

Towards the next generation of image guidance for endoscopic procedures

CVPR Workshop on 3D Computer Vision in Medical Environments

June 16th 2019

Mathias Unberath, PhD

Assistant Research Professor

Department of Computer Science

Johns Hopkins University



Xingtong Liu
Graduate Student
Department of Computer Science



Ayushi Sinha, PhD
Assistant Research Scientist
Computational Sensing and Robotics



Russell H Taylor, PhD
John C. Malone Professor
Department of Computer Science



Gregory Hager, PhD
Mandell Bellmore Professor
Department of Computer Science



Masaru Ishii, MD
Associate Professor
Department of Otolaryngology

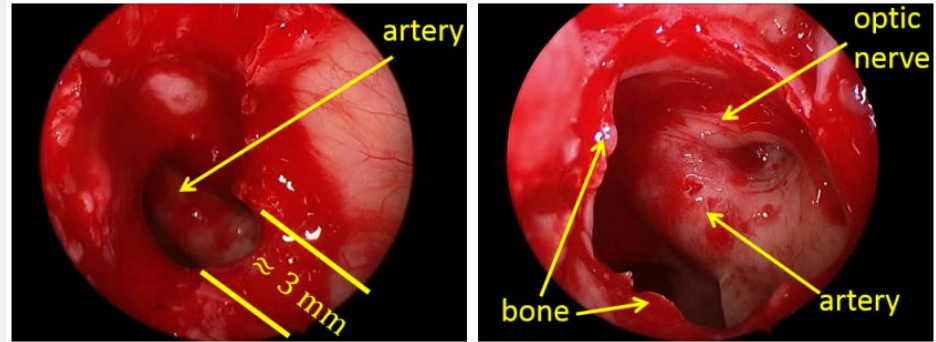
Some Background: Clinical and Technical

Navigating Sinus Surgery



Endoscopic Sinus Surgery

- Functional sinus surgery
 - Close proximity to critical structures
 - Surgical navigation desired



Challenges of Conventional Navigation

- Patient-specific 3D model of anatomy
 - Pre-operative (potentially outdated)
 - Obtained from CT scan (usually)
- Intra-operative registration: Optical tracking
 - CT to marker (via surface digitization)
 - Endoscope / tool to anatomy
 - Line of sight constraints
 - Visualization on model

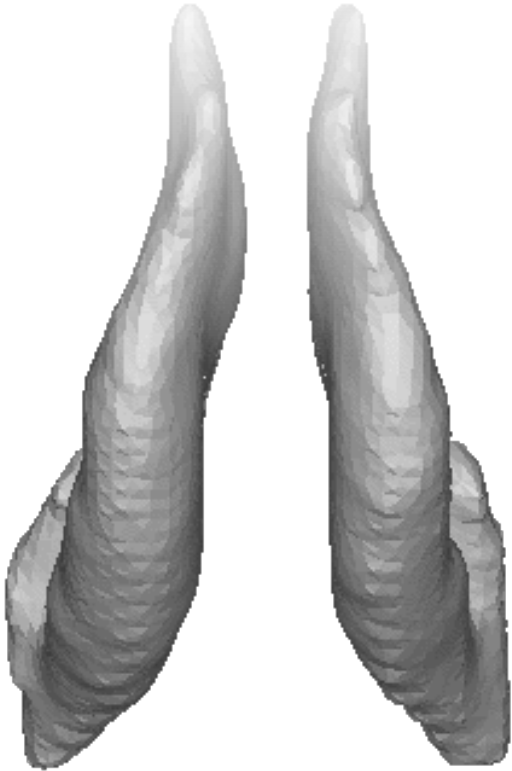


Challenges of Conventional Navigation

- Patient-specific 3D model of anatomy
 - Pre-operative (potentially outdated)
 - Obtained from CT scan (usually)
- Intra-operative registration: Optical tracking
 - CT to marker (via surface digitization)
 - Endoscope / tool to anatomy
 - Line of sight constraints
 - Visualization on model
- Observations
 - Complex setups **increase procedure time**
 - Disruptive workflows **promote frustration**
 - **Where to innovate?**



Step 1: Navigating in the Absence of CT



- Patient-specific 3D model of anatomy
 - ~~Pre-operative (potentially outdated)~~
 - ~~Obtained from CT scan (usually)~~
 - **Population-derived atlas of sinus anatomy**
- Intra-operative registration: Optical tracking
 - ~~CT to marker (via surface digitization)~~
 - **Model to video registration**
 - Endoscope / tool to anatomy
 - Line of sight constraints
 - Visualization on model

Step 2: Navigating Without Prior Information

~~Pre-operative specific 3D model of anatomy~~

~~Pre-operative (potentially outdated)~~

~~Obtained from CT scan (usually)~~

Reconstructed from endoscopy sequence

~~Pre-operative registration: Optical tracking~~

~~CT to marker (via surface digitization)~~

~~Endoscope / tool to anatomy~~

~~Line of sight constraints~~

~~Visualization on model~~

Everything relative to endoscopy

Towards Next-generation Image Guidance

Navigating in the Absence of CT



Building the Population-based Model

- Build statistical shape models
 - **Principal component analysis**
 - Capture anatomical variation
- Given shapes, $\mathbf{V} = [\mathbf{v}_1 \ \mathbf{v}_2 \ \dots \ \mathbf{v}_{n_s}]^T$ with correspondences, we can compute:

$$\text{Mean: } \bar{\mathbf{V}} = \frac{1}{n_s} \sum_{i=1}^{n_s} \mathbf{V}_i$$

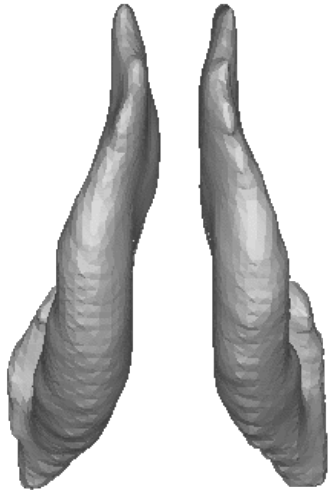
$$\text{Variance: } \Sigma = \frac{1}{n_s} \sum_{i=1}^{n_s} (\mathbf{V}_i - \bar{\mathbf{V}})(\mathbf{V}_i - \bar{\mathbf{V}})^T$$

$$\Sigma = [\mathbf{m}_1 \ \dots \ \mathbf{m}_{n_s}] \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_{n_s} \end{bmatrix} [\mathbf{m}_1 \ \dots \ \mathbf{m}_{n_s}]^T$$

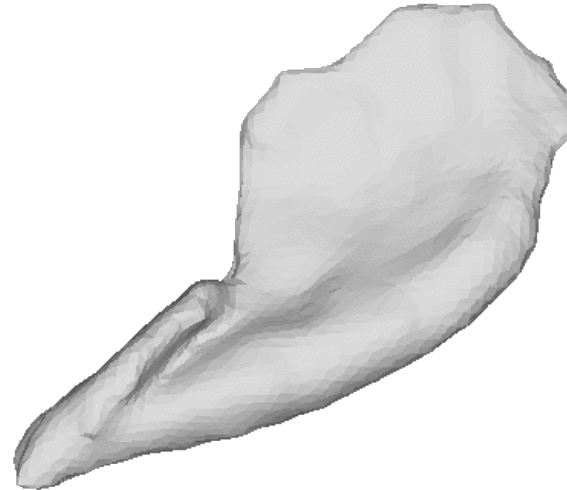
Building the Population-based Model

- Build statistical shape models
 - Principal component analysis
 - **Capture anatomical variation (middle turbinate)**

mean +0.00 std. dev.



mean +0.00 std. dev.



Estimating Patient Anatomy

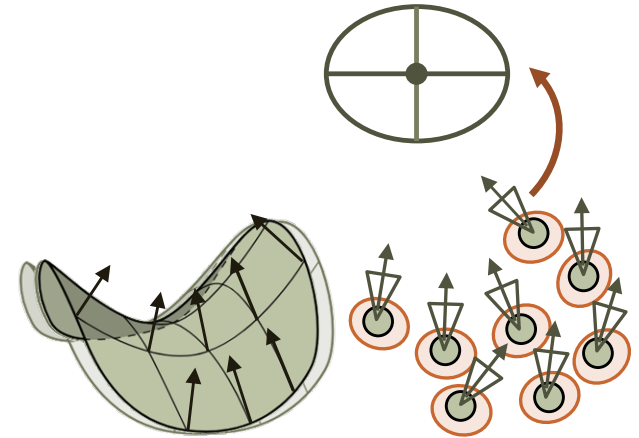
- Deformable registration
 - Optimize shape model model parameters
 - Align with endoscopic video
- Given a new shape \mathbf{V}^* , we can compute:

$$\text{Weights: } s_i = \mathbf{w}_i^T (\mathbf{V}^* - \bar{\mathbf{V}}) \quad \text{Estimated shape: } \tilde{\mathbf{V}}^* = \bar{\mathbf{V}} + \sum_{i=1}^{n_m} s_i \mathbf{w}_i$$
$$\mathbf{w}_i = \sqrt{\lambda_i} \mathbf{m}_i$$

Estimating Patient Anatomy

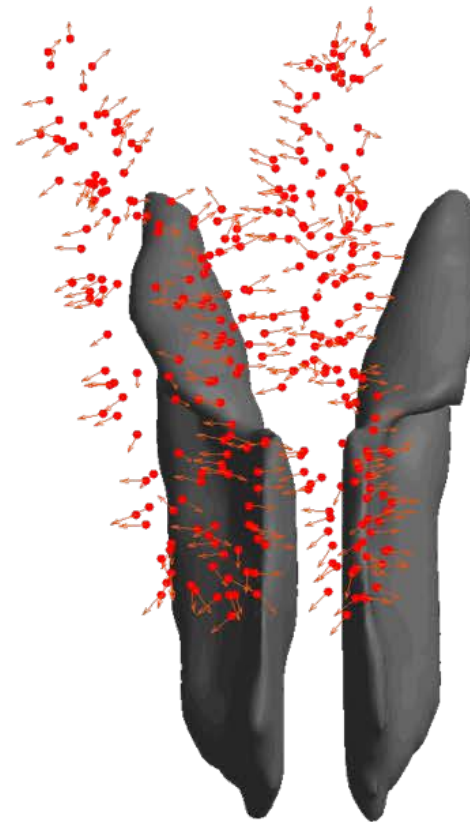
- Deformable registration
 - Optimize shape model model parameters
 - Align with endoscopic video
- Simultaneously, align rigidly

