a}(h, \gamma(\mathbf{d}), \mathcal{F}_{color}(\mathbf{x})), \quad (14)$$

$$\text{where } \mathcal{F}_{color}(\mathbf{x}) = \{F_i^{(2)}(\mathbf{x}_{0 \rightarrow i})\}_{i=1}^N. \quad (15)$$

The density  $\sigma^b$  and color  $\mathbf{c}^b$  of the base branch is decoded solely by  $h$  without viewing vector input since this branch is designed to model the underlying true geometry and diffuse colors. As for the adaptive branch, both its density  $\sigma^a$  and color  $\mathbf{c}^a$  are dependent on view direction  $\gamma(\mathbf{d})$  since it is designed to model the inconsistencies. Note that we also use image feature as input since we find it is helpful experimentally possibly due to its robustness, i.e. deep features with large enough context information. Specifically, we use the features  $\{F_i^{(2)}(\mathbf{x}_{0 \rightarrow i})\}_{i=1}^N$  from the last layer of the FPN network.

**2) Uncertainty-aware Adaptation:** Since the uncertainty measure reflects the geometric reliability and photometric inconsistency of each point, it is suitable to be used as weight to balance and control the contribution of the two branches. Specifically, radiance fields predicted by the two branches are spatially weighted summed according to uncertainty score  $U(\mathbf{r})$  as follows:

$$\mathbf{c} = \mathbf{c}^b \cdot (1 - U(\mathbf{r})) + \mathbf{c}^a \cdot U(\mathbf{r}), \quad (16)$$

$$\sigma = \sigma^b \cdot U(\mathbf{r}) + \sigma^a \cdot (1 - U(\mathbf{r})). \quad (17)$$

The final image is volumetric rendered following Eq. (3) and Eq. (4). Under this setting, in color prediction, the adaptive branch  $\mathbf{c}^a$  contributes more to regions with higher uncertainty, capturing view-dependent photometric effects that vary significantly with different viewing angles. The base branch  $\mathbf{c}^b$  is more influential in regions with higher confidence, where the appearance is stable and predictable. Conversely, in density prediction, the base branch  $\sigma^b$  is dominant in regions with higher uncertainty. This helps maintain geometric reliability and multi-view consistency, avoiding shape-radiance ambiguity which can lead to degenerated reconstructions. For regions with higher confidence and fewer directional variations, the adaptive branch  $\sigma^a$  could contribute more, providing additional textures and details to enhance the reconstruction quality (See Tab. V).

**Discussion.** Following this adaptation, the base branch provides reliable and accurate geometry information to reduce the uncertainty and improve the robustness for these rays (See Fig. 6). In color rendering, sampled rays with higher uncertainty, prone to complex photometric effects, are weighted towards the adaptive branch to compensate for more view-dependent details (See Fig. 9). With the uncertainty map modulating the base and adaptive branches, our UC-NeRF is enabled to construct the neural radiance field with spatial uncertainty, facilitating more stable and controllable learning to solve the shape-radiance ambiguity. With the balance modulated in the two branches, the differential uncertainty-aware adaptation is capable of synchronizing the uncertainty learned by multi-view stereo with the neural rendering process.

### D. Distillation from Monocular Geometry Priors

To further improve the geometry consistency, we exploit the geometry priors from monocular images to guide the training of UC-NeRF for scale-aware depth learning. We employ a two-fold approach: 1) Scale distillation from SfM. 2) Uncertainty-guided monocular depth distillation.

**1) Scale Distillation:** To preserve the scale consistency among sparse views, We first incorporate the sparse depth from SfM, which intrinsically preserves real-world scale. This sparse depth is used to guide our method in learning the correct scale of the scene through a scale distillation Loss. Similar to Eq. (9), to alleviate the negative influence of inaccurate depth, we adopt the weight  $\exp(-(\omega/\bar{\omega})^2)$  to suppress the supervision from unreliable depth estimations (i.e. with larger reprojection errors):

$$\mathcal{L}_{scale} = \sum_{\mathbf{r} \in \mathcal{R}} \exp(-(\omega/\bar{\omega})^2) \|\hat{D}(\mathbf{r}) - D_{sfm}(\mathbf{r})\|_1, \quad (18)$$



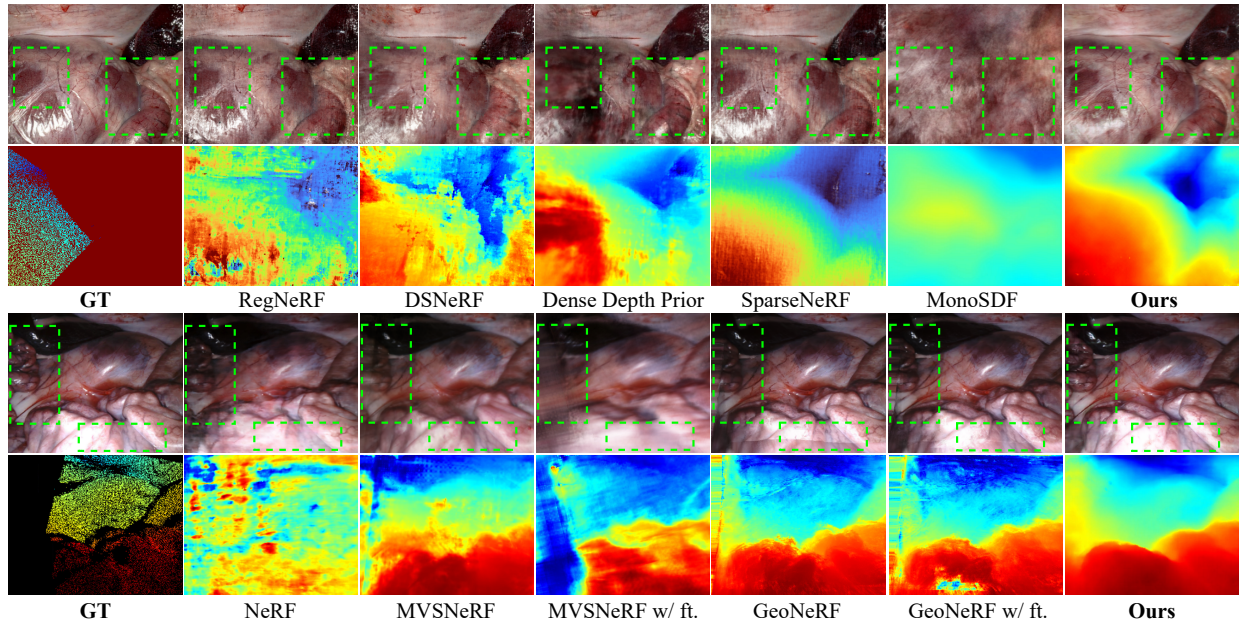


Fig. 5. **Qualitative Comparisons of rendered color and depth.** Given sparse input views, existing approaches show rendering results with blur and artifacts, suffering from photometric inconsistency. Our UC-NeRF can generate fine-grained details and consistent depth.

where  $\mathcal{R}$  denotes the sampled set of pixels in the region where SfM depth is available. We minimize the  $L_1$  loss between the rendered depth  $\hat{D}(\mathbf{r})$  and the SfM sparse depth  $D_{sfm}$ , with SfM reprojection error  $\omega$  to weight the loss.

