



# EndoFlow-SLAM: Real-Time Endoscopic SLAM with Flow-Constrained Gaussian Splatting

Taoyu Wu<sup>1</sup>, Yiyi Miao<sup>1</sup>, Zhuoxiao Li<sup>1,2</sup>, Haocheng Zhao<sup>1</sup>, KangDang<sup>1</sup>, Jionglong Su<sup>1</sup>, Limin Yu<sup>1</sup>, and Haoang Li<sup>2</sup>

<sup>1</sup>Xi'an Jiaotong Liverpool University <sup>2</sup>The Hong Kong University of Science and Technology (Guangzhou)



**Abstract:** Efficient three-dimensional reconstruction and real-time visualization are critical in surgical scenarios such as endoscopy. In recent years, 3D Gaussian Splatting (3DGS) has demonstrated remarkable performance in efficient 3D reconstruction and rendering. Most 3DGS-based Simultaneous Localization and Mapping (SLAM) methods only rely on the appearance constraints for optimizing both 3DGS and camera poses. However, in endoscopic scenarios, the challenges include photometric inconsistencies caused by non-Lambertian surfaces and dynamic motion from breathing affects the performance of SLAM systems. To address these issues, we additionally introduce optical flow loss as a geometric constraint, which effectively constrains both the 3D structure of the scene and the camera motion. Furthermore, we propose a depth regularization strategy to mitigate the problem of photometric inconsistencies and ensure the validity of 3DGS depth rendering in endoscopic scenes. In addition, to improve scene representation in the SLAM system, we improve the 3DGS refinement strategy by focusing on view-points corresponding to Keyframes with suboptimal rendering quality frames, achieving better rendering results. Extensive experiments on the C3VD static dataset and the StereoMIS dynamic dataset demonstrate that our method outperforms existing state-of-the-art methods in novel view synthesis and pose estimation, exhibiting high performance in both static and slightly dynamic surgical scenes.

## Background

- Importance in Surgery**  
Accurate camera pose estimation and tissue reconstruction are crucial for minimally invasive surgeries. They improve spatial awareness between organs and instruments, aiding surgical navigation and postoperative evaluation.
- Challenges in Endoscopy**  
Limited fields of view, dynamic tissue deformations, and complex lighting conditions reduce accuracy. Traditional visual SLAM often assumes rigid environments and relies on sparse reconstructions, which are unsuitable for endoscopic scenarios with weak textures.
- Limitations of Dense SLAM**  
Dense SLAM enables real-time reconstructions but typically depends on RGB-D data. However, in endoscopic environments, limited instrument mobility and access constraints hinder effective data acquisition, leading to incomplete and occluded scene representations.
- Advances with NeRF and 3DGS**  
Recent SLAM systems integrate Neural Radiance Fields (NeRF) and 3D Gaussian Splatting (3DGS). NeRF offers high-precision view synthesis but is computationally heavy, while 3DGS provides faster rendering with realistic images, achieving effective tracking and dense reconstruction.

## Methodology

### Preliminary: 3D Gaussian Splatting

3DGS [10] is a differentiable rendering framework that models a 3D scene using a collection of Gaussian primitives. The  $i$ -th Gaussian is characterized by a center  $\mu_i$ , an opacity  $\alpha_i$ , and a covariance matrix  $\Sigma_i$  and can be expressed as:

$$G_i(\mathbf{X}) = \alpha_i \cdot \exp\left\{-\frac{1}{2}(\mathbf{X} - \mu_i)^T \Sigma_i^{-1}(\mathbf{X} - \mu_i)\right\}, \quad (1)$$

where  $\mathbf{X}$  represents an arbitrary point in 3D space. In the rendering process, these Gaussians are projected onto the 2D image plane along camera rays, and their contributions are composited using an alpha-blending scheme.