Can be solved with the
Generalized Deformable Most Likely Oriented Point (GD-IMLOP)
algorithm



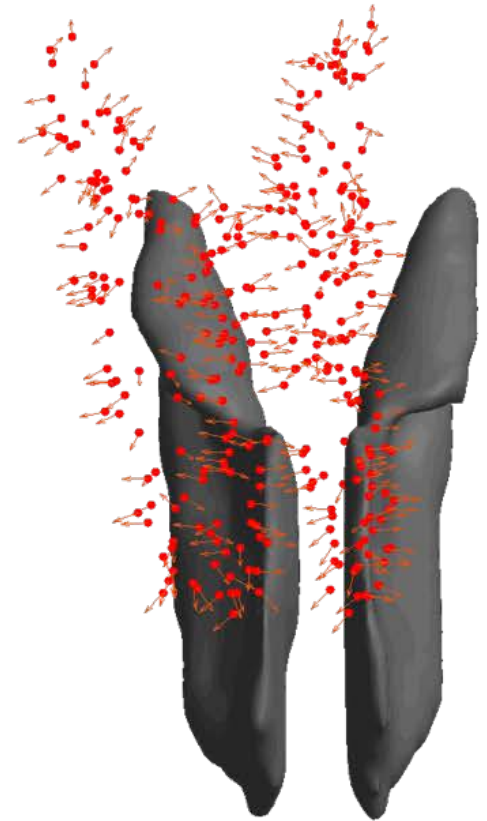
Estimating Patient Anatomy

- Deformable registration
 - Optimize shape model model parameters
 - Align with endoscopic video
- Simultaneous deformable and rigid alignment to unseen shape V^*
- Great!



Estimating Patient Anatomy

- Deformable registration
 - Optimize shape model model parameters
 - Align with endoscopic video
- Simultaneous deformable and rigid alignment to unseen shape V^*
- Great!
- **But wait ...**
Where do we get the new shape from?
How does this link to endoscopy?



Estimating Patient Anatomy

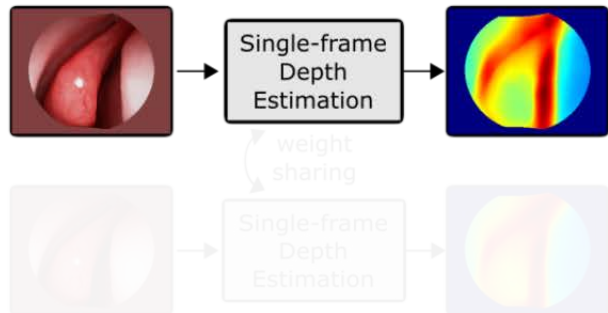
- Deformable registration
 - Optimize shape model model parameters
 - Align with endoscopic video
- Estimating unseen shapes V^* from endoscopic video



... some AI maybe?



Training Phase



This is what we are after here

Endoscopic image in → Depth map out

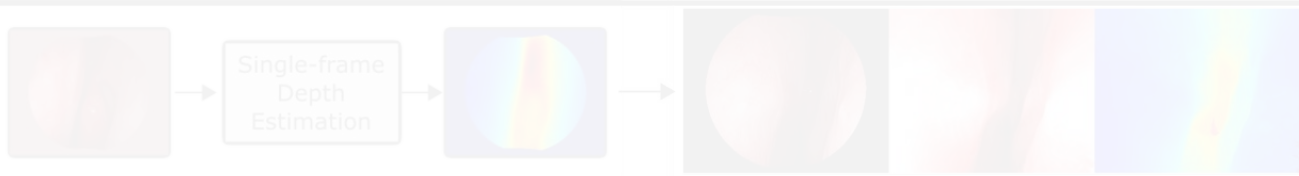
ConvNets are trained via backpropagation

→ Need informative gradients

→ Consequently, need informative loss

→ **How to supervise learning?**

Application Phase

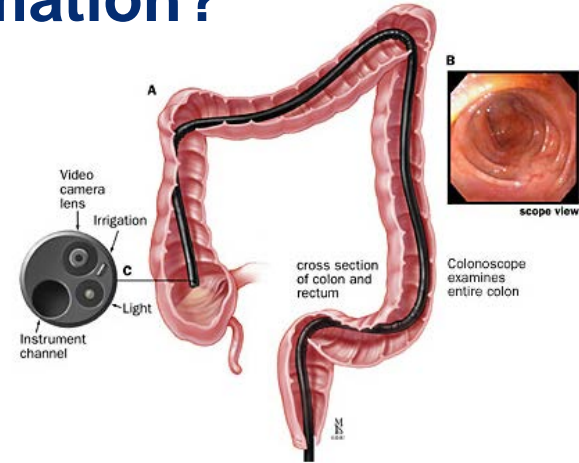
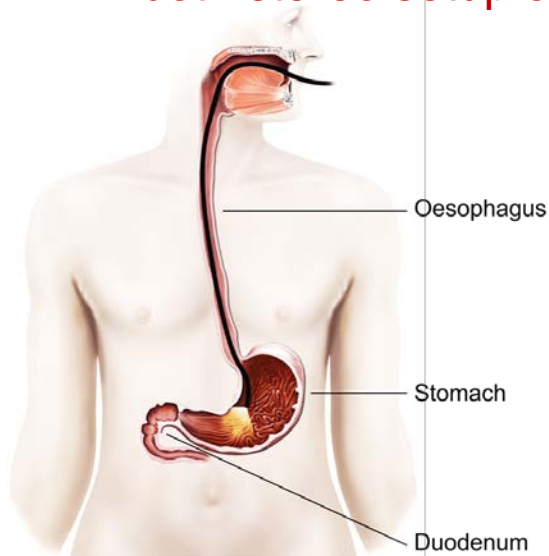




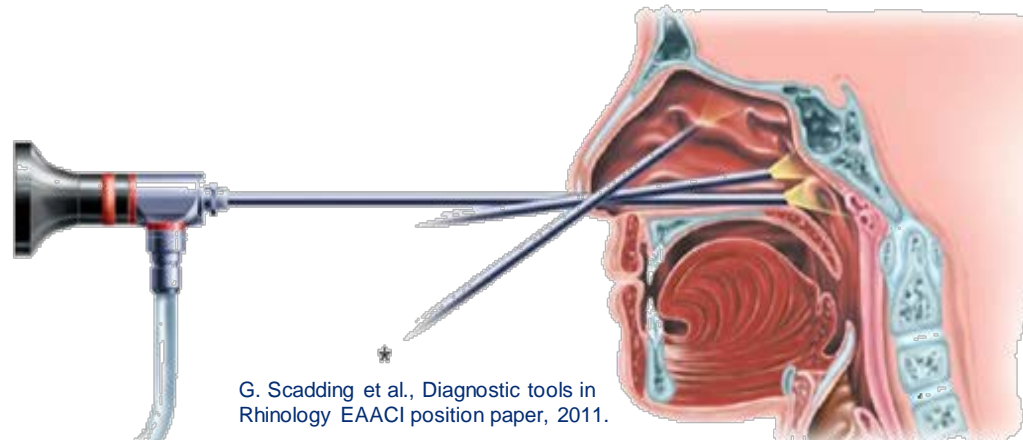
How to supervise monocular depth estimation?

Remembering the application: **Endoscopy**

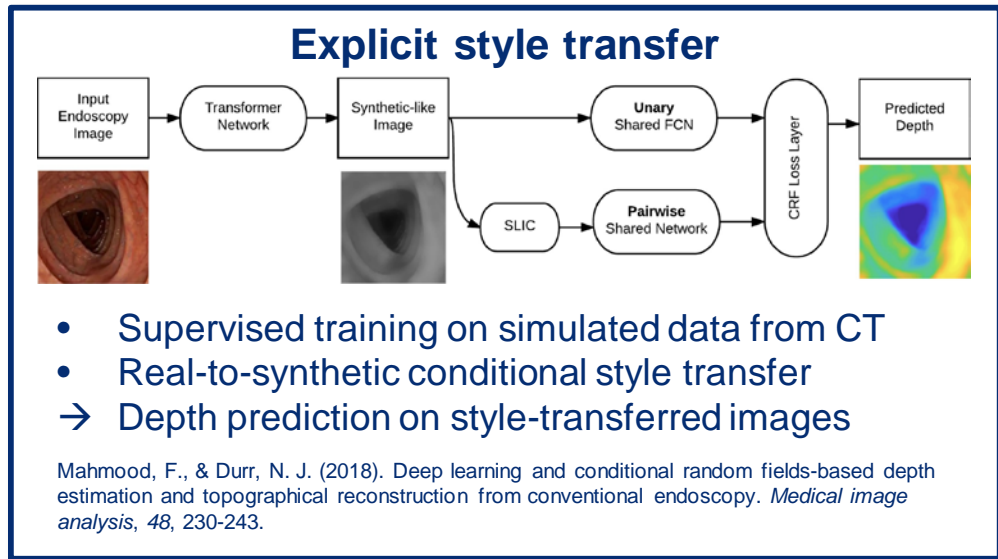
- **Miniaturized equipment** to inspect difficult to access anatomy
- **Prohibitively disruptive to install dedicated hardware, both stereo setup or depth sensing**



