

Efficient Deformable Tissue Reconstruction via Orthogonal Neural Plane

Chen Yang, Kailing Wang, Yuehao Wang, Qi Dou, Xiaokang Yang, Wei Shen

Abstract—Intraoperative imaging techniques for reconstructing deformable tissues in vivo are pivotal for advanced surgical systems. Existing methods either compromise on rendering quality or are excessively computationally intensive, often demanding dozens of hours to perform, which significantly hinders their practical application. In this paper, we introduce Fast Orthogonal Plane (Forplane), a novel, efficient framework based on neural radiance fields (NeRF) for the reconstruction of deformable tissues. We conceptualize surgical procedures as 4D volumes, and break them down into static and dynamic fields comprised of orthogonal neural planes. This factorization discretizes the four-dimensional space, leading to a decreased memory usage and faster optimization. A spatiotemporal importance sampling scheme is introduced to improve performance in regions with tool occlusion as well as large motions and accelerate training. An efficient ray marching method is applied to skip sampling among empty regions, significantly improving inference speed. Forplane accommodates both binocular and monocular endoscopy videos, demonstrating its extensive applicability and flexibility. Our experiments, carried out on two in vivo datasets, the EndoNeRF and Hamlyn datasets, demonstrate the effectiveness of our framework. In all cases, Forplane substantially accelerates both the optimization process (by over 100 times) and the inference process (by over 15 times) while maintaining or even improving the quality across a variety of non-rigid deformations. This significant performance improvement promises to be a valuable asset for future intraoperative surgical applications. The code of our project is now available at <https://github.com/Loping151/ForPlane>.

Index Terms—Endoscopy, tissues reconstruction, optical imaging.

I. INTRODUCTION

RECONSTRUCTING deformable tissues from endoscopy videos in surgery is a crucial and promising field in medical image computing. High-quality tissue reconstruction can assist surgeons in avoiding critical structures, *e.g.*, blood

vessels and nerves, as well as improve the observation of lesions such as tumors and enlarged lymph nodes [3], [38]. Moreover, tissue reconstruction significantly contributes to creating a virtual surgical training environment [15]. Such detailed reconstructions not only facilitate skill acquisition and expedite the learning curve for endoscopists but also generate substantial training data. This data is crucial for developing virtual reality (VR) and augmented reality (AR) training modules in surgical education, as well as for enhancing the learning algorithms of surgical robots [20], [21]. However, achieving precise identification and visualization of these structures and lesions remains a challenge for existing methods. Moreover, the extensive training and inference time required by these methods hinders their practical application during surgery. This paper aims to address these challenges and seize the opportunities by focusing on high-quality real-time reconstruction of deformable tissues in both monocular and binocular endoscopy videos, providing valuable insights for future intraoperative applications.

Previous methods [23]–[25], [42] primarily focus on static surgical scenes and disregard the presence of diverse surgical instruments in endoscopy videos. These methods rely on simultaneous localization and mapping (SLAM) techniques, generating individual meshes for each frame to reconstruct surgical procedures. However, they fail to account for non-rigid transformations in deformable tissues and potential occlusions caused by surgical instruments. As a result, their reconstruction performance lacks fidelity, contains gaps, and does not realistically represent the surgical environment. These limitations significantly hinder their applicability in data simulation and intraoperative assistance.

Recently, the integration of implicit scene representations and differentiable volume rendering has exhibited remarkable efficacy in capturing complex scenarios. A notable example is neural radiance fields (NeRF) [26], a seminal work in neural rendering that introduced neural implicit fields for continuous scene representations. NeRF has demonstrated exceptional performance in tasks such as high-quality view synthesis and 3D reconstruction across diverse contexts [48]. Furthermore, NeRF and its variants have showcased diverse applications in medical imaging, encompassing but not restricted to medical imaging segmentation [13], deformable tissue reconstruction [46], [51], radiograph measurement reconstruction [6], [16], and organ shape completion [37], [52].

A notable advancement in reconstructing deformable tissue is EndoNeRF [46], a recent approach that has demonstrated

This work was supported in part by the National Key R&D Program of China 2022YFF1202600, in part by the National Natural Science Foundation of China under Grant 62176159, in part by the Natural Science Foundation of Shanghai 21ZR1432200, and in part by the Shanghai Municipal Science and Technology Major Project 2021SHZDZX0102. (Corresponding author: Wei Shen.)

C. Yang, K. Wang, X. Yang and W. Shen are with the MoE Key Laboratory of Artificial Intelligence, AI Institute, Shanghai Jiao Tong University, Shanghai 200240, China (e-mail: {ycyangchen, wangkailing151, xkyang, wei.shen}@sjtu.edu.cn).

Y. Wang and Q. Dou are with Dept. of Computer Science and Engineering, The Chinese University of Hong Kong, Hong Kong 999077 (e-mail: yhwang21@cse.cuhk.edu.hk; qidou@cuhk.edu.hk)

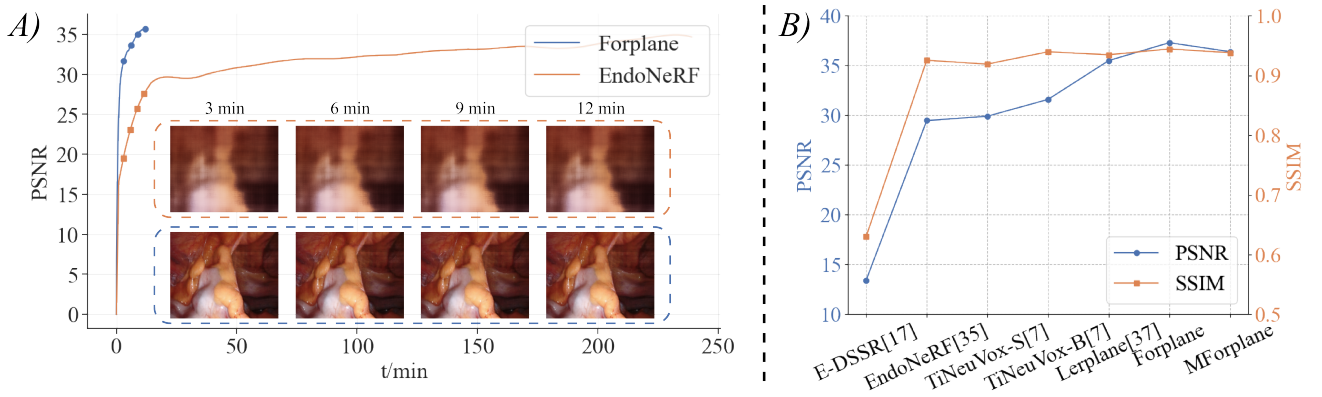


Fig. 1. A) The speed-quality comparison between our Forplane and EndoNeRF [46]. We present the performance of the two methods in terms of PSNR and provide insights into their training times. Furthermore, we showcase the reconstruction results obtained from both methods at various training time intervals, specifically at 3, 6, 9, and 12 minutes. These results clearly demonstrate the significant convergence speed exhibited by Forplane. B) The metrics (PSNR, SSIM) of 7 different methods on the *pushing tissues* scene from EndoNeRF dataset [46]. Both metrics exhibit a positive correlation, indicating that higher values shows better performance.

exceptional capabilities in 3D reconstruction and deformation tracking of surgical scenes within the context of robotic surgery. EndoNeRF [46] employs a canonical neural radiance field in conjunction with a time-dependent neural displacement field to effectively model deformable tissues using binocular captures within a single viewpoint configuration. However, despite its impressive achievements in reconstructing deformable tissues, EndoNeRF [46] faces computational challenges due to intensive optimization processes. The optimization process for EndoNeRF [46] typically requires tens of hours to complete, as each pixel generated necessitates a substantial number of neural network calls. This computational bottleneck significantly restricts the widespread adoption of such methods in surgical procedures. Additionally, its performance heavily relies on precise binocular depth estimation, further impeding its broader applications.

To overcome the challenges mentioned above, we propose a novel approach called Forplane (Fast Orthogonal Plane), which delivers efficient and high-quality deformable tissue reconstruction and offers rapid training and inference with endoscopy videos. Forplane integrates space-time decomposition with neural radiance fields to achieve rapid and accurate reconstruction of deformable tissues during surgical procedures. By treating surgical procedures as 4D volumes, with time as an orthogonal axis to spatial coordinates, Forplane discretizes the continuous space into static and dynamic fields. Static fields are represented by spatial planes, while dynamic fields are represented by space-time planes. We utilize limited resolution features to represent these planes, allowing for bilinear interpolation when querying features at any spatio-temporal point. This design markedly reduces computational costs compared to MLP-reliant methods [32], [33], [46], cutting down complexity from $O(N^4)$ to $O(N^2)$. Additionally, the static field facilitates information sharing across neighboring timesteps, addressing the limitations imposed by restricted viewpoints. Furthermore, drawing inspiration from the observation that certain tissues exhibit more frequent deformations, we develop an importance sampling method, which strategically higher sampling probability towards tissues that are either

occluded by surgical tools or exhibit more extensive motion range.

A preliminary version of Forplane was introduced in Lerplane [51], where its ability to achieve rapid reconstruction of deformable tissues on binocular endoscopy videos was demonstrated. **1) Advancements in Rendering and Speed:** In this extended version, we have made notable advancements on rendering quality, inference speed and adaptability. Specifically, we develop an efficient ray marching algorithm to guide the rendering procedure, significantly improving inference speed and rendering quality. **2) MForplane for Monocular Videos:** Considering that many surgical procedures employ monocular devices, we present MForplane, a version of Forplane adapted for monocular videos. MForplane performs effectively on monocular sequences alone with minimal performance drop, though this is a more challenging task. **3) Comprehensive Evaluation with Hamlyn Datasets:** In addition to validating our results with the EndoNeRF dataset [46], we have expanded our evaluation to include the public Hamlyn dataset [28], [44]. This dataset offers more complex sequences that encompass various challenging factors such as intracorporeal scenes with weak textures, deformations, motion blur, reflections, surgical tools, and occlusions. The inclusion of this dataset allows us to assess the performance of Forplane in more diverse and realistic surgical scenarios. Fig. 1 shows that Forplane excels in its significantly faster optimization, demonstrating superior quantitative and qualitative performance in 3D reconstruction and deformation tracking within surgical scenes in comparison to preceding methods. Fig. 2 demonstrates the superior reconstruction quality of Forplane. These results herald significant potential for future intraoperative applications.

The remainder of this paper is organized as follows. Section II provides a concise review of the literature pertinent to our domain. Section III delineates the principal components of our Forplane framework. Section IV details the experimental results, validating the efficacy of our approach. In Section V, we elucidate the mechanisms underlying Forplane's capacity for efficient reconstruction of deformable tissues and outlines future research directions. Finally, Section VI synthesizes our

findings and contributions into a conclusion.

