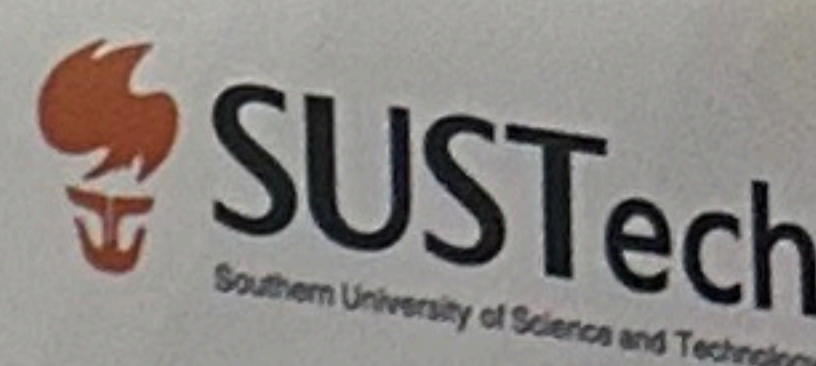


ReSurgSAM2:

Referring Segment Anything in Surgical Video via Credible Long-term Tracking

Haofeng Liu¹, Mingqi Gao², Xuxiao Luo¹, Ziyue Wang¹, Guanyi Qin¹, Junde Wu³, Yueming Jin^{1†}
¹National University of Singapore ²Southern University of Science and Technology ³University of Oxford
 haofeng.liu@u.nus.edu, ymjn@nus.edu.sg



Introduction

- Precise segmentation of surgical instruments and tissues is critical for automation, guidance, and training.
- Current methods:** generate collective masks without **interactivity** and lack hands-free **text-driven** interaction, real-time use, and long-term stability.
- Key question:** how to achieve **accurate, real-time** referring segmentation with robust **long-term** tracking in surgical videos?