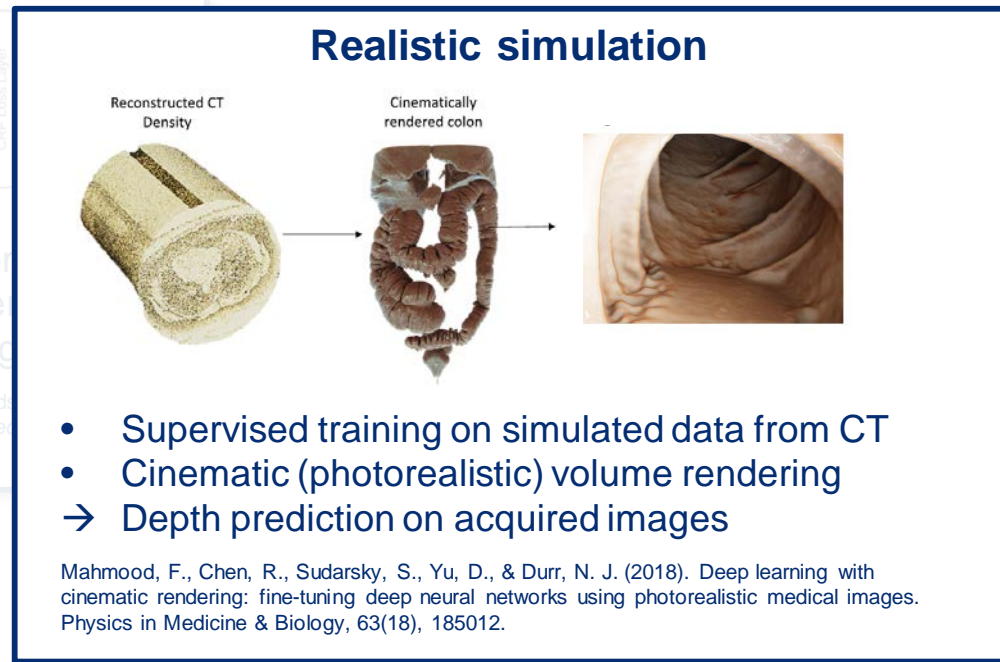
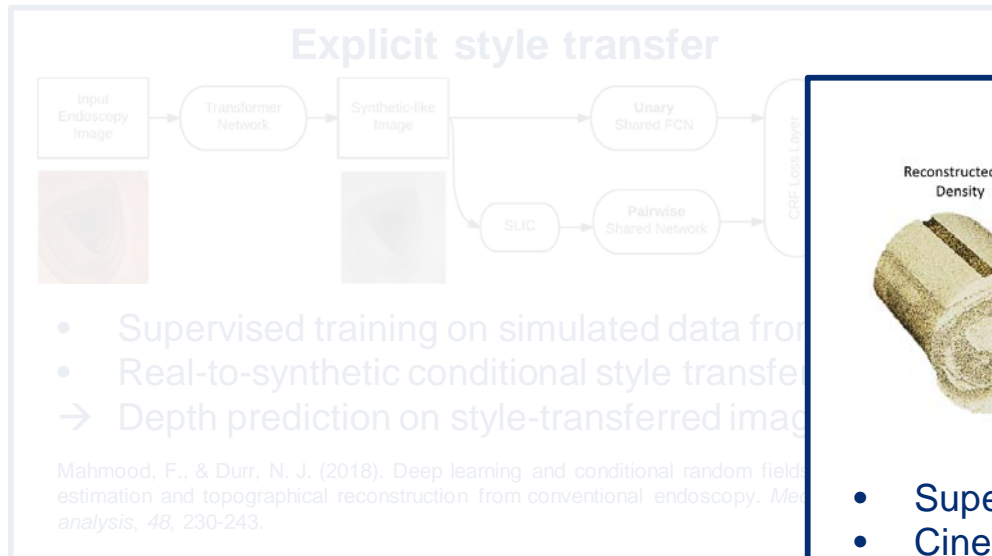
<http://www.alfasurgerycenter.com/procedures.html>



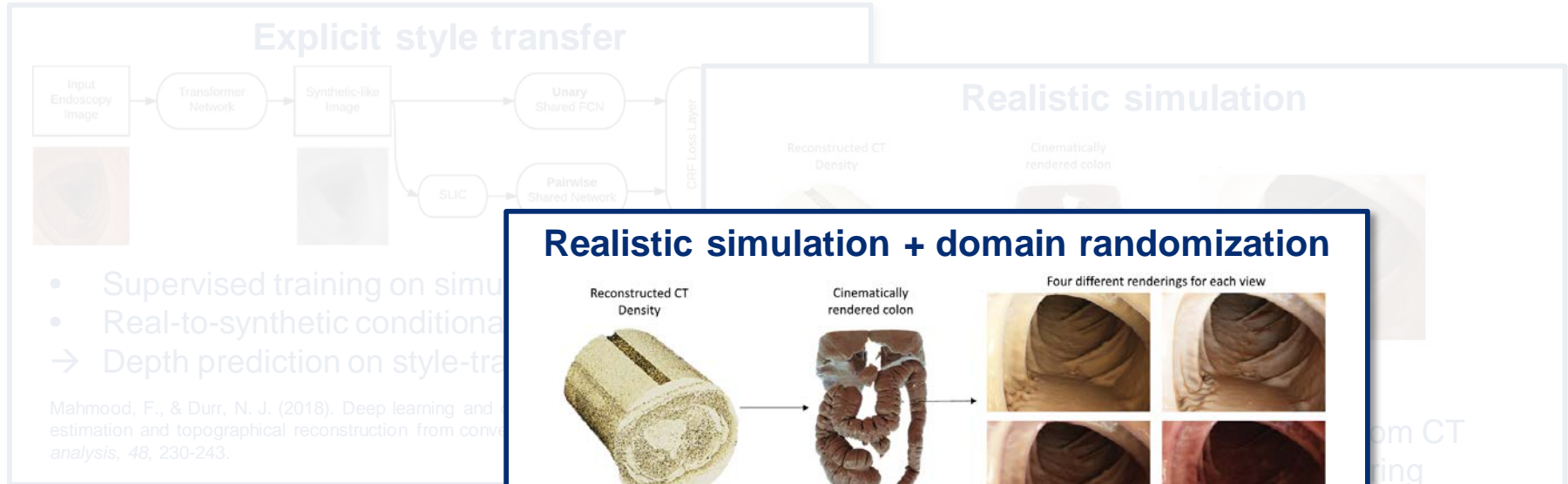
How to supervise monocular depth estimation?



How to supervise monocular depth estimation?



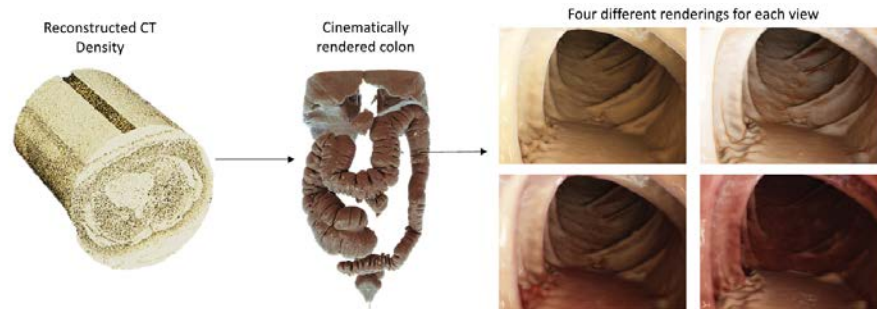
How to supervise monocular depth estimation?



- Supervised training on simulated data
 - Real-to-synthetic conditional GAN
- Depth prediction on style-transferred images

Mahmood, F., & Durr, N. J. (2018). Deep learning with cinematic rendering: fine-tuning deep neural networks using photorealistic medical images. *Physics in Medicine & Biology*, 48, 230-243.

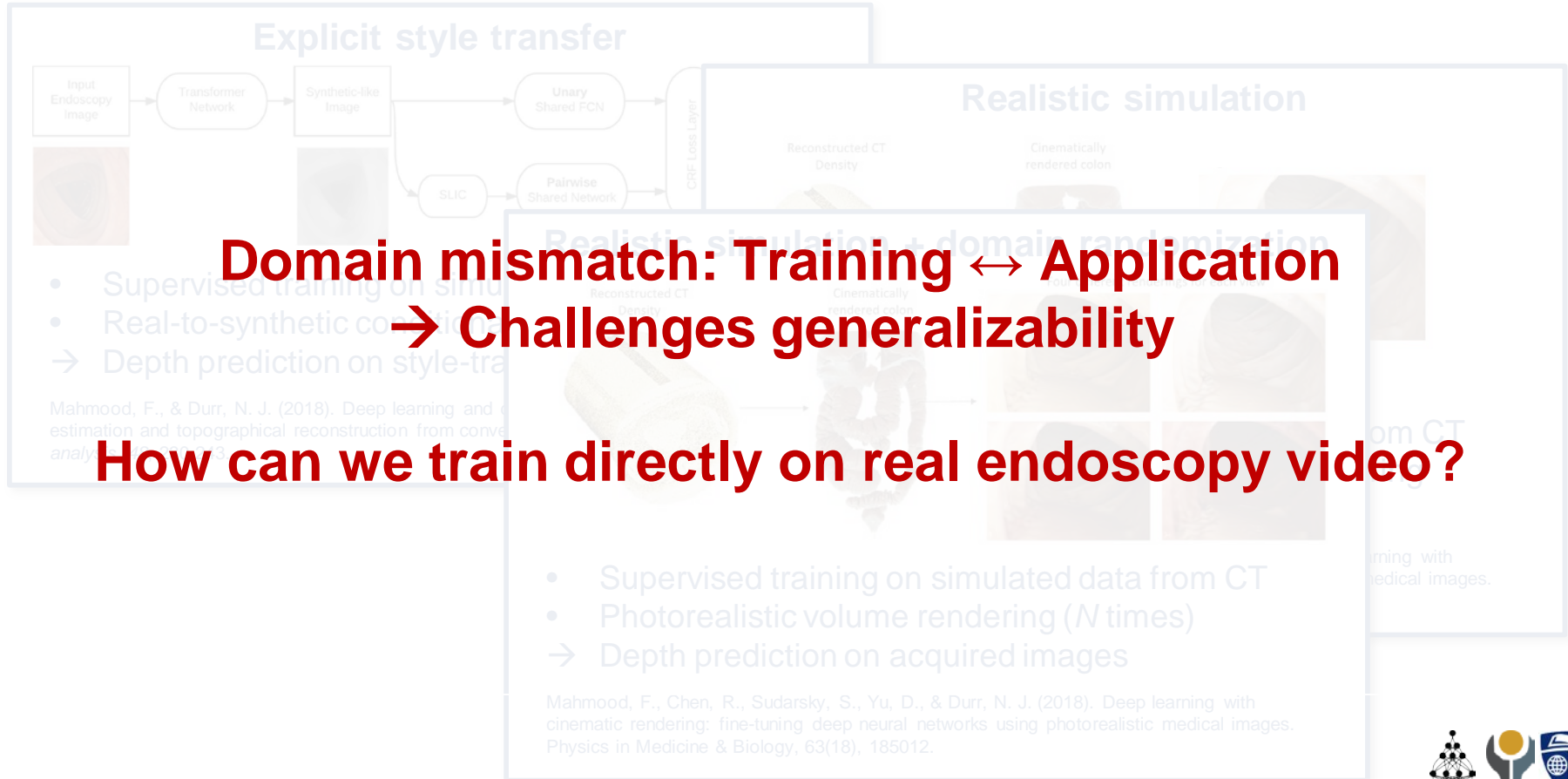
Realistic simulation + domain randomization



