

ENDO-G²T: GEOMETRY-GUIDED & TEMPORALLY AWARE TIME-EMBEDDED 4DGS FOR ENDOSCOPIC SCENES

Yangle Liu¹ Fengze Li^{1,2} Kan Liu^{1,2} Jieming Ma^{2*}

¹ University of Liverpool

² Xi'an Jiaotong–Liverpool University

ABSTRACT

Endoscopic (endo) video exhibits strong view-dependent effects such as specularities, wet reflections, and occlusions. Pure photometric supervision misaligns with geometry and triggers early geometric drift, where erroneous shapes are reinforced during densification and become hard to correct. We ask how to anchor geometry early for 4D Gaussian splatting (4DGS) while maintaining temporal consistency and efficiency in dynamic endoscopic scenes. Thus, we present Endo-G²T, a geometry-guided and temporally aware training scheme for time-embedded 4DGS. First, geo-guided prior distillation converts confidence-gated monocular depth into supervision with scale-invariant depth and depth-gradient losses, using a warm-up-to-cap schedule to inject priors softly and avoid early overfitting. Second, a time-embedded Gaussian field represents dynamics in XYZT with a rotor-like rotation parameterization, yielding temporally coherent geometry with lightweight regularization that favors smooth motion and crisp opacity boundaries. Third, keyframe-constrained streaming improves efficiency and long-horizon stability through keyframe-focused optimization under a max-points budget, while non-keyframes advance with lightweight updates. Across EndoNeRF and StereoMIS-P1 datasets, Endo-G²T achieves state-of-the-art results among monocular reconstruction baselines.

Index Terms— Endoscopy, Monocular geometry prior distillation, 4D Gaussian splatting (4DGS), Temporal consistency, Keyframe-constrained streaming

1. INTRODUCTION

Endoscopic imaging is challenging for 3D and 4D reconstruction [1]. Specularities and wet reflections break Lambertian assumptions; tissues deform nonrigidly with occasional topology changes; narrow field of view yields sparse baselines; and instruments cause frequent, structured occlusions. Classical endoscopic SLAM and multiview stereo use illumination models, shading/reflectance priors, tool masking, nonrigid registration, and stereo fusion [2], yet degrade under strong highlights, large elastic motion, low texture, and

long sequences with drift [3]. Failures include depth bias near glossy mucosa, instability at instrument boundaries, and cumulative pose error in extended procedures [4].

Neural rendering reframes reconstruction by learning radiance and geometry from images. NeRF methods model a continuous, view-dependent radiance field with high-fidelity novel views but require long optimization and dense sampling [5]. 3D Gaussian splatting (3DGS) [6] speeds up rendering by projecting Gaussian primitives to screen space with differentiable splatting, achieving real-time quality. Its dynamic extension, 4D Gaussian splatting (4DGS) [7], introduces time-varying primitives and a time-embedded motion space, enabling fast dynamic reconstruction and much faster convergence for videos.

Motivated by these advances, several works adapt neural rendering to endoscopy. EndoNeRF-style methods add tool-aware sampling, learned deformation fields, and stereo or photometric cues to regularize depth near tissue surfaces [8, 9]. Endoscopic Gaussian-splatting variants, such as Endo-4DGS and successors, initialize from monocular pseudo-depth, use confidence-guided learning with surface-normal and depth regularization, and employ lightweight deformation heads for temporal modeling [10, 11]. Despite these gains, two issues persist: early geometric drift when pseudo-depth is biased or injected too strongly, which is then reinforced during densification; and temporal decoherence with uncontrolled point growth on long, fast, and occluded sequences, which reduces stability and frame rate.

These limitations imply two requirements: an appearance-agnostic geometric anchor early in training, and a schedule that preserves temporal consistency and efficiency. Recent visual-geometry transformers satisfy the first need by producing pixel-aligned monocular depth with confidence in a single pass, well suited for early distillation [12]. Their streaming variant adds causal attention and cached memory for incremental inference over long sequences, inspiring our system design [13]. On the representation side, higher-dimensional Gaussian fields with rotor-based rotations yield smooth, stable trajectories in time-embedded spaces [14], and streaming Gaussian pipelines show that keyframe selection and point-budget control effectively curb long-horizon drift while sustaining throughput [15].

* Corresponding author.