### Non-keyframe Optimization

Existing 3DGS-based SLAM systems [25, 16, 19] typically employ photometric ( $\mathcal{L}_{rgb}$ ) and depth ( $\mathcal{L}_{depth}$ ) error minimization for camera pose estimation. However, this paradigm faces scale ambiguity challenges in endoscopic scenarios due to monocular depth estimation limitations [28, 15]. To address this limitation, we propose a modified optimization framework incorporating scale-invariant loss  $\mathcal{L}_{scale}$  [21] and depth gradient regularization  $\mathcal{L}_{depth}^{reg}$ . The depth regularization term is defined as the weighted gradient difference between the estimated and ground truth depth maps:  $\mathcal{L}_{depth}^{reg} = \frac{1}{N} \sum_{i=1}^N (w_h \cdot |\nabla_h d_i| + w_v \cdot |\nabla_v d_i|)$ , where  $d_i$  denotes the depth value at pixel  $i$ ,  $\nabla_h$  and  $\nabla_v$  represent horizontal and vertical gradients respectively,  $w_h$  and  $w_v$  are corresponding weight factors, and  $N$  represents the total number of image pixels. The estimated camera pose  $\hat{T}_{t+1}$  at time  $t+1$  can be optimized by minimizing the photometric, scale-invariant and depth regularization losses:

$$\hat{T}_{t+1} = \underset{T_{t+1}}{\operatorname{argmin}} \lambda_1 \mathcal{L}_{rgb} + \lambda_2 \mathcal{L}_{depth}^{reg} + \lambda_3 \mathcal{L}_{scale}, \quad (2)$$