- Supervised training on simulated data from CT
 - Photorealistic volume rendering (N times)
- Depth prediction on acquired images

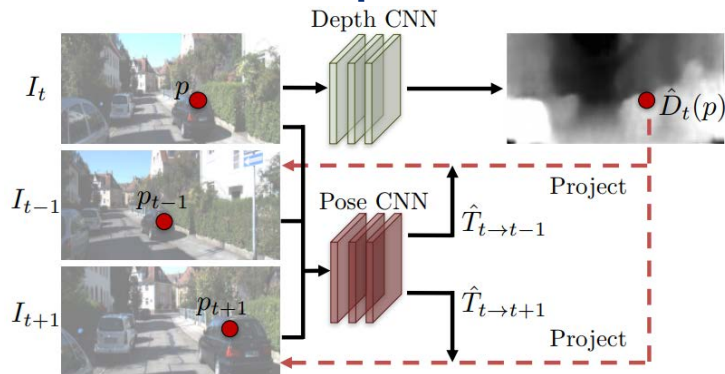
Mahmood, F., Chen, R., Sudarsky, S., Yu, D., & Durr, N. J. (2018). Deep learning with cinematic rendering: fine-tuning deep neural networks using photorealistic medical images. *Physics in Medicine & Biology*, 63(18), 185012.

How to supervise monocular depth estimation?



How to supervise monocular depth estimation?

Self-supervision



- Predict depth on target, synthesize neighbor views
- Photometric reconstruction loss for training
- Self-supervision, directly on acquired video

Zhou, T., Brown, M., Snavely, N., & Lowe, D. G. (2017). Unsupervised learning of depth and ego-motion from video. In Proceedings of the IEEE CVPR (pp. 1851-1858).

→ Depth prediction **Does this work for endoscopy?**

Mahmood, F., Chen, R., Sudarsky, S., Yu, D., & Durr, N. J. (2018). Deep learning with cinematic rendering: fine-tuning deep neural networks using photorealistic medical images. *Physics in Medicine & Biology*, 63(18), 185012.



How to supervise monocular depth estimation?

Classical – Structure from Motion



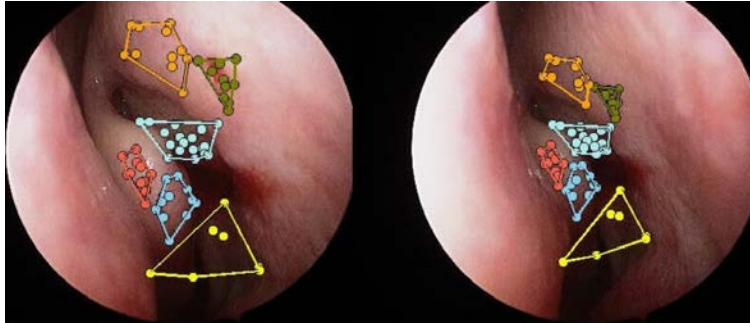
- Feature matching
 - Triangulation and bundle adjustment
- Reconstruction from acquired images

Snaveley, N., Seitz, S. M., & Szeliski, R. (2006, July). Photo tourism: exploring photo collections in 3D. In ACM transactions on graphics (TOG) (Vol. 25, No. 3, pp. 835-846). ACM.

Does this work for endoscopy?

How to supervise monocular depth estimation?

Classical – Structure from Motion

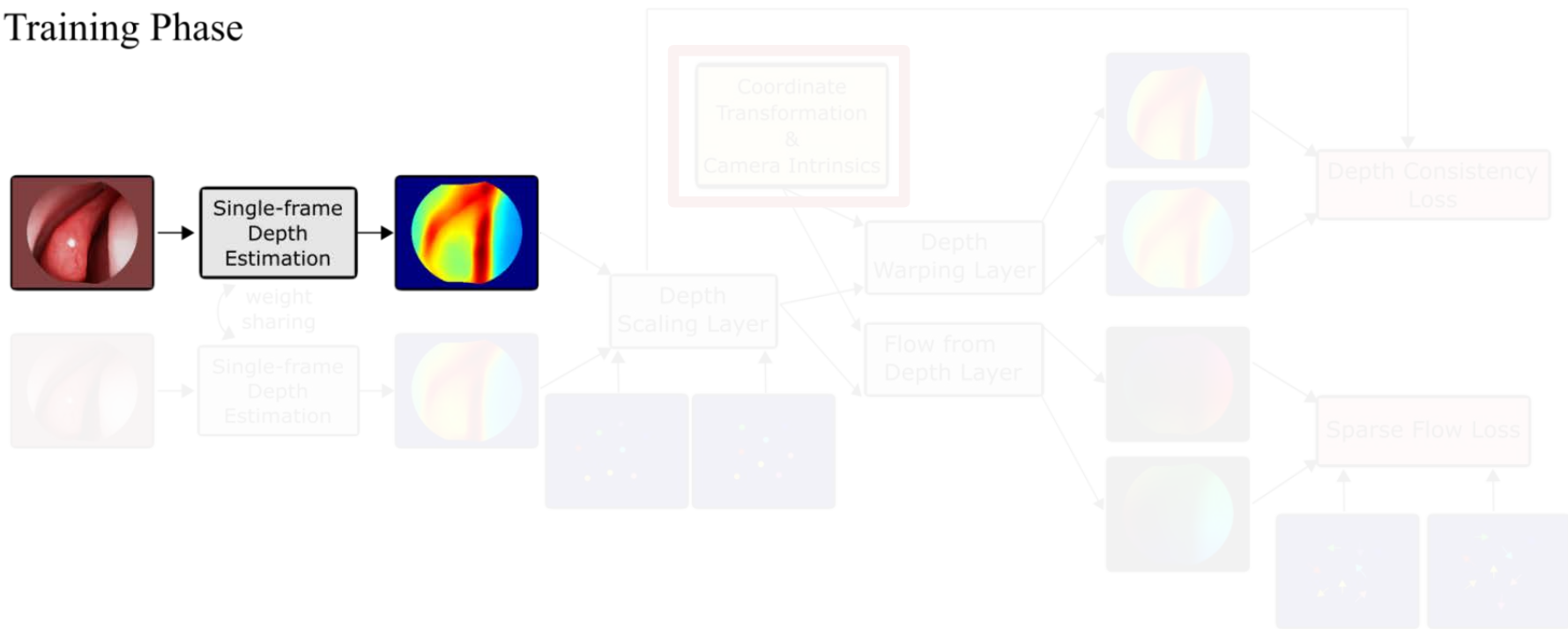


- SURF feature matching, hierarchical refinement
 - Triangulation and bundle adjustment
- Reconstruction from acquired images (**sparse**)

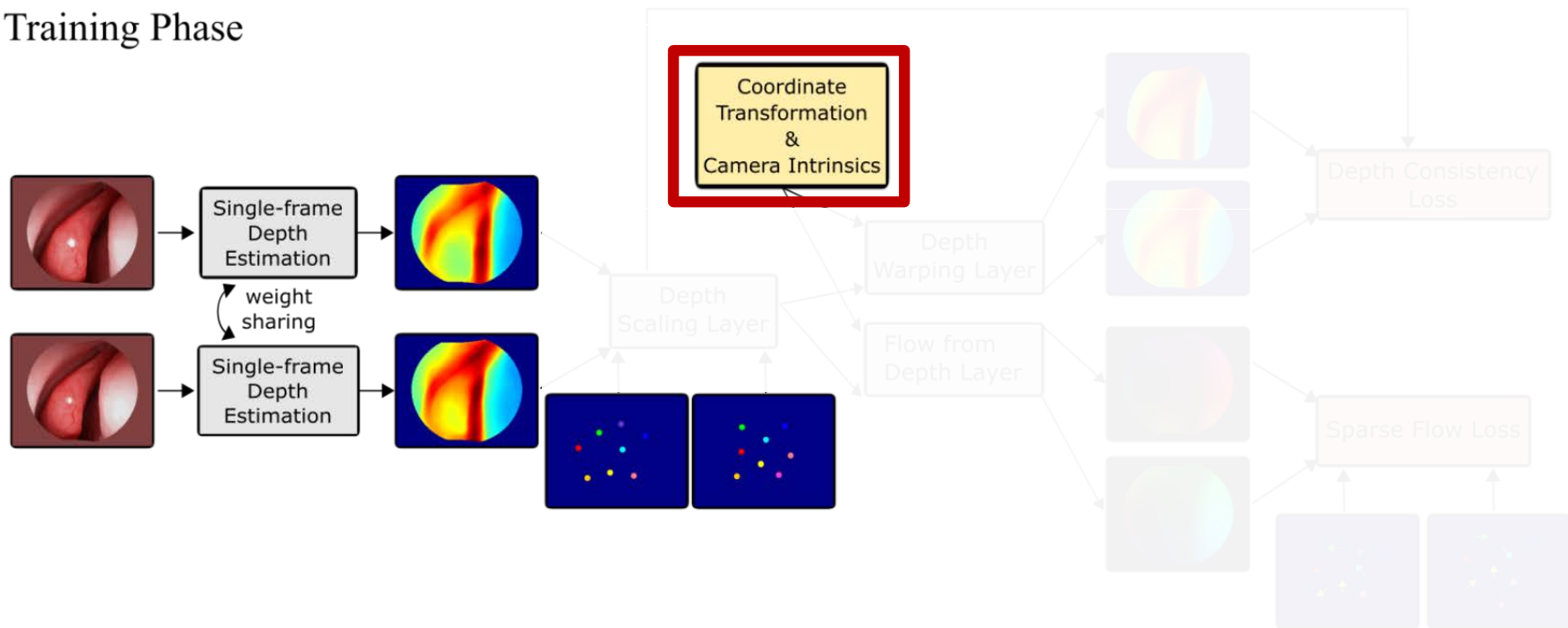
Leonard, S., Reiter, A., Sinha, A., Ishii, M., Taylor, R. H., & Hager, G. D. (2016, March). Image-based navigation for functional endoscopic sinus surgery using structure from motion. In Medical Imaging 2016: Image Processing (Vol. 9784, p. 97840V).