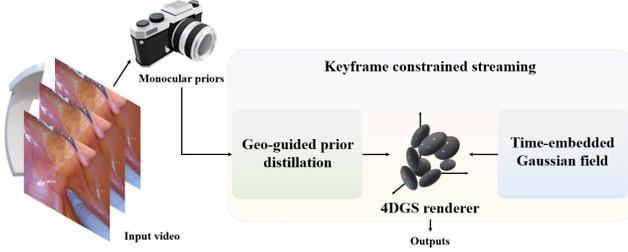


Fig. 1. Overview of ENDO-G²T.

Motivated by these observations, we propose ENDO-G²T, a geometry guided and temporally aware training scheme for time-embedded 4DGS. Our design directly targets the above failure modes with three modules that align with the identified requirements.

Thus, the key contributions are as follows. **(i)** ENDO-G²T introduces geo-guided prior distillation (GPD) that anchors early geometry via confidence-gated monocular priors with scale-invariant and depth-gradient supervision under a warm-up-to-cap schedule; **(ii)** it employs a time-embedded gaussian field (TEGF) with a rotor-like rotation parameterization and lightweight regularization to improve temporal coherence and reduce floaters; **(iii)** it adopts keyframe-constrained streaming (KCS) with a max-points budget to curb point growth, sustain throughput, and mitigate long-horizon drift; **(iv)** across EndoNeRF and StereoMIS-P1 datasets, ENDO-G²T achieves state-of-the-art monocular 4DGS results at high frame rates.

2. ENDO-G²T

ENDO-G²T comprises three modules: GPD that anchors geometry early with confidence-gated monocular depth under a warm-up-to-cap schedule, a TEGF that models dynamics in XYZT with a rotor-like rotation and lightweight regularization, and keyframe constrained streaming that preserves efficiency and long-horizon stability via keyframe-focused optimization and a global max-points budget.

2.1. Geo-guided prior distillation

This module injects appearance-agnostic supervision early by distilling pixel-aligned monocular geometry into the rendered depth. For each training view we use the observed RGB image I , the current rendering \hat{I} and rendered depth \hat{D} , together with external priors (D^*, C^*) where D^* is the monocular depth and $C^* \in [0, 1]$ is its confidence. Supervision is restricted to a valid pixel set selected by (i) a frame-adaptive confidence threshold and (ii) a physically plausible depth range; optional instrument masks further exclude tool regions. Prior-based losses are activated only when the fraction of valid pixels is sufficiently large (at least 10%) to avoid spurious supervision. Depths are normalized per frame and

compared in the log domain to achieve scale invariance, and a depth-gradient term aligns edge structure. Both prior terms are introduced with a warm-up-to-cap schedule so that the influence of priors grows gradually and remains bounded, mitigating early overfitting to noisy estimates.

Photometric reconstruction combines an L_1 term and SSIM,

$$\mathcal{L}_{\text{photo}} = (1 - \lambda_{\text{dssim}}) \|I - \hat{I}\|_1 + \lambda_{\text{dssim}} (1 - \text{SSIM}(I, \hat{I})). \quad (1)$$

where $\lambda_{\text{dssim}} \in [0, 1]$ balances the two components and SSIM is computed in a standard patchwise manner. We also define the valid pixel set by confidence, depth range, and an optional instrument mask:

$$\Omega_v = \left\{ p \text{ pixel} : C^*(p) \geq \tau, D_{\min} \leq D^*(p) \leq D_{\max}, \mathcal{M}_{\text{inst}}(p) = 0 \right\}. \quad (2)$$

where $\tau = \max\{0.01, 0.5 \max_q C^*(q)\}$ is an adaptive confidence threshold, D_{\min}, D_{\max} bound plausible depths, and $\mathcal{M}_{\text{inst}} \in \{0, 1\}^{H \times W}$ is an optional instrument mask. Prior-based terms are enabled only if $|\Omega_v|/(HW) \geq 0.1$.