2) *Uncertainty-guided Mono-Depth Distillation:* To optimize the region where SfM is sparse or unavailable, we leverage a monocular depth estimation model, Dense Prediction Transformer (DPT) [44], to guide our UC-NeRF. We introduce an uncertainty-guided mono-depth distillation, taking the uncertainty map as a reference to employ different losses spatially for depth supervision to sampled patches.

We first sample image patches in the region with high uncertainty denoted as  $\Omega_h$ . Since the constraint of SfM sparse point is not enough for the high uncertainty region  $\Omega_h$ , we apply a scale-invariant depth gradient loss  $\mathcal{L}_{grad}$  to supervise the gradient difference between rendered depth  $\hat{D}$  with the monocular depth  $D_{dpt}$ .

$$\mathcal{L}_{grad} = \sum_{\mathcal{M} \in \Omega_h} \sum_{\mathbf{r} \in \mathcal{M}} \|\nabla(D_{dpt}(\mathbf{r}) - (\hat{D}(\mathbf{r}) \cdot s + q))\|_1, \quad (19)$$

where  $s$  and  $q$  represent the scale and shift computed by linear least squares to convert the patches to the same scale following [45]. Conversely, low uncertainty regions  $\Omega_l$  contain more view-consistent correspondences with reliable depth, where the edge-aware smooth loss is employed as a regularization term to refine the continuity of the rendered depth.

$$\mathcal{L}_{reg} = \sum_{\mathcal{M} \in \Omega_l} \sum_{\mathbf{r} \in \mathcal{M}} \exp(-\beta \nabla D_{dpt}(\mathbf{r})) \|\nabla \hat{D}(\mathbf{r})\|_1, \quad (20)$$

where  $\beta$  denotes a hyperparameter to control the smooth extent. The exponential term serves as a weight that decreases as the depth gradient from DPT increases to preserve the edges, thereby promoting smoothness and continuity in low-uncertainty regions.

Through this uncertainty-guided mono-depth distillation, our method exploits the monocular geometry priors in the scene and provides balanced supervision with the uncertainty map as the reference. Our method not only incorporates detailed geometry and scale cues from single-view depth prediction but also aligns its understanding with the global, scale-aware information from SfM. It allows our model to effectively deal with various situations in the scene, enhancing the stability and performance of depth rendering.

3) *Training Loss:* The total training loss for UC-NeRF is formulated by:

$$\mathcal{L} = \lambda_1 \mathcal{L}_{rgb} + \lambda_2 \mathcal{L}_{con} + \lambda_3 \mathcal{L}_{scale} + \lambda_4 \mathcal{L}_{grad} + \lambda_5 \mathcal{L}_{reg}, \quad (21)$$

where  $\lambda_1, \lambda_2, \lambda_3, \lambda_4, \lambda_5$  denote the loss weights. In practice, we set  $\lambda_1 = 10, \lambda_2 = 0.5, \lambda_3 = 0.5, \lambda_4 = 0.5, \lambda_5 = 0.05$ .

## IV. EXPERIMENTS

### A. Implementation Details

1) *Experimental Setup:* The code of our method is based on PyTorch, running on NVIDIA GeForce RTX 3090 GPU. We process rays for scale distillation in a batch size of 1024, and patches for guided patch sampling with size of  $6 \times 6$  in batch size of 50. We use Adam Optimizer for our network, with a learning rate of  $6 \times 10^{-4}$  with a cosine decay scheduler. For the sampling points on a ray, we adopt 90 points during training and inference for efficiency. For the consistency learner, we utilize a pre-trained Cas-MVSNet model [41], with the depth hypothesis planes decreasing as (48, 32, 8) in three stages. All the experiments and comparisons use the image size of  $320 \times 256$ . To compare with NeRF-based methods which require no pre-training, we train the model from scratch. For NeRF with pre-training methods, we utilize their released pre-trained model and train the fine-tuned model on each scene.



TABLE I

**QUANTITATIVE COMPARISON OF NOVEL VIEW SYNTHESIS PERFORMANCE ON SCARED DATASET.**  $\mathbb{P}$  DENOTES FOR NERF WITH PRE-TRAINING METHODS.  $\mathbb{P}$  W/ FT. DENOTES FOR THE PRE-TRAINED MODEL AFTER FINE-TUNING.  $\mathbb{G}$  DENOTES FOR NERF WITH GEOMETRY GUIDANCE.