Yes(-ish).
So let's use this, then!

Training Phase



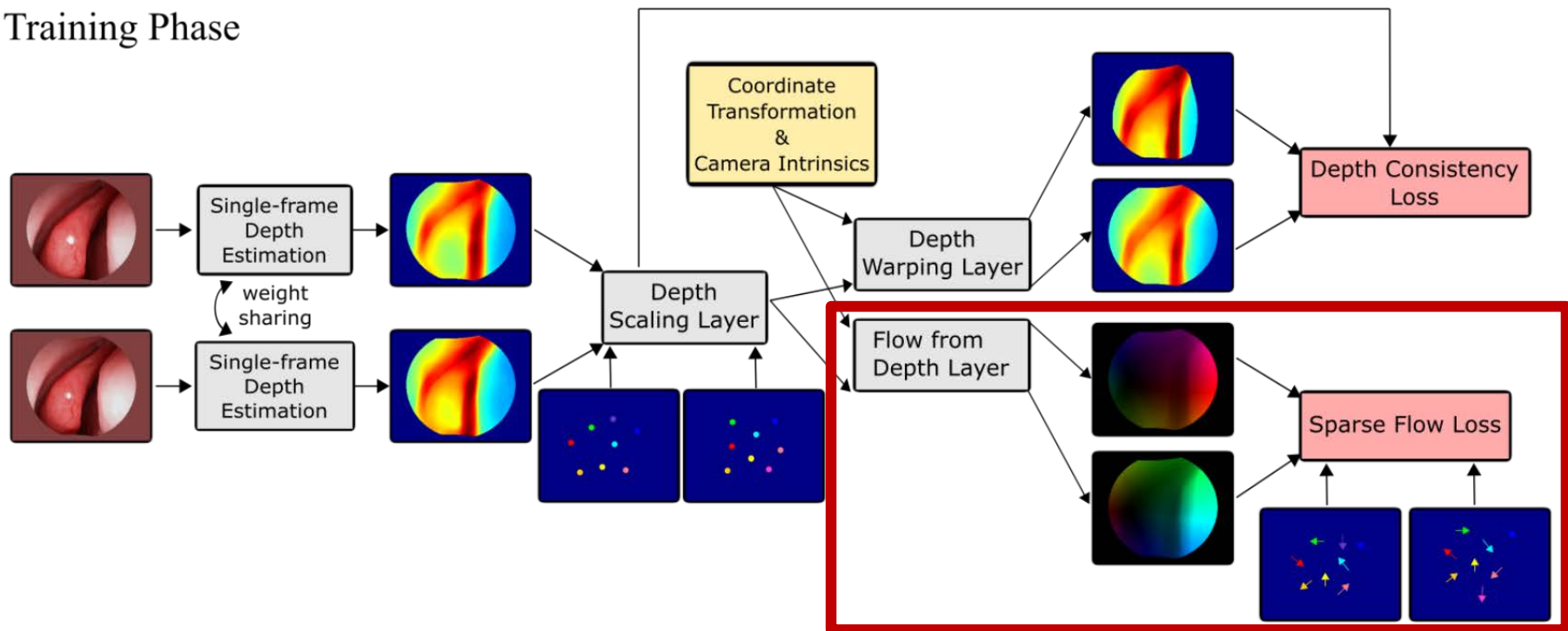
Training Phase



Structure from motion (SfM)-based self-supervision

- Run SfM on short video sequence (15 to 30 frames)
- Siamese network → Process multiple frames

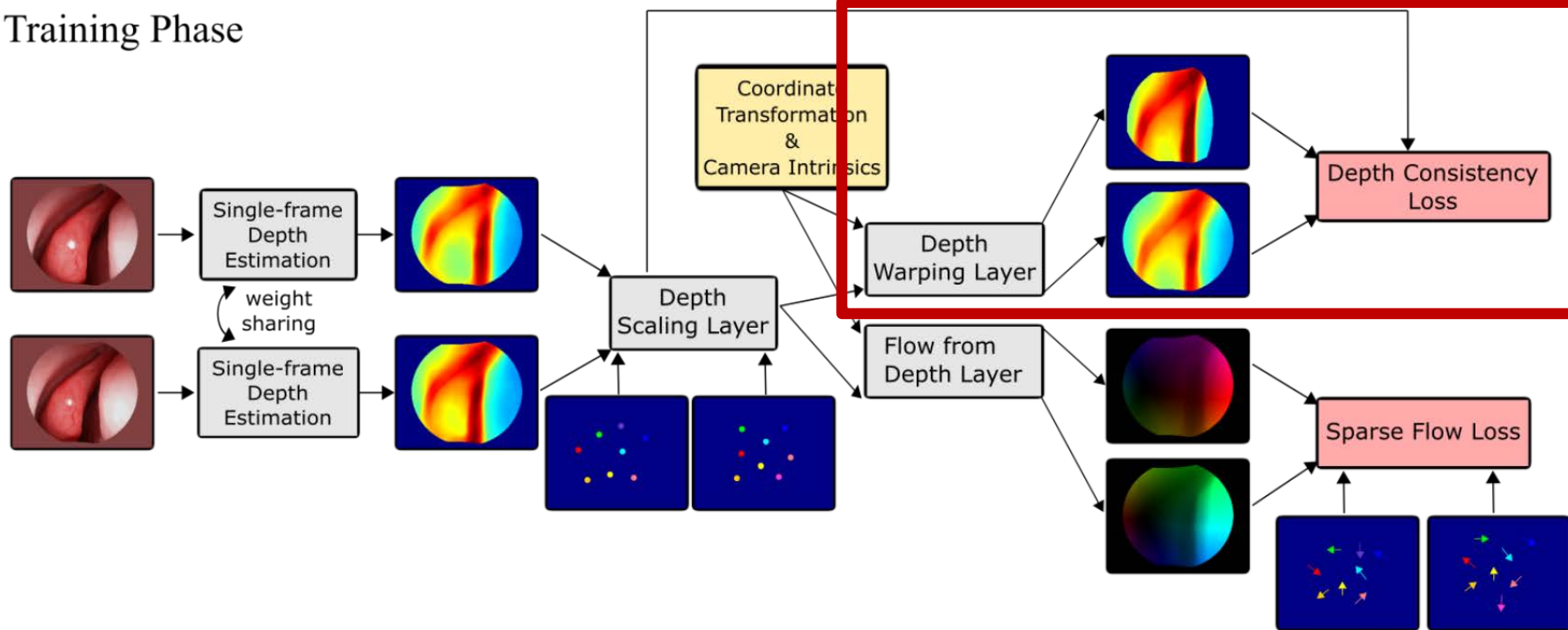
Training Phase



Sparse Flow Loss

- True 2D optical flow from 3D reconstructed points (SfM)
- Estimated optical flow from depth prediction

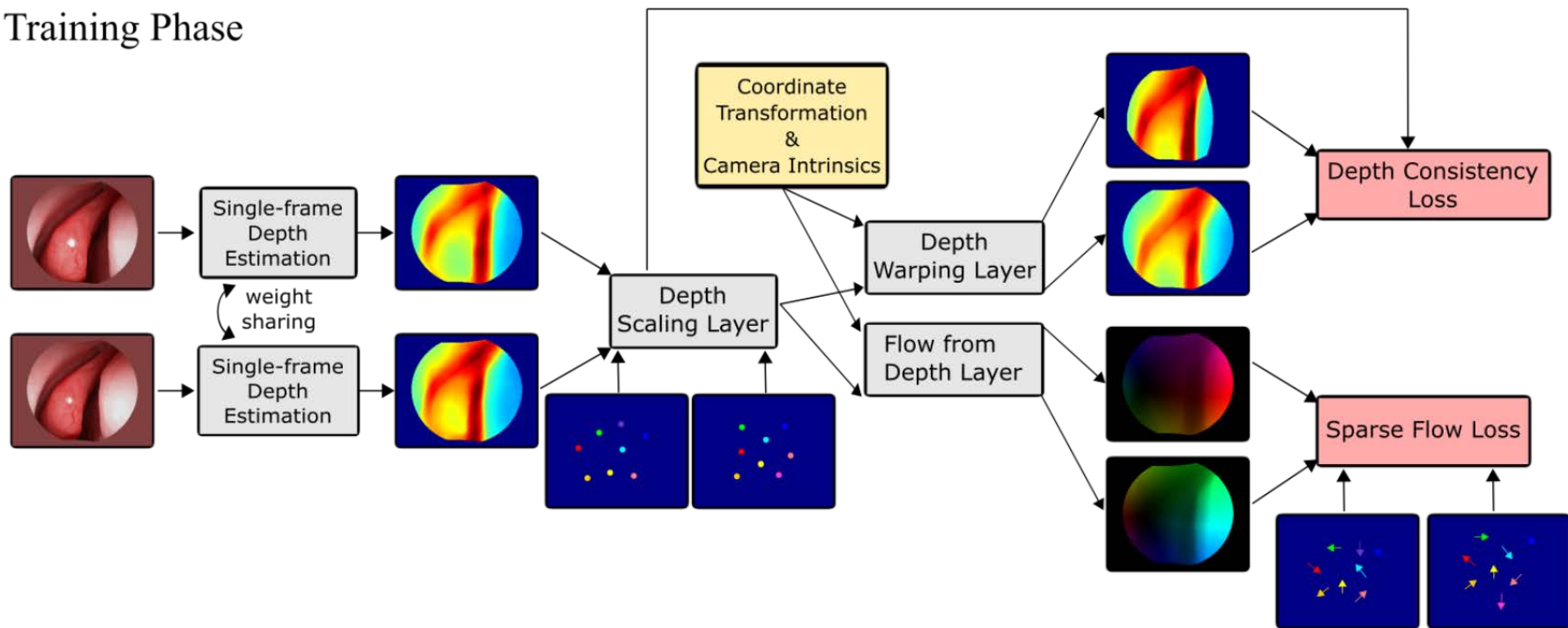
Training Phase



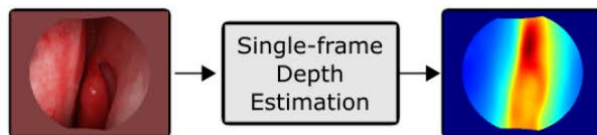
Depth Consistency Loss

- Differentiable warping operation to warp estimated depth into neighbor frame
- Enforces consistency among predictions