II. RELATED WORK

A. Surgical Scene Reconstruction

Numerous endeavors have been undertaken to reconstruct deformable tissues within surgical scenes.

1) *Non-implicit Representations*: There have been numerous noteworthy non-implicit approaches in the realm of surgical scene reconstruction [41]. Earlier SLAM-based studies such as [43], [54], [55] leveraged depth estimation from stereo videos and fuse depth maps in 3D space for reconstruction. However, these methods either neglected the presence of surgical tools or oversimplified the scenes by presuming them to be static. Later advancements, such as SuPer [18] and EDSSR [22], proposed frameworks for stereo depth estimation with tool masking, and SurfelWarp [12] performed single-view 3D deformable reconstruction. These approaches relied heavily on deformation tracking via a sparse warp field, which compromised their efficacy when encountering deformations that exceed simple non-topological changes.

2) *Novel Implicit Representations*: Implicit representations like NeRF [26] have made remarkable contributions to medical imaging. In contrast to non-implicit representations, which encode discrete features or signal values directly, implicit representations use generator functions that associate input coordinates with their respective values within the input space [27]. Recently, EndoNeRF [46] emerged as a promising solution. It utilized NeRF with tool-guided ray casting, stereo depth-cueing ray marching, and stereo depth-supervised optimization, yielding high-quality non-rigidity reconstruction. However, EndoNeRF optimized an entire spatial temporal field, resulting in significant time and resource consumption. Lerplane [51] factorizes the scene into explicit 2D planes of static and dynamic fields to accelerate optimization, but the inference speed is insufficient to meet clinical demands. Therefore, for effective inter-operative application, more intensive efforts must be dedicated towards speeding up the process without sacrificing the reconstruction performance.

B. Implicit Representations in Medical Imaging

Implicit representations have yield significant contributions in other realms of medical imaging. Some methods leverage NeRF to perform visualization of medical images. For instance, Li et al. [16] used NeRF to acquire 3D ultrasound reconstructed spine image volume to assess spine deformity. MedNeRF [6] rendered CT projections given a few or even a single-view X-ray utilizing NeRF and GAN. NeAT [37] used a hybrid explicit-implicit neural representation for tomographic image reconstruction, showing better quality than traditional methods. CoIL [45] leveraged coordinate-based neural representations for estimating high-fidelity measurement fields in the context of sparse-view CT. EndoNeRF [46], applied dynamic nerf [33] to perform deformable surgical scene reconstruction, showing promising performance on stereo 3D reconstruction of deformable tissues in robotic surgery. Some methods use implicit functions to do shape reconstruction. ImplicitAtlas [52] proposed a data-efficient shape model based

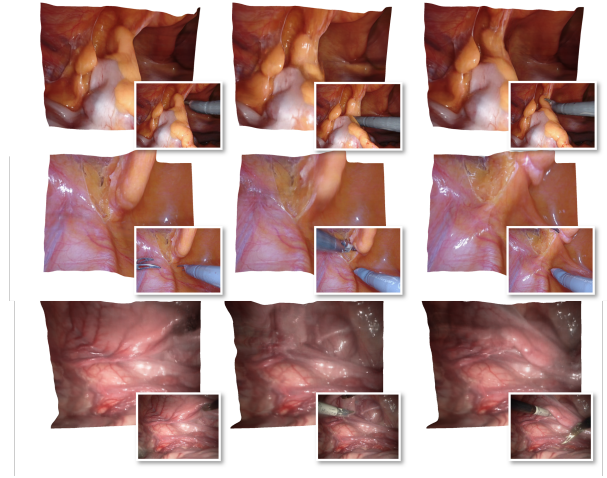


Fig. 2. Reconstruction results of deformable tissues. We show the deformable tissues reconstructed by Forplane among surgical procedures with the corresponding captured images in the lower-right corner.

on templates for organ shape reconstruction and interpolation. Fang et al. [9] introduced a curvature-enhanced implicit function network for high-quality tooth model generation from CBCT images. Raju et al. [34] presented deep implicit statistical shape models (DISSMs) for the 3D delineation of medical images. For segmentation, IOSNET [13] used neural fields to create continuous segmentation maps which converge fast and are memory-efficient. TiAVox [56] used a time-aware attenuation voxel approach for sparse-view 4D DSA reconstruction. NeSVoR [49] modeled the underlying volume as a continuous function to perform slice-to-volume reconstruction. Reed et al. devised a reconstruction pipeline that utilizes implicit neural representations in conjunction with a novel parametric motion field warping technique to perform limited-view 4D-CT reconstruction of rapidly deforming scenes [36]. Schmidt et al. introduced RING [40] which estimates flow efficiently and KINFlow [39] which enables a prior-free estimation of deformation, both using implicit neural representation and graph-based model.

III. MATERIALS AND METHODS

We focus on the task of efficiently reconstructing deformable tissues from both monocular and binocular endoscopy videos. Mathematically, we represent a surgical procedure as a 4D volume denoted by \mathbf{V}_{4d} , with dimensions $H \times W \times D \times T$. Here, H, W, D represents the 3D space of the scene, and the T represents the time dim, assumed to be orthogonal to the 3D space. Unlike previous methods that treat the surgical procedure as independent static 3D volumes per time step $\{\mathbf{V}_{3d}^1, \mathbf{V}_{3d}^2, \dots, \mathbf{V}_{3d}^T\}$ or use a time-dependent displacement field \mathbf{G}_ϕ to model the tissue deformations, we factorizes the 4D volumes into a static field \mathbf{V}_{3d}^s and a dynamic field \mathbf{V}_{3d}^d composed of 2D neural planes. This factorization significantly enhances convergence speed and improves the representational capacity.

In this section, we begin with reviewing the key techniques used in NeRF [26] and EndoNeRF [46] (Sec. III-A). Subsequently, we introduce our novel and efficient orthogonal

plane representation for surgical scenes (Sec. III-B). To reconstruct deformable tissues at any time step, we first employ a novel spatiotemporal sampling algorithm to identify high-priority tissue pixels and generate corresponding rays (Sec. III-C). Next, we design one efficient ray marching method to generate discrete samples along the selected rays (Sec. III-D). These samples are then used to retrieve corresponding features from orthogonal planes with linear interpolation. The obtained features, along with the spatiotemporal information of the samples, are fed into a lightweight MLP that predicts radiance and density for each sample (Sec. III-E). Finally, we apply standard volume rendering to render the accumulated color and depth along the selected rays. To improve the spatiotemporal smoothness and decomposition, we design various regularization strategies on $\{\mathbf{V}_{3d}^s, \mathbf{V}_{3d}^d\}$ as well as the lightweight MLP (Sec. III-F). We present an optimization strategy for Forplane that diminishes its reliance on binocular depth so that it can work with monocular scopes (Sec. III-G). An illustration of the overall framework is presented in Fig. 3.

A. NeRF and EndoNeRF

We provide a brief overview of the pipeline of NeRF [26] and EndoNeRF [46]. Given a set of posed images $\{\mathbf{I}_i\}_{i=1}^N$, NeRF represents a 3D scene using volume density and directional emitted radiance for each point in space with a coordinate-based neural network θ . To render one image \mathbf{I}_i , a ray casting process is performed for each pixel, where a ray $\mathbf{r}(t) = \mathbf{x}_o + t\mathbf{d}$ is projected through the camera pose. Here, \mathbf{x}_o is the camera origin, \mathbf{d} is the ray direction, and t denotes the distance of a point along the ray from the origin. NeRF uses a positional encoding $\gamma(\cdot)$ to map the coordinates \mathbf{x} and viewing direction \mathbf{d} into a higher dimensional space:

$$\gamma(\mathbf{x}) = \left[\sin(\mathbf{x}), \cos(\mathbf{x}), \dots, \sin(2^{L-1}\mathbf{x}), \cos(2^{L-1}\mathbf{x}) \right]^T. \quad (1)$$

The process of casting a ray through the scene and generating discrete samples along the ray is termed ray marching. The color $\hat{\mathbf{C}}(\mathbf{r})$ for a specific ray $\mathbf{r}(t)$ is computed using volume rendering, which involves integrating the weighted volumetric radiance within the near and far bounds t_n and t_f of the ray:

$$\hat{\mathbf{C}}(\mathbf{r}) = \int_{t_n}^{t_f} w(t) \cdot \underbrace{\mathbf{c}(\mathbf{r}(t), \mathbf{d})}_{\text{radiance}} dt \quad (2)$$

The integration weights $w(t)$ for volume rendering are given by:

$$w(t) = \underbrace{\exp\left(-\int_{t_n}^t \sigma(\mathbf{r}(s)) ds\right)}_{\text{visibility of } \mathbf{r}(t) \text{ from } \mathbf{o}} \cdot \underbrace{\sigma(\mathbf{r}(t))}_{\text{density at } \mathbf{r}(t)} \quad (3)$$

The training of θ is supervised by an L2 photometric reconstruction loss.

EndoNeRF [46] employs a canonical radiance field $F_\Theta(\mathbf{x}, \mathbf{d})$ and a time-dependent displacement field $G_\Phi(\mathbf{x}, \tau)$ to represent the surgical scene. The time-dependent displacement field maps the input space-time coordinates (\mathbf{x}, τ) to the displacement between point \mathbf{x} at time step τ and its

corresponding point in the canonical field. For any given time step τ , the radiance and density at \mathbf{x} can be obtained by querying $F_\Theta(\mathbf{x} + G_\Phi(\mathbf{x}, \tau), \mathbf{d})$. EndoNeRF employs the same positional encoding algorithm to map the input coordinates and time step into Fourier features before feeding them into the networks.

B. Fast Orthogonal Plane Representation

The surgical procedure consists of a series of consecutive frames, each representing a separate scene. Notwithstanding the mutable nature of tissue transformations over time, a significant portion of the tissue structure demonstrates continuity across consecutive frames. Observing this, we propose a method that efficiently capitalizes on the time-invariant components of the tissue structures to construct a static field \mathbf{V}_{3d}^s . The static field is designed to encapsulate invariant tissue structures across these frames, thereby enabling the reuse of static components without necessitating data duplication for each temporal step. This design significantly reduces the memory requirement by obviating the need to repeatedly store information pertaining to the static aspects of the scene in every frame. In addressing the time-aware deformations, we have formulated a dynamic field, denoted as \mathbf{V}_{3d}^d . This field is uniquely tasked with capturing the deviations from the static field. This focused approach allows the dynamic field to circumvent the need for a comprehensive reconstruction of the entire scene, thereby substantially diminishing the computational burden. To build the static field, we adopt the orthogonal space planes (*i.e.*, G_{XY} , G_{YZ} and G_{XZ}) to represent the static components among the surgical procedure. The use of orthogonal space planes has been proven effective and compact in various static scene reconstruction methods [5], [29]. As for the \mathbf{V}_{3d}^d , it is designed to capture the time-dependent appearance. Since the time dimension is orthogonal to the 3D space, it is natural to employ space-time planes (*i.e.*, G_{TX} , G_{TY} and G_{TZ}) to represent the dynamic field. Each space plane has dimensions of $N \times N \times D$, and each space-time plane has dimensions of $N \times M \times D$, where N and M denote the spatial and temporal resolutions, respectively, and D represents the size of the feature stored within the plane.