Settings	Methods	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	Abs Rel $\downarrow$	Sq Rel $\downarrow$	RMSE $\downarrow$	RMSE log $\downarrow$	$\delta < 1.25 \uparrow$
-	NeRF [37]	24.27 $\pm$ 1.95	0.712 $\pm$ 0.10	0.196 $\pm$ 0.07	0.324 $\pm$ 0.12	2.776 $\pm$ 2.14	6.118 $\pm$ 2.10	0.453 $\pm$ 0.14	0.470 $\pm$ 0.12
-	ActiveNeRF [40]	23.29 $\pm$ 1.56	0.696 $\pm$ 0.11	0.215 $\pm$ 0.08	0.326 $\pm$ 0.13	2.927 $\pm$ 2.54	6.274 $\pm$ 2.52	0.433 $\pm$ 0.14	0.465 $\pm$ 0.12
$\mathbb{P}$	PixelNeRF [10]	15.24 $\pm$ 1.57	0.432 $\pm$ 0.08	0.308 $\pm$ 0.04	0.327 $\pm$ 0.08	2.371 $\pm$ 1.43	5.794 $\pm$ 1.83	0.363 $\pm$ 0.06	0.531 $\pm$ 0.09
	MVSNeRF [11]	17.35 $\pm$ 1.63	0.582 $\pm$ 0.07	0.408 $\pm$ 0.04	0.266 $\pm$ 0.05	2.020 $\pm$ 1.01	5.401 $\pm$ 1.87	0.318 $\pm$ 0.07	0.536 $\pm$ 0.08
	GeoNeRF [5]	23.50 $\pm$ 1.30	0.792 $\pm$ 0.07	0.152 $\pm$ 0.05	0.096 $\pm$ 0.05	0.196 $\pm$ 0.13	1.628 $\pm$ 0.83	0.093 $\pm$ 0.04	0.951 $\pm$ 0.04
$\mathbb{P}$ w/ ft.	PixelNeRF [10]	21.45 $\pm$ 1.87	0.673 $\pm$ 0.09	0.288 $\pm$ 0.08	0.271 $\pm$ 0.11	1.479 $\pm$ 1.24	4.890 $\pm$ 1.63	0.301 $\pm$ 0.08	0.817 $\pm$ 0.05
	MVSNeRF [11]	23.09 $\pm$ 1.69	0.751 $\pm$ 0.07	0.248 $\pm$ 0.06	0.246 $\pm$ 0.04	1.357 $\pm$ 1.12	4.021 $\pm$ 1.39	0.280 $\pm$ 0.06	0.854 $\pm$ 0.05
	GeoNeRF [5]	24.59 $\pm$ 1.48	0.825 $\pm$ 0.10	0.126 $\pm$ 0.09	0.119 $\pm$ 0.10	0.513 $\pm$ 0.15	2.747 $\pm$ 1.17	0.155 $\pm$ 0.05	0.889 $\pm$ 0.04
$\mathbb{G}$	NerfingMVS [13]	22.70 $\pm$ 1.76	0.598 $\pm$ 0.05	0.285 $\pm$ 0.05	0.228 $\pm$ 0.11	1.275 $\pm$ 1.61	4.251 $\pm$ 2.06	0.290 $\pm$ 0.08	0.598 $\pm$ 0.14
	RegNeRF [14]	25.18 $\pm$ 1.41	0.746 $\pm$ 0.11	0.152 $\pm$ 0.07	0.226 $\pm$ 0.06	1.923 $\pm$ 2.41	4.287 $\pm$ 2.44	0.276 $\pm$ 0.06	0.593 $\pm$ 0.15
	MonoSDF [4]	18.49 $\pm$ 0.49	0.471 $\pm$ 0.06	0.382 $\pm$ 0.02	0.246 $\pm$ 0.11	1.677 $\pm$ 1.40	4.257 $\pm$ 1.98	0.256 $\pm$ 0.09	0.628 $\pm$ 0.15
	DS-NeRF [15]	24.93 $\pm$ 1.55	0.714 $\pm$ 0.10	0.227 $\pm$ 0.08	0.329 $\pm$ 0.15	3.396 $\pm$ 3.79	6.506 $\pm$ 3.23	0.375 $\pm$ 0.12	0.476 $\pm$ 0.12
	Dense Prior [12]	23.43 $\pm$ 1.42	0.816 $\pm$ 0.05	0.147 $\pm$ 0.04	0.088 $\pm$ 0.03	0.149 $\pm$ 0.10	1.458 $\pm$ 0.66	0.087 $\pm$ 0.03	0.932 $\pm$ 0.04
	SparseNeRF [35]	25.92 $\pm$ 2.06	0.776 $\pm$ 0.10	0.150 $\pm$ 0.06	0.222 $\pm$ 0.02	1.100 $\pm$ 0.35	4.095 $\pm$ 1.20	0.239 $\pm$ 0.03	0.589 $\pm$ 0.04
	<b>Ours</b>	26.40 $\pm$ 1.39	0.855 $\pm$ 0.05	0.107 $\pm$ 0.03	0.053 $\pm$ 0.03	0.121 $\pm$ 0.03	1.304 $\pm$ 0.68	0.074 $\pm$ 0.03	0.965 $\pm$ 0.04

TABLE II

**QUANTITATIVE COMPARISON OF NOVEL VIEW SYNTHESIS PERFORMANCE ON HAMLYN DATASET.**  $\mathbb{P}$  DENOTES FOR NERF WITH PRE-TRAINING METHODS.  $\mathbb{P}$  W/ FT. DENOTES FOR THE PRE-TRAINED MODEL AFTER FINE-TUNING.  $\mathbb{G}$  DENOTES FOR NERF WITH GEOMETRY GUIDANCE.