To compare geometry without enforcing an absolute scale, depths are min to max normalized per frame in code. We then measure a scale invariant discrepancy in the log domain on Ω_v ,

$$\mathcal{L}_{\text{SILog}} = 10 \sqrt{\text{Var}_{p \in \Omega_v}(g(p)) + \beta \text{Mean}_{p \in \Omega_v}(g(p))^2}. \quad (3)$$

$$g(p) = \log(\tilde{D}(p) + \epsilon) - \log(\tilde{D}^*(p) + \epsilon). \quad (4)$$

where \tilde{D} and \tilde{D}^* are the per frame normalized versions of \hat{D} and D^* , $\epsilon > 0$ prevents singularities, and $\beta > 0$ controls the penalty on global log scale bias. To sharpen geometry at boundaries we align first order depth gradients on the valid set,

$$\mathcal{L}_{\text{grad}} = \frac{1}{|\Omega_v|} \sum_{p \in \Omega_v} \left(\|\nabla_x \hat{D}(p) - \nabla_x D^*(p)\|_1 + \|\nabla_y \hat{D}(p) - \nabla_y D^*(p)\|_1 \right). \quad (5)$$

where ∇_x and ∇_y are forward differences along the horizontal and vertical axes in pixels.

To avoid overfitting to prior noise at the beginning of training, the prior weights follow a warm up to cap schedule,

$$\lambda_{\text{SI}}(t) = \lambda_{\text{SI},0} \min\left(1, \frac{t}{T_{\text{warm}}}\right) w_{\text{max}}, \quad (6)$$

$$\lambda_{\nabla}(t) = \lambda_{\nabla,0} \min\left(1, \frac{t}{T_{\text{warm}}}\right) w_{\text{max}}.$$

where t is the global iteration, $T_{\text{warm}} > 0$ is the warm up length, $\lambda_{\text{SI},0}, \lambda_{\nabla,0} > 0$ are base coefficients, and $w_{\text{max}} \in$

$(0, 1]$ caps the effective strength of priors. The module objective is

$$\mathcal{L}_{\text{geo}} = \mathcal{L}_{\text{photo}} + \lambda_{\text{SI}}(t) \mathcal{L}_{\text{SILog}} + \lambda_{\nabla}(t) \mathcal{L}_{\text{grad}}. \quad (7)$$

All symbols in (1) to (7) are defined above and finite value checks are applied before accumulation to ensure numerical stability.

2.2. Time-embedded Gaussian field

This module represents dynamics by lifting Gaussian primitives into a time-embedded space and evolving their parameters smoothly. For each primitive we maintain a 3D center, a diagonal scale, a rotation, an opacity, and spherical-harmonic color coefficients. The covariance used for splatting is

$$\Sigma_i(t) = R_i(t) S_i^2(t) R_i^\top(t), \quad (8)$$

where $\mu_i(t) \in \mathbb{R}^3$ is the center, $S_i(t)$ contains axis scales, $R_i(t) \in \text{SO}(3)$ is the rotation, and thus $\Sigma_i(t) \in \mathbb{R}^{3 \times 3}$ is positive definite by construction; $\alpha_i(t) \in (0, 1)$ is opacity and $\mathbf{c}_i(t) \in \mathbb{R}^{(n+1)^2 \times 3}$ holds spherical-harmonic coefficients up to degree n . Positions later advance with a per-primitive velocity and rotations are updated by a minimal rotor operator; the corresponding variables $v_i(t)$, $\rho_i(t)$ and the frame interval $\Delta t > 0$ are introduced in the subsequent evolution equations.

To bias toward decisive visibility and coherent motion with minimal overhead we include two lightweight stabilizers. Opacity entropy encourages $\alpha_i(t)$ to concentrate near zero or one,

$$\mathcal{L}_{\text{ent}} = -\frac{1}{N} \sum_{i=1}^N \left(\alpha_i(t) \log \alpha_i(t) + (1 - \alpha_i(t)) \log (1 - \alpha_i(t)) \right). \quad (9)$$

where N is the number of active primitives at time t . Local velocity coherence smooths motion in a joint space and time neighborhood,

$$\mathcal{L}_{\text{vel}} = \frac{1}{N} \sum_{i=1}^N \frac{1}{|\mathcal{N}_k(i, t)|} \sum_{j \in \mathcal{N}_k(i, t)} \left\| \underbrace{\mu_i(t) - \mu_i(t - \Delta t)}_{v_i(t)} - \underbrace{\mu_j(t) - \mu_j(t - \Delta t)}_{v_j(t)} \right\|_1. \quad (10)$$

where $\mathcal{N}_k(i, t)$ denotes the k nearest neighbors of primitive i in a joint position–time metric, where $k \in \mathbb{N}$ is a fixed neighborhood size. When N is large, the inner sum is evaluated on a subsample to respect memory limits. All symbols in (8) to (10) are defined above.