To render one tissue pixel p_{ij} in a specific time step τ , we first cast a ray $\mathbf{r}(t)$ from \mathbf{x}_o to the pixel. We then sample spatial-temporal points along the ray, obtaining their 4D coordinates. We acquire a feature vector for a point $\mathbf{P}(x, y, z, \tau)$ by projecting it onto each plane and using bilinear interpolation \mathcal{B} to query features from the six feature planes:

$$\mathbf{v}(x, y, z, \tau) = \mathcal{B}(G_{XY}, x, y) \odot \mathcal{B}(G_{YZ}, y, z) \cdots \mathcal{B}(G_{YT}, y, \tau) \odot \mathcal{B}(G_{ZT}, z, \tau), \quad (4)$$

where \odot represents element-wise multiplication, inspired by [4], [10], [11]. The fused feature vector \mathbf{v} is then passed to a tiny MLP θ , which predicts the color $\mathbf{c}(\mathbf{r}(t), \mathbf{d})$ and density $\sigma(\mathbf{r}(t))$ of the point. Finally, we leverage the Eq. 2 to get the predicted color $\hat{\mathbf{C}}$. Inspired by the hybrid representation of static fields [29], we build $\{\mathbf{V}_{3d}^s, \mathbf{V}_{3d}^d\}$ with multi-resolution planes.

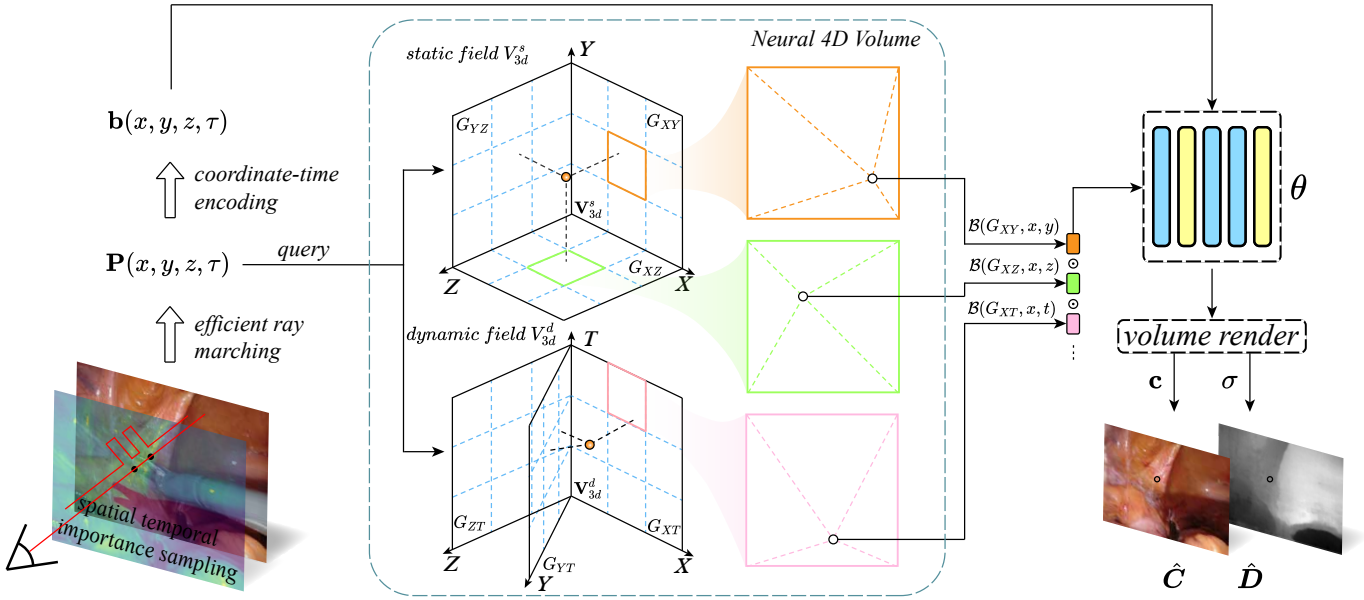


Fig. 3. The rendering pipeline of Forplane. Forplane begins with spatiotemporal importance sampling (Sec. III-C), which assigns higher sampling probabilities to tissue pixels occluded by tools or with extensive motion. Following the selection of the pixel (ray) to be rendered, an efficient ray marching algorithm (Sec. III-D) samples spatial points along the ray near the tissue surface. The coordinates and time embeddings of sampled points are used to query their corresponding features from multi-scale orthogonal feature planes (e.g., G_{XY}) via linear interpolation (Sec. III-B). For simplicity, single-resolution planes exemplify these multi-scale neural planes. The features from different planes and scales are fused using element-wise multiplication in Forplane. To enhance temporal information, the fused plane features are concatenated with positional encoded coordinate and time embeddings (Sec. III-E and Fig. 4), and then passed through a lightweight MLP θ . Finally, a volume rendering scheme is applied to generate predicted color and depth values for each selected ray.

The factorization of surgical procedure brings three main benefits: 1) Significant acceleration: Existing methods for reconstructing surgical procedures using pure implicit representations require traversing all possible positions in space-time, resulting in high computational and time complexity. In contrast, querying a Forplane is fast and efficient, involving only several bilinear interpolations among feature planes and a vector matrix product. This reduces the computational cost from $O(N^4)$ to $O(N^2)$. Moreover, this factorization lightens the burden on the MLP, enabling the use of lighter MLP architectures, comprising just 2 fully connected layers, each with 64 channels. In contrast, EndoNeRF [46] employs a complex MLP structure, processing the positional encoding of the input location $\gamma(x)$ through 8 fully connected ReLU layers, each with 256 channels. 2) Sharing learned scene priors: The static field in Forplane facilitates information sharing across different time steps, improving reconstruction performance by incorporating learned scene priors. 3) Modeling arbitrary deformation: Compared with displacement field used by [7], [8], [19], [32], [46] which struggles with changes in scene topology, the dynamic field can easily model complex tissue deformations, allowing for more accurate reconstruction in surgical scenes.

C. Spatiotemporal Importance Sampling

Tool occlusion during robotic surgery poses significant challenges in accurately reconstructing occluded tissues due to their infrequent representation in the training set. As a result, different pixels encounter varying levels of difficulty in the learning process, exacerbating the complexity of tissue reconstruction. Furthermore, the presence of stationary tissues

over time leads to repeated training on these pixels, yielding minimal impact on convergence and diminishing overall efficiency [50]. To tackle these challenges, we devise a novel spatiotemporal importance sampling strategy. It prioritizes tissue pixels that have been occluded by tools or exhibit extensive motion ranges by assigning them higher sampling probabilities.

In particular, we utilize binary masks $\{M_i\}_{i=1}^T$ and temporal differences among frames to generate sampling weight maps $\{W_i\}_{i=1}^T$. These weight maps represent the sampling probabilities for each pixel/ray, drawing inspiration from EndoNeRF [46]. One sampling weight map W_i can be determined by:

$$W_i = \min(\max_{\substack{i-n < j \\ < i+n}} (\|I_i \odot M_i - I_j \odot M_j\|_1) / 3, \alpha) \odot \Omega_i, \\ \Omega_i = \beta (T M_i / \sum_{i=1}^T M_i), \quad (5)$$

where α is a lower-bound to avoid zero weight among unchanged pixels, Ω_i specifies higher importance scaling for those tissue areas with higher occlusion frequencies, and β is a hyper-parameter for balancing augmentation among frequently occluded areas and time-variant areas. By unitizing spatiotemporal importance sampling, Forplane concentrates on tissue areas and speeds up training, improving the rendering quality of occluded areas and prioritizing tissue areas with higher occlusion frequencies and temporal variability.

D. Efficient Ray Marching

The process of volume rendering greatly benefits from the precise sampling of spatiotemporal points, particularly around

tissue regions. Both EndoNeRF [46] and the original version Lerplane [51] introduce specialized sampling schemes to enhance point accuracy. EndoNeRF’s approach involves a stereo depth-cueing ray marching module that utilizes stereo depth data to guide point sampling near tissue surfaces. However, this method is limited in scenarios lacking accurate stereo depth information. Conversely, our original version Lerplane [51] employs a sample-net for surface-focused sampling. While effective, this approach requires extensive sampling across the pipeline, leading to reduced inference speeds.

To address this issue, we design an indicator grid that directs the Forplane’s focus towards the tissue surface. This method is inspired by the grid representation method [5], [11], [17], [29] used in static scene reconstruction, and it serves to optimize ray marching. The indicator grid functions as a binary map which is cached and updated throughout the training process, signifying empty areas. This low cost grid allows for the early termination of the marching process, based on the transmittance along the ray.

For optimization efficiency, we integrate the indicator grid update with Forplane’s standard training. We start with the assumption that all space is dense, and the indicator grid is fully occupied. As we introduce a series of spatial temporal points to Forplane, we pass a small batch sampled from them to the tiny MLP θ to obtain the corresponding density, updating the indicator grid accordingly. Initially, the optimization process is somewhat slow. However, after several iterations, the indicator grid becomes capable of identifying the density distribution within the 4D volume. This capability significantly reduces the need for unnecessary point sampling, which in turn greatly enhances the speed of training.

Our low cost indicator grid allows for efficient ray marching during both the optimization and inference stage. Leveraging the pre-learned scene context distribution allows us to ignore points with minimal contribution to rendering results and sample points near the surface more accurately. This change improves rendering quality, especially for deformable tissues. Compared to the original version [51], our novel efficient ray marching enables faster training and substantially improves rendering speed, approximately 5 times faster (*1.73 fps* vs. *0.38 fps*), representing a significant advancement for intraoperative use.

E. Coordinate-Time Encoding

In Section III-A, it is illustrated that previous methods utilize a positional encoding function $\gamma(\cdot)$ to map coordinates and viewing directions into a higher-dimensional space. The positional encoding of the viewing direction is exploited to capture the view-dependent appearance, which operates on the assumption of a stationary environmental light source. However, this modeling method proves ineffective within the context of an in vivo sequence. The primary reason for this inefficacy is the dynamic nature of the light source, which moves in conjunction with the camera, bound by specific constraints. Consequently, the modeling of the view-dependent appearance becomes arbitrary, which leads to inaccurate scene reconstruction.

