# FLex: Joint Pose and Dynamic Radiance Fields Optimization for Stereo Endoscopic Videos

Florian Philipp Stilz[*,1,2], Mert Asim Karaoglu[*,1,2], Felix Tristram[*,1], Nassir Navab[1], Benjamin Busam[1], and Alexander Ladikos[2]

[1] Technical University Munich
[2] ImFusion GmbH

**Abstract.** Reconstruction of endoscopic scenes is an important asset for various medical applications, from post-surgery analysis to educational training. Neural rendering has recently shown promising results in endoscopic reconstruction with deforming tissue. However, the setup has been restricted to a static endoscope, limited deformation, or required an external tracking device to retrieve camera pose information of the endoscopic camera. With FLex we adress the challenging setup of a moving endoscope within a highly dynamic environment of deforming tissue. We propose an implicit scene separation into multiple overlapping 4D neural radiance fields (NeRFs) and a progressive optimization scheme jointly optimizing for reconstruction and camera poses from scratch. This improves the ease-of-use and allows to scale reconstruction capabilities in time to process surgical videos of 5,000 frames and more; an improvement of more than ten times compared to the state of the art while being agnostic to external tracking information. Extensive evaluations on the StereoMIS dataset show that FLex significantly improves the quality of novel view synthesis while maintaining competitive pose accuracy.

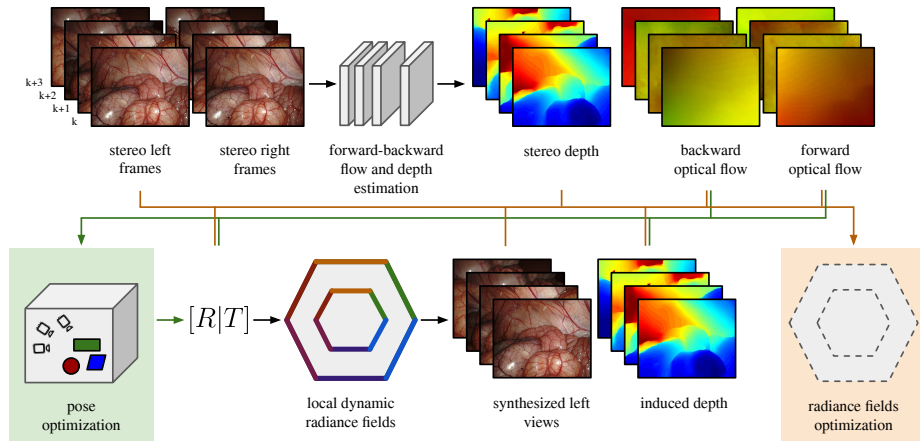**Keywords:** 3D Reconstruction · Neural Rendering · Robotic Surgery

## 1 Introduction

Visually and geometrically accurate reconstructions of surgical scenes are crucial components for various computer vision and AR/VR applications such as post surgical longitudinal assessment [8], surgical training [6], as well as for data generation for other learning-based computer vision and robotics applications [9]. However, endoscopic videos present a range of visual and practical difficulties for contemporary reconstruction methods, including strong non-homeomorphic deformations, prolonged recording times and the challenge of determining camera positions. Often, this leads to a dependency on external tools for acquisition, diminishing the ease-of-use of the reconstruction framework.

---

* The authors contributed equally.

Contact author: Florian Philipp Stilz (*florian.stilz@tum.de*).

**Fig. 1.** Joint pose and radiance fields optimization. $k$ indexes the frames along the temporal dimension. Orange and green arrows show the flow of inputs and outputs in the optimization processes of pose and radiance fields.

Prior methods widely explore usage of explicit representations such as sparse and dense point clouds [19,16] in visual odometry and simultaneous localization and mapping (SLAM) frameworks. Even though these approaches often provide efficient solutions for combined camera tracking and reconstruction, their incomplete geometry modeling results in limitations when rendering views from new camera poses. Our method is similar regarding joint optimization of localization and reconstruction, however, we instead use a dynamic neural radiance fields (NeRF) [11] based architecture to reconstruct the scene together with capturing the dynamics, thereby enabling high quality time dependent novel view synthesis.