Settings	Methods	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	Abs Rel $\downarrow$	Sq Rel $\downarrow$	RMSE $\downarrow$	RMSE log $\downarrow$	$\delta < 1.25 \uparrow$
-	NeRF [37]	23.86 $\pm$ 4.71	0.716 $\pm$ 0.15	0.318 $\pm$ 0.19	0.688 $\pm$ 0.19	36.64 $\pm$ 15.2	26.09 $\pm$ 6.18	0.596 $\pm$ 0.11	0.486 $\pm$ 0.08
-	ActiveNeRF [40]	18.81 $\pm$ 3.56	0.536 $\pm$ 0.14	0.488 $\pm$ 0.12	0.770 $\pm$ 0.18	41.72 $\pm$ 16.4	32.94 $\pm$ 6.15	0.917 $\pm$ 0.25	0.327 $\pm$ 0.07
$\mathbb{P}$	PixelNeRF [10]	14.21 $\pm$ 3.47	0.621 $\pm$ 0.09	0.524 $\pm$ 0.09	0.690 $\pm$ 0.15	34.06 $\pm$ 8.97	25.77 $\pm$ 4.16	0.547 $\pm$ 0.11	0.439 $\pm$ 0.09
	MVSNeRF [11]	16.11 $\pm$ 2.04	0.642 $\pm$ 0.09	0.513 $\pm$ 0.14	0.681 $\pm$ 0.13	35.50 $\pm$ 7.57	24.34 $\pm$ 3.28	0.539 $\pm$ 0.08	0.450 $\pm$ 0.14
	GeoNeRF [5]	23.31 $\pm$ 3.09	0.802 $\pm$ 0.09	0.241 $\pm$ 0.07	0.596 $\pm$ 0.18	30.92 $\pm$ 11.9	17.84 $\pm$ 6.54	0.498 $\pm$ 0.13	0.682 $\pm$ 0.26
$\mathbb{P}$ w/ ft.	PixelNeRF [10]	21.45 $\pm$ 3.74	0.659 $\pm$ 0.12	0.470 $\pm$ 0.10	0.663 $\pm$ 0.14	32.91 $\pm$ 3.30	22.37 $\pm$ 3.02	0.524 $\pm$ 0.07	0.510 $\pm$ 0.06
	MVSNeRF [11]	23.19 $\pm$ 3.49	0.792 $\pm$ 0.06	0.310 $\pm$ 0.06	0.622 $\pm$ 0.11	29.42 $\pm$ 7.89	19.86 $\pm$ 1.87	0.501 $\pm$ 0.05	0.541 $\pm$ 0.03
	GeoNeRF [5]	23.53 $\pm$ 3.34	0.783 $\pm$ 0.11	0.268 $\pm$ 0.11	0.621 $\pm$ 0.16	31.12 $\pm$ 12.1	19.70 $\pm$ 5.76	0.545 $\pm$ 0.12	0.617 $\pm$ 0.22
$\mathbb{G}$	NerfingMVS [13]	19.36 $\pm$ 3.78	0.643 $\pm$ 0.19	0.345 $\pm$ 0.19	0.792 $\pm$ 0.18	32.71 $\pm$ 11.7	20.31 $\pm$ 5.13	0.510 $\pm$ 0.39	0.629 $\pm$ 0.28
	RegNeRF [14]	26.38 $\pm$ 4.33	0.780 $\pm$ 0.15	0.274 $\pm$ 0.21	0.703 $\pm$ 0.21	38.56 $\pm$ 13.9	24.45 $\pm$ 6.35	0.594 $\pm$ 0.12	0.515 $\pm$ 0.21
	MonoSDF [4]	20.20 $\pm$ 3.69	0.701 $\pm$ 0.08	0.473 $\pm$ 0.07	0.628 $\pm$ 0.20	28.77 $\pm$ 15.0	19.48 $\pm$ 7.37	0.472 $\pm$ 0.10	0.506 $\pm$ 0.17
	DS-NeRF [15]	24.12 $\pm$ 7.82	0.687 $\pm$ 0.24	0.351 $\pm$ 0.29	0.692 $\pm$ 0.11	35.43 $\pm$ 7.47	23.86 $\pm$ 4.07	0.741 $\pm$ 0.49	0.463 $\pm$ 0.14
	Dense Prior [12]	23.79 $\pm$ 4.52	0.710 $\pm$ 0.18	0.257 $\pm$ 0.21	0.616 $\pm$ 0.17	30.45 $\pm$ 11.2	19.43 $\pm$ 4.61	0.502 $\pm$ 0.29	0.636 $\pm$ 0.20
	SparseNeRF [35]	26.72 $\pm$ 3.64	0.812 $\pm$ 0.13	0.293 $\pm$ 0.22	0.610 $\pm$ 0.18	30.61 $\pm$ 11.8	18.50 $\pm$ 4.62	0.485 $\pm$ 0.08	0.585 $\pm$ 0.22
	<b>Ours</b>	26.87 $\pm$ 1.89	0.819 $\pm$ 0.07	0.215 $\pm$ 0.06	0.530 $\pm$ 0.15	26.37 $\pm$ 9.52	16.15 $\pm$ 3.59	0.450 $\pm$ 0.07	0.741 $\pm$ 0.19

**2) Datasets:** We train our method on the SCARED Dataset [16] and Hamlyn Dataset [17]–[19], which contain challenging endoscopic scenes with weak textures, reflections, and occlusions. Following the preprocessing in [12], after filtering out the frames with motion blur or flaws, we obtained 9 and 6 scenes from the SCARED and Hamlyn datasets respectively. Specifically, we first use a sliding window approach and extract the sharpest image from every N frame. Images with severe occlusions are also removed to avoid introducing noise into the training process. Next, we downsample the dataset to enhance computational efficiency and avoid redundancy from very similar frames. Specifically, we manually select frames to ensure each frame provides unique coverage of the surgical scene, with an overlap of approximately 60-80 % between consecutive frames following [12]. This step leads to reduced data size, maintaining a balance between comprehensive scene coverage and training efficiency.

For Hamlyn Dataset, we follow the volumetric reconstruction part of Endo-Depth-and-Motion [46] to specifically focus

on static scenes with larger camera translation to validate our method. Note that the collected scenes only contain monocular images, taken from the left camera view for both SCARED and Hamlyn Dataset. We use COLMAP [47] to preprocess the datasets to get the camera poses and sparse pointclouds following [1], [15]. We obtain the monocular depth from DPT [44] as priors to guide our UC-NeRF. Similar to the data split in [1], [13], [14], we collect 20 images covering a local area for each scene and hold out 1/2 of these views for testing and the remaining for training. During training, we take 7 closest neighboring views as source views for the consistency learner to extract features and 3D neural volumes.

**3) Metrics:** We report the quantitative comparisons of appearance and geometry to show the synthesis performance, including the mean and standard deviation across the test novel views. We further highlight the table cell to show **best**, **second best**, and **third best**. We adopt the PSNR, SSIM [48] and LPIPS [49] for assessing image synthesis quality. To evaluate the depth results, we scale the predicted depth maps

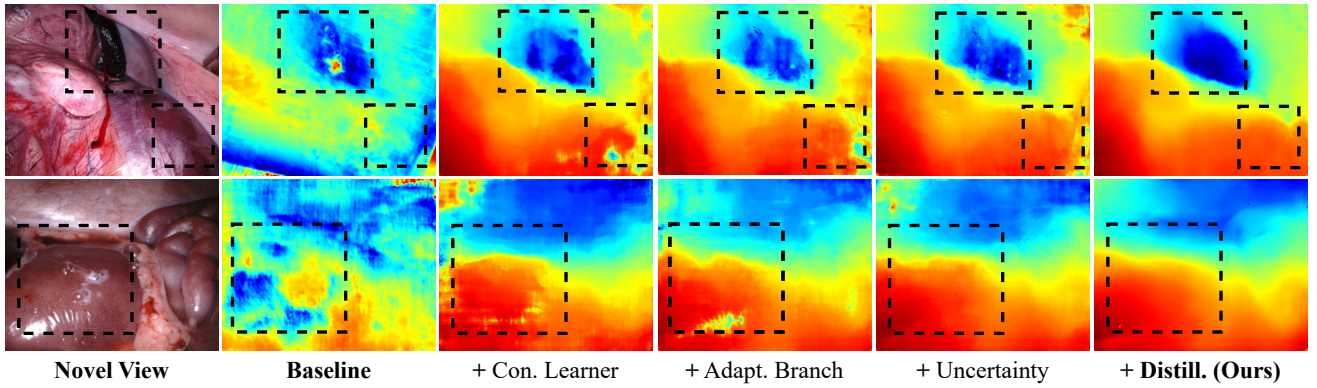


Fig. 6. Visualization of the ablation study on different components. Using all the components achieves the best depth estimation.

with median scaling [50]. Following [13], we adopt the error and accuracy metrics for depth evaluation, including Abs Rel, Sq Rel, RMSE in mm, RMSE log in log mm and  $\delta < t$  in %.