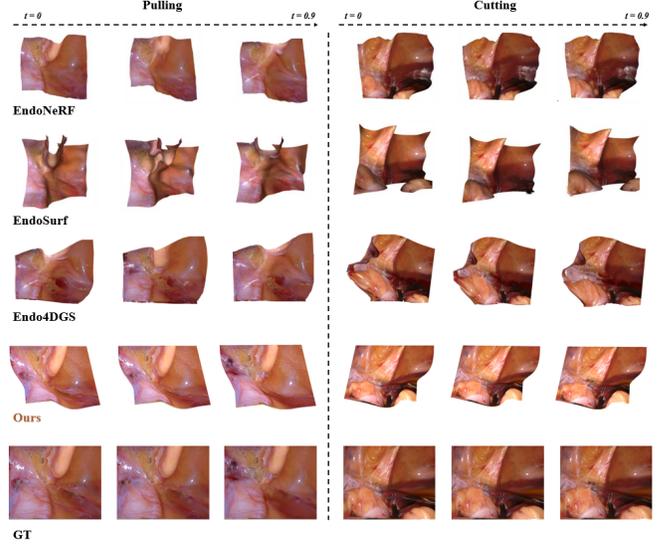


Fig. 2. Comparison on EndoNeRF Cutting/Pulling: Endo-G²T produces sharper tissue boundaries and fewer floaters than other baselines.

2.3. Keyframe constrained streaming

This module maintains stability and throughput on long sequences by interleaving full refinement with lightweight updates. The video is partitioned by a stride $w \in \mathbb{N}$ into keyframes \mathcal{K} and candidate frames \mathcal{C} , where

$$\begin{aligned} \mathcal{K} &= \{ f \in \{1, \dots, F\} : f \equiv 1 \pmod{w} \}, \\ \mathcal{C} &= \{1, \dots, F\} \setminus \mathcal{K}. \end{aligned} \quad (11)$$

Keyframes run full optimization with densification and pruning, while candidate frames apply lightweight image-space updates that keep the model size fixed. To prevent uncontrolled growth, the active set of Gaussians obeys a global budget at every time t :

$$\begin{aligned} |\mathcal{G}_t| &\leq G_{\text{max}}, \\ &\text{for all } t \text{ associated with frames } 1, \dots, F. \end{aligned} \quad (12)$$

At keyframes the budget is enforced by balancing additions and removals, which curbs long-horizon drift and preserves frame rate. For throughput reporting, a raster-only mode disables input and output so that measured frames per second reflect GPU rasterization alone.

3. EXPERIMENTS

We evaluate ENDO-G²T on EndoNeRF (cutting, pulling) [8] and StereoMIS [16]. For EndoNeRF, we use a 7:1 train–validation split and report PSNR, SSIM, and LPIPS on held-out views. For StereoMIS, following the Endo-4DGS [10] protocol, we reconstruct frames 800–1000 from

Table 1. Quantitative comparison on EndoNeRF (cutting, pulling) and StereoMIS-P1 (frames 800–1000). Best results are in bold. FPS is measured in raster-only mode.

Models	EndoNeRF–Cutting			EndoNeRF–Pulling			StereoMIS–P1 (800–1000)			FPS \uparrow
	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	
EndoNeRF [8]	35.840	0.942	0.057	35.430	0.939	0.064	30.500	0.880	0.078	0.2
EndoSurf [9]	34.893	0.952	0.107	34.910	0.955	0.124	31.221	0.890	0.071	0.04
Endo-4DGS [10]	36.165	0.959	0.039	37.014	0.960	0.041	32.188	0.898	0.066	100
ST-Endo4DGS [11]	39.290	0.973	0.016	38.280	0.966	0.024	32.900	0.905	0.060	123
Endo-G²T (ours)	40.080	0.982	0.007	38.290	0.970	0.016	33.580	0.914	0.056	148

Table 2. Ablation on KCS. “KF” denotes keyframe scheduling, “w” is the keyframe stride, “KF-only” optimizes keyframes with candidates frozen, “w/o KF (w=1)” disables keyframe scheduling, and “budget” is the global max-points cap G_{\max} .