EndoNeRF [23] is the first in the line of works that adapt a dynamic NeRF archicture [13] for endoscopic scenes. EndoSurf [27] builds on top of the previous work and substitutes the representation from volumetric density to a signed-distance function (SDF) [22]. LerPlane [26] and its follow-up work ForPlane [25] utilize explicit data structures as in [2,3,1] for faster optimization and higher rendering quality. While these works present great results and a promising research direction, unlike our method, they rely on external, reliable measurement or computation of the camera poses which are difficult to obtain in endoscopic environments. While EndoNeRF [23], EndoSurf [27], and LerPlane [26] only test their methods on the EndoNeRF dataset [23], which consists of a static camera and has around 150 frames per sequence, ForPlane [25] also tests on the Hamlyn dataset [20,12], containing around 301 frames per sequence, utilizing the camera poses estimated using Endo-Depth-and-Motion [14]. In our experiments, we extend these investigations to recordings with moving cameras along with tissue deformations with significantly longer durations of up to 5,000 frames.

More recently, various works investigated pose optimization within NeRF setups integrating core components of visual odometry (VO) and SLAM pipelines, with experiments on natural scenes. NeRF-- [24] and BARF [7], propose solutions for jointly learning the poses as a part of NeRF optimization. LocalRF [10] further extends this idea for larger scale scenes. However, unlike in the case of endoscopic scenes, these methods employ radiance fields with static scene assumptions, and, as we show in our experiments, tend to have a performance drop when confronted with highly dynamic content.

In this work, we present FLex; a NeRF-based architecture capable of high quality novel view synthesis from pose-free surgical videos, containing strong deformations. Employing a progressive optimization scheme [10], FLex utilizes local dynamic radiance fields to jointly optimize for pose and scene representation. In the context of NeRF applications on endoscopic scenes, along with a concurrent work following a different approach, BASED [17], we believe FLex is the first to investigate joint pose optimization. In addition to architectural differences, in this work we jointly target efficient scaling to longer sequences. Extensively evaluated on the StereoMIS [4] dataset, our method shows state-of-the-art results in novel view synthesis. Furthermore, it achieves competitive tracking accuracy compared to recent methods designed specifically for this task [4].
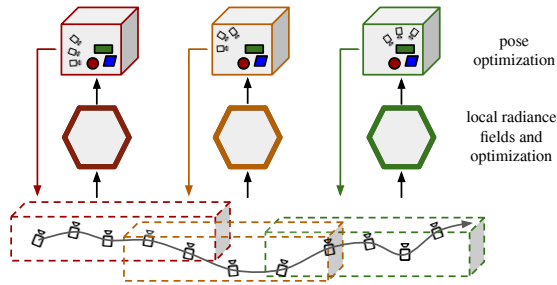
To summarize our contributions are the following:

- We present a novel NeRF architecture for the challenging task of 4D reconstruction in highly dynamic endoscopic scenes in the absence of camera pose information, achieved by joint optimization in a progressive strategy.
- Dissecting the scene into multiple smaller 4D models with overlaps allows to efficiently scale reconstruction in time. We conduct experiments to showcase this on up to 5,000 frames, more than ten times the size used in previous works. Together with our first contribution this improves the applicability to real world use-cases.
- Experimental results on the StereoMIS dataset reveal a clear advantage in terms of novel view synthesis with competitive accuracy in camera poses.

## 2   Method

### 2.1   Overview

Given a rectified stereo-endoscopic video, our goal is to reconstruct the 4D scene accurately without prior camera pose information. For this we propose a new method **FLex**, standing for **F**low-optimized **L**ocal **Hex**planes, depicted in Fig. 1, and combines advancements from recent NeRF literature to build multiple smaller dynamic models that are progressively optimized. In contrast to prior work [27,23,25], we do not have one unified representation of the scene but multiple smaller overlapping ones. The representation of local models allows us to represent larger scenes, both temporally and spatially, accurately without incurring prohibitive memory growth while maintaining a high feature grid resolution.

**Fig. 2.** Joint progressive pose and local dynamic radiance fields optimization. Spatial extents clustered within the bounding boxes of different colors represent the spatio-temporal domain of the corresponding local radiance fields. The arrow on the camera trajectory shows the temporal direction.

Furthermore, adopting a progressive optimization scheme enables the optimization of poses from scratch. For further regularization, we add supervision through optical flow and stereo depth.
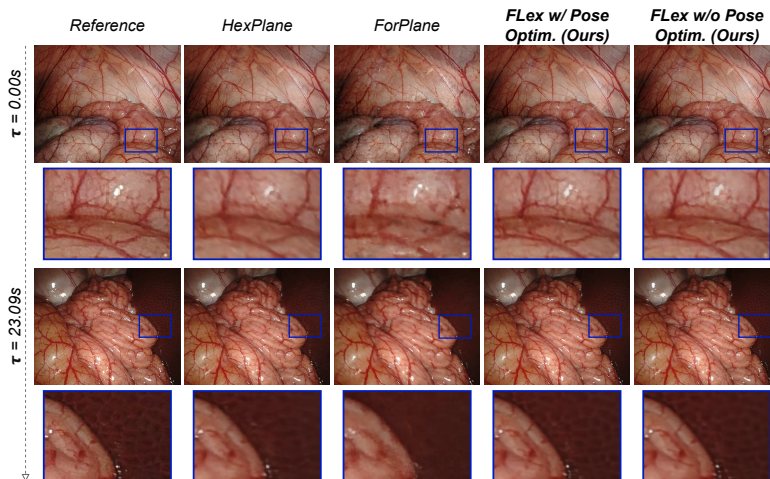
### 2.2    4D Scene Representation