### B. Comparison Study

We report the quantitative results in Tab. I and Tab. II, and qualitative comparison in Fig. 5. We first compare with NeRF [1] and ActiveNeRF [40] to validate the performance of our multi-view uncertainty estimation in dealing with endoscopic sparse views. Our method achieves consistent improvement on both color and geometry rendering compared to ActiveNeRF [40] that incorporates the uncertainty by modeling the rendered radiance value as a Gaussian distribution.

1) *NeRF w/ Pre-training*: We compare both the pre-trained generalizable model and fine-tuned model on the SCARED dataset and the Hamlyn dataset. PixelNeRF [10], MVSNeRF [11], GeoNeRF [5] have trained generalizable model on large datasets (DTU [51], LLFF [28], IBRNet [33]) with large cost in time and computation. However, due to the domain gap between general scene and surgical scene, directly adapting the pre-trained generalizable NeRF model to the surgical scene shows unsatisfying results. As shown in Fig. 5, MVSNeRF [11] and GeoNeRF [5] produce blurs and flaws in the rendered color, and depth with floating artifacts and noises on the surgical tissue. After fine-tuning each scene, while the RGB performance has been improved, the geometry performance degenerates due to shape-radiance ambiguity, especially in the region with changing illumination and reflective surface.

2) *NeRF w/ Geometry Guidance*: To enhance the geometry learning of the radiance field, NerfingMVS [13], DS-NeRF [15], Dense Depth Priors [12] take advantage of the sparse depth priors from COLMAP during the pre-processing stage, as constraints for ray sampling range or supervision to rendered depth. However, as presented in Fig. 5, they fail to tackle the severe photometric inconsistencies under sparse surgical views, causing deficiencies in rendered color and depth.

While leveraging multiple monocular geometric priors, MonoSDF [4] exhibits blurred rendered color, lack of detail, and overly smooth surfaces in depth. RegNeRF [14] addresses input sparse views by introducing a patch-based regularizer that improves performance. Although its color reconstruction outperforms other baselines, the degeneration in depth indicates suffering from shape-radiance ambiguity. We

further compare our method with the recent SparseNeRF [35] which regularizes the rendered depth with depth ranking and continuity constraints. Despite its improvements in color rendering, it shows limited depth promotion due to its unified optimization for regions differing in uncertainty. Our UC-NeRF demonstrates robustness to outliers with uncertainty guidance, resulting in accurate depth while preserving view-dependent effects. It achieves consistent improvement in both color and depth rendering results, with higher mean and lower standard deviation compared to baselines.

### C. Ablation Study

1) *Different Components of UC-NeRF*: We conduct an ablation study by gradually adding more components to analyze the effectiveness of different components in our UC-NeRF, presenting the qualitative results in Fig. 6 and quantitative comparisons in Tab. IV. Specifically, we utilize MVSNeRF [11] as the baseline method for comparison. The notation “+ Con. Learner” indicates the addition of our proposed consistency learner. The notation “+ Adapt. Branch” refers to further including the adaptive branch and simply summing up both branches’ prediction. The notation “+ Uncertainty” denotes the proposed uncertainty-aware adaptation used to fuse the two branches (i.e. weighted sum). Finally, the notation “+ Distill.” denotes the inclusion of distillation from monocular geometry priors. Adding the consistency learner helps extract view-consistent correspondences from sparse views with guidance from SfM, facilitating both the geometry and color rendering. Then, the adaptive branch enhances NeRF with the high-level feature, however, causing degeneration in geometry due to the photometric inconsistency (e.g. the hole in the reflective surface in example 2). To address this problem, the uncertainty-aware adaptation merges the output from the base-adaptive branch using the uncertainty map, adaptively handling the photometric inconsistency. Finally, the distillation from monocular geometry priors further enhances the performance.

2) *Effect of deep feature*: We validate the effectiveness of the high-level image features in recovering the view-dependent effect at the adaptive branch. As shown in Fig. 9, integrated with the deep features, the adaptive branch is capable of capturing high-frequency details and view-dependent effects, e.g. reflective surfaces, tiny veins, and blood.



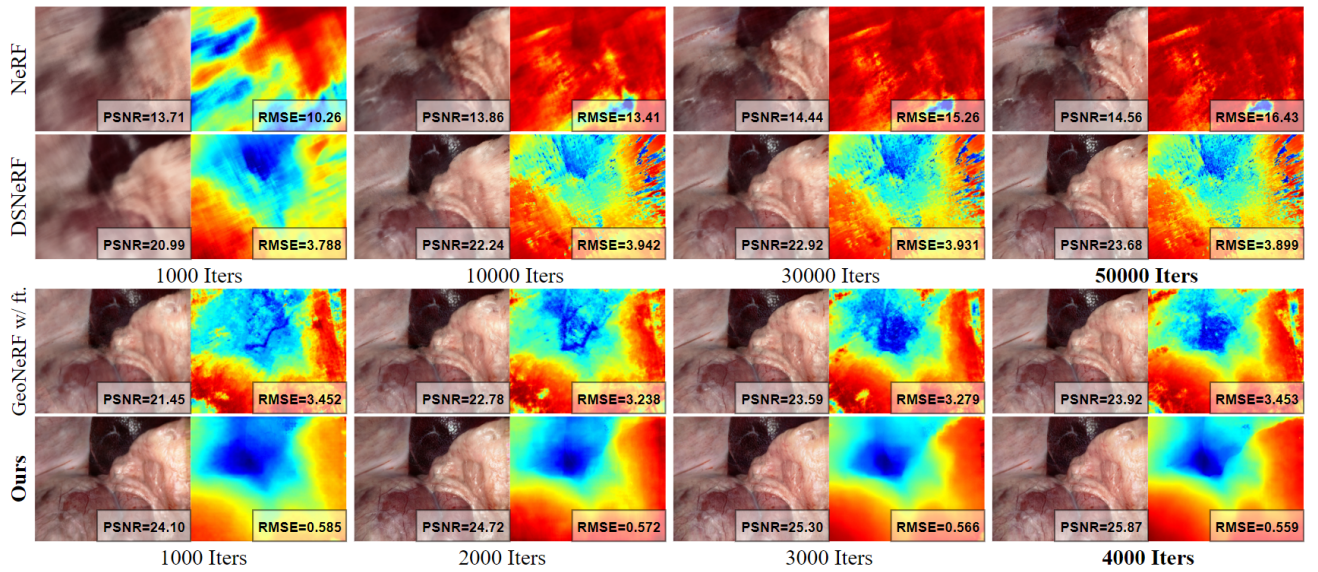


Fig. 7. Visualization for the ablation study on training efficiency. Our UC-NeRF achieves the most efficient and effective training.