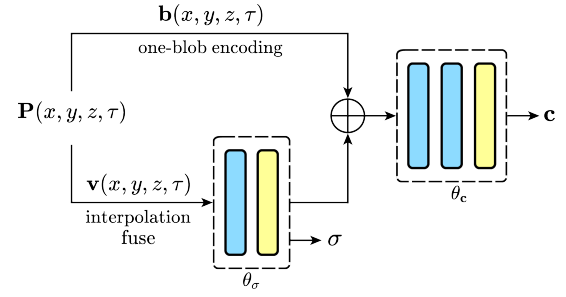


Fig. 4. The procedure of Coordinate-Time Encoding. We pass positional encoded coordinate and time embeddings to enhance the temporal information. The tiny MLP (θ) consists of the sigma net θ_σ and the color net θ_c . The blue and yellow boxes denote MLP layers with and without ReLU activation.

Algorithm 1: The Rendering Procedure

Input : A specific ray $\mathbf{r}(t) = (\mathbf{x}_o + t\mathbf{d})$, 4 levels of the static field and the dynamic field $\{\mathbf{V}_{3d}^s, \mathbf{V}_{3d}^d\}^4$, the indicator grid \mathbf{G}_o and the tiny MLP θ

Output : The expected color $\hat{\mathbf{C}}$ and depth $\hat{\mathbf{D}}$ of $\mathbf{r}(t)$

- 1 $\mathcal{P}_u \leftarrow$ Uniformly sampled points
- 2 $\mathcal{P} \leftarrow \mathbf{G}_o(\mathcal{P}_u)$ to get sample points around the surface
- 3 **for** $\mathcal{P}(x, y, z, \tau)$ in \mathcal{P} **do**
- 4 $\mathbf{v}(x, y, z, \tau) \leftarrow \mathbf{1}$ to initialize
- 5 **for** G, ij in
- 6 $\text{zip}(\{\mathbf{V}_{3d}^s, \mathbf{V}_{3d}^d\}^4, \{(xy, xz, xt, yz, y\tau, z\tau)\})$ **do**
- 7 $\mathbf{v}(x, y, z, \tau) \leftarrow \mathbf{v}(x, y, z, \tau) \odot \mathcal{B}(G, i, j)$,
- 8 \mathcal{B} is the bi-linear interpolation on G
- 9 **end**
- 10 $\text{features} \leftarrow \text{concat}(\mathbf{b}(x, y, z, \tau), \mathbf{v}(x, y, z, \tau))$
- 11 $(\mathbf{c}_i, \sigma_i) \leftarrow \theta(\text{features})$
- 12 **end**
- 13 $\hat{\mathbf{C}} \leftarrow \int_{t_n}^{t_f} w(t) \cdot \mathbf{c}(\mathbf{r}(t), \mathbf{d}) dt$, $\hat{\mathbf{D}} \leftarrow \int_{t_n}^{t_f} w(t) \cdot t dt$,
- 14 $w(t) = \exp\left(-\int_{t_n}^t \sigma(\mathbf{r}(s) ds)\right) \sigma(\mathbf{r}(t))$
- 15 **return** $\hat{\mathbf{C}}, \hat{\mathbf{D}}$;

Rather than inputting positionally-encoded viewing information, we propose an enhancement on spatiotemporal information, as shown in Fig. 4. This enhancement involves concatenating interpolated features in Eq. 4 with positionally encoded coordinates and temporal embeddings. This new approach ensures a more precise representation of the scene and a more accurate reflection of the dynamic conditions inherent in an in vivo sequence. In this work, we use one-blob [30] encoding separately on each of the four coordinate values. The formulation of one-blob encoding is:

$$\mathbf{b}(s) = g_{s, \xi^2} \left(\frac{i - 0.5}{k} \right), i = 1, 2, 3 \dots k, \quad (6)$$

where g_{s, ξ^2} is a Gaussian kernel with mean value s and variance ξ^2 . The encoding function maps $s \in \mathbb{R}$ into a higher dimensional space \mathbb{R}^k . The encoding along with the fused features \mathbf{v} from feature planes is input to the MLP θ , which predicts σ and \mathbf{c} of each point. Then we utilize Eq. 2 to render the expected color $\hat{\mathbf{C}}$ and depth $\hat{\mathbf{D}}$ of one specific ray.

Algorithm 1 demonstrates our rendering procedure, starting with a ray $\mathbf{r}(t) = (\mathbf{x}_o + t\mathbf{d})$ selected through spatiotemporal importance sampling. We uniformly sample points along $\mathbf{r}(t)$, forming a point-set \mathcal{P}_u . This set is passed to the indicator

grid \mathbf{G}_o to obtain \mathcal{P} , which is distributed around the target tissue. We initialize each point in \mathcal{P} with 1 (the multiplicative identity). Subsequently, we apply bi-linear interpolation to extract features from different feature planes (e.g., G_{XY}) using corresponding dimension (e.g., xy) and update \mathbf{v} with element-wise multiplication accordingly. These features, combined with positional encoded coordinate and time embeddings $\mathbf{b}(x, y, z, \tau)$, are processed by the tiny MLP θ to yield color \mathbf{c}_i and density σ_i . Finally, volume rendering (Eq. 2) is applied to calculate the color $\hat{\mathbf{C}}$ and depth \hat{D} of $\mathbf{r}(t)$ by integrating radiance and density along the ray.

F. Optimization

Reconstructing surgical procedures rapidly and accurately presents significant challenges due to the inherent uncertainty and limited information. To expedite the optimization of the tiny MLP and orthogonal planes, we incorporate not only supervision from color and depth but also spatiotemporal continuity constraints. Furthermore, we devise a disentangle loss to aid in the disentanglement of static and dynamic fields. The following sections comprehensively demonstrate the utilization of the employed loss function in our approach.

1) Color Loss: We utilize captured images to optimize both the tiny MLP θ and the neural plane representation $\phi = \{\mathbf{V}_{3d}^s, \mathbf{V}_{3d}^d\}$ concurrently. We define the color loss in the following manner:

$$\mathcal{L}_{\text{rgb}}(\theta, \phi) = \sum_i \mathbb{E}_{\mathbf{r} \sim I_i} \left[\left\| \hat{\mathbf{C}}(\mathbf{r}) - \mathbf{C}_i^{\text{gt}}(\mathbf{r}) \right\|_2^2 \right], \quad (7)$$

where $\mathbf{C}_i^{\text{gt}}(\mathbf{r})$ represents the ground truth color of the ray \mathbf{r} passing through a pixel in image I_i .

2) Depth Loss: Following the methodology laid out in EndoNeRF [46], we also employ depth information generated by stereo matching to assist in the optimization of the neural scene. We define the depth loss as follows:

$$\mathcal{L}_{\text{depth}}(\theta, \phi) = \sum_i \mathbb{E}_{\mathbf{r} \sim D_i} \left[H_\delta \left(\hat{D}(\mathbf{r}) - D_i^{\text{gt}}(\mathbf{r}) \right) \right], \quad (8)$$

where the H_δ represent standard huber loss and $D_i^{\text{gt}}(\mathbf{r})$ represents the depth predicted by the stereo matching algorithm of the ray \mathbf{r} .

3) Total Variation Loss: Inspired by [11] and [31], we also leverage the total variation regularization to help optimizing our method. Specifically, the \mathcal{L}_{TV} is defined as:

$$\mathcal{L}_{TV} = \sum_{G \in \mathbf{V}_{3d}^s} \sum_{g \in G} \left\| \Delta_h^2(G, g) + \Delta_w^2(G, g) \right\|, \quad (9)$$

where $\mathbf{V}_{3d}^s = \{G_{XY}, G_{YZ}, G_{XZ}\}$ denotes all space planes in Forplane. $\Delta_h^2(G, g)$ represents the squared difference between the g th value in plane $G := (i, j)$ and the g th value in plane $G := (i+1, j)$ normalized by the resolution, and analogously for $\Delta_w^2(G, g)$.

4) Time Smoothness Loss: To robustly reconstruct deformable tissue under limited view, we further incorporate time smoothness regularization for all space-time planes. This time smoothness term, akin to total variation loss, aims to

ensure similarity between adjacent frames. The formulation is as follows:

$$\mathcal{L}_{TS} = \sum_{G \in \mathbf{V}_{3d}^d} \sum_{g \in G} \left\| \Delta_t^2(G, g) \right\|, \quad (10)$$

where $\mathbf{V}_{3d}^d = \{G_{XT}, G_{YT}, G_{ZT}\}$ is the set of the space-time planes in Forplane. $\Delta_t^2(G, g)$ represents the difference along the time axis.

5) Disentangle Loss: As discussed in Sec. I, surgical scenes pose a unique challenge due to the limited viewpoints they offer. This limitation necessitates sharing information across non-sequential timesteps to enhance the reconstruction quality. Our methodology employs a static-dynamic structure to model these surgical scenes. However, this structure creates an ambiguity for the model in discerning which areas should be modeled in the static field and which ones in the dynamic field. To tackle this issue, we commence by initializing the values in \mathbf{V}_{3d}^d to 1, which do not influence the features in \mathbf{V}_{3d}^s during element-wise multiplication. Subsequently, we optimize the model using a specialized disentangle loss function. The \mathcal{L}_{DE} is described as follows:

$$\mathcal{L}_{DE} = \sum_{G \in \mathbf{V}_{3d}^d} \sum_{g \in G} \|1 - g\|. \quad (11)$$

With this regularization, the features of the space-time planes will tend to remain the initial value. This strategy allows us to maximize the use of the static field in modeling static scenes while also appropriately accounting for dynamic elements.

Total Loss: The total loss optimized during each iteration is defined as:

$$\mathcal{L}_{\text{total}}(\theta, \phi) = \mathcal{L}_{\text{rgb}} + \lambda_d \mathcal{L}_D + \lambda_{tv} \mathcal{L}_{TV} + \lambda_{ts} \mathcal{L}_{TS} + \lambda_{de} \mathcal{L}_{DE}. \quad (12)$$

In all experiments conducted, we set the parameters $\lambda_d = 1, \lambda_{tv} = 0.001, \lambda_{ts} = 0.05, \lambda_{de} = 0.001$, unless otherwise specified in the experimental design.

