

Endo-FAST3r: Endoscopic Foundation model Adaptation for Structure from motion

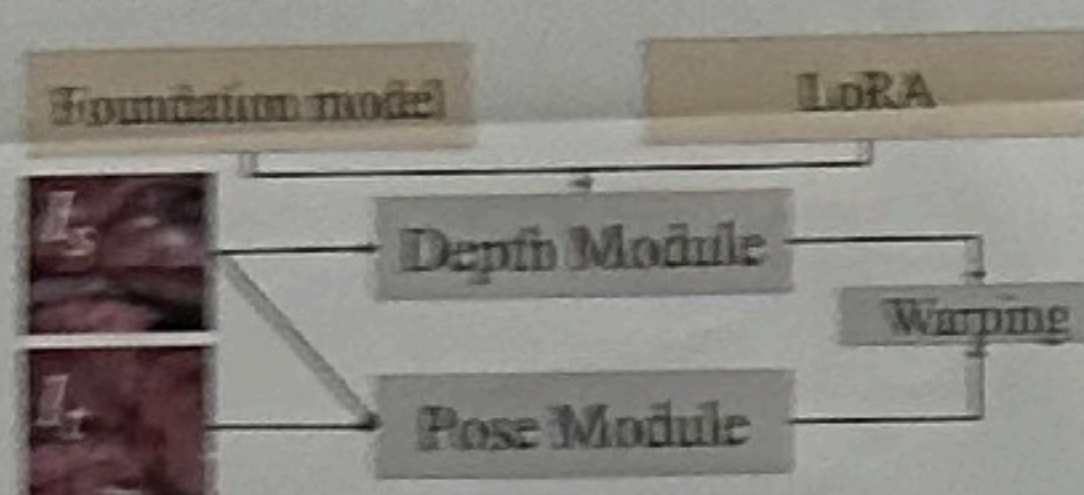
Mona Sheikh Zeinoddin^{1,2}, Mobarak I. Hoque^{1,4}, Zafer Tandogdu^{3,6}, Greg L. Shaw³, Matthew J. Clarkson^{1,4}, Evangelos B. Mazomenos^{1,4}, Danail Stoyanov^{1,5}

¹Hawkes Institute, University College London, London, UK ²Institute of Health Informatics, University College London, London, UK ³Dept. of Urology, University College London Hospitals, London, UK ⁴Dept. of Medical Physics & Biomedical Engineering, University College London, UK ⁵Dept. of Computer Science, University College London, UK ⁶Division of Surgery and Interventional Science, University College London, UK



Motivation & Background

- Camera Pose & Depth estimation is essential to achieve 3D scene understanding in robotic-assisted surgery.
- One of the most widely used frameworks to perform these two tasks is the self-supervised reprojection loss [2] pipeline.