### Keyframe Optimization with Flow Constraints

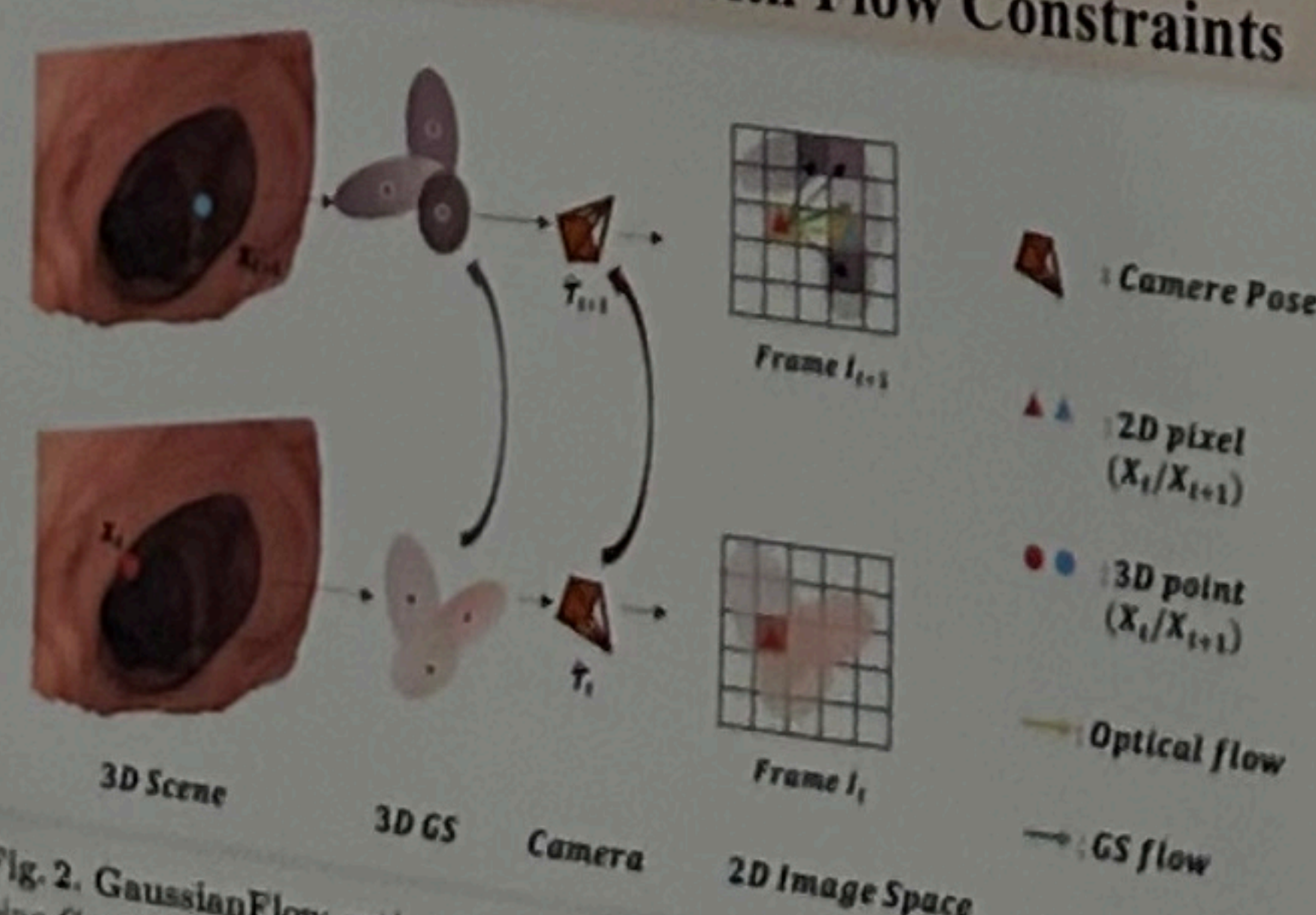


Fig. 2. GaussianFlow estimation. At time  $t$ , each pixel  $x_t$  results from  $K$  overlapping Gaussians. At time  $t+1$ , each  $K$  Gaussian will have a corresponding Gaussian flow (Black arrow). By accumulating these Gaussian flows, we obtain the overall Gaussian flow. Our goal is to minimize the difference between GS flow and Optical flow by optimizing both the camera pose  $\hat{T}_{t+1}$  and 3DGS primitive  $\hat{G}$ .

## Keyframe Optimization and Global Refinement

### Keyframe-Oriented Local Bundle Adjustment

**Keyframe-Oriented Local Bundle Adjustment.** To further improve the expressiveness of the scene, we also introduce flow loss as a geometric constraint in Mapping. After the camera tracking module, we perform 3D Gaussian map representation optimization only when frame  $I_{t+1}$  is designated as a Keyframe. We follow the Keyframe management strategy from [19]. We perform BA on all frames within the Keyframe window. By minimizing the loss function  $\mathcal{L}_{rgb}$ ,  $\mathcal{L}_{depth}$  and  $\mathcal{L}_{flow}$  in the objective function (3), we simultaneously optimize the estimated camera pose  $\hat{T}_{t+1}$  and 3DGS primitive  $\hat{G}$ .