NeRFs[11] implicitly model a 3D scene utilizing differentiable volume rendering to predict pixel colors. They can be adapted to a 4D scene representation by adding the timestep $k$ as an additional input to the model. HexPlane [1] models a dynamic scene using an explicit 4D feature grid paired with an implicit MLP. The grid is constructed from several planes, where each plane represents a combination of two dimensions, yielding six planes in total. During ray-casting, the corresponding features on each plane are extracted for the spatio-temporal locations and combined by multiplication and concatenation before being fed into smaller MLPs and similarly to NeRF [11] rendered with volumetric rendering.

### 2.3    Progressive Optimization

Endoscopic videos contain two main challenges for NeRF architectures: They are dependent on external tools for accurate pose estimation and can constitute arbitrarily long sequences. In order to tackle these two problems in a robust and efficient way, our joint pose and radiance fields optimization scheme utilizes the concepts of progressive optimization and dynamic allocation of local HexPlane models as visualized in Fig. 2 and inspired by LocalRF [10].

In the scope of progressive optimization, we start with the first five frames of the sequence, then we consecutively add one frame at a time, initializing it's camera pose parameters with the prior camera pose. When the appended frame increases the number of frames above a preset threshold, $t_k$, or the distance between the optimized position of the camera and the center of the current local model is larger than a distance threshold, $t_d$, we instantiate a new local model,

**Fig. 3.** Qualitative results on a 1,000 frame scene with breathing deformations and camera motion.

setting this new frame to be its origin. To ensure consistency across local models, we assign the last thirty frames of the previous model to be overlapping with the new local model. To secure a coherent trajectory during progressive optimization, we consistently sample rays from the last four appended frames. When a new local model is initiated, the weights of the previous one are frozen and offloaded from the GPU to prevent unnecessary memory usage. During inference, if a pose corresponds to the spatial and temporal extent of multiple local models, each model's contribution is aggregated into the ray-casting formulation with blending weights linearly set on the overlapping regions based on the proximity to the centers of the local models. Before a new local model is initialized, the last model goes into a refinement phase where the pose and model parameters are optimized with batches of samples uniformly picked along its entire span.

### 2.4   Training Objectives

We employ the common photometric loss $\mathcal{L}_{rgb}$ as defined in Eq. (1) with ground-truth $C(\mathbf{r})$ and predicted $\hat{C}(\mathbf{r})$ pixel values for ray $\mathbf{r}$ within the set of rays $\mathcal{R}$:

$$\mathcal{L}_{rgb} = \frac{1}{|\mathcal{R}|} \sum_{\mathbf{r} \in \mathcal{R}} \left\| \hat{C}(\mathbf{r}) - C(\mathbf{r}) \right\|_2^2 \tag{1}$$

Additionally, we also use the depth supervision loss along with the line-of-sight prior as introduced by Rematas et al. [15]. We denote them together as $\mathcal{L}_z$. The line-of-sight prior regularizes the density values along a ray to be concentrated on the actual surface, thereby, together with the depth loss, improving the capture of the scene geometry.

Furthermore, our method is optimized via an optical flow loss $\mathcal{L}_f$ in both temporal directions as described in Eq. (3). The estimated optical flow $\hat{F}_{k\to k\pm 1}$ is induced via finding the surface point for a given ray at time $k$ in 3D with the help of the predicted depth using the projection from 2D to 3D $\pi_{3D}$ and then transforming the point to the adjacent timestep $k \pm 1$ using the relative camera extrinsics $[R|T]_{k\to k\pm 1}$. The resulting 3D point is then projected from 3D to 2D via $\pi_{2D}$ using the known camera intrinsics and compared to the initial pixel location $\mathbf{p}(\mathbf{r})$ at time k, see Eq. (2).