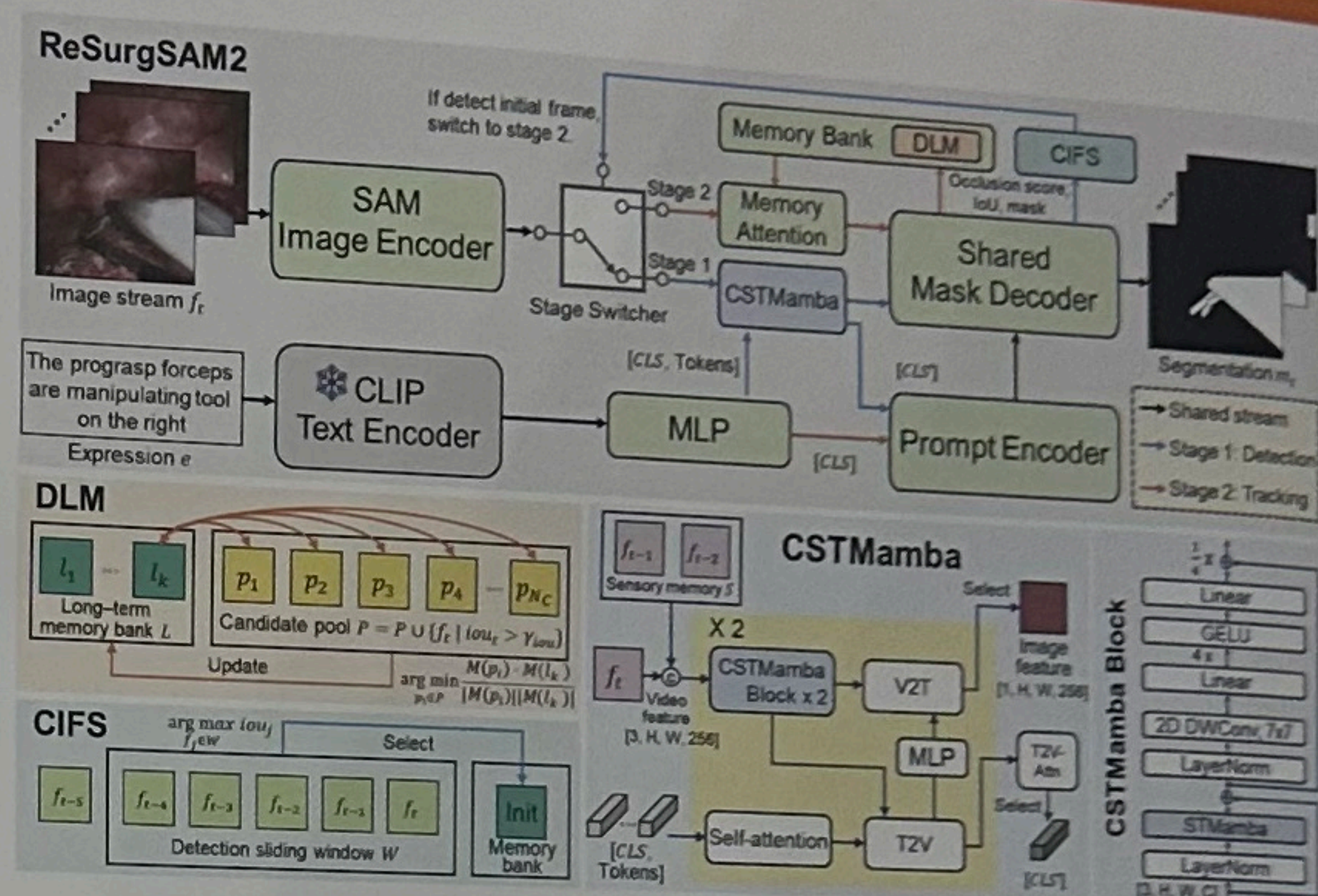
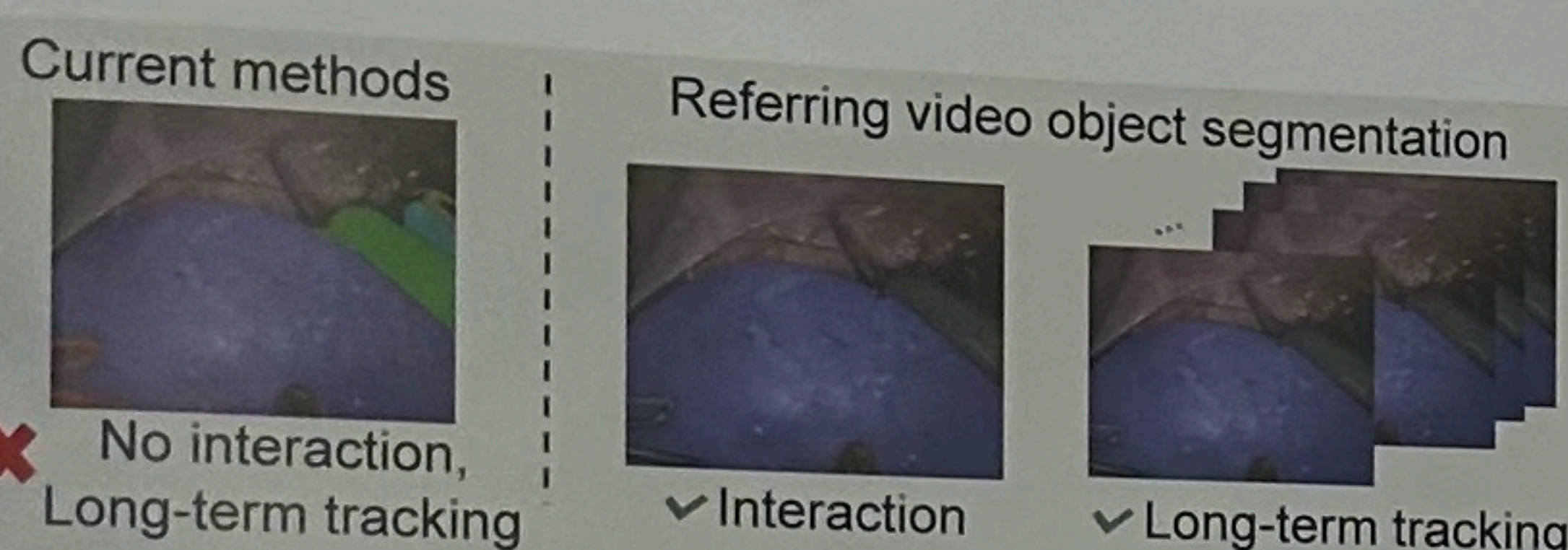