Training Phase



Application Phase



Dataset and Architecture

- Endoscopic video (no tools) of 6 consenting patients
 - 8 minutes of video total; rectified, and downsampled to 256 x 320 pixels
 - Different endoscopes for every patient
 - 4 patients with corresponding CT data (ground truth, disregarding erectile tissue)

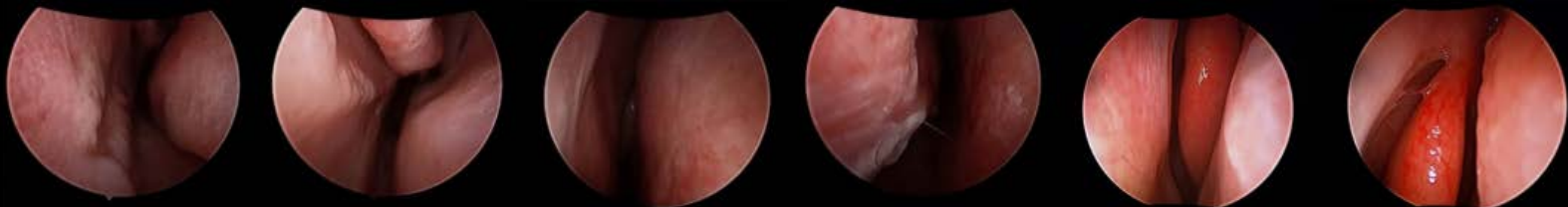


Dataset and Architecture

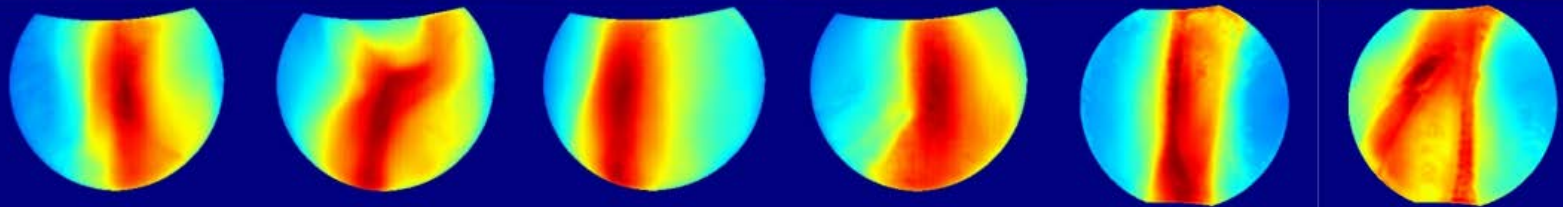
- Endoscopic video (no tools) of 6 consenting patients
 - 8 minutes of video total; rectified, and downsampled to 256 x 320 pixels
 - Different endoscopes for every patient
 - 4 patients with corresponding CT data (ground truth, disregarding erectile tissue)
- Depth estimation architecture
 - U-Net (8 M params): Easy to train on sparse signals but overfits heavily
 - FC-DenseNet-57 (1.5 M params): Generalizes well but hard to train from scratch
 - Teacher-Student approach
 - Teacher self-supervised learning
 - Teacher supervises student
 - Student self-supervised learning
 - Code available on GitHub: [lppllpp1920/EndoscopyDepthEstimation-Pytorch](https://github.com/lppllpp1920/EndoscopyDepthEstimation-Pytorch)



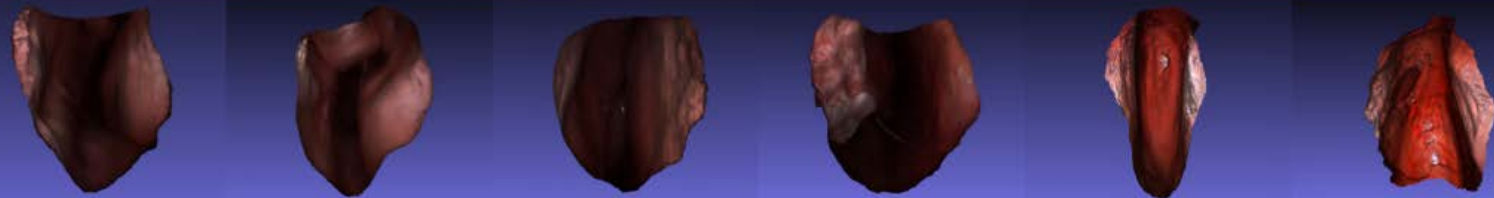
Input Video



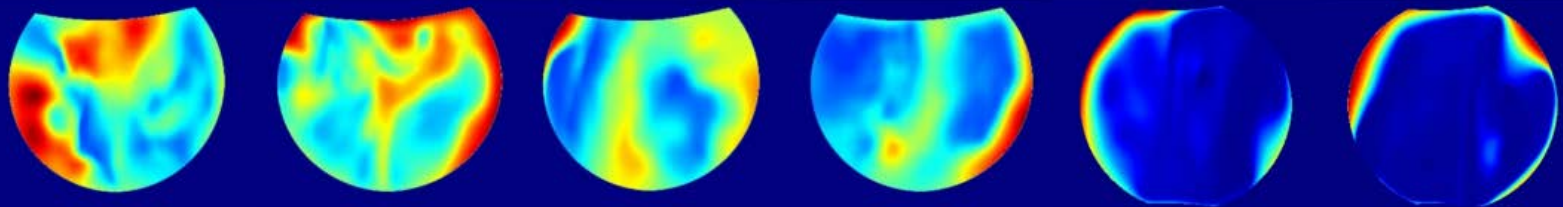
Our depth



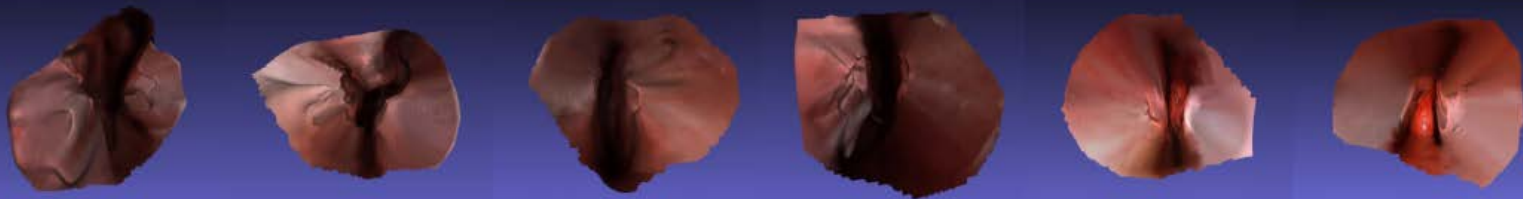
Our recon.



SfmLearner



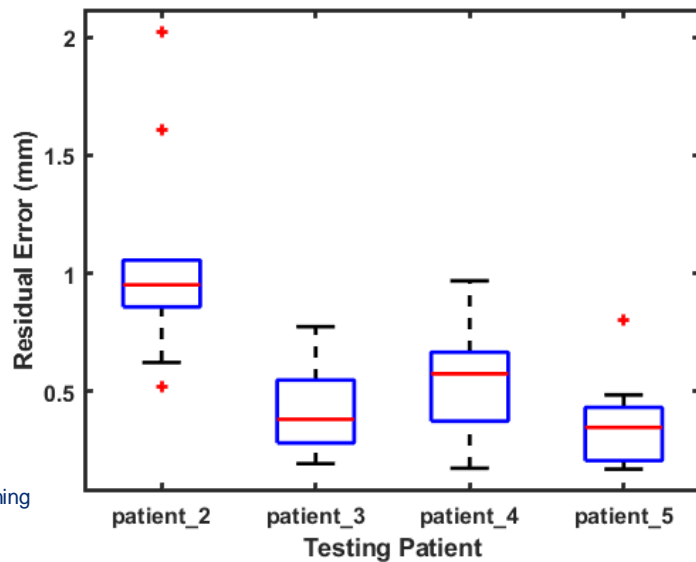
SfmLearner
recon.



Quantitative Results

- Leave-one-out training
- Randomly sample 20 frames per left-out patient
 - Estimate depth
 - Register to patient CT surface via GD-IMLOP (no shape deformation)
 - Compute residual error

- Sub-millimeter accuracy in most cases!
 - SfmLearner: > 10 mm
 - Deep (dark) regions exhibit high variation
→ Outliers
 - CT is imperfect ground truth (erectile tissue)

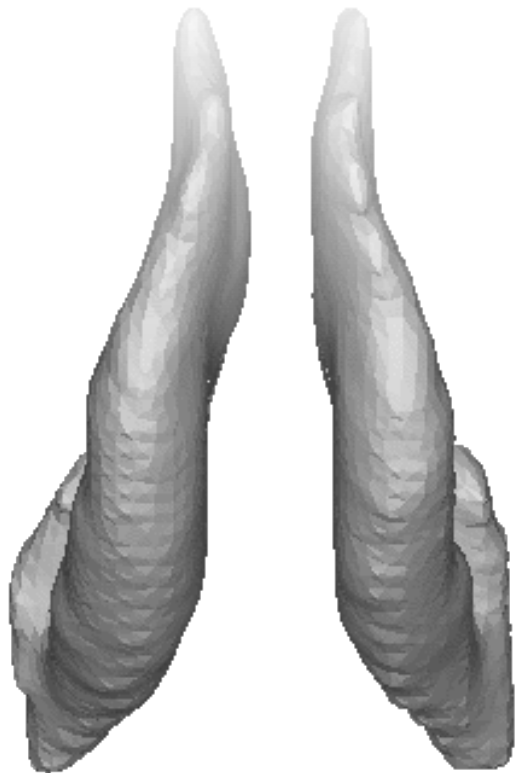


Towards Next-generation Image Guidance

Navigating Without Prior Information



Estimating Patient-specific Anatomy



Potential sources of patient-specific models

- CT scans
- Statistical shape model
- ...