G. Forplane with Monocular Inputs

In endoscopic procedures, monocular cameras are commonly used due to their compact size, despite some laparoscopes having stereo scope cameras. In the context of these scopes, the reliance on stereo depth in the original method [51] poses limitations. To address this, our enhanced version introduces an innovative optimization method tailored for monocular scopes. This method integrates data-driven monocular geometric cues into our training pipeline more effectively. In particular, we employ an off-the-shelf monocular depth predictor [2] to generate a depth map $\{\hat{D}_i\}_{i=1}^N$ for each input RGB image $\{I_i\}_{i=1}^N$. Note that estimating absolute scale in surgical scenes is challenging; hence, depth should be treated as a relative cue. Yet, this relative depth information is beneficial even over larger distances in the image.

In light of these considerations, we introduce a monocular depth loss $\mathcal{L}_{\text{Mono}}$ to enforce consistency between the esti-

mated depth $\hat{D}(\mathbf{r})$ and the predicted monocular depth $\bar{D}(\mathbf{r})$:

$$\mathcal{L}_{Mono}(\theta, \phi) = \sum_i \mathbb{E}_{\mathbf{r} \sim \bar{D}_i} \left\| (\eta \hat{D}(\mathbf{r}) + \epsilon) - \bar{D}(\mathbf{r}) \right\|^2,$$

$$\text{where } \eta = \left(\hat{D}(\mathbf{r})^T \hat{D}(\mathbf{r}) \right)^{-1} \hat{D}(\mathbf{r})^T \bar{D}(\mathbf{r}),$$

$$\epsilon = \bar{D}(\mathbf{r}) - \eta \hat{D}(\mathbf{r}). \quad (13)$$

The η and ϵ are the closed-form solution obtained via a least-squares criterion. These parameters are calculated separately for each batch, as the predicted depth maps can vary in scale and shift across different batches. For monocular endoscopic sequences, we replace \mathcal{L}_D with \mathcal{L}_{Mono} during training.

IV. EXPERIMENTS

A. Datasets

1) *EndoNeRF Dataset*: We assess the performance of Forplane using the EndoNeRF dataset [46]. This dataset is a collection of typical robotic surgery videos captured via stereo cameras from a single viewpoint during in-house DaVinci robotic prostatectomy procedures. Specifically designed to capture challenging surgical scenarios characterized by non-rigid deformation and tool occlusion, the dataset serves as an ideal platform for evaluation. The dataset is composed of six video clips (in total 807 frames) with resolution of 512×640 and lasting 4-8 seconds at a frame rate of 15 frames per second. The surgical challenges depicted in each case varies. Specifically, two cases (*thin structure*, *traction*) show traction on thin structures. Another two cases (*pushing tissues*, *pulling tissues*) depict significant manipulation of tissue through pushing or pulling. The final two cases (*cutting twice*, *tearing tissues*) capture the process of tissue cutting. These scenarios effectively illustrate the challenges inherent in managing soft tissue deformation and tool occlusion during surgery.

2) *Hamlyn Dataset*: We evaluate our method using the public Hamlyn dataset [28], [44], which includes both phantom heart and in-vivo sequences captured during da Vinci surgical robot procedures. The rectified images, stereo depth, and camera calibration information are from [35]. To generate instrument masks, we utilize the widely-used vision foundation model, Segment Anything [14], which enables semi-automatic segmentation of surgical instruments.

The Hamlyn dataset presents a rigorous evaluation scenario as it contains sequences that depict intracorporeal scenes with various challenges, such as weak textures, deformations, reflections, surgical tool occlusion, and illumination variations. We select seven specific sequences from the Hamlyn dataset (Sequence: rectified01, rectified06, rectified08, and rectified09), each comprising 301 frames with a resolution of 480×640 . These sequences span approximately 10 seconds and feature scenarios involving surgical tool occlusion, extensive tissue mobilization, large deformations, and even internal tissue exposure. To train and evaluate our method effectively, we divide the frames of each sequence into two sets: a training set of 151 frames and an evaluation set comprising the remaining frames. This division allows us to train our method adequately while providing a robust set for evaluation.



Fig. 5. Visualization of datasets. The left column is from EndoNeRF dataset and the right is from Hamlyn dataset. From left to right is original image, tool mask and the binocular depth with tool mask.

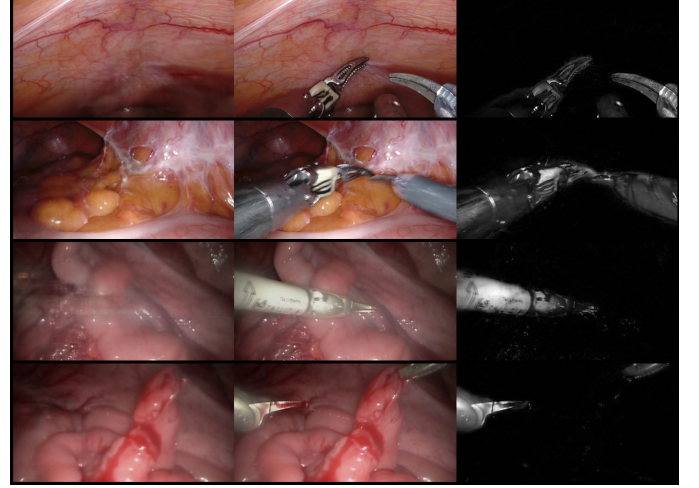


Fig. 6. Reconstruction results with difference maps. From left to right, we present the image rendered by Forplane, the corresponding ground truth image, and the difference map between them. These difference maps serve as a clear visual representation of our strong ability to reconstruct deformable tissues from tool-occluded endoscopy videos.

B. Baselines

We compare our approach with three state-of-the-art methods: EDSSR [22], EndoNeRF [46] and TiNeuVox [8]. Each of these methods represents a different type of deformable tissue reconstruction. EDSSR [22] leverages tissue deformation fields and volumetric fusion. EndoNeRF [46] utilizes dynamic neural radiance fields in MLPs to represent deformable surgical scenes and optimizes shapes and deformations in a learning-based manner. We follow the standard setting in [46] with 200k optimizing iterations. TiNeuVox [8] integrates optimizable time-aware voxel features for faster optimization. Note that [8] is not specified for surgical procedures, so we modify the original TiNeuVox model with tool mask-guided ray casting and stereo depth-cueing ray marching which are proposed by [46]. These methods are equipped with same hyper-parameters as [46] and enables the TiNeuVox to reconstruct the dynamic tissues. We report two versions of TiNeuVox: TiNeuVox-S and TiNeuVox-B, representing the small and base versions, respectively. We also report Lerplane's performance [51], the previous version of Forplane. For fair comparison, we train Lerplane, TiNeuVox and Forplane with 32k iterations. We further report Forplane with only 9k iterations, showing significant optimization speed.

C. Qualitative and Quantitative Results

We show the quantitative results on EndoNeRF dataset in Fig. 7 and all quantitative results on two datasets are summarized in Table I and Table II. All experiments were conducted on an Ubuntu 20.04 system equipped with one NVIDIA GeForce RTX 3090 GPU. Four metrics are used

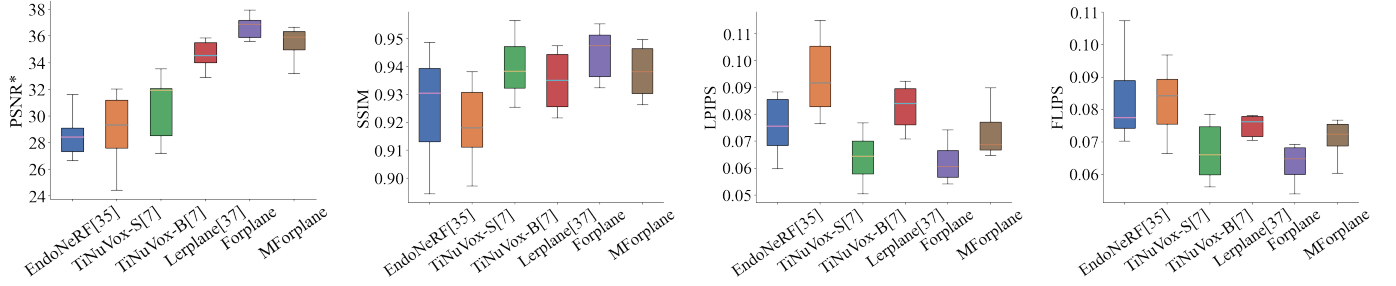


Fig. 7. The Quantitative Comparison of Similarity Metrics for Reconstructed Volumes. These figures present a quantitative comparison of similarity metrics between the input slices and the corresponding slices extracted from the reconstructed volumes. Four different similarity metrics, namely PSNR*, SSIM, LPIPS, FLIPS, are evaluated. Each figure depicts the performance of six different methods.

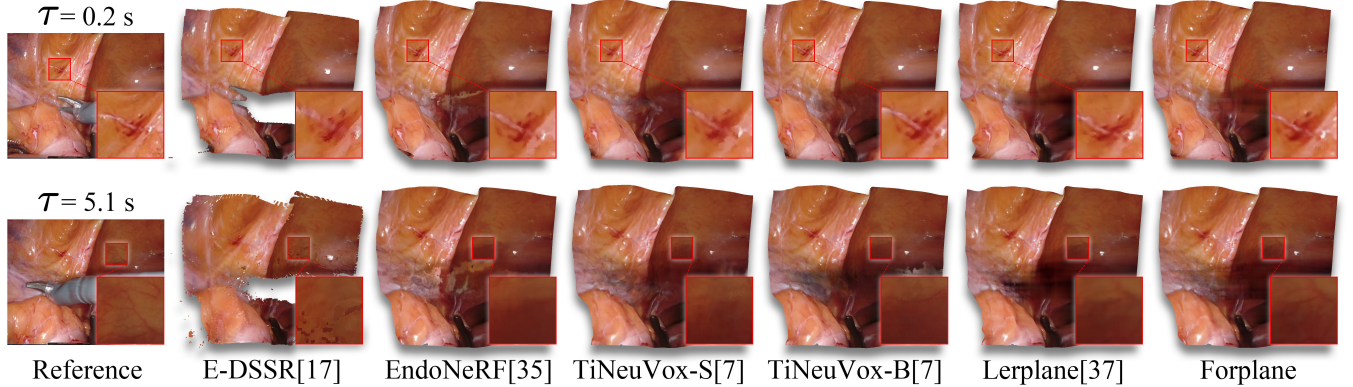


Fig. 8. Results on the scene *cutting twice* from EndoNeRF dataset [46]. We present a comparison of tissue reconstructions from various methods. Forplane consistently shows more detailed reconstructions of the surgical procedure (e.g. capillaries and peritoneum) compared to other methods.