Variants	EndoNeRF–Cutting			EndoNeRF–Pulling			StereoMIS–P1 (800–1000)		
	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
Endo-G ² T (w/ KF, w/o budget)	38.912	0.971	0.012	37.041	0.962	0.020	32.217	0.902	0.061
Endo-G ² T (KF-only, w=5)	39.146	0.975	0.010	37.498	0.965	0.019	32.803	0.907	0.059
Endo-G ² T (w/o KF, w=1)	39.867	0.979	0.008	38.061	0.968	0.016	33.012	0.910	0.057
Endo-G ² T (w/ KF, w=10)	39.722	0.978	0.009	38.003	0.967	0.016	33.164	0.911	0.057
Endo-G ² T (w/ KF, w=3)	39.954	0.981	0.007	38.241	0.970	0.016	33.471	0.914	0.057
Endo-G²T (w/ KCS, w/ budget, w=5)	40.080	0.982	0.007	38.290	0.970	0.016	33.580	0.914	0.056

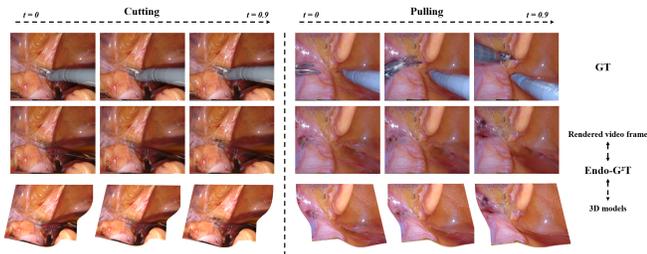


Fig. 3. Endo-G²T on EndoNeRF Cutting/Pulling from $t=0$ to $t=0.9$: GT, our renders, and reconstructed surfaces, showing coherent geometry and crisp opacity.

scene P1 in a monocular setting to ensure fair comparison. All models are trained in PyTorch with Adam on a single NVIDIA RTX 4090, using a learning rate of 1.6×10^{-3} , mixed photometric supervision (L1+SSIM), and our warm-up-to-cap geometry-prior schedule. Inference FPS is measured in a raster-only mode that disables I/O to reflect pure GPU rasterization.

Table 1 and Figures 2–3 jointly show quantitative and qualitative gains. On EndoNeRF–Cutting, Endo-G²T achieves 40.080 PSNR, 0.982 SSIM, and 0.007 LPIPS, improving over ST-Endo4DGS by 0.790 PSNR (2.0% relative), +0.009 SSIM, and a 56.3% LPIPS reduction. On EndoNeRF–Pulling, our method reaches 38.290/0.970/0.016, surpassing Endo-4DGS by +3.45% PSNR, +0.010 SSIM, and 61.0% lower LPIPS. On StereoMIS–P1, we obtain 33.580/0.914/0.056 with a +2.1% PSNR gain over ST-Endo4DGS and a 6.7% LPIPS drop. The qualitative frames in Figure 2 show sharper tissue boundaries and fewer floaters, while Figure 3 illustrates temporally

coherent geometry and crisp opacity across time. Raster-only throughput rises to 148 FPS, a 20.3% increase over ST-Endo4DGS, indicating that accuracy gains come without sacrificing speed.

For ablations, the core modeling components (GPD and TEGF) remain fixed to ensure backbone fairness; we vary only the system schedule in Table 2. Removing the global budget while keeping keyframes degrades Cutting PSNR from 40.080 to 38.912 (2.9%) and raises LPIPS from 0.007 to 0.012 (71% higher), reflecting point explosion and instability. Freezing candidates (KF-only) saves compute but trails the full model across datasets, and disabling keyframes (w=1) recovers some accuracy yet loses the periodic re-anchoring benefits. Varying the stride reveals a sweet spot at $w=5$ yielding slightly lower PSNR/SSIM or higher LPIPS. The full KCS configuration (budget + $w=5$) is consistently best, supporting the proposed efficiency–stability trade-off.

4. CONCLUSION

We presented ENDO-G²T, a geometry-guided and temporally aware training scheme for time-embedded 4D Gaussian splatting in endoscopy. The method unifies geo-guided prior distillation, a time-embedded Gaussian field with rotor-based evolution, and keyframe-constrained streaming with a max-points budget. Experiments on EndoNeRF and StereoMIS–P1 show state-of-the-art accuracy with high throughput, and ablations verify the importance of the scheduling strategy. Future work will study cross-view prior calibration, uncertainty-aware fusion, and adaptive keyframe selection for clinical deployment.