Can we build a patient-specific, dense 3D model

- intra-operatively and
- on-the-fly?

Estimating Patient-specific Anatomy

Multiple sources of patient-specific models

CT scans

Statistical shape model

...

How to build a patient-specific, dense 3D model

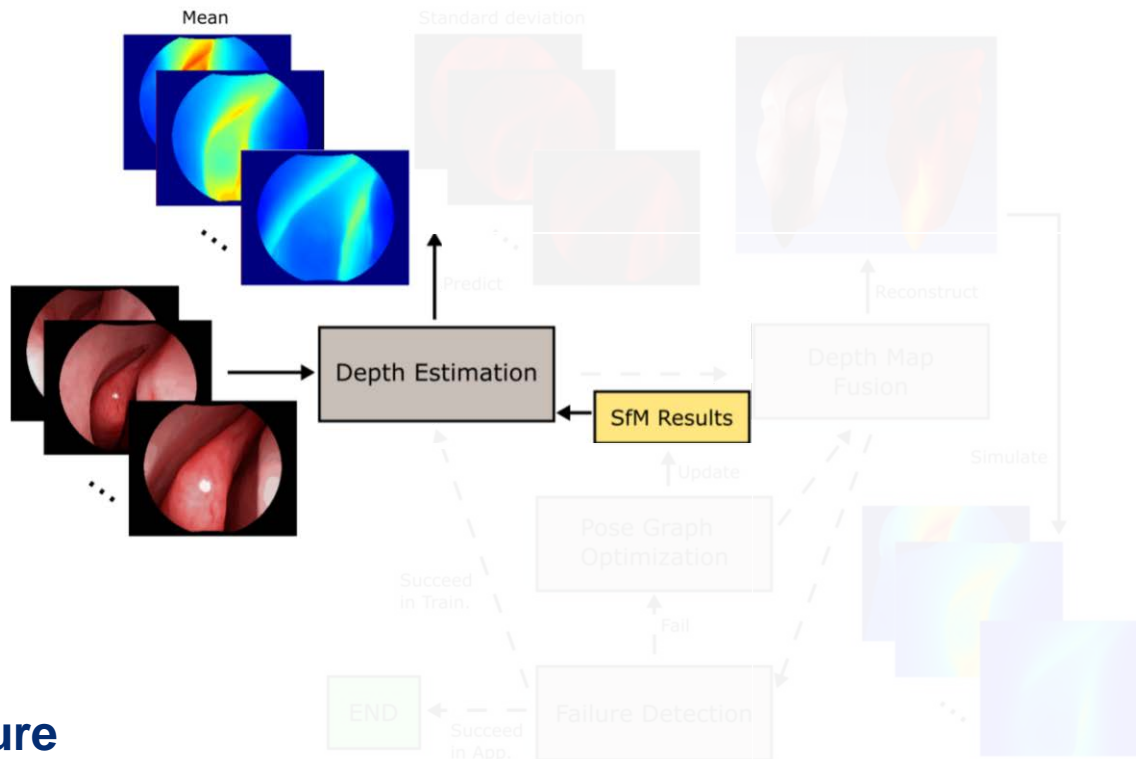
intra-operatively and

on-the-fly?

How do we benefit two ways

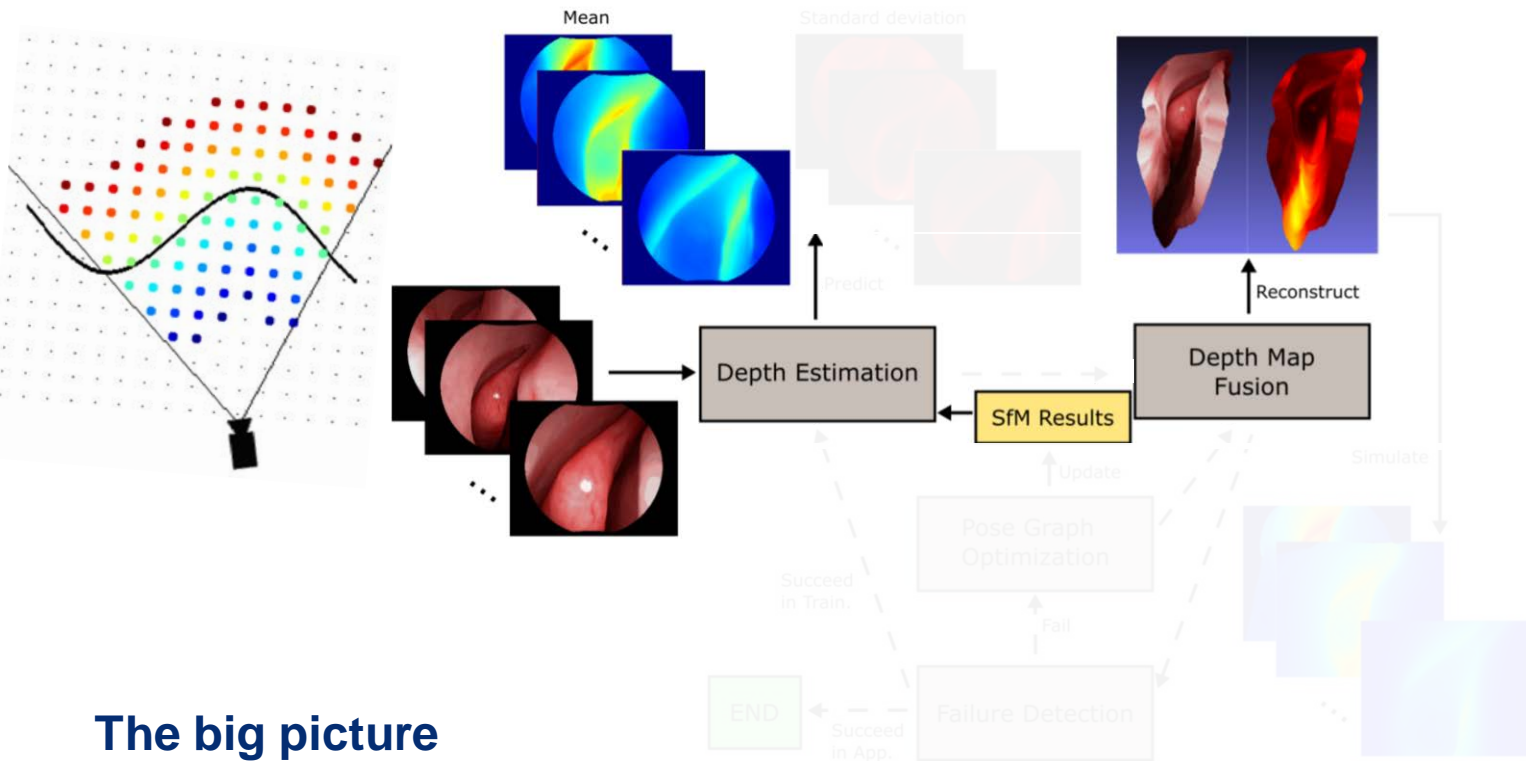
Bootstrapping for **dense depth supervision**