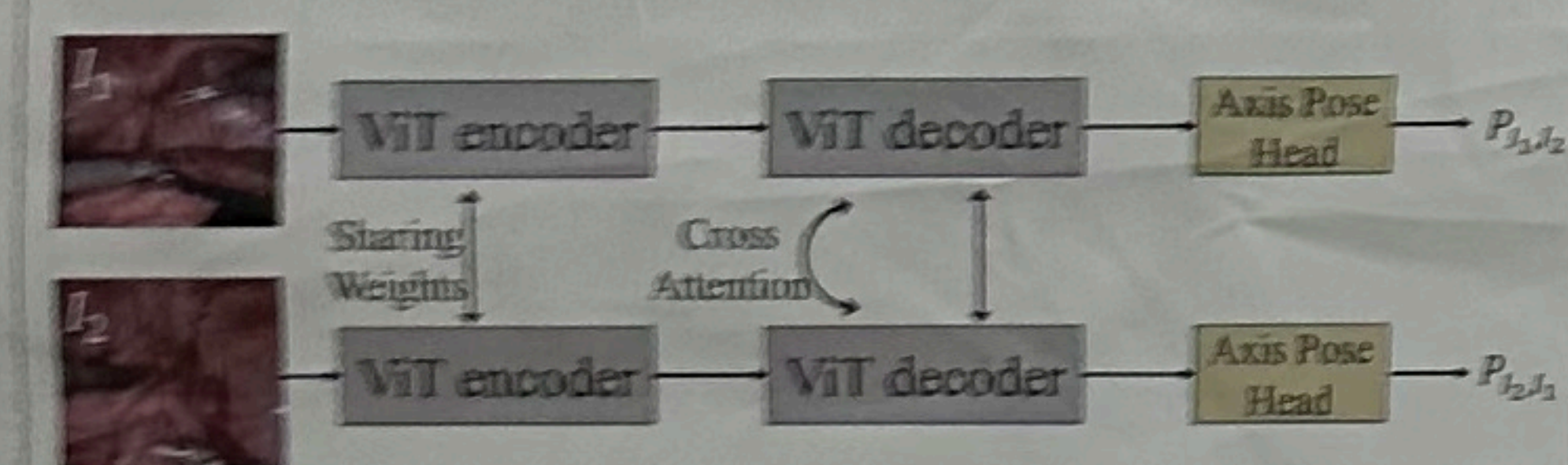


- Previous works have investigated the use of foundation models in the depth module via Low-Rank Adaptation (LoRA) [1]-based techniques.
- Major limitations of current literature:
- No work has investigated the use of foundation models in the pose module.
- The limiting low rank update space of LoRA-based approaches.