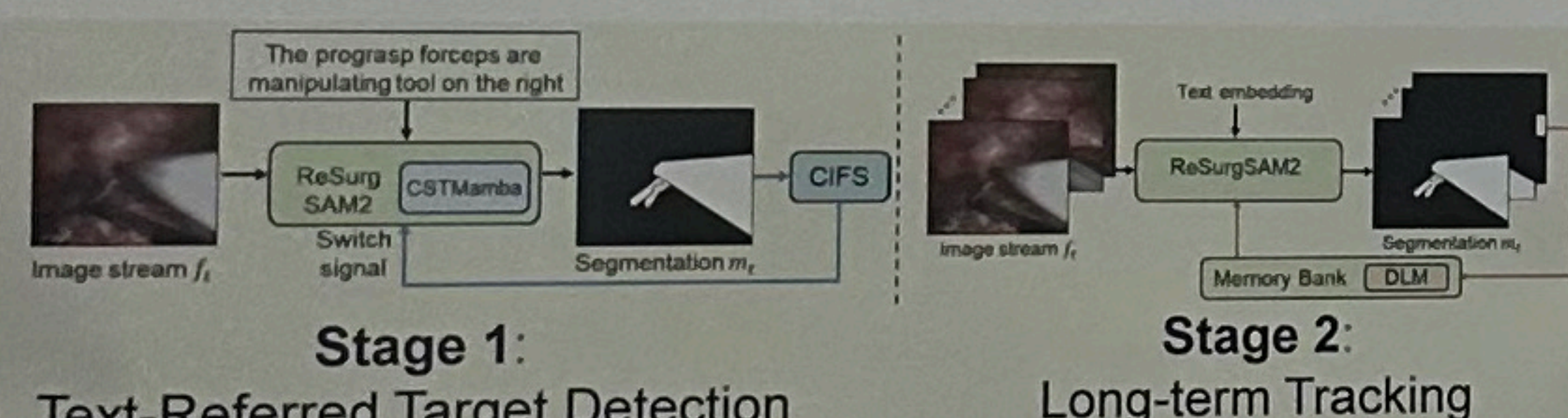