Uncertainty of depth estimates



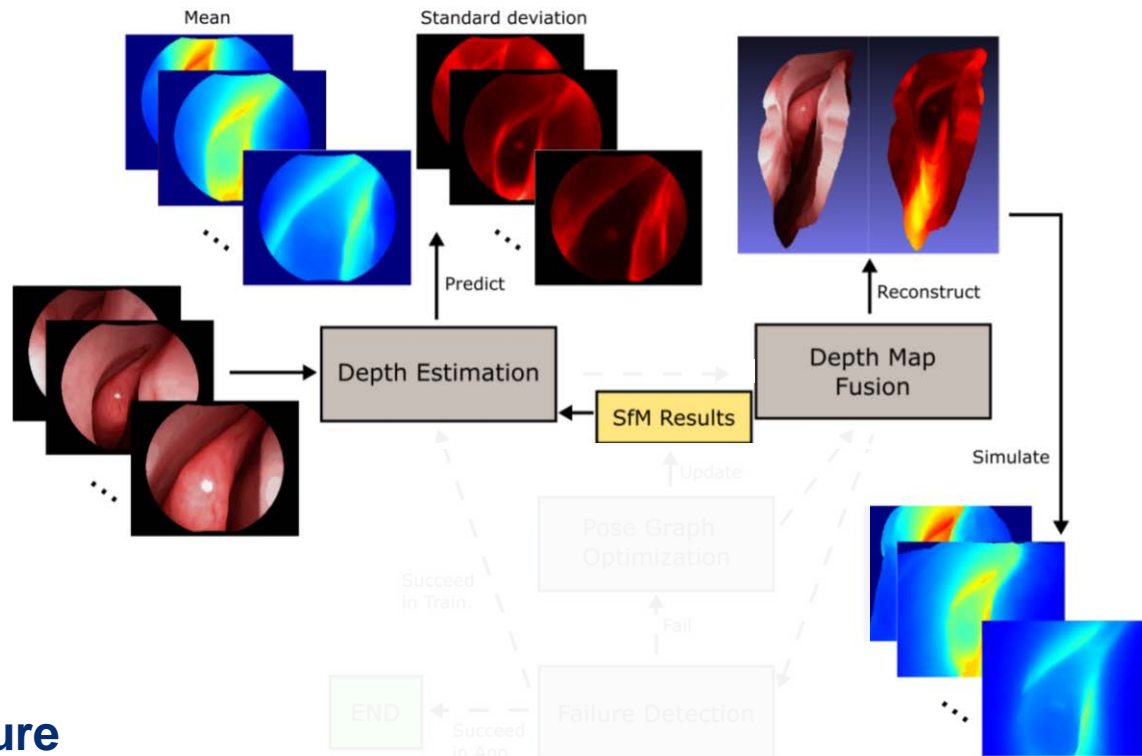
The big picture

1. Self-supervised training of depth estimation (now on long video sequences)



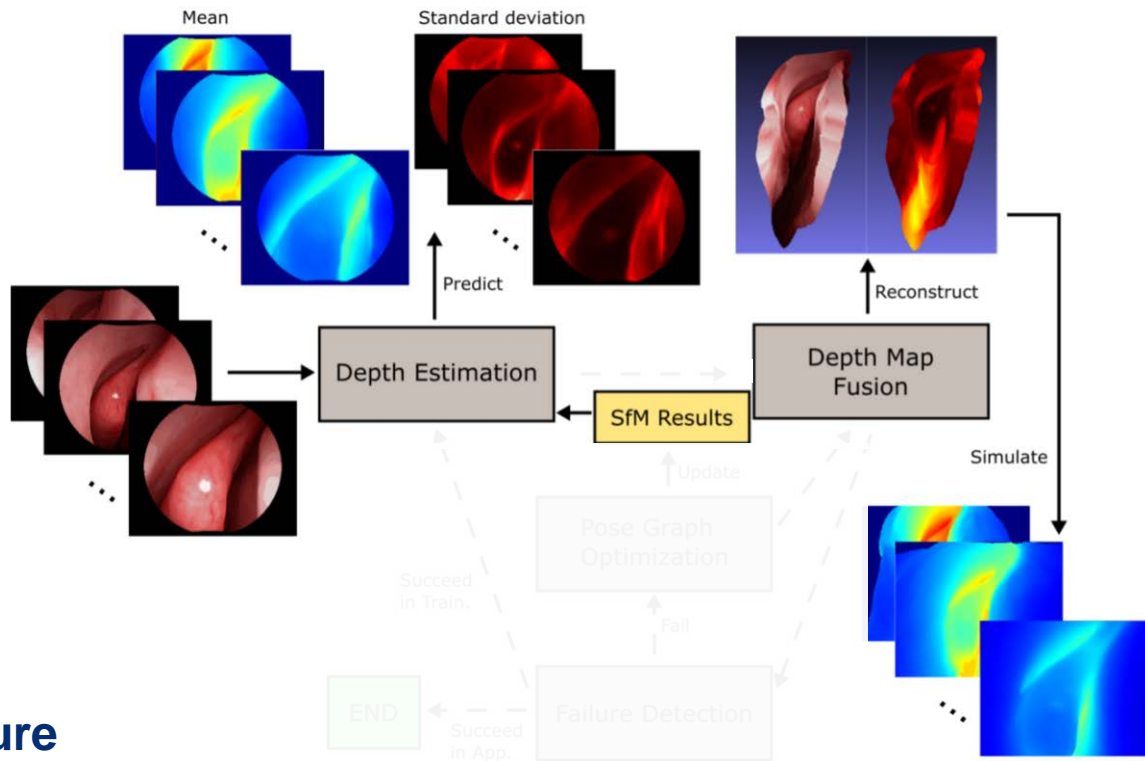
The big picture

1. Self-supervised training of depth estimation (now on long video sequences)
2. Volumetric fusion (truncated signed distance function) → Mean, STD



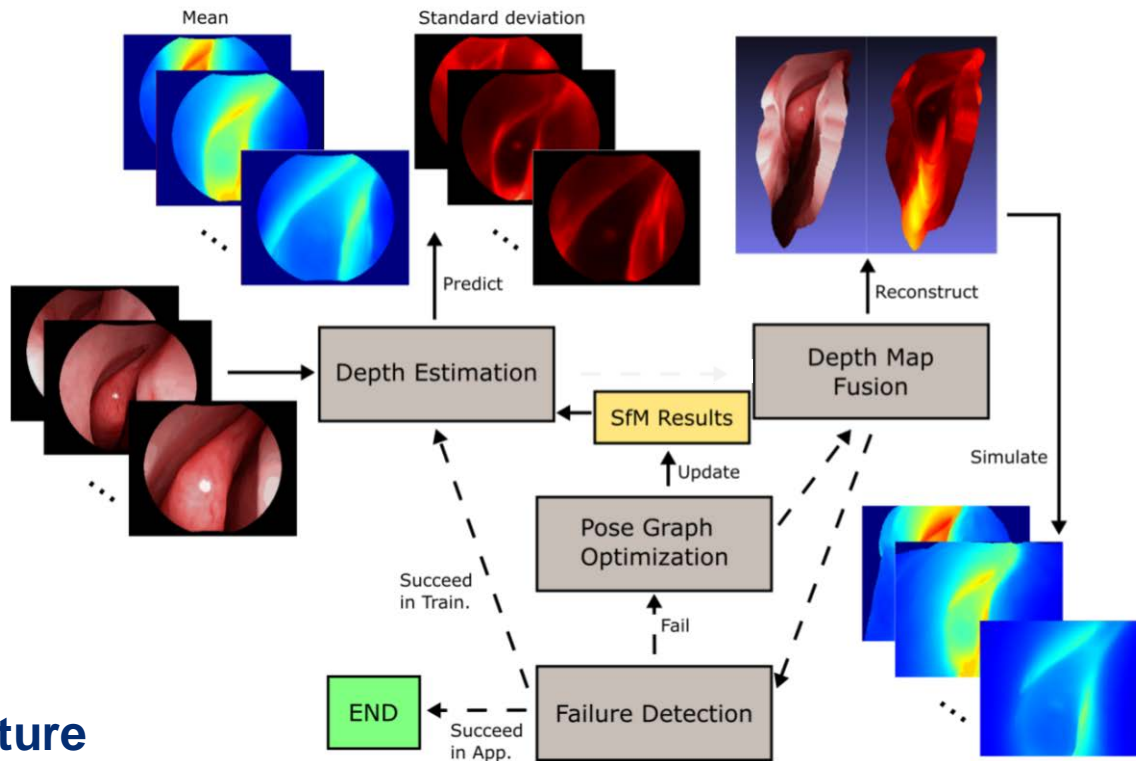
The big picture

1. Self-supervised training of depth estimation (now on long video sequences)
2. Volumetric fusion (truncated signed distance function) → Mean, STD
3. Bootstrapping → Dense supervision of mean depth and uncertainty



The big picture

1. Self-supervised training of depth estimation (now on long video sequences)
 2. Volumetric fusion (truncated SVD fusion) → Mean, STD
 3. Bootstrapping → Dense supervision of mean depth and uncertainty
- But wait, there's more!**



More big picture

- SfM results can be incorrect (few points etc.) → **Fusion will be wrong**
- Consistency between simulated and estimated depth → **Failure detection**
- If close → Pose graph refinement; If far off → Re-run SfM







Results and Observations

- Again, leave-one-out and GD-IMPLOP to patient CT
- Sub-millimeter errors
- Error seems higher → **Misleading**
 - Reconstruction is of ~ 1 minute video not just a single frame
 - Registration has larger residual, but average is over much larger region

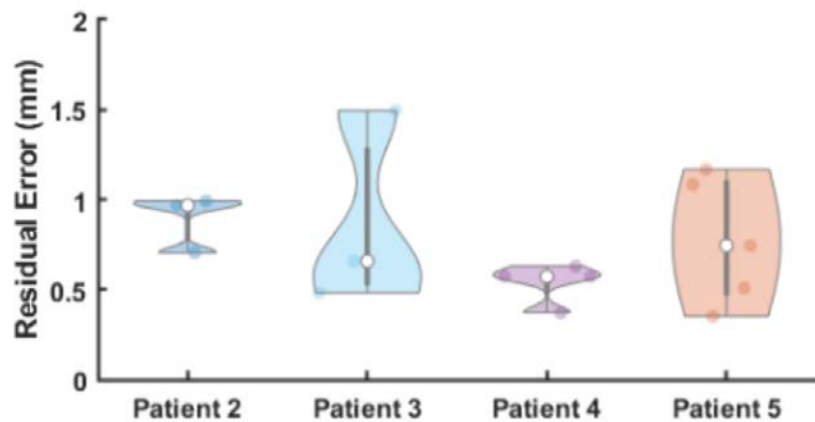
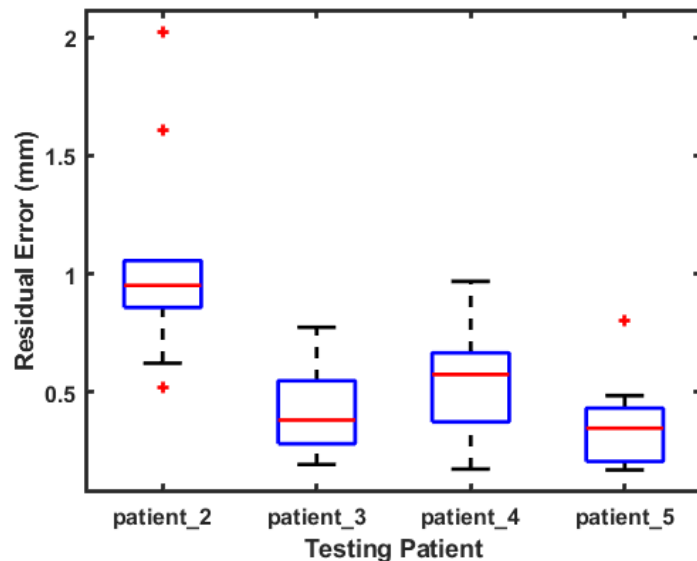


Image Guidance for Endoscopic Procedures

Concluding Remarks – Accounting for Anatomical Change



Where do we go from here?


Quantitative endoscopy

- Longitudinal monitoring of anatomical change
- E.g. for monitoring polyp behavior after steroid injection

The fairly untapped supreme discipline...

Monitoring anatomical change during surgery

- How to deal with tools?
- Blood, gore, and all other sorts of unseen variation?

A close-up, high-resolution photograph of a human eye, showing the iris and pupil in detail. The eye is the central focus on the left side of the slide.

**Thank you.
Questions?**