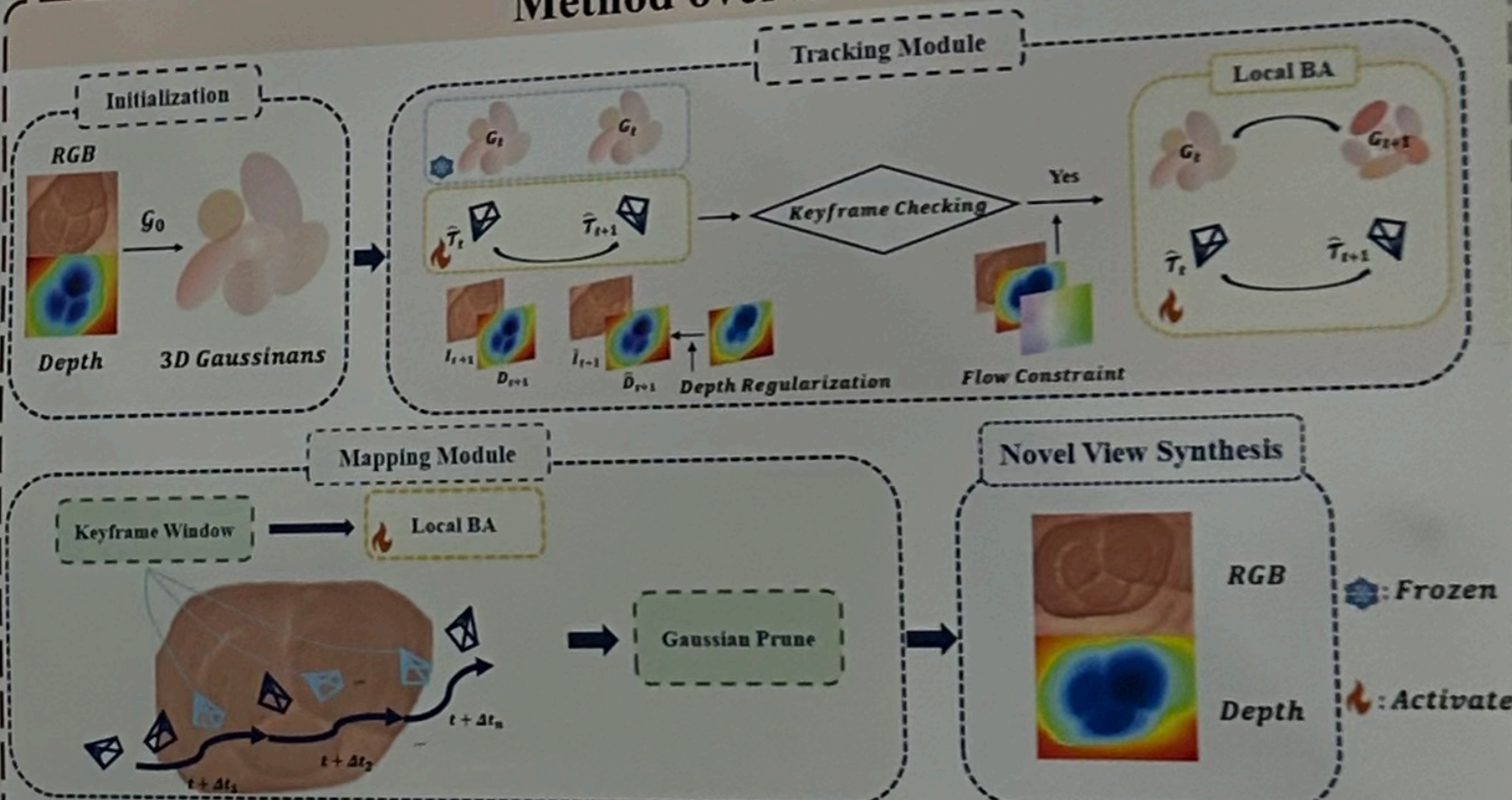
### Global Refinement

For each selected viewpoint, the optimization of the Gaussian primitive is performed using the following loss function:

$$\mathcal{L}_{refine} = (1 - \lambda_{ssim}) \mathcal{L}_{rgb} + \lambda_{ssim} (1 - \text{SSIM}) + |\nabla \hat{D} - \nabla D|_2^2, \quad (4)$$

where  $\lambda_{ssim}$  controls the relative weight of the Structural Similarity Index (SSIM). The SSIM term ensures structural preservation, maintaining the integrity of the scene's overall structure while minimizing pixel-level discrepancies.

## Method overview



**Method overview.** Given the first RGB-D image, we initialize the 3D Gaussians and subsequently perform camera tracking and mapping iteratively. Tracking handles non-keyframes with depth regularization for scale consistency (optimizing camera pose only), while keyframes incorporate optical flow as a geometric constraint in Local BA to simultaneously optimize poses and 3DGS primitives. Mapping applies BA with optical flow constraints to optimize poses and 3DGS primitives across the Keyframe window.