TABLE III

COMPUTATION COST OF DIFFERENT METHODS ON THE SCARED DATASET [16]. WE COMPARE THE TIME AND GPU MEMORY REQUIRED FOR TRAINING, AND RENDER TIME DURING INFERENCE.

Methods	Training		Inference
	GPU Hour ↓	GPU Mem. ↓	Render Time ↓
NeRF [1]	34.48h	3.705G	1.855s
ActiveNeRF [40]	35.08h	4.365G	2.655s
RegNeRF [14]	25.84h	21.94G	17.75s
DSNeRF [15]	362.5h	10.21G	1.846s
Dense Prior [12]	35.71h	5.681G	1.834s
SparseNeRF [35]	49.52h	21.94G	19.24s
GeoNeRF w/ ft. [5] <sup>†</sup>	2.529h	10.95G	7.796s
<b>UC-NeRF (Ours)</b>	<b>2.201h</b>	<b>7.638G</b>	<b>0.828s</b>

<sup>†</sup>Pretrained on DTU [51], LLFF [28], IBNet [33].

**3) Efficiency:** To evaluate the efficiency of our UC-NeRF, we visualize the average metrics changing with iterations of “scene 09” in the SCARED Dataset in Fig. 8 and present the qualitative result in Fig. 7. It shows that both the PSNR curve and RMSE curve of our method converge faster in the beginning iterations, while other methods converge slowly and yield low performance. As depicted in Fig. 7, only after 1000 iterations of training, our method has achieved PSNR 24.10 and RMSE 0.586, which outperforms NeRF [1] and DSNeRF [15] in both geometry and color rendering after 50000 iters training. Even after generic training on multiple large datasets, GeoNeRF [5] still yields unsatisfying performance in surgical datasets which have large domain gap to general scenes. To fully compare the training and inference efficiency, we report the time (GPU Hour) and GPU memory for training, and render time during inference in Tab. III. Thanks to our generalizable training across multiple surgical scenes, our UC-NeRF demonstrates superior efficiency in training time. With few hypothesis planes in the consistency learner and a reduction of sampled points during NeRF rendering, the inference time surpasses other state-of-the-art baselines.

**4) Robustness to Number of Source Views:** We conduct an ablation study to validate our method’s robustness to different

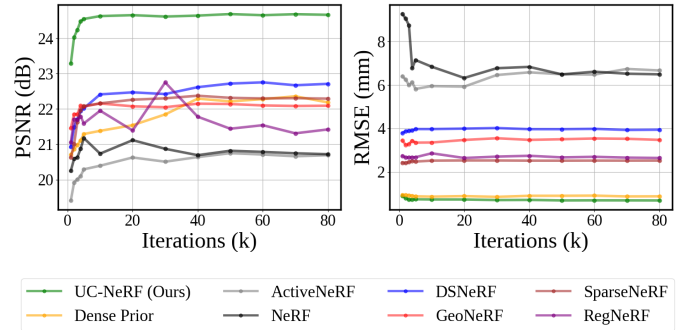


Fig. 8. Training efficiency. We visualize the charts of PSNR-iteration and RMSE-iteration for training “scene 09” in the SCARED Dataset [16].

input source views for the consistency learner. During training, the input source views are selected as the  $N$  closest views to the target view. As shown in Tab. VI, our method is robust to the number of input views and achieves the best performance with 7 input views. The slight performance drop from using 7 views to 9 is probably because using too many source views may introduce excessive information into training, leading to information redundancy with distractions or noises for rendering the target view.

**5) Robustness to Different Training Data Size:** We conduct an ablation study to validate the robustness to training data sizes (i.e. total number of training images). Specifically, after filtering out the low-quality images, we evenly sample different numbers of images (i.e. 50, 30, 10) as training images and validate on the same test set. As shown in Tab. IV, it is observed that our sparse set shows competitive performance with the dense training set while being much more efficient. With the training size increases, it is observed that the performance of rendered images becomes better, because the source views would contain more contextual information to the target view.

This experiment helps validate that our method can tackle surgical datasets of different sizes, showing robustness to the endoscopic sparse views and demonstrating the potential for



TABLE IV  
ABLATION STUDY ON DIFFERENT TRAINING DATA SIZE.

Training Size	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	Abs Rel $\downarrow$	Sq Rel $\downarrow$	RMSE $\downarrow$	RMSE log $\downarrow$	$\delta < 1.25 \uparrow$	GPU Hour $\downarrow$
50	26.73 $\pm$ 1.83	0.860 $\pm$ 0.06	0.109 $\pm$ 0.04	0.055 $\pm$ 0.04	0.128 $\pm$ 0.04	1.327 $\pm$ 0.63	0.082 $\pm$ 0.05	0.967 $\pm$ 0.05	10.83h
30	26.49 $\pm$ 1.44	0.857 $\pm$ 0.06	0.107 $\pm$ 0.05	0.056 $\pm$ 0.03	0.126 $\pm$ 0.05	1.317 $\pm$ 0.55	0.077 $\pm$ 0.03	0.965 $\pm$ 0.04	6.423h
<b>10 (Ours)</b>	26.40 $\pm$ 1.39	0.855 $\pm$ 0.05	0.107 $\pm$ 0.03	0.053 $\pm$ 0.03	0.121 $\pm$ 0.04	1.304 $\pm$ 0.68	0.074 $\pm$ 0.03	0.965 $\pm$ 0.04	2.201h

TABLE V  
ABLATION STUDY ON THE INFLUENCE OF DIFFERENT COMPONENTS.

Components	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	RMSE $\downarrow$
Baseline	22.21 $\pm$ 1.68	0.652 $\pm$ 0.08	0.250 $\pm$ 0.07	5.069 $\pm$ 1.88
+ Con. Learner	25.18 $\pm$ 1.60	0.808 $\pm$ 0.07	0.163 $\pm$ 0.07	1.875 $\pm$ 0.82
+ Adapt. Branch	25.78 $\pm$ 1.69	0.830 $\pm$ 0.06	0.119 $\pm$ 0.06	1.895 $\pm$ 1.27
+ Uncertainty	26.12 $\pm$ 1.45	0.843 $\pm$ 0.06	0.115 $\pm$ 0.03	1.524 $\pm$ 0.81
<b>+ Distill. (Ours)</b>	26.40 $\pm$ 1.39	0.855 $\pm$ 0.05	0.107 $\pm$ 0.03	1.304 $\pm$ 0.68

TABLE VI  
ABLATION STUDY ON THE NUMBER OF SOURCE VIEWS.