$$\hat{\mathcal{F}}_{k\to k\pm 1}(\mathbf{r}) = \mathbf{p}(\mathbf{r}) - \pi_{2D}\left([R|T]_{k\to k\pm 1}\,\pi_{3D}(\mathbf{r}, \hat{D})\right) \tag{2}$$

$$\mathcal{L}_f = \frac{1}{|\mathcal{R}|}\sum_{\mathbf{r}\in\mathcal{R}}\left\|\hat{\mathcal{F}}_{k\to k\pm 1}(\mathbf{r}) - \mathcal{F}_{k\to k\pm 1}(\mathbf{r})\right\|_1 \tag{3}$$

All the aforementioned losses are aggregated into our final loss function $\mathcal{L}$. It is essential to highlight that we employ all three loss terms to optimize our method FLex. The camera extrinsics are only optimized by the optical flow loss $\mathcal{L}_f$.

Note that the optical flow loss $\mathcal{L}_f$ is entirely removed after 20% of the refinement iterations to ensure an early pose convergence.

$$\mathcal{L} = \mathcal{L}_{rgb} + \lambda_z\mathcal{L}_z + \lambda_f\mathcal{L}_f \tag{4}$$

## 3    Experiments

### 3.1    Dataset and Evaluation Metrics

We systematically assess the efficacy of our approach using the publicly available StereoMIS [4] dataset, recorded using a stereo endoscope of a da Vinci Xi robot; ground-truth camera trajectories are measured using the forward kinematics. In total we extract five sequences for general comparison, each 1,000 frames long equivalent to ca. 29 seconds per clip. Two of these scenes exhibit concurrent breathing motion and camera movement, while another two showcase pronounced non-rigid deformations induced by surgical tools amidst subtle changes in camera perspective. The remaining scene presents a nearly static environment accompanied by rapid camera movements. Furthermore, we create two additional longer sequences to study the method's behavior given a larger temporal and spatial extent. One sequence contains deformations induced by surgical tools and consists of 5,000 frames, while the other scene incorporates more rapid camera motion and comprises 4,000 frames. For measuring our method quantitatively, we follow the evaluation of [23] by making use of PSNR, SSIM, and LPIPS metrics (denoted with subscripts "a" and "v" for AlexNet [5] and VGG [18] backbones) and L1-distance metric, indicated in mm, to assess the captured geometry by comparing the induced and stereo-estimated depth images. We evaluate the estimated camera poses using root-mean-squared absolute trajectory error (ATE-RMSE), relative translational and rotational pose errors (RPE-Trans and RPE-Rot).

**Table 1.** View synthesis quality on StereoMIS dataset. The metrics are computed as an average for five 1,000 frame sequences. L1-Distance is computed between the synthesized and the ground truth depth images in mm. The best result for each metric is marked in bold.

| Model | PSNR ↑ | SSIM ↑ | LPIPS$_a$ ↓ | LPIPS$_v$ ↓ | L1-Distance ↓ |
|---|---|---|---|---|---|
| EndoNeRF [23] | 21.99 | 0.590 | 0.496 | 0.514 | – |
| EndoSurf [27] | 25.18 | 0.622 | 0.528 | 0.529 | 8.105 |
| ForPlane [25] | 30.35 | 0.783 | 0.208 | 0.301 | 23.717 |
| LocalRF$^{†2}$ [10] | 27.41 | 0.781 | 0.245 | 0.288 | 4.576 |
| HexPlane$^{†1}$ [1] | 30.85 | 0.819 | 0.211 | 0.273 | 1.532 |
| FLex w/o Pose Optim. (Ours) | **31.10** | **0.836** | 0.200 | **0.244** | 1.456 |
| FLex w/ Pose Optim. (Ours) | 30.62 | 0.818 | **0.179** | 0.245 | **1.273** |

## 3.2   Implementation Details

In our local models, we set the dimension of the spatial feature grids to 512 for $(x, y, z)$, and the temporal dimension is set to half of the image sequence length of a scene. The feature dimension is 72 in total for both density and color. For a fair comparison, we ensure equal capacity for all methods using explicit data structures [1,10,25]. Additionally, we adopt a coarse-to-fine approach as in HexPlane [1] to start with a lower grid resolution and increase over time to the settings mentioned above. Except for our method with pose optimization and LocalRF [10], we use Robust-Pose Estimation [4] to estimate the camera poses. In the experiment tables, we use HexPlane$^{†1}$ to denote an improved version with scene contraction, depth loss, and optical flow loss; and LocalRF$^{†2}$ to depict a version of it which substitutes the monocular with the stereo depth estimation.