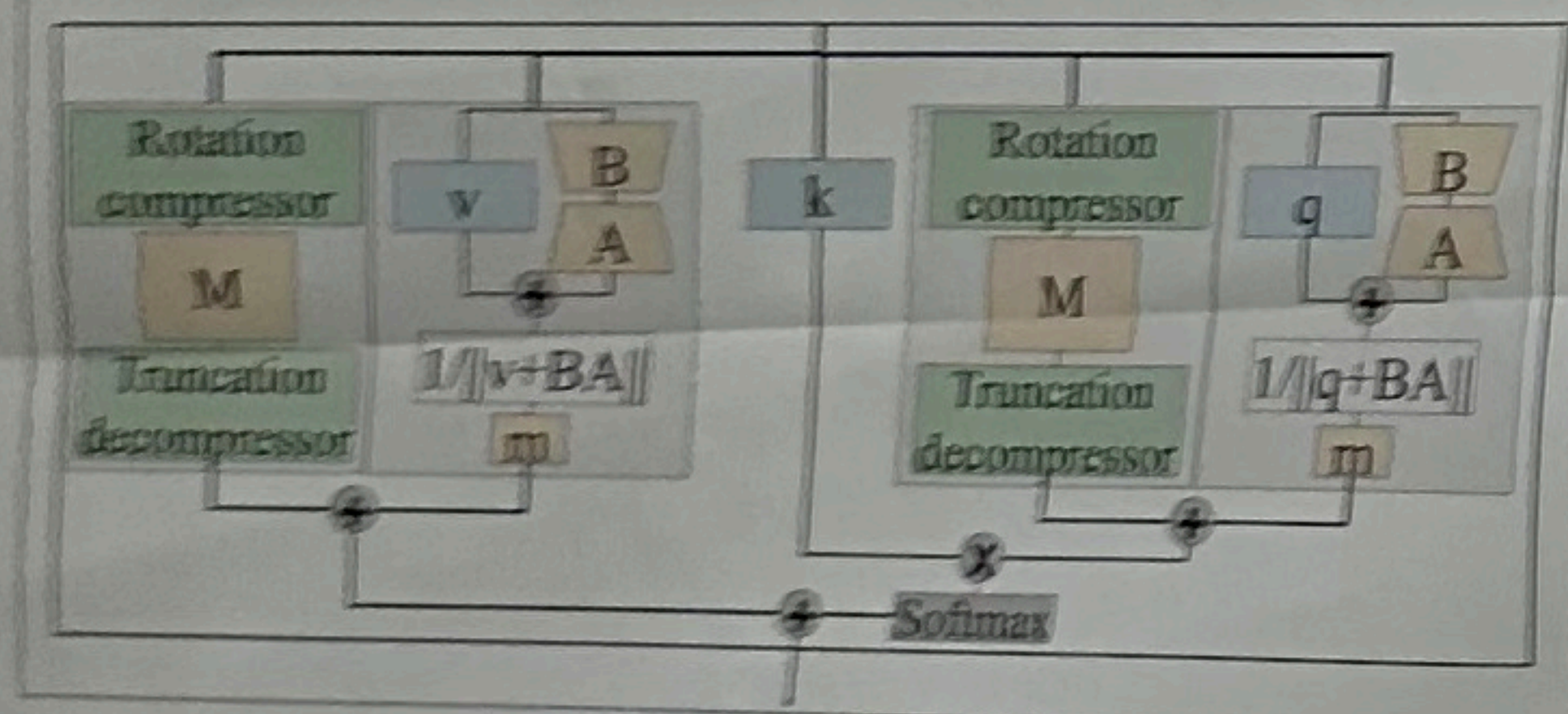
Endo-FAST3r Contributions

- In this work, Endo-FAST3r: Endoscopic Foundation model Adaptation for Structure from motion, we introduce:

- Reloc3rX: Extending the foundation model Reloc3r [3] by designing the Axis Pose Head.



- DoMoRA: Enabling both low- and full-rank updates.



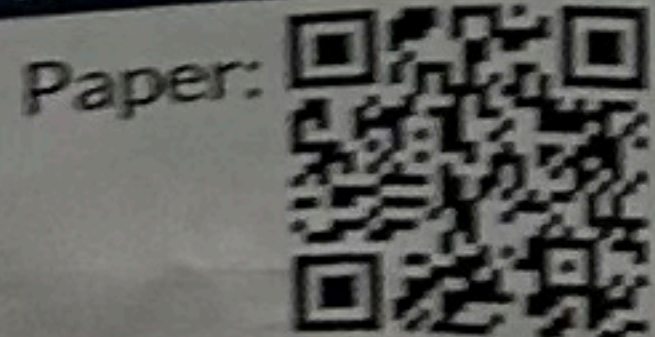
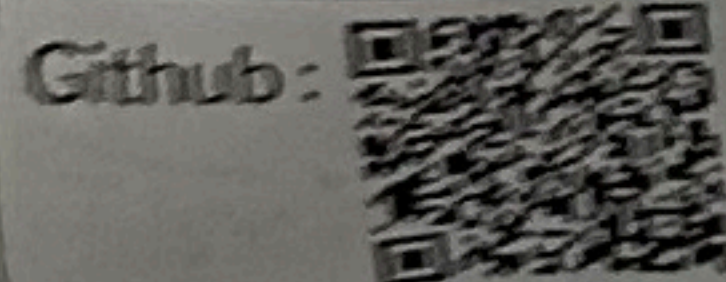
References

- [1] Hu, Edward S, et al. "LoRA: Low-Rank Adaptation of Large Language Models." In ICLR (2022)
- [2] Shao, Shuwei, et al. "Self-supervised monocular depth and ego-motion estimation in endoscopy: Appearance flow to the rescue." In ICRA (2021)
- [3] Dong, Siyan, et al. "Reloc3r." In CVPR (2025)