for evaluation, *i.e.* PSNR, SSIM [47], LPIPS [53], FLIP [1]. These metrics are applied to the whole image including the masked empty areas. These areas may diminish the discernible differences between predicted and ground truth images. To address this, we introduce a modified metric called masked PSNR (PSNR*), which focuses solely on measuring the PSNR of tissue pixels, allowing for a more accurate comparison. Apparently, Forplane outperforms other state-of-the-art methods across all evaluation metrics, demonstrating significantly superior performance. E-DSSR [22] achieves the fastest rendering speed due to its neural computing-light rendering procedure; however, its reconstruction quality is relatively low compared to other methods. EndoNeRF [46] shows promising reconstruction quality but requires approximately 14 hours for per-scene training. In contrast, our Forplane completes one scene learning in just 10 minutes while achieving better performance across all evaluation metrics. TiNeuVox-S [8] has a similar training time as Forplane-32k but exhibits significantly lower performance compared to our approach. TiNeuVox-B [8], with a larger parameter quantity than TiNeuVox-S, exhibits stronger fitting ability. While better reconstruction quality is achieved on datasets with better quality (EndoNeRF dataset), the improvement diminishes on more challenging datasets (Hamlyn dataset). Furthermore, TiNeuVox-B accommodates a greater number of parameters, leading to increased computational time for both optimization and inference relative to TiNeuVox-S. In comparison, Forplane achieves high-quality reconstruction on

both datasets with significantly faster rendering speed even with only 3 minutes optimization. Our previous method [51] exhibits performance akin to rapid optimization speeds, yet it lags in terms of rendering speed. Conversely, Forplane not only delivers superior overall quality but also boasts a rendering process that is roughly five times swifter than that of the previous version [51]. This highlights the remarkable efficiency of our ray marching approach. Furthermore, qualitative comparisons are provided in Fig. 8 and the fine detailed reconstruction results are shown in Fig. 2 and Fig. 6. The shared static field in Forplane enables better utilization of information throughout the surgical procedure, resulting in finer and more accurate details compared to other methods. Notably, Forplane achieves these improvements with less training time and faster rendering speed.

Experimental results indicate that Forplane excels in computational speed and reconstruction quality, positioning it as a promising tool for future clinical and intraoperative applications in surgery. Through its pioneering method that bolsters optimization without sacrificing the quality of reconstruction, Forplane stands poised to transform the landscape of surgical procedures.

D. Ablation Study

In this section, we perform a series of experiments to validate the effectiveness of our proposed methods. The experiments are conducted on the EndoNeRF dataset [46] with

TABLE I

MEAN VALUES OF QUANTITATIVE METRICS FOR DIFFERENT MODELS ON THE TWO DATASETS (STANDARD DEVIATION IN PARENTHESES).
 \uparrow INDICATES THAT HIGHER VALUES INDICATE HIGHER ACCURACY, AND VICE VERSA.

Methods	PSNR \uparrow	PSNR* \uparrow	SSIM \uparrow	LPIPS \downarrow	FLIP \downarrow
EndoNeRF Dataset [46]					
E-DSSR [22]	13.398 (1.270)	12.997 (1.232)	0.630 (0.057)	0.423 (0.047)	0.426 (0.056)
EndoNeRF [46]	29.477 (1.886)	28.560 (1.782)	0.926 (0.021)	0.080 (0.019)	0.083 (0.014)
TiNeuVox-S [8]	29.913 (3.069)	28.980 (2.862)	0.919 (0.015)	0.094 (0.015)	0.083 (0.011)
TiNeuVox-B [8]	31.601 (2.855)	30.668 (2.697)	0.940 (0.012)	0.064 (0.010)	0.067 (0.009)
Lerplane [51]	35.504 (1.161)	34.575 (1.129)	0.935 (0.011)	0.083 (0.009)	0.075 (0.010)
Forplane-9k	33.374 (1.074)	32.435 (1.165)	0.907 (0.014)	0.127 (0.019)	0.093 (0.007)
Forplane-32k	37.306 (1.406)	36.367 (1.521)	0.945 (0.010)	0.062 (0.008)	0.063 (0.006)
Hamlyn Dataset [28], [44]					
E-DSSR [22]	18.150 (2.571)	17.337 (2.402)	0.640 (0.060)	0.393 (0.066)	0.259 (0.065)
EndoNeRF [46]	34.879 (1.784)	34.066 (1.717)	0.951 (0.011)	0.071 (0.017)	0.070 (0.012)
TiNeuVox-S [8]	35.277 (1.682)	34.464 (1.501)	0.953 (0.014)	0.085 (0.029)	0.067 (0.016)
TiNeuVox-B [8]	33.764 (2.047)	32.951 (1.974)	0.942 (0.020)	0.146 (0.061)	0.078 (0.022)
Lerplane [51]	32.455 (2.247)	31.629 (2.064)	0.935 (0.021)	0.124 (0.041)	0.098 (0.028)
Forplane-9k	35.301 (2.241)	34.475 (2.076)	0.945 (0.018)	0.093 (0.035)	0.074 (0.020)
Forplane-32k	37.474 (2.401)	36.647 (2.232)	0.960 (0.014)	0.058 (0.025)	0.059 (0.017)

TABLE II

TRAINING TIME AND TEST SPEED FOR DIFFERENT METHODS

Methods	Train Time	Test Speed
E-DSSR [22]	13 mins	28.0 fps
EndoNeRF [46]	>10 hours	0.11 fps
TiNeuVox-S [8]	12 mins	0.56 fps
TiNeuVox-B [8]	90 mins	0.18 fps
Lerplane [51]	10 mins	0.38 fps
Forplane-9k	3 mins	<u>1.73</u> fps
Forplane-32k	10 mins	<u>1.73</u> fps

9k iterations, and the reported metric values are averaged.

Sampling Strategy We compare our proposed spatiotemporal importance sampling strategy with two other sampling strategies: Naive Sampling, which avoids tool masks and assigns equal weights to all pixels, and EndoNeRF Sampling, which assigns higher probabilities to highly occluded areas as described in [46].

Encoding Method We compare our coordinate-time encoding method with two methods: Dummy Encoding, which replaces all the coordinate-time encoded parameters with a constant value to preserve the parameter quantity of the tiny MLP, and Direction Encoding, where one-blob positional encoding is applied to the view direction and the encoded parameters are passed to the tiny MLP, as described in [26].

Space-time Disentangle As detailed in Sec. III-F.5, we use a specialized initialization for \mathbf{V}_{3d}^d and pair it with \mathcal{L}_{DE} to enhance the decomposition between \mathbf{V}_{3d}^d and \mathbf{V}_{3d}^s . For comparison, we also trained a version of Forplane without these adjustments, termed w/o Disentangle.

Ray Marching In this improved version, we introduce an efficient ray marching algorithm designed to accurately generate spatial-temporal points. To evaluate the effectiveness of this algorithm, we conducted a comparative analysis against two existing methods: the sample-net method (termed Sample-Net) and the stereo depth-cueing ray marching method (termed Depth-Cueing). Sample-Net, as proposed in the original version of our work [51], utilizes a simplified representation of the 4D scene, aiming to enhance the accuracy of point sampling. In contrast, the Depth-Cueing method, introduced in EndoNeRF [46], employs Gaussian transfer functions informed by stereo depth data to guide the sampling of points near tissue

TABLE III

MEAN VALUES OF EVALUATION METRICS FOR DIFFERENT ABLATION MODELS ON THE ENDONeRF [46] DATASET.

Methods	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	FLIP \downarrow
Naive Sampling	33.146	0.906	0.135	0.094
EndoNeRF Sampling	33.019	0.906	0.135	0.093
Dummy Encoding	33.055	0.905	0.129	0.094
Direction Encoding	33.062	0.904	0.133	0.095
W/o Disentangle	18.015	0.869	0.302	0.349
Full Model	33.374	0.907	0.127	0.093

TABLE IV

EVALUATION METRICS FOR DIFFERENT RAY MARCHING STRATEGY ON THE ENDONeRF [46] DATASET.

Methods	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	Test Speed \uparrow
Depth-Cueing	31.325	0.897	0.132	1.95 fps
Sample-Net	32.889	0.901	0.126	0.38 fps
Ours	33.374	0.907	0.127	1.73 fps

surfaces, aiming to optimize the rendering of fine details.

The results of the ablation study are summarized in Table III and Table IV. Our spatiotemporal importance sampling achieves superior quality compared to the other sampling strategies, confirming its effectiveness. Regarding encoding methods, our coordinate-time encoding outperforms the others, while the performance of Dummy and Direction Encoding is similar, validating of our hypothesis in Section III-E. Forplane w/o Disentangle shows notable degradation, primarily from optimization ambiguities. Lacking clear guidance on decomposition, the network struggles to distinguish between static and dynamic elements given the restricted viewpoints. The results presented in Table IV indicate that the depth-cueing method yields suboptimal performance, primarily due to the reliance on noisy and incomplete stereo depth data, as depicted in Fig. 5. Additionally, its inapplicability for temporal interpolation further limits its utility due to the absence of stereo depth in certain contexts. Conversely, while the sample-net method is versatile, it leads to increased computational demands during inference, negatively impacting test speed. Our efficient ray marching method effectively mitigates these issues, striking a balance between rendering speed and reconstruction quality. It

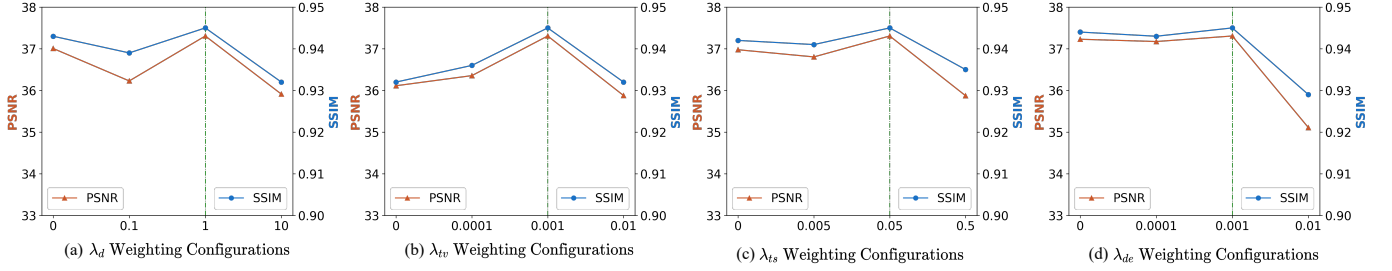


Fig. 9. Comparative reconstruction results with varied weighting parameters. This figure illustrates the performance impact of different weighting configurations on the reconstruction process, with the horizontal axis represented in a logarithmic scale. Fig. 9 (a) shows the performance of Forplane with different λ_d . Fig. 9 (b) shows λ_{tv} , Figure (c) shows λ_{ts} and Figure (d) shows λ_{de} .

TABLE V

THE PERFORMANCE BETWEEN FORPLANE AND MFORPLANE. THE TIME ROW REPRESENTS FOR OPTIMIZATION TIME.