Furthermore, we set $\lambda_z = 0.01$ and $\lambda_f = 1.0$ as illustrated in Eq. (4) and $t_k = 100$ and $t_d = 1.0$. Overall, we train our method without pose optimization for 100 iterations per frame with a batch size of 4,096 rays, which takes approximately 7 hours on up to 40 GB of an Nvidia Tesla A100 and FLex with pose optimization for an additional 100 iterations per frame during the prior progressive optimization which yields ca. 20 hours on the same hardware configuration. In comparison, HexPlane, for the exact same settings, takes ca. 6 hours to train. We do not mask the tools in any experiment. In addition, we use RAFT [21] for estimating both optical flow from frame-to-frame $\mathcal{F}_{k \to k \pm 1}(\mathbf{r})$ for both directions and for obtaining stereo depth $D$, which we both use as pseudo-ground-truth for model optimization.

## 3.3   Quantitative and Qualitative Results

We conduct a comprehensive comparison of our method, both with and without pose optimization, against the latest state-of-the-art (SoTA) NeRF methods designed for endoscopy [25,23,27] and two additional baselines [1,10]. The results in Table 1, summarizing the average results across all 5 scenes, demonstrate that our method without pose optimization consistently outperforms all baselines and notably surpasses the current endoscopic SoTA, ForPlane. In addition our method with pose optimization (FLex w/ Pose Optim.) also manages

**Table 2.** Ablation study on longer StereoMIS sequences. L1-distance is computed between the synthesized and the estimated stereo depth images in mm. The best results are marked in bold.

| Model | Frame # | PSNR ↑ | SSIM ↑ | LPIPS$_a$ ↓ | LPIPS$_v$ ↓ | L1-Distance ↓ |
|---|---|---|---|---|---|---|
| HexPlane[†1] [1] | 4,000 | 24.79 | 0.614 | 0.545 | 0.510 | 4.856 |
| FLex w/o Pose Optim. (Ours) | 4,000 | **26.09** | **0.661** | **0.498** | **0.469** | **3.567** |
| HexPlane[†1] [1] | 5,000 | 28.55 | 0.718 | 0.453 | 0.471 | 1.902 |
| FLex w/o Pose Optim. (Ours) | 5,000 | **29.97** | **0.773** | **0.386** | **0.413** | **1.704** |

to outperform ForPlane in all metrics and is competitve to ours without pose optimization (FLex w/o Pose Optim.) and the improved HexPlane. These quantitative findings are substantiated by our qualitative results presented in Fig. 3, highlighting that FLex, with and without prior poses, achieves superior image quality compared to the most competitive baselines [1,25].

### 3.4   Impact of Sequence Length on View Synthesis Quality

We evaluate our model's effect on longer sequences. Ensuring a reasonable hardware memory cap of 16 GB vRAM, we set the spatial grid sizes to 128. Additionally, we set the maximum temporal dimensions for our method and HexPlane to 50 per local model and 100, respectively, and we train the models for 100 iterations per frame. As displayed in Table 2, our method preserves its reconstruction quality scaling to the longer sequences and achieves a higher performance difference to HexPlane in comparison to the difference on shorter sequences.

### 3.5   Pose Accuracy

We compare our method against a SoTA method in visual odometry for endoscopic scenes, Robust-Pose Estimation [4], on 3 sequences each with 1,000 frames. As highlighted in Table 3, FLex performs competitively achieving close results to the baseline. However, we emphasize that this task is not the main focus of our work and can be improved using robust optimization and globally consistent methods in the future.

**Table 3.** Average Pose accuracy on StereoMIS dataset. ATE-RMSE and RPE-Trans are in mm, RPE-Rot is in degrees. The best results are marked in bold.

| Model | ATE-RMSE ↓ | RPE-Trans ↓ | RPE-Rot ↓ |
|---|---|---|---|
| Robust-Pose Estimation [4] | $\mathbf{2.164} \pm 2.68e-1$ | $\mathbf{0.073} \pm 3e-5$ | $\mathbf{0.043} \pm 2e-6$ |
| LocalRF[†2] [10] | $7.704 \pm 1.506$ | $0.160 \pm 8e-4$ | $0.119 \pm 2e-5$ |
| FLex w/ Pose Optim. (Ours) | $2.565 \pm 1.6e-1$ | $0.127 \pm 9e-4$ | $0.102 \pm 4e-6$ |

## 4   Conclusion

In this work, we present FLex, a novel method for reconstructing pose-free, long surgical videos with challenging tissue deformations and camera motion. Our approach successfully eliminates the reliance on prior poses by jointly optimizing for reconstruction and camera trajectory. FLex improves upon the scalability of dynamic NeRFs for larger scenes thus becoming more applicable to realistic surgical recordings, while improving over current methods on the StereoMIS dataset in terms of novel view synthesis with competitive pose accuracy. We believe that FLex can pave the way towards more easily accessible, realistic and reliable 4D endoscopy reconstructions to improve post surgical analysis and medical education.