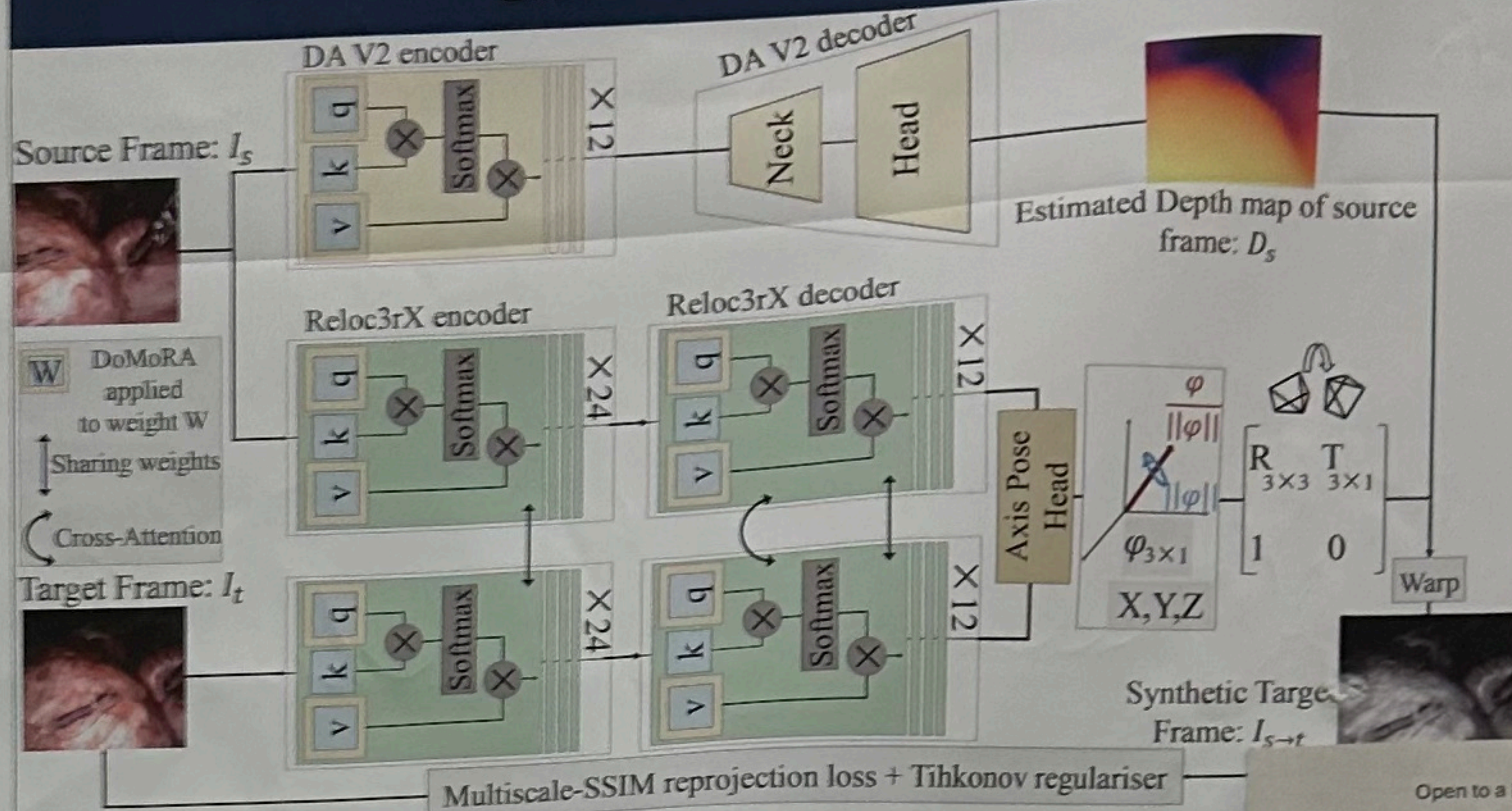
Let's discuss

- What is the biggest barrier of using foundation models in the surgical field?
- How does the self-supervised reprojection loss pipeline work?
- What is the main difference of LoRA and DoMoRA?

Interested in our work?



