



**ICLR**

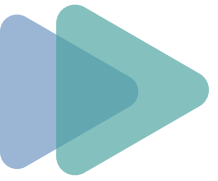
# HFSTI-Net: Hierarchical Frequency-spatial-temporal Interactions for Video Polyp Segmentation

**Yuanqin He<sup>1</sup>, Guilian Chen<sup>1</sup>, Yuhua Zhang<sup>1</sup>, Huisi Wu<sup>1\*</sup>, Jing Qin<sup>2</sup>**

<sup>1</sup>College of Computer Science and Software Engineering, Shenzhen University

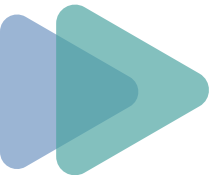
<sup>2</sup>Centre for Smart Health, School of Nursing, The Hong Kong Polytechnic University

2400101064@mails.szu.edu.cn, hswu@szu.edu.cn

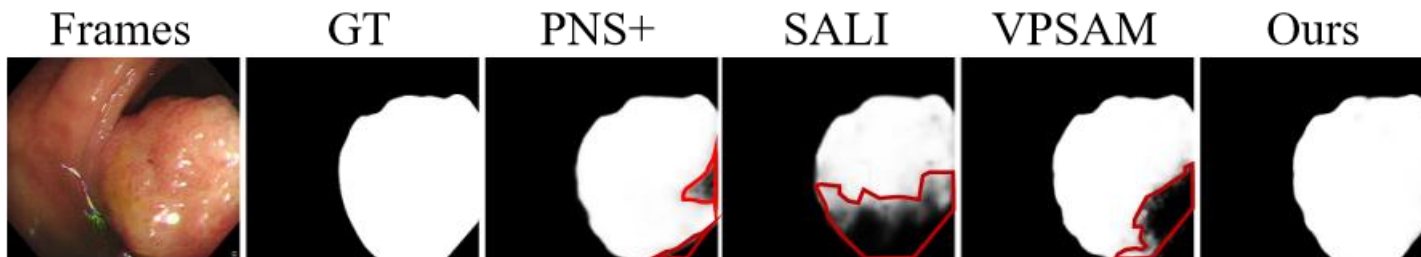


## Background

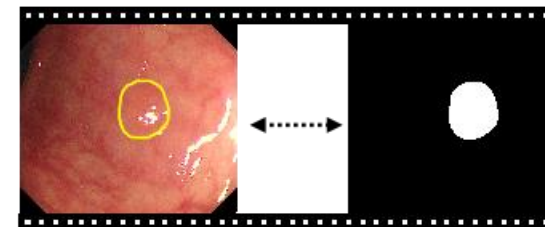
- Colorectal cancer is the third most common cancer worldwide.
- Colonoscopy with polyp removal is the most effective prevention strategy.
- Diagnostic accuracy depends heavily on endoscopist experience and attention.
- Accurate, real-time video polyp segmentation is essential to assist detection, reduce missed cases, and improve consistency across clinicians.