## References

1. Cao, A., Johnson, J.: Hexplane: A fast representation for dynamic scenes. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 130–141 (2023)
2. Chen, A., Xu, Z., Geiger, A., Yu, J., Su, H.: Tensorf: Tensorial radiance fields. In: European Conference on Computer Vision. pp. 333–350. Springer (2022)
3. Fridovich-Keil, S., Meanti, G., Warburg, F.R., Recht, B., Kanazawa, A.: K-planes: Explicit radiance fields in space, time, and appearance. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 12479–12488 (2023)
4. Hayoz, M., Hahne, C., Gallardo, M., Candinas, D., Kurmann, T., Allan, M., Sznitman, R.: Learning how to robustly estimate camera pose in endoscopic videos. International journal of computer assisted radiology and surgery pp. 1–8 (2023)
5. Krizhevsky, A., Sutskever, I., Hinton, G.E.: Imagenet classification with deep convolutional neural networks. Advances in neural information processing systems **25** (2012)
6. Lange, T., Indelicato, D.J., Rosen, J.M.: Virtual reality in surgical training. Surgical oncology clinics of North America **9**(1), 61–79 (2000)
7. Lin, C.H., Ma, W.C., Torralba, A., Lucey, S.: Barf: Bundle-adjusting neural radiance fields. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 5741–5751 (2021)
8. Liu, X., Stiber, M., Huang, J., Ishii, M., Hager, G.D., Taylor, R.H., Unberath, M.: Reconstructing sinus anatomy from endoscopic video–towards a radiation-free approach for quantitative longitudinal assessment. In: Medical Image Computing and Computer Assisted Intervention–MICCAI 2020: 23rd International Conference, Lima, Peru, October 4–8, 2020, Proceedings, Part III 23. pp. 3–13. Springer (2020)
9. Long, Y., Cao, J., Deguet, A., Taylor, R.H., Dou, Q.: Integrating artificial intelligence and augmented reality in robotic surgery: An initial dvrk study using a surgical education scenario. In: 2022 International Symposium on Medical Robotics (ISMR). pp. 1–8. IEEE (2022)
10. Meuleman, A., Liu, Y.L., Gao, C., Huang, J.B., Kim, C., Kim, M.H., Kopf, J.: Progressively optimized local radiance fields for robust view synthesis. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 16539–16548 (2023)

11. Mildenhall, B., Srinivasan, P.P., Tancik, M., Barron, J.T., Ramamoorthi, R., Ng, R.: Nerf: Representing scenes as neural radiance fields for view synthesis. Communications of the ACM **65**(1), 99–106 (2021)
12. Mountney, P., Stoyanov, D., Yang, G.Z.: Three-dimensional tissue deformation recovery and tracking. IEEE Signal Processing Magazine **27**(4), 14–24 (2010)
13. Pumarola, A., Corona, E., Pons-Moll, G., Moreno-Noguer, F.: D-nerf: Neural radiance fields for dynamic scenes. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 10318–10327 (2021)
14. Recasens, D., Lamarca, J., Fácil, J.M., Montiel, J., Civera, J.: Endo-depth-and-motion: Reconstruction and tracking in endoscopic videos using depth networks and photometric constraints. IEEE Robotics and Automation Letters **6**(4), 7225–7232 (2021)
15. Rematas, K., Liu, A., Srinivasan, P.P., Barron, J.T., Tagliasacchi, A., Funkhouser, T., Ferrari, V.: Urban radiance fields. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 12932–12942 (2022)
16. Rodriguez, J.J.G., Montiel, J., Tardos, J.D.: Nr-slam: Non-rigid monocular slam. arXiv preprint arXiv:2308.04036 (2023)
17. Saha, S., Liu, S., Lin, S., Lu, J., Yip, M.: Based: Bundle-adjusting surgical endoscopic dynamic video reconstruction using neural radiance fields. arXiv preprint arXiv:2309.15329 (2023)
18. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. arXiv preprint arXiv:1409.1556 (2014)
19. Song, J., Wang, J., Zhao, L., Huang, S., Dissanayake, G.: Dynamic reconstruction of deformable soft-tissue with stereo scope in minimal invasive surgery. IEEE Robotics and Automation Letters **3**(1), 155–162 (2017)
20. Stoyanov, D., Mylonas, G.P., Deligianni, F., Darzi, A., Yang, G.Z.: Soft-tissue motion tracking and structure estimation for robotic assisted mis procedures. In: International Conference on Medical Image Computing and Computer-Assisted Intervention. pp. 139–146. Springer (2005)
21. Teed, Z., Deng, J.: Raft: Recurrent all-pairs field transforms for optical flow pp. 402–419 (2020)
22. Wang, P., Liu, L., Liu, Y., Theobalt, C., Komura, T., Wang, W.: Neus: Learning neural implicit surfaces by volume rendering for multi-view reconstruction. arXiv preprint arXiv:2106.10689 (2021)
23. Wang, Y., Long, Y., Fan, S.H., Dou, Q.: Neural rendering for stereo 3d reconstruction of deformable tissues in robotic surgery. In: International Conference on Medical Image Computing and Computer-Assisted Intervention. pp. 431–441. Springer (2022)
24. Wang, Z., Wu, S., Xie, W., Chen, M., Prisacariu, V.A.: Nerf–: Neural radiance fields without known camera parameters. arXiv preprint arXiv:2102.07064 (2021)
25. Yang, C., Wang, K., Wang, Y., Dou, Q., Yang, X., Shen, W.: Efficient deformable tissue reconstruction via orthogonal neural plane. arXiv preprint arXiv:2312.15253 (2023)
26. Yang, C., Wang, K., Wang, Y., Yang, X., Shen, W.: Neural lerplane representations for fast 4d reconstruction of deformable tissues. arXiv preprint arXiv:2305.19906 (2023)
27. Zha, R., Cheng, X., Li, H., Harandi, M., Ge, Z.: Endosurf: Neural surface reconstruction of deformable tissues with stereo endoscope videos. In: International Conference on Medical Image Computing and Computer-Assisted Intervention. pp. 13–23. Springer (2023)