Endo-FAST3r Architecture



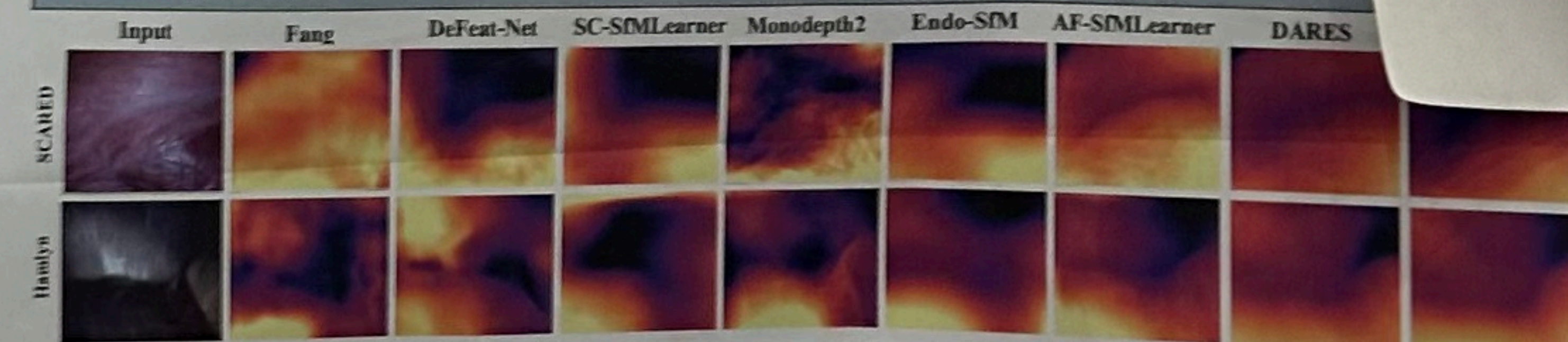
Results & Conclusion

- Trained on the SCARED dataset and Evaluation performed on the rigid SCA and non-rigid StereoMIS datasets.

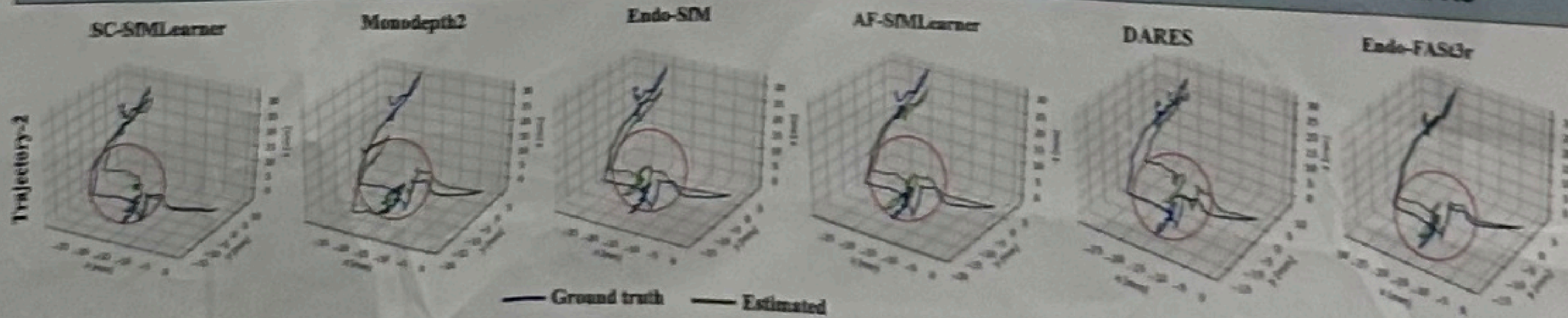
Comparison on rigid scenes with benchmark methods in depth estimation and pose estimation