Figure 1. Overview of ReSurgSAM2

Method

Two-stage design: Stage 1 detects text-referred targets; Stage 2 enables long-term tracking.



- CSTMamba** – Efficient cross-modal spatio-temporal fusion with depth-wise convolution and text-vision fusion.
- CIFS** – Credible initial frame selection reduces error accumulation.
- DLM** – Diversity-driven memory builds a hybrid short/long-term memory, improving long-sequence stability.

Qualitative Analysis

- Complex Scenes** – segments the specified instrument.
- Clearer Boundaries** – Accurate instrument/tissue segmentation.
- Stable Tracking** – Maintains consistency during occlusion, motion.

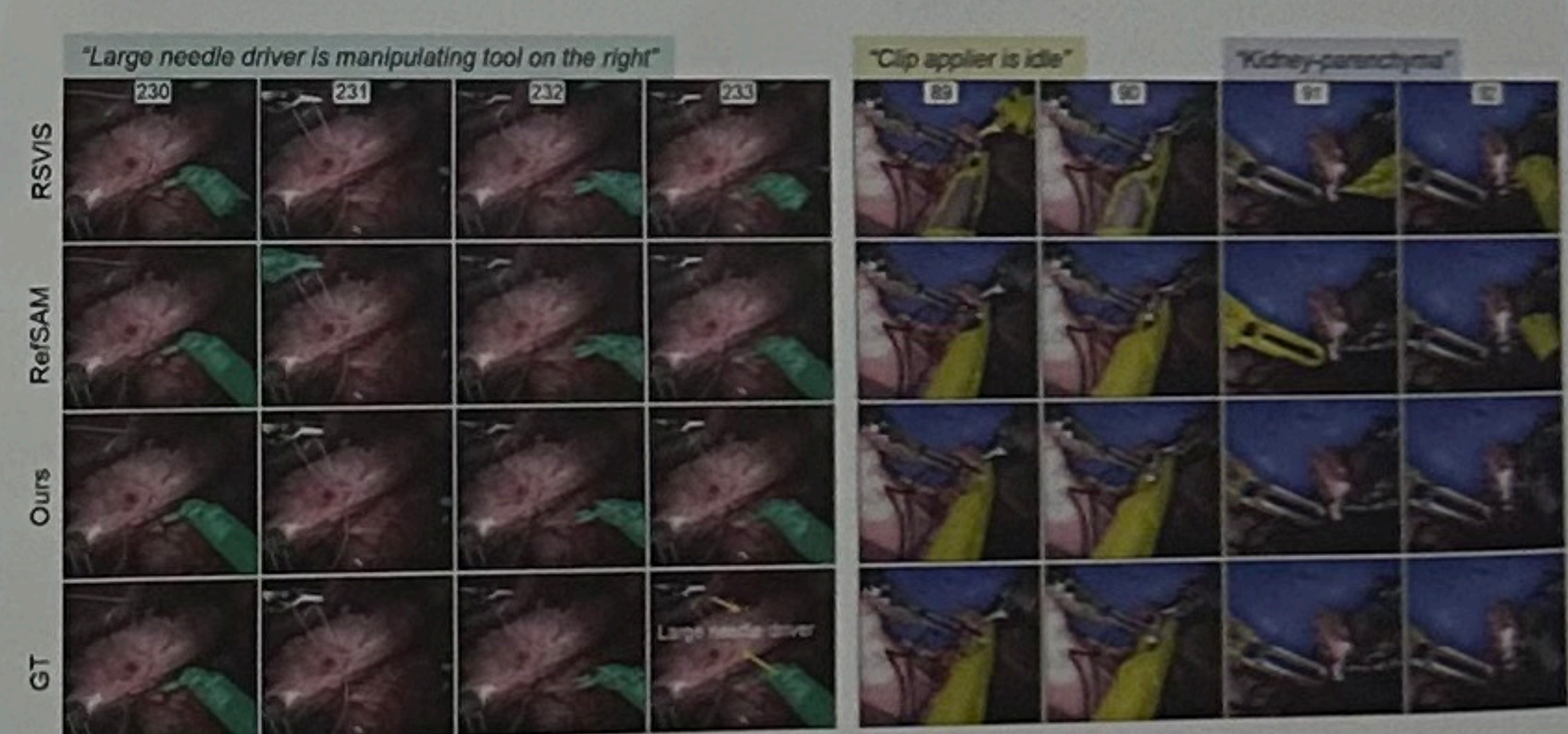


Figure 2. Visual comparison between ReSurgSAM2 and the state-of-the-art.

Experiment

Dataset: Ref-EndoVis17 and Ref-EndoVis18

Table 1: Dataset statistics for Ref-EndoVis17/EndoVis18 datasets.

Dataset	Training				Testing			
	Sequence	Frame	Object	Pair	Sequence	Frame	Object	Pair
Ref-EndoVis17(tool)	7	2100	20	4873	3	900	10	2265
Ref-EndoVis18(tool)	11	1639	34	3787	4	596	15	1384
Ref-EndoVis18(tissue)	11	1639	25	2995	4	596	7	807

Metric: J (region accuracy), F (boundary accuracy), J&F, FPS.

Comparison Experiment

- State-of-the-art Accuracy.**
- Robust Long-term Tracking.**
- Real-time Performance** – Runs efficiently at 61.2 FPS.

Table 2: Quantitative comparison with state-of-the-art methods.

Method	Setting	Ref-EndoVis17(tool)			Ref-EndoVis18(tool)			Ref-EndoVis18(tissue)			FPS
		J&F	J	F	J&F	J	F	J&F	J	F	
ReferFormer [27]	Offline	62.41	62.28	62.55	71.09	70.96	71.23	61.84	69.9	53.78	42.3
MUTR [28]	Offline	60.97	60.76	61.18	67.56	67.79	67.33	63.53	71.48	55.58	32.3
RSVIS [24]	Online	61.22	61.37	61.07	68.35	68.55	68.15	65.69	72.91	58.47	22.1
OnlineRefer [26]	Online	60.32	60.29	60.34	72.19	71.88	72.50	70.56	77.58	63.55	25.6
RefSAM [11]	Online	63.56	63.77	63.35	72.86	73.40	72.31	71.90	77.66	66.14	25.4
ReSurgSAM2	Online	77.73	77.77	77.69	80.62	80.94	80.31	75.09	80.03	69.25	61.2
		+14.17			+7.76			+3.19			+18.9

Table 3: Component Contribution Analysis

Stage 2	CSTMamba	CIFS	DLM	J&F	J	F	FPS
✓				61.15	61.46	60.84	70.1
✓	✓			63.79	63.77	63.82	68.2
✓	✓	✓		68.56	68.51	68.61	67.5
✓	✓	✓	✓	74.70	74.67	74.72	63.1
✓	✓	✓	✓	77.73	77.77	77.69	61.2

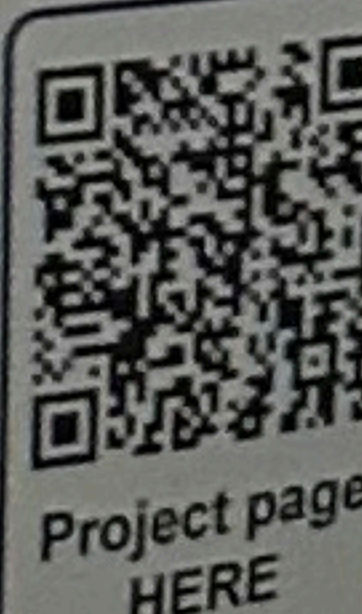
Table 4: Memory Bank Comparison

Method	J&F	J	F
Vanilla	74.70	74.67	74.72
Extended	74.68	74.64	74.72
Interval	75.32	75.27	75.37
DLM	77.73	77.77	77.69

- Alation:** separation (+2.64), fusion (+4.77), initialization (+6.14), memory (+3.03)
- Memory Bank:** DLM significantly improves long-term tracking stability compared with different memory variants.

Conclusion

- ReSurgSAM2 enables hands-free, text-driven segmentation with real-time and long-term tracking.
- It supports intraoperative guidance, analytics, and surgical training.



Project page
HERE