# FLex: Joint Pose and Dynamic Radiance Fields Optimization for Stereo Endoscopic Videos - Supplementary Material

Florian Philipp Stilz[*,1,2], Mert Asim Karaoglu[*,1,2], Felix Tristram[*,1], Nassir Navab[1], Benjamin Busam[1], and Alexander Ladikos[2]

[1] Technical University Munich
[2] ImFusion GmbH

**Table 1.** Per sequence comparison of pose accuracy on StereoMIS dataset. LocalRF[†2] is the vanilla LocalRF, but optimized via stereo depth. The ATE-RMSE and the RPE-Trans are in mm, and the RPE-Rot is in degrees.

| Model | Deform. | Camera Motion | Tool | ATE-RMSE ↓ | RPE-Trans ↓ | RPE-Rot ↓ |
|---|---|---|---|---|---|---|
| Robust Pose Estimation | ✓ | ✓ | ✗ | 2.407 | **0.068** | **0.043** |
| | ✗ | ✓ | ✗ | 2.640 | **0.080** | **0.054** |
| | ✓ | ✓ | ✗ | **1.444** | **0.071** | **0.032** |
| | Average | | | **$2.164 \pm 2.68e-1$** | **$0.073 \pm 3e-5$** | **$0.043 \pm 2e-6$** |
| LocalRF[†2] | ✓ | ✓ | ✗ | 8.210 | 0.155 | 0.136 |
| | ✗ | ✓ | ✗ | 8.888 | 0.198 | 0.143 |
| | ✓ | ✓ | ✗ | 6.013 | 0.128 | 0.079 |
| | Average | | | $7.704 \pm 1.506$ | $0.436 \pm 8e-4$ | $0.119 \pm 2e-5$ |
| Ours w/ Pose Optim. | ✓ | ✓ | ✗ | **2.106** | 0.099 | 0.090 |
| | ✗ | ✓ | ✗ | **2.509** | 0.113 | 0.123 |
| | ✓ | ✓ | ✗ | 3.081 | 0.168 | 0.093 |
| | Average | | | $2.565 \pm 1.6e-1$ | $0.127 \pm 9e-4$ | $0.102 \pm 4e-6$ |

**Table 2.** Per sequence comparison of novel synthesis quality on StereoMIS dataset. L1 Distance is computed between the synthesized and the ground truth depth images in mm. HexPlane[†1] stands for an optimized baseline with scene contraction, depth loss, and optical flow loss, while LocalRF[†2] is the vanilla LocalRF, but optimized via stereo depth. Blue indicates the best result and red the second best for each metric, respectively.