	Method	AbsRel↓	SqRel↓	RMSE↓	$\delta \uparrow$	ATE-T1 ↓	ATE-T2 ↓	Total	Tra
SCARED	DeFeat-Net	0.077	0.792	6.688	0.941	0.1765	0.0995	14.8	14.8
	SC-SfMLearner	0.068	0.645	5.988	0.957	0.0767	0.0509	14.8	14.8
	Monodepth2	0.069	0.577	5.546	0.948	0.0769	0.0554	14.8	14.8
	Endo-SfM	0.062	0.606	5.726	0.957	0.0759	0.0500	14.8	14.8
	AF-SfMLearner	0.059	0.435	4.925	0.974	0.0757	0.0501	14.8	14.8
	Yang et al.	0.062	0.558	5.585	0.962	0.0723	0.0474	2.0	2.0
	Zero-Shot DA V2	0.091	1.056	7.601	0.916	-	-	-	-
	Zero-Shot Reloc3r	-	-	-	-	0.0938	0.0735	-	-
	DARES	0.052	0.356	4.483	0.980	0.0752	0.0498	24.9	24.9
	EndoFAST3r (Ours)	0.051	0.354	4.480	0.998	0.0702	0.0438	24.9	24.9
Humanly	Endo Depth & Motion	0.185	5.424	16.100	0.732	-	-	14.8	14.8
	AF-SfMLearner	0.168	4.440	13.870	0.770	-	-	14.8	14.8
	EndoFAST3r (Ours)	0.166	4.529	13.718	0.778	-	-	24.9	24.9

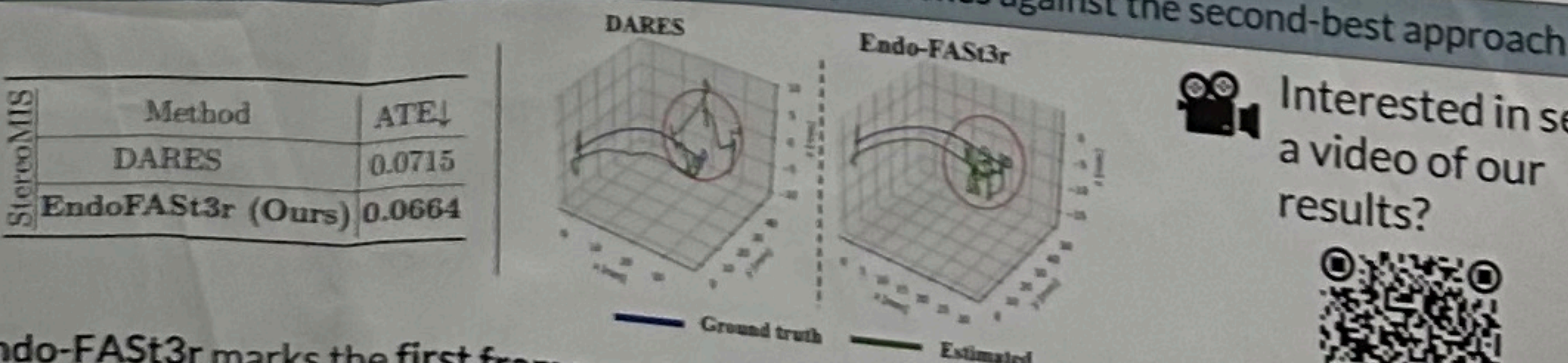
Qualitative comparison of our depth estimation on rigid scenes against benchmarks



Qualitative comparison of our ego-motion estimation on rigid scenes against benchmarks

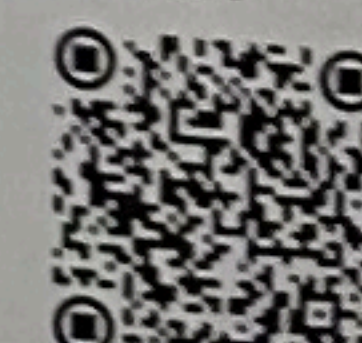


Comparison of our ego-motion estimation on non-rigid scenes against the second-best approach



- Endo-FAST3r marks the first framework to use foundation models for pose estimation in surgical environments, and it does so with NO ground truth data.
- Endo-FAST3r surpasses all SOTA methods, reaching an improvement of 9.34% in camera pose estimation and 2% in depth estimation over the nearest competitor.

Interested in seeing a video of our results?



Mona Sheikh Zeinoddin
mona.zeinoddin.22@ucl.ac.uk

UCL
Hawkes Institute

Engineering and
Physical Sciences
Research Council

EPSRC

UCL
Health Informatics

ROBOT VISION
SURGICAL ROBOT VISION
research group