Methods	PSNR↑	PSNR*↑	SSIM↑	LPIPS↓	FLIP↓
EndoNeRF Dataset [46]					
Forplane	37.306	36.367	0.945	0.062	0.063
MForplane	36.403	35.464	0.938	0.073	0.071
Hamlyn Dataset [28], [44]					
Forplane	37.474	36.647	0.960	0.058	0.059
MForplane	37.034	36.208	0.959	0.060	0.063

outperforms the depth-cueing method in terms of PSNR and SSIM, and also surpasses the sample-net method in inference speed, thus demonstrating its potential for real-time application scenarios.

E. Hyperparameters

As delineated in Section III-F, we set the λ parameters to a specific weighting configuration across all experiments. This subsection is dedicated to conducting a comprehensive analysis to discern the effects of various hyperparameters on Forplane’s performance. To this end, we train Forplane using an array of λ weightings on the EndoNeRF dataset [46], extending across 32k iterations. The results, depicted in Fig. 9, suggest that Forplane exhibits considerable robustness to variations in loss weighting configurations, only excessively large weightings could potentially impair performance.

F. Monocular Reconstruction

We show the performance of MForplane and Forplane on two datasets in Table. V. By incorporating depth priors from monocular depth estimation networks, MForplane achieves fast reconstruction of dynamic tissues in monocular endoscopic videos. However, when compared to Forplane, we observe a slight decrease in the reconstruction quality of MForplane (e.g. 1.81% in PSNR, 0.4% in SSIM). Furthermore, there is a slight increase in the training time (from 10 minutes to 12 minutes), primarily attributed to the additional computational workload involved in calculating the monocular depth loss. Nevertheless, it is important to note that the ability to achieve fast reconstruction of dynamic tissues using one monocular endoscopic video holds significant promise for various applications in medical imaging. Therefore, the minor reduction in quality is deemed acceptable given the enormous potential and practicality of MForplane.

TABLE VI

MEAN VALUES OF EVALUATION METRICS FOR DIFFERENT MLP STRUCTURE ON THE ENDONERF [46] DATASET.

Methods	PSNR↑	SSIM↑	Train Time↓	Test Speed↑
Large MLP	34.285	0.913	23 mins	0.75 fps
Tiny MLP	33.374	0.907	3 mins	1.73 fps

G. The Choice of MLP Structure

As mentioned in Section III-B, our fast orthogonal plane representation reduces the computational load on the MLP, allowing for a more efficient MLP architecture in complex surgical reconstruction tasks. To validate this approach, we conducted an experimental comparison where our tiny MLP (2 fully-connected ReLU layers with 64 channels) was replaced with the more elaborate MLP structure (8 fully-connected ReLU layers with 256 channels) utilized by EndoNeRF [46]. The experiments are conducted on the EndoNeRF dataset with 9k iterations, and the reported metric values are averaged. The experimental results, detailed in Table VI, indicate that the use of a larger MLP configuration does not significantly enhance reconstruction quality but results in a substantial increase in training duration for the same iteration count. These outcomes strongly corroborate our proposition that a streamlined MLP architecture contributes to enhanced optimization and inference efficiency, underscoring the effectiveness of our proposed method.

H. Disentangle Static and Dynamic Field

We propose an effective method for visualizing the static and dynamic fields in surgical scenes, showcasing the efficacy of our factorization approach. Our approach utilizes element-wise multiplication for feature aggregation across different feature planes, as outlined in Eq. 4. By setting all values in the dynamic field to 1 and performing inference, we obtain the network’s output as the rendered result of the static field. Similarly, setting all values in the static field to 1 provides the rendered result of the dynamic field. The time-space decomposition is clearly illustrated in Fig. 10, offering a compelling demonstration of the effectiveness of our proposed structure. This property offers valuable insights into tissue examination and procedural understanding during surgery, with potential applications in virtual surgery and intraoperative utilization.

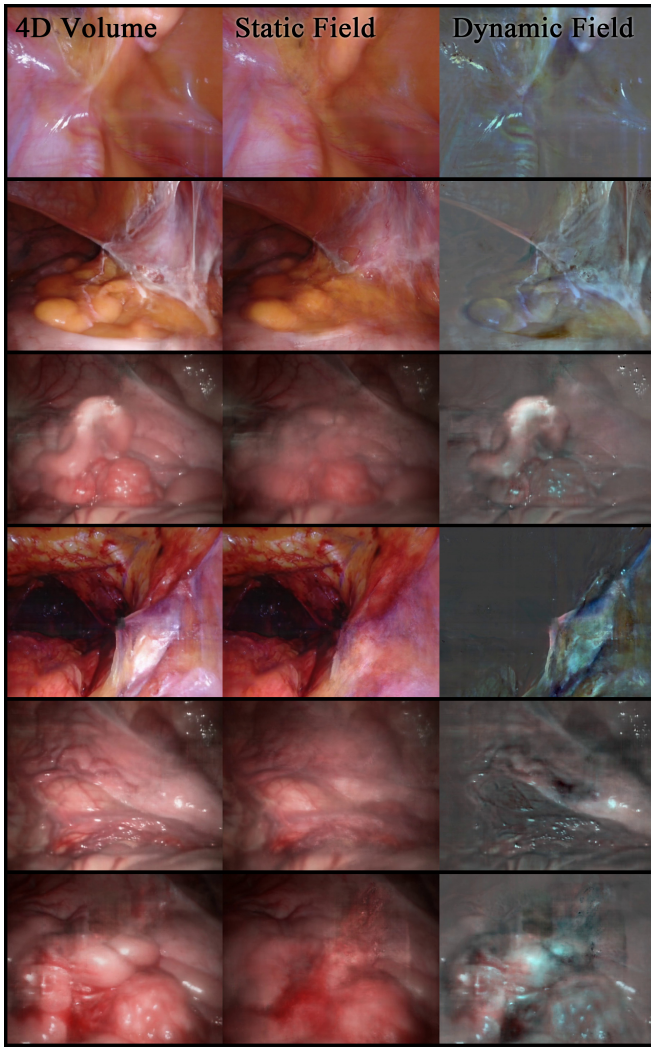


Fig. 10. Visualization of static and dynamic fields. The visualizations provided demonstrate the decomposition of static and dynamic fields in surgical procedures, showcasing the effectiveness of utilizing scene context within Forplane. The first image of each row shows the result of full 4D volume. The second image corresponds to the static field, while the third image portrays the dynamic field, which is re-normalized to $[0, 1]$ for optimal visualization.

V. DISCUSSION

Forplane exhibits SOTA performance in rapidly reconstructing high-quality deformable tissues within both monocular and binocular endoscopic videos. This significant advancement is attributed to a series of critical observations and methodological innovations: 1) Combination of Implicit and Explicit Representations: While implicit representations excel in detailed scene restoration, the computational demands are substantial. Conversely, explicit representations, though computationally efficient, often lack detail. Forplane integrates the strengths of both implicit and explicit representations, thereby achieving efficient training and inference, alongside high-quality reconstruction. 2) Selective Focus on Significant Regions: Acknowledging the varying significance of different tissues, Forplane strategically focuses on critical regions. This targeted approach accelerates the optimization process by allocating computational resources to areas of the scene with

the most impact on overall reconstruction quality. 3) Utilization of Temporal Information: Given the limited viewpoint in endoscopic environments, Forplane effectively leverages temporal data. This utilization of temporal information is key in enhancing the accuracy of the reconstruction, particularly in dynamically changing surgical scenes.

Despite its notable achievements, Forplane has areas that warrant further development. First, the MForplane variant, designed to operate without stereo depth, still depends on manually generated masks for surgical instrument identification. This process is both challenging and labor-intensive, highlighting the need for advancements in automated instrument labeling and tracking. Secondly, the rendering speed of Forplane, although significantly faster than its predecessor (approximately 5 times quicker at 1.73 fps on a single RTX3090 GPU), does not yet meet the demands for real-time intraoperative assistance. This limitation underscores the necessity for further improvements in rendering speed. Lastly, the functionality of Forplane is presently focused on the reconstruction of deformable tissues. To fully support comprehensive surgical training for robotic systems, a more detailed analysis of surgical procedures is required. This includes not only deformable tissue reconstruction but also the intricate dynamics of surgical instruments, environmental lighting conditions, and camera motion.

VI. CONCLUSION

In this paper, we propose a fast dynamic reconstruction framework, Forplane, targeted on deformable tissues during surgery. Forplane represents surgical procedures with orthogonal neural planes and incorporates advancements in pixel sampling, spatial sampling, and information enhancement to improve rendering quality, accelerate optimization and inference. Evaluations on two in vivo datasets demonstrate the superior performance of Forplane, requiring less training time and offering faster inference speed compared to other methods. In addition, Forplane effectively handles both monocular and binocular endoscopy videos, maintaining nearly identical high-quality reconstruction performance. We firmly believe that Forplane holds substantial promise for intraoperative applications among surgery.

REFERENCES

- [1] P. Andersson, J. Nilsson, T. Akenine-Möller, M. Oskarsson, K. Åström, and M. D. Fairchild. Flip: A difference evaluator for alternating images. *PACMCGIT*, 2020.
- [2] S. F. Bhat, R. Birkel, D. Wofk, P. Wonka, and M. Müller. Zoedepth: Zero-shot transfer by combining relative and metric depth. *arXiv preprint arXiv:2302.12288*, 2023.
- [3] Lorenzo Bianchi, Umberto Barbaresi, Laura Cencenelli, Barbara Bortolani, Caterina Gaudiano, Francesco Chessa, Andrea Angiolini, Simone Lodi, Angelo Porreca, Federico Mineo Bianchi, Carlo Casablanca, Amelio Ercolino, Alessandro Bertaccini, Rita Golfieri, Emanuela Marcelli, and Riccardo Schiavina. The impact of 3d digital reconstruction on the surgical planning of partial nephrectomy: A case-control study. still time for a novel surgical trend? *Clinical Genitourinary Cancer*, 18(6):e669–e678, 2020.
- [4] A. Cao and J. Johnson. Hexplane: A fast representation for dynamic scenes. In *CVPR*, 2023.
- [5] A. Chen, Z. Xu, A. Geiger, J. Yu, and H. Su. Tensorf: Tensorial radiance fields. In *ECCV*, 2022.