## Quantitative Results

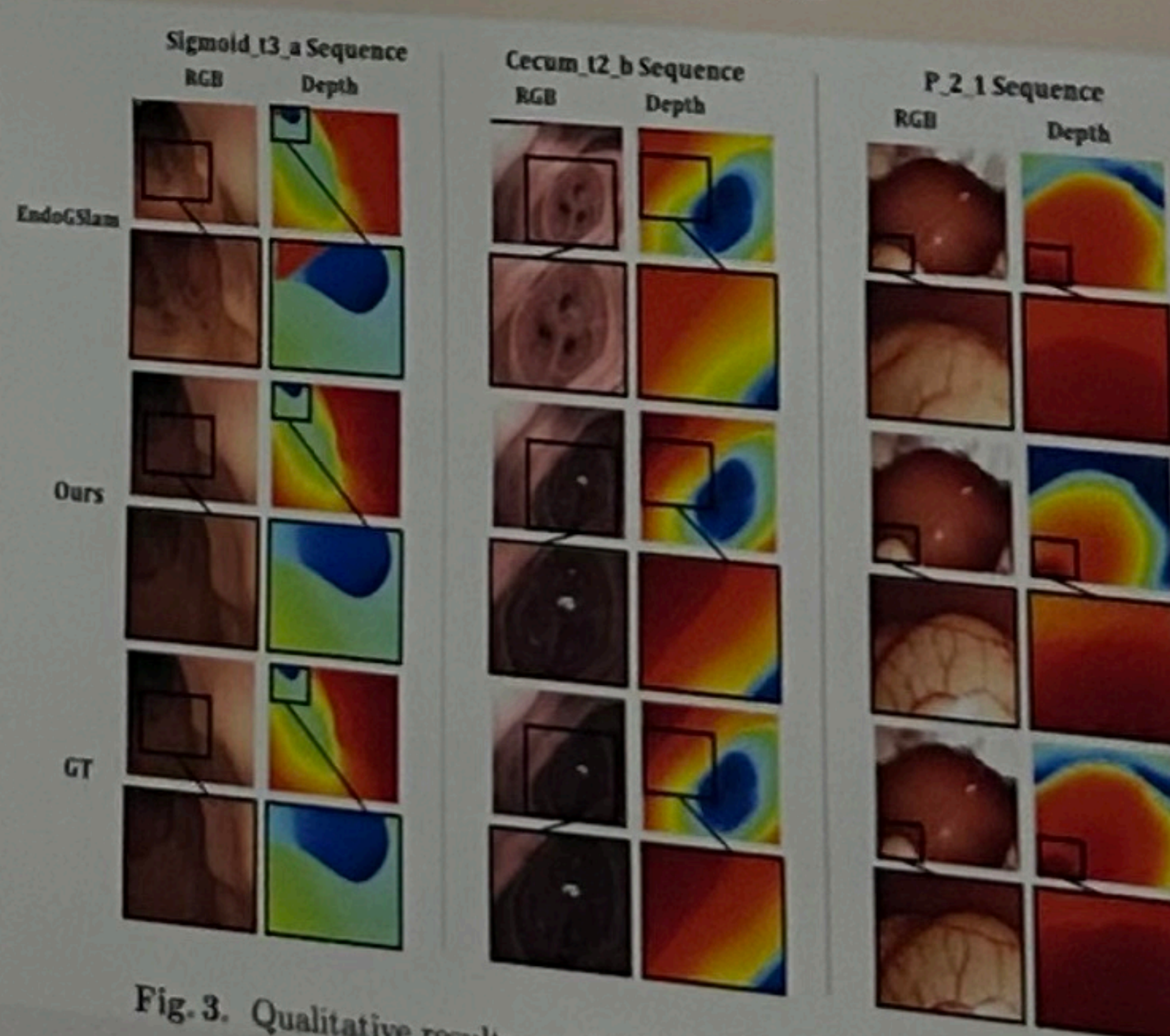


Fig. 3. Qualitative results on C3VD and StereoMIS Dataset.

## Qualitative Results

Table 1. Quantitative results on the C3VD dataset.

Methods	PSNR↑	SSIM↑	LPIPS↓	RMSE (mm)↓	ATE (mm)↓
NICE-SLAM [29]	22.07	0.73	0.33	1.88	0.48
Endo-Depth [22]	18.13	0.64	0.33	5.10	1.25
EndoGSLAM-H [25]	22.16	0.77	0.22	2.17	0.34
Ours	25.18	0.82	0.27	1.54	0.23

Table 2. Quantitative results on the StereoMIS dataset.

Methods	PSNR↑	SSIM↑	LPIPS↓	ATE (mm)↓
NICE-SLAM [29]	13.07	0.49	0.61	38.24
ESLAM [8]	18.70	0.54	0.57	16.73
EndoGSLAM-H [25]	16.67	0.52	0.45	18.82
Ours	21.96	0.59	0.27	

## Conclusion

In this paper, we introduce EndoFlow-SLAM, a SLAM framework based on 3DGS. This framework enables accurate camera tracking and high-quality novel view synthesis in endoscopic scenes. By incorporating optical flow as a geometric constraint, EndoFlow-SLAM is better at handling the dynamic changes caused by breathing in real-world endoscopic environments. Extensive experimental results demonstrate that, compared to traditional SLAM methods and those based on 3DGS, EndoFlow-SLAM achieves superior tracking and rendering performance. Future work will focus on the construction and optimization of object-level 3DGS to better adapt to intense