Number of Views	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	RMSE $\downarrow$
3 Views	26.13 $\pm$ 1.24	0.841 $\pm$ 0.05	0.116 $\pm$ 0.04	1.474 $\pm$ 0.71
5 Views	26.20 $\pm$ 1.22	0.845 $\pm$ 0.05	0.115 $\pm$ 0.04	1.393 $\pm$ 0.64
<b>7 Views</b>	26.40 $\pm$ 1.39	0.855 $\pm$ 0.05	0.107 $\pm$ 0.03	1.304 $\pm$ 0.68
9 Views	26.15 $\pm$ 1.35	0.847 $\pm$ 0.06	0.118 $\pm$ 0.03	1.321 $\pm$ 0.71

real-world surgery where data acquisition opportunities are naturally restricted.

## V. DISCUSSION

The objective of this paper is to improve NeRF's ability to handle the shape-ambiguity problem caused by challenges in surgical scenes, such as endoscopic sparse views and photometric inconsistencies. Our method is the first to address the challenging sparse view NeRF problem in surgical scenes. The proposed UC-NeRF is both effective and efficient in surgical novel view synthesis because of three main reasons. First, the view-specific uncertainty information is estimated through the consistency learner guided by the sparse SfM depth to measure the extent of photometric inconsistency across multi-view inputs (See Fig. 4). Moreover, our method is capable of generic training across multiple surgical scenes thanks to the conditional inputs from the consistency learner. Second, unlike other few-shot NeRF methods to inject the unified constraints, we take benefits of the learnt uncertainty information to explicitly model the regions in the designed dual-branch NeRF, i.e. fusing the base and adaptive branch with different weight to generate view-dependent appearance and consistent geometry. This design is further demonstrated in Fig. 9 to visualize the different modeling effects decomposed in the two branches. Third, we further introduce the distillation from monocular geometry priors to improve the accuracy and robustness of the rendered depth (see Fig. 6). Specifically, the scale distillation improves the scale consistency through the sparse SfM depth constraints. The uncertainty-guided mono-depth distillation employ the supervision and regularization adaptively in spatial following the uncertainty estimation.

Despite the aforementioned strengths, our method is limited to static or semi-static surgical scene, which hinders

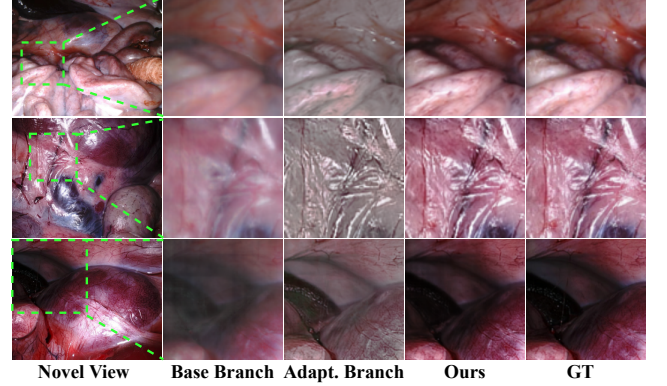


Fig. 9. **View synthesis decomposition in dual-branch NeRF.** We decompose the synthesized images in the base branch, adaptation branch, and the full model of UC-NeRF. It shows the ability of the high-level features to model view-dependent effects.

its application to the highly dynamic surgical scenes. In the future work, we will try to integrate time dimension to enable the spatiotemporal NeRF, allowing both time-varying and view-free rendering for color and depth. More efficient neural representations are also considered to integrate with our method to improve the efficiency for real-time rendering, which benefits the downstream tasks like surgical navigation, 3D reconstruction, surgical skill learning, etc.

## VI. CONCLUSIONS

In this paper, the UC-NeRF network addresses the difficulties posed by surgical sparse views and photometric inconsistencies, achieving robust and efficient novel view synthesis in minimally invasive surgery. By leveraging a consistent learner and an uncertainty map, UC-NeRF enhances geometric correspondence and significantly reduces shape-radiance ambiguity. The introduction of an uncertainty-aware conditional NeRF also refines the learning process of view-dependent appearances, ensuring a balance between geometric precision and photorealistic rendering. Our experimental evaluation shows that our method outperforms other few-shot NeRF methods in both efficiency and effectiveness.

## REFERENCES

- [1] B. Mildenhall *et al.*, "Nerf: Representing scenes as neural radiance fields for view synthesis," *Commun. ACM*, vol. 65, no. 1, pp. 99–106, 2021.
- [2] B. Tobias *et al.*, "Stitching and surface reconstruction from endoscopic image sequences: a review of applications and methods," *IEEE JBHI*, vol. 20, no. 1, pp. 304–321, 2014.
- [3] M. Lena *et al.*, "Optical techniques for 3d surface reconstruction in computer-assisted laparoscopic surgery," *Med. Image. Anal.*, vol. 17, no. 8, pp. 974–996, 2013.
- [4] Z. Yu *et al.*, "Monosdf: Exploring monocular geometric cues for neural implicit surface reconstruction," *NeurIPS*, 2022.
- [5] M. Johari *et al.*, "Geonerf: Generalizing nerf with geometry priors," in *CVPR*, 2022, pp. 18 365–18 375.