- [6] A. Corona-Figueroa, J. Frawley, S. Bond-Taylor, S. Bethapudi, H. PH Shum, and C. G. Willcocks. Mednerf: Medical neural radiance fields for reconstructing 3d-aware ct-projections from a single x-ray. In *EMBC*, 2022.
- [7] Y. Du, Y. Zhang, H. Yu, J. B. Tenenbaum, and J. Wu. Neural radiance flow for 4d view synthesis and video processing. In *ICCV*, 2021.
- [8] J. Fang, T. Yi, X. Wang, L. Xie, X. Zhang, W. Liu, M. Nießner, and Q. Tian. Fast dynamic radiance fields with time-aware neural voxels. In *SIGGRAPH Asia*, 2022.
- [9] Y. Fang, Z. Cui, L. Ma, L. Mei, B. Zhang, Y. Zhao, Z. Jiang, Y. Zhan, Y. Pan, M. Zhu, et al. Curvature-enhanced implicit function network for high-quality tooth model generation from cbct images. In *MICCAI*, 2022.
- [10] S. Fridovich-Keil, G. Meanti, F. R. Warburg, B. Recht, and A. Kanazawa. K-planes: Explicit radiance fields in space, time, and appearance. In *CVPR*, 2023.
- [11] S. Fridovich-Keil, A. Yu, M. Tancik, Q. Chen, B. Recht, and A. Kanazawa. Plenoxels: Radiance fields without neural networks. In *CVPR*, 2022.
- [12] Wei Gao and Russ Tedrake. Surfelpwarp: Efficient non-volumetric single view dynamic reconstruction, 2019.
- [13] M. O. Khan and Y. Fang. Implicit neural representations for medical imaging segmentation. In *MICCAI*, 2022.
- [14] A. Kirillov, E. Mintun, N. Ravi, H. Mao, C. Rolland, L. Gustafson, T. Xiao, S. Whitehead, A. C. Berg, W. Lo, et al. Segment anything. *arXiv preprint arXiv:2304.02643*, 2023.
- [15] Tim Lange, Daniel J. Indelicato, and Joseph M. Rosen. Virtual reality in surgical training. *Surgical Oncology Clinics of North America*, 9(1):61–79, 2000. Surgical Techniques and Outcomes.
- [16] H. Li, H. Chen, W. Jing, Y. Li, and R. Zheng. 3d ultrasound spine imaging with application of neural radiance field method. In *IJUS*, 2021.
- [17] R. Li, M. Tancik, and A. Kanazawa. Nerfacc: A general nerf acceleration toolbox. *arXiv preprint arXiv:2210.04847*, 2022.
- [18] Y. Li, F. Richter, J. Lu, E. K. Funk, R. K. Orosco, J. Zhu, and M. C. Yip. Super: A surgical perception framework for endoscopic tissue manipulation with surgical robotics. *RA-L*, 2020.
- [19] Z. Li, S. Niklaus, N. Snavely, and O. Wang. Neural scene flow fields for space-time view synthesis of dynamic scenes. In *CVPR*, 2021.
- [20] Yonghao Long, Jianfeng Cao, Anton Deguet, Russell H. Taylor, and Qi Dou. Integrating artificial intelligence and augmented reality in robotic surgery: An initial dvrc study using a surgical education scenario. In *2022 International Symposium on Medical Robotics (ISMR)*, pages 1–8, 2022.
- [21] Yonghao Long, Chengkun Li, and Qi Dou. Robotic surgery remote mentoring via ar with 3d scene streaming and hand interaction, 2022.
- [22] Y. Long, Z. Li, C. H. Yee, C. F. Ng, R. H. Taylor, M. Unberath, and Q. Dou. E-dssr: efficient dynamic surgical scene reconstruction with transformer-based stereoscopic depth perception. In *MICCAI*, 2021.
- [23] G. Luegmair, D. D. Mehta, J. B. Kobler, and M. Döllinger. Three-dimensional optical reconstruction of vocal fold kinematics using high-speed video with a laser projection system. *TMI*, 2015.
- [24] N. Mahmoud, T. Collins, A. Hostettler, L. Soler, C. Doignon, and J. M. Martinez Montiel. Live tracking and dense reconstruction for handheld monocular endoscopy. *TMI*, 2018.
- [25] L. Maier-Hein, A. Groch, A. Bartoli, S. Bodenstedt, G. Boissonnat, P.-L. Chang, N. T. Clancy, D. S. Elson, S. Haase, E. Heim, et al. Comparative validation of single-shot optical techniques for laparoscopic 3-d surface reconstruction. *TMI*, 2014.
- [26] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM*, 2021.
- [27] Amirali Molaei, Amirhossein Aminimehr, Armin Tavakoli, Amirhossein Kazerouni, Bobby Azad, Reza Azad, and Dorit Merhof. Implicit neural representation in medical imaging: A comparative survey. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2381–2391, 2023.
- [28] P. Mountney, D. Stoyanov, and G. Yang. Three-dimensional tissue deformation recovery and tracking. *IEEE Signal Processing Magazine*, 2010.
- [29] T. Müller, A. Evans, C. Schied, and A. Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ToG*, 2022.
- [30] T. Müller, B. McWilliams, F. Rousselle, M. Gross, and J. Novák. Neural importance sampling. *ToG*, 2019.
- [31] M. Niemeyer, J. T. Barron, B. Mildenhall, M. S. Sajjadi, A. Geiger, and N. Radwan. Regnerf: Regularizing neural radiance fields for view synthesis from sparse inputs. In *CVPR*, 2022.
- [32] K. Park, U. Sinha, J. T. Barron, S. Bouaziz, D. B. Goldman, S. M. Seitz, and R. Martin-Brualla. Nerfies: Deformable neural radiance fields. In *CVPR*, 2021.
- [33] A. Pumarola, E. Corona, G. Pons-Moll, and F. Moreno-Noguer. D-nerf: Neural radiance fields for dynamic scenes. In *CVPR*, 2021.
- [34] A. Raju, S. Miao, D. Jin, L. Lu, J. Huang, and A. P. Harrison. Deep implicit statistical shape models for 3d medical image delineation. In *AAAI*, 2022.
- [35] D. Recasens, J. Lamarca, J. M. Fácil, J. M. M. Montiel, and J. Civera. Endo-depth-and-motion: reconstruction and tracking in endoscopic videos using depth networks and photometric constraints. *RA-L*, 2021.
- [36] Albert W. Reed, Hyojin Kim, Rushil Anirudh, K. Aditya Mohan, Kyle Champley, Jingu Kang, and Suren Jayasuriya. Dynamic ct reconstruction from limited views with implicit neural representations and parametric motion fields, 2021.
- [37] D. Rückert, Y. Wang, R. Li, R. Idoughi, and W. Heidrich. Neat: Neural adaptive tomography. *TOG*, 2022.
- [38] Riccardo Schiavina, Lorenzo Bianchi, Marco Borghesi, Francesco Chessa, Laura Cerenelli, Emanuela Marcelli, and Eugenio Brunocilla. Three-dimensional digital reconstruction of renal model to guide pre-operative planning of robot-assisted partial nephrectomy. *International Journal of Urology*, 26(9):931–932, 2019.
- [39] Adam Schmidt, Omid Mohareri, Simon DiMaio, and Septimiu E Salcudean. Fast graph refinement and implicit neural representation for tissue tracking. In *2022 International Conference on Robotics and Automation (ICRA)*, pages 1281–1288. IEEE, 2022.
- [40] Adam Schmidt, Omid Mohareri, Simon DiMaio, and Septimiu E Salcudean. Recurrent implicit neural graph for deformable tracking in endoscopic videos. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 478–488. Springer, 2022.
- [41] Adam Schmidt, Omid Mohareri, Simon DiMaio, Michael Yip, and Septimiu E. Salcudean. Tracking and mapping in medical computer vision: A review, 2023.
- [42] M. Semmler, S. Kniesburges, V. Birk, A. Ziethe, R. Patel, and M. Döllinger. 3d reconstruction of human laryngeal dynamics based on endoscopic high-speed recordings. *TMI*, 2016.
- [43] J. Song, J. Wang, L. Zhao, S. Huang, and G. Dissanayake. Dynamic reconstruction of deformable soft-tissue with stereo scope in minimal invasive surgery. *RA-L*, 2017.
- [44] D. Stoyanov, G. P. Mylonas, F. Deligianni, A. Darzi, and G. Yang. Soft-tissue motion tracking and structure estimation for robotic assisted mis procedures. In *MICCAI*, 2005.
- [45] Y. Sun, J. Liu, M. Xie, B. Wohlberg, and U. S. Kamilov. Coil: Coordinate-based internal learning for tomographic imaging. *IEEE Transactions on Computational Imaging*, 2021.
- [46] Y. Wang, Y. Long, S. H. Fan, and Q. Dou. Neural rendering for stereo 3d reconstruction of deformable tissues in robotic surgery. In *MICCAI*, 2022.
- [47] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli. Image quality assessment: From error visibility to structural similarity. *TIP*, 2004.
- [48] Yiheng Xie, Towaki Takikawa, Shunsuke Saito, Or Litany, Shiqin Yan, Numair Khan, Federico Tombari, James Tompkin, Vincent Sitzmann, and Srinath Sridhar. Neural fields in visual computing and beyond. In *Computer Graphics Forum*. Wiley Online Library, 2022.
- [49] Junshen Xu, Daniel Moyer, Borjan Gagoski, Juan Eugenio Iglesias, P. Ellen Grant, Polina Golland, and Elfar Adalsteinsson. Nvsor: Implicit neural representation for slice-to-volume reconstruction in mri. *IEEE Transactions on Medical Imaging*, 42(6):1707–1719, 2023.
- [50] Chen Yang, Peihao Li, Zanwei Zhou, Shanxin Yuan, Bingbing Liu, Xiaokang Yang, Weichao Qiu, and Wei Shen. Nerfvs: Neural radiance fields for free view synthesis via geometry scaffolds. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16549–16558, 2023.
- [51] C. Yang, K. Wang, Y. Wang, X. Yang, and W. Shen. Neural lerplane representations for fast 4d reconstruction of deformable tissues. *MICCAI*, 2023.
- [52] J. Yang, U. Wickramasinghe, B. Ni, and P. Fua. Implicitatlas: learning deformable shape templates in medical imaging. In *CVPR*, 2022.
- [53] R. Zhang, P. Isola, A. A. Efros, E. Shechtman, and O. Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*, 2018.
- [54] H. Zhou and J. Jagadeesan. Real-time dense reconstruction of tissue surface from stereo optical video. *TMI*, 2019.
- [55] H. Zhou and J. Jagadeesan. Emdq-slam: Real-time high-resolution reconstruction of soft tissue surface from stereo laparoscopy videos. In *MICCAI*, 2021.
- [56] Zhenghong Zhou, Huangxuan Zhao, Jiemin Fang, Dongqiao Xiang, Lei Chen, Lingxia Wu, Feihong Wu, Wenyu Liu, Chuansheng Zheng, and Xinggao Wang. Tiavox: Time-aware attenuation voxels for sparse-view 4d dsa reconstruction, 2023.