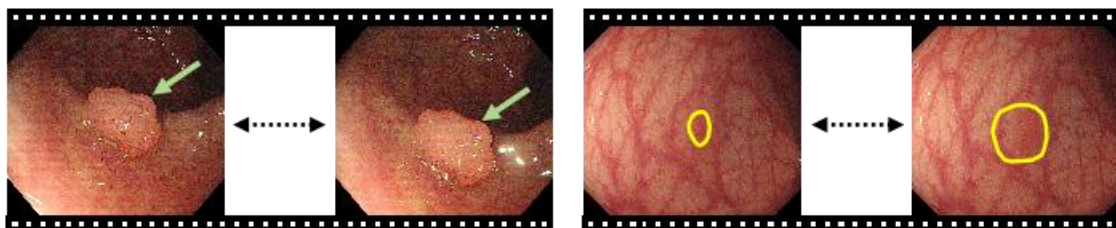
# Challenges



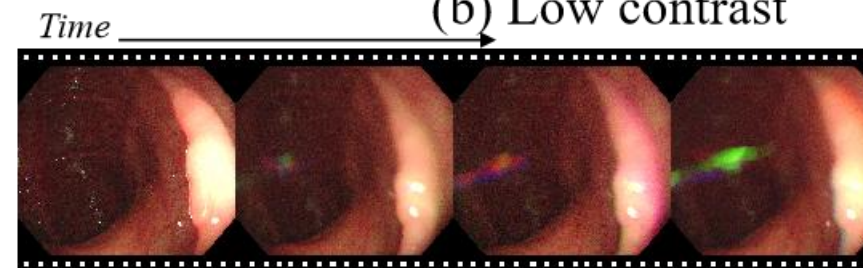
(a) Shape collapse



(b) Low contrast



(c) Significant variations



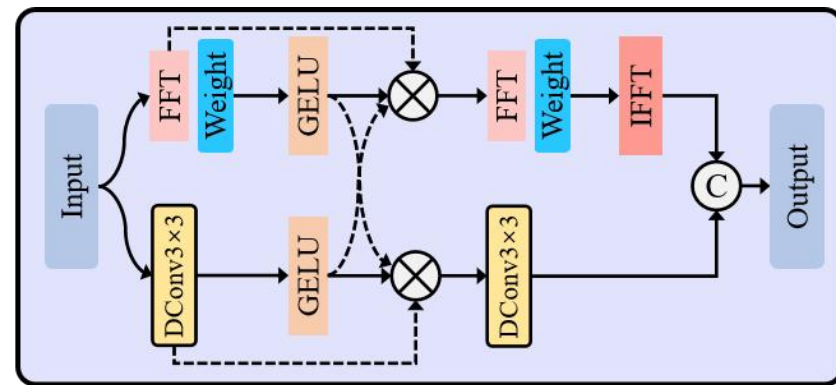
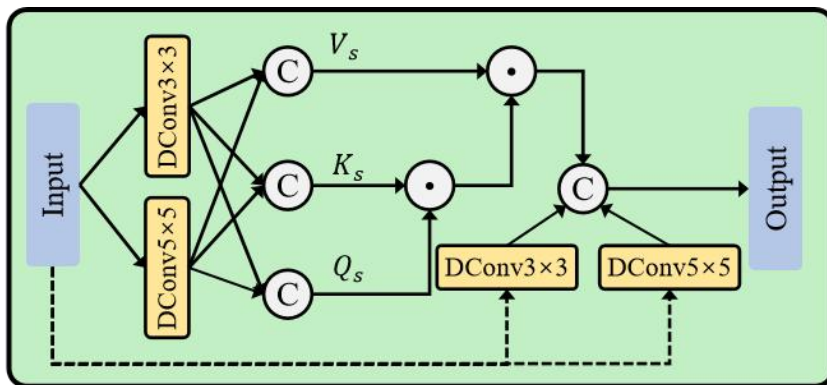
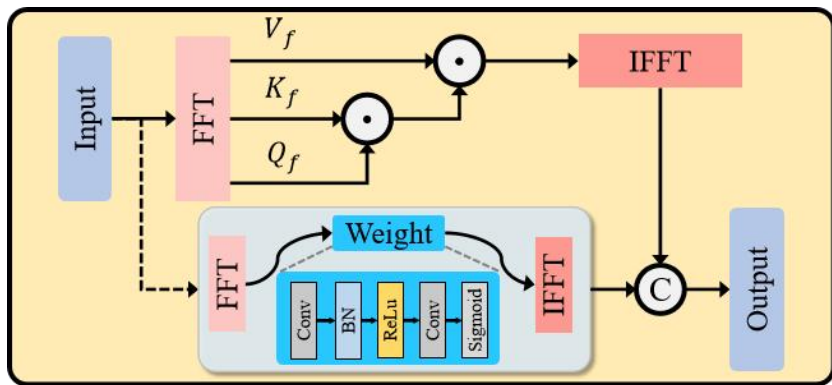
(d) Consecutive low-quality frames

- Challenge 1: Shape Collapse Phenomenon  
Low contrast causes poor boundary delineation and incomplete segmentation.
- Challenge 2: Episodic Amnesia in Videos  
Motion and blur break temporal tracking, causing inconsistency and target loss.



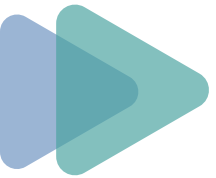


# Method

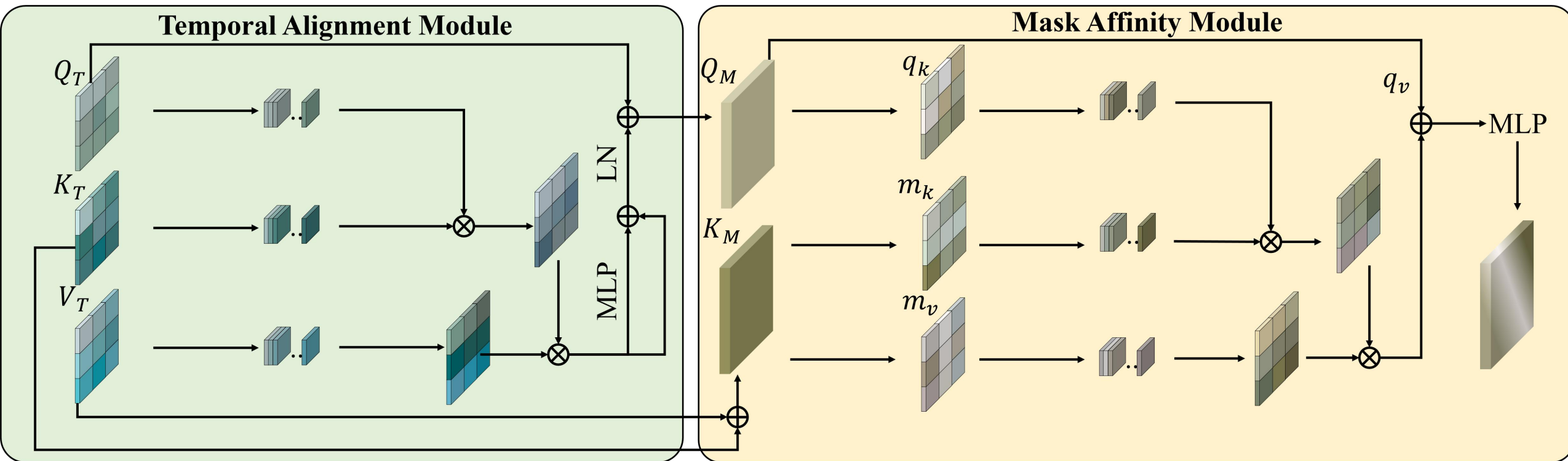


(a) Frequency Filter Block (FFB)    (b) Spatial Refinement Block (SRB)    (c) Interwoven Fusion Block (IFB)

- FFB applies Fast Fourier Transform to capture global semantic distributions.
- SRB preserves fine-grained edge details and local structures, enhancing boundary localization.
- IFB uses gated attention to align and fuse spectral–spatial features for accurate polyp boundary delineation.



# Method



- TAM retrieves relevant context via cross-attention with memory bank.
- MAM aligns current predictions with past masks, maintaining polyp identity under blur, occlusion, or motion.