- [6] M. Shivali *et al.*, “Augmented reality in surgical navigation: A review of evaluation and validation metrics,” *Appl. Sci.*, vol. 13, no. 3, p. 1629, 2023.
- [7] T. Rui *et al.*, “Augmented reality technology for preoperative planning and intraoperative navigation during hepatobiliary surgery: A review of current methods,” *HBPD INT.*, vol. 17, no. 2, pp. 101–112, 2018.
- [8] P. Veronica *et al.*, “Envisors: Enhanced vision system for robotic surgery. a user-defined safety volume tracking to minimize the risk of intraoperative bleeding,” *Front. Robot. AI*, vol. 4, p. 15, 2017.
- [9] C. Long *et al.*, “Slam-based dense surface reconstruction in monocular minimally invasive surgery and its application to augmented reality,” *Comput. Methods Programs. Biomed.*, vol. 158, pp. 135–146, 2018.
- [10] A. Yu *et al.*, “pixelnerf: Neural radiance fields from one or few images,” in *CVPR*, 2021, pp. 4578–4587.
- [11] A. Chen *et al.*, “Mvsnerf: Fast generalizable radiance field reconstruction from multi-view stereo,” in *ICCV*, 2021, pp. 14 124–14 133.
- [12] B. Roessle *et al.*, “Dense depth priors for neural radiance fields from sparse input views,” in *CVPR*, 2022, pp. 12 892–12 901.
- [13] Y. Wei *et al.*, “Nerfingmvs: Guided optimization of neural radiance fields for indoor multi-view stereo,” in *ICCV*, 2021, pp. 5610–5619.
- [14] M. Niemeyer *et al.*, “Regnerf: Regularizing neural radiance fields for view synthesis from sparse inputs,” in *CVPR*, 2022, pp. 5480–5490.
- [15] K. Deng *et al.*, “Depth-supervised nerf: Fewer views and faster training for free,” in *CVPR*, 2022, pp. 12 882–12 891.
- [16] M. Allan *et al.*, “Stereo correspondence and reconstruction of endoscopic data challenge,” *arXiv:2101.01133*, 2021.
- [17] P. Mountney *et al.*, “Three-dimensional tissue deformation recovery and tracking,” *IEEE Signal Process. Mag.*, vol. 27, no. 4, pp. 14–24, 2010.
- [18] D. Stoyanov *et al.*, “Real-time stereo reconstruction in robotically assisted minimally invasive surgery,” in *MICCAI*, 2010, pp. 275–282.
- [19] P. Pratt *et al.*, “Dynamic guidance for robotic surgery using image-constrained biomechanical models,” in *MICCAI*, 2010, pp. 77–85.
- [20] C. McMillan *et al.*, “Unstructured lumigraph rendering,”
- [21] S. Sinha *et al.*, “Piecewise planar stereo for image-based rendering,” in *ICCV*, 2009, pp. 1881–1888.
- [22] G. Chaurasia *et al.*, “Depth synthesis and local warps for plausible image-based navigation,” *ACM TOG*, vol. 32, no. 3, pp. 1–12, 2013.
- [23] M. LEVOY *et al.*, “Light field rendering,” in *Computer graphics proceedings, annual conference series*, 1996, pp. 31–42.
- [24] D. WOOD, “Surface light fields for 3d photography,” *SIGGRAPH2000, Jul.*, pp. 287–296, 2000.
- [25] A. Chen *et al.*, “Deep surface light fields,” *ACM Comput. Graph.*, vol. 1, no. 1, pp. 1–17, 2018.
- [26] P. Srinivasan *et al.*, “Learning to synthesize a 4d rgb-d light field from a single image,” in *ICCV*, 2017, pp. 2243–2251.
- [27] I. Choi *et al.*, “Extreme view synthesis,” in *ICCV*, 2019, pp. 7781–7790.
- [28] B. Mildenhall *et al.*, “Local light field fusion: Practical view synthesis with prescriptive sampling guidelines,” *ACM TOG*, vol. 38, no. 4, pp. 1–14, 2019.
- [29] P. Srinivasan *et al.*, “Pushing the boundaries of view extrapolation with multiplane images,” in *CVPR*, 2019, pp. 175–184.
- [30] Y. Wang *et al.*, “Neural rendering for stereo 3d reconstruction of deformable tissues in robotic surgery,” in *MICCAI*, 2022, pp. 431–441.
- [31] X. Sun *et al.*, “Dynamic surface reconstruction in robot-assisted minimally invasive surgery based on neural radiance fields,” *IJCAI*, 2023.
- [32] R. Zha *et al.*, “Endosurf: Neural surface reconstruction of deformable tissues with stereo endoscope videos,” in *MICCAI*, 2023, pp. 13–23.
- [33] Q. Wang *et al.*, “Ibrnet: Learning multi-view image-based rendering,” in *CVPR*, 2021, pp. 4690–4699.
- [34] J. Schonberger *et al.*, “Structure-from-motion revisited,” in *ICCV*, 2016, pp. 4104–4113.
- [35] G. Wang *et al.*, “Sparsenerf: Distilling depth ranking for few-shot novel view synthesis,” in *ICCV*, 2023.
- [36] S. Hu *et al.*, “Consistentnerf: Enhancing neural radiance fields with 3d consistency for sparse view synthesis,” *arXiv:2305.11031*, 2023.
- [37] M. Ricardo *et al.*, “Nerf in the wild: Neural radiance fields for unconstrained photo collections,” in *CVPR*, 2021, pp. 7210–7219.
- [38] J. Shen *et al.*, “Stochastic neural radiance fields: Quantifying uncertainty in implicit 3d representations,” in *3DV*, 2021, pp. 972–981.
- [39] —, “Conditional-flow nerf: Accurate 3d modelling with reliable uncertainty quantification,” in *ECCV*, 2022, pp. 540–557.
- [40] X. Pan *et al.*, “Activenerf: Learning where to see with uncertainty estimation,” in *ECCV*, 2022, pp. 230–246.
- [41] X. Gu *et al.*, “Cascade cost volume for high-resolution multi-view stereo and stereo matching,” in *CVPR*, 2020, pp. 2495–2504.
- [42] T. Lin *et al.*, “Feature pyramid networks for object detection,” in *CVPR*, 2017, pp. 2117–2125.
- [43] Y. Yao *et al.*, “MVSNet: Depth inference for unstructured multi-view stereo,” in *ECCV*, 2018, pp. 767–783.
- [44] R. Ranftl *et al.*, “Vision transformers for dense prediction,” in *ICCV*, 2021, pp. 12 179–12 188.
- [45] R. René *et al.*, “Towards robust monocular depth estimation: Mixing datasets for zero-shot cross-dataset transfer,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 44, no. 3, pp. 1623–1637, 2020.
- [46] D. Recasens *et al.*, “Endo-depth-and-motion: Reconstruction and tracking in endoscopic videos using depth networks and photometric constraints,” *IEEE RAL*, vol. 6, no. 4, pp. 7225–7232, 2021.
- [47] J.L. Schönberger *et al.*, “Structure-from-motion revisited,” in *CVPR*, 2016.
- [48] Z. Wang *et al.*, “Image quality assessment: from error visibility to structural similarity,” *IEEE Trans. Image Process.*, vol. 13, no. 4, pp. 600–612, 2004.
- [49] R. Zhang *et al.*, “The unreasonable effectiveness of deep features as a perceptual metric,” in *CVPR*, 2018, pp. 586–595.
- [50] T. Zhou *et al.*, “Unsupervised learning of depth and ego-motion from video,” in *ICCV*, 2017, pp. 1851–1858.
- [51] R. Jensen *et al.*, “Large scale multi-view stereopsis evaluation,” in *CVPR*, 2014, pp. 406–413.