| Model | Deform. | Camera Motion | Tool | PSNR ↑ | SSIM ↑ | LPIPS$_a$ ↓ | LPIPS$_v$ ↓ | L1 Distance ↓ |
|---|---|---|---|---|---|---|---|---|
| EndoNeRF | ✓ | ✓ | ✗ | 25.28 | 0.628 | 0.507 | 0.500 | — |
|  | ✗ | ✓ | ✗ | 11.22 | 0.546 | 0.517 | 0.573 | — |
|  | ✓ | ✗ | ✓ | 23.87 | 0.569 | 0.501 | 0.517 | — |
|  | ✓ | ✓ | ✗ | 26.56 | 0.636 | 0.409 | 0.456 | — |
|  | ✓ | ✗ | ✓ | 23.03 | 0.569 | 0.545 | 0.524 | — |
| EndoSurf | ✓ | ✓ | ✗ | 25.14 | 0.619 | 0.516 | 0.533 | 1.143 |
|  | ✗ | ✓ | ✗ | 29.81 | 0.766 | 0.513 | 0.552 | 1.895 |
|  | ✓ | ✗ | ✓ | 23.42 | 0.582 | 0.493 | 0.505 | 14.025 |
|  | ✓ | ✓ | ✗ | 25.39 | 0.609 | 0.482 | 0.505 | 17.147 |
|  | ✓ | ✗ | ✓ | 22.14 | 0.533 | 0.618 | 0.567 | 6.318 |
| LerPlane | ✓ | ✓ | ✗ | 29.47 | 0.766 | 0.182 | 0.292 | 8.175 |
|  | ✗ | ✓ | ✗ | 36.17 | 0.900 | 0.139 | 0.273 | 20.875 |
|  | ✓ | ✗ | ✓ | 27.28 | 0.710 | 0.272 | 0.355 | 30.702 |
|  | ✓ | ✓ | ✗ | 32.70 | 0.850 | 0.149 | 0.214 | 30.325 |
|  | ✓ | ✗ | ✓ | 26.11 | 0.690 | 0.296 | 0.373 | 28.511 |
| LocalRF[†2] | ✓ | ✓ | ✗ | 29.02 | 0.818 | 0.177 | 0.233 | 3.926 |
|  | ✗ | ✓ | ✗ | 35.07 | 0.895 | 0.172 | 0.235 | 2.701 |
|  | ✓ | ✗ | ✓ | 22.54 | 0.701 | 0.318 | 0.352 | 6.424 |
|  | ✓ | ✓ | ✗ | 31.22 | 0.841 | 0.166 | 0.209 | 4.418 |
|  | ✓ | ✗ | ✓ | 19.21 | 0.649 | 0.394 | 0.409 | 5.411 |
| HexPlane[†1] | ✓ | ✓ | ✗ | 31.50 | 0.854 | 0.170 | 0.233 | 1.098 |
|  | ✗ | ✓ | ✗ | 36.70 | 0.916 | 0.178 | 0.251 | 2.497 |
|  | ✓ | ✗ | ✓ | 27.13 | 0.745 | 0.258 | 0.317 | 1.397 |
|  | ✓ | ✓ | ✗ | 33.04 | 0.872 | 0.162 | 0.206 | 1.154 |
|  | ✓ | ✗ | ✓ | 25.90 | 0.710 | 0.287 | 0.359 | 1.516 |
| Ours w/o Pose Optim. | ✓ | ✓ | ✗ | 31.91 | 0.875 | 0.155 | 0.201 | 1.234 |
|  | ✗ | ✓ | ✗ | 37.03 | 0.917 | 0.173 | 0.242 | 1.708 |
|  | ✓ | ✗ | ✓ | 27.21 | 0.773 | 0.242 | 0.274 | 1.533 |
|  | ✓ | ✓ | ✗ | 33.71 | 0.889 | 0.152 | 0.184 | 1.052 |
|  | ✓ | ✗ | ✓ | 25.65 | 0.724 | 0.279 | 0.318 | 1.752 |
| Ours w/ Pose Optim. | ✓ | ✓ | ✗ | 31.53 | 0.851 | 0.139 | 0.196 | 1.052 |
|  | ✗ | ✓ | ✗ | 35.25 | 0.885 | 0.180 | 0.262 | 1.093 |
|  | ✓ | ✗ | ✓ | 28.07 | 0.801 | 0.187 | 0.246 | 1.453 |
|  | ✓ | ✓ | ✗ | 31.97 | 0.816 | 0.157 | 0.216 | 1.062 |
|  | ✓ | ✗ | ✓ | 26.30 | 0.736 | 0.230 | 0.304 | 1.703 |