5. REFERENCES

- [1] Zhuoyue Yang, Ju Dai, and Junjun Pan, “3d reconstruction from endoscopy images: A survey,” *Computers in biology and medicine*, vol. 175, pp. 108546, 2024.
- [2] Oscar G Grasa, Ernesto Bernal, Santiago Casado, Ismael Gil, and JMM Montiel, “Visual slam for handheld monocular endoscope,” *IEEE transactions on medical imaging*, vol. 33, no. 1, pp. 135–146, 2013.
- [3] Ying Zhou, Shiquan He, Hao Wang, Fan Huang, Mei Liu, Qiang Li, and Zhiwei Wang, “Improved self-supervised monocular endoscopic depth estimation based on pose alignment-friendly dynamic view selection,” in *2024 IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*. IEEE, 2024, pp. 3005–3012.
- [4] Liang Qiu and Hongliang Ren, “Endoscope navigation with slam-based registration to computed tomography for transoral surgery,” *International Journal of Intelligent Robotics and Applications*, vol. 4, no. 2, pp. 252–263, 2020.
- [5] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng, “Nerf: Representing scenes as neural radiance fields for view synthesis,” *Communications of the ACM*, vol. 65, no. 1, pp. 99–106, 2021.
- [6] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkuehler, and George Drettakis, “3d gaussian splatting for real-time radiance field rendering,” *ACM Trans. Graph.*, vol. 42, no. 4, July 2023.
- [7] Guanjun Wu, Taoran Yi, Jiemin Fang, Lingxi Xie, Xiaopeng Zhang, Wei Wei, Wenyu Liu, Qi Tian, and Xinggang Wang, “4d gaussian splatting for real-time dynamic scene rendering,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2024, pp. 20310–20320.
- [8] Yuehao Wang, Yonghao Long, Siu Hin Fan, and Qi Dou, “Neural rendering for stereo 3d reconstruction of deformable tissues in robotic surgery,” in *International conference on medical image computing and computer-assisted intervention*. Springer, 2022, pp. 431–441.
- [9] Ruyi Zha, Xuelian Cheng, Hongdong Li, Mehrtash Harandi, and Zongyuan Ge, “Endosurf: Neural surface reconstruction of deformable tissues with stereo endoscope videos,” in *International conference on medical image computing and computer-assisted intervention*. Springer, 2023, pp. 13–23.
- [10] Yiming Huang, Beilei Cui, Long Bai, Ziqi Guo, Mengya Xu, Mobarakol Islam, and Hongliang Ren, “Endo-4dgs: Endoscopic monocular scene reconstruction with 4d gaussian splatting,” in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2024, pp. 197–207.
- [11] Fengze Li, Jishuai He, Jieming Ma, and Zhijing Wu, “Real-time spatio-temporal reconstruction of dynamic endoscopic scenes with 4d gaussian splatting,” in *2025 IEEE 22nd International Symposium on Biomedical Imaging (ISBI)*. IEEE, 2025, pp. 1–5.
- [12] Jianyuan Wang, Minghao Chen, Nikita Karaev, Andrea Vedaldi, Christian Rupprecht, and David Novotny, “Vggt: Visual geometry grounded transformer,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2025, pp. 5294–5306.
- [13] Dong Zhuo, Wenzhao Zheng, Jiahe Guo, Yuqi Wu, Jie Zhou, and Jiwen Lu, “Streaming 4d visual geometry transformer,” *arXiv preprint arXiv:2507.11539*, 2025.
- [14] Yuanxing Duan, Fangyin Wei, Qiyu Dai, Yuhang He, Wenzheng Chen, and Baoquan Chen, “4d-rotor gaussian splatting: towards efficient novel view synthesis for dynamic scenes,” in *ACM SIGGRAPH 2024 Conference Papers*, 2024, pp. 1–11.
- [15] Jinbo Yan, Rui Peng, Zhiyan Wang, Luyang Tang, Jiayu Yang, Jie Liang, Jiahao Wu, and Ronggang Wang, “Instant gaussian stream: Fast and generalizable streaming of dynamic scene reconstruction via gaussian splatting,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2025, pp. 16520–16531.
- [16] Michel Hayoz, Christopher Hahne, Mathias Gallardo, Daniel Candinas, Thomas Kurmann, Maximilian Allan, and Raphael Sznitman, “Learning how to robustly estimate camera pose in endoscopic videos,” *International journal of computer assisted radiology and surgery*, vol. 18, no. 7, pp. 1185–1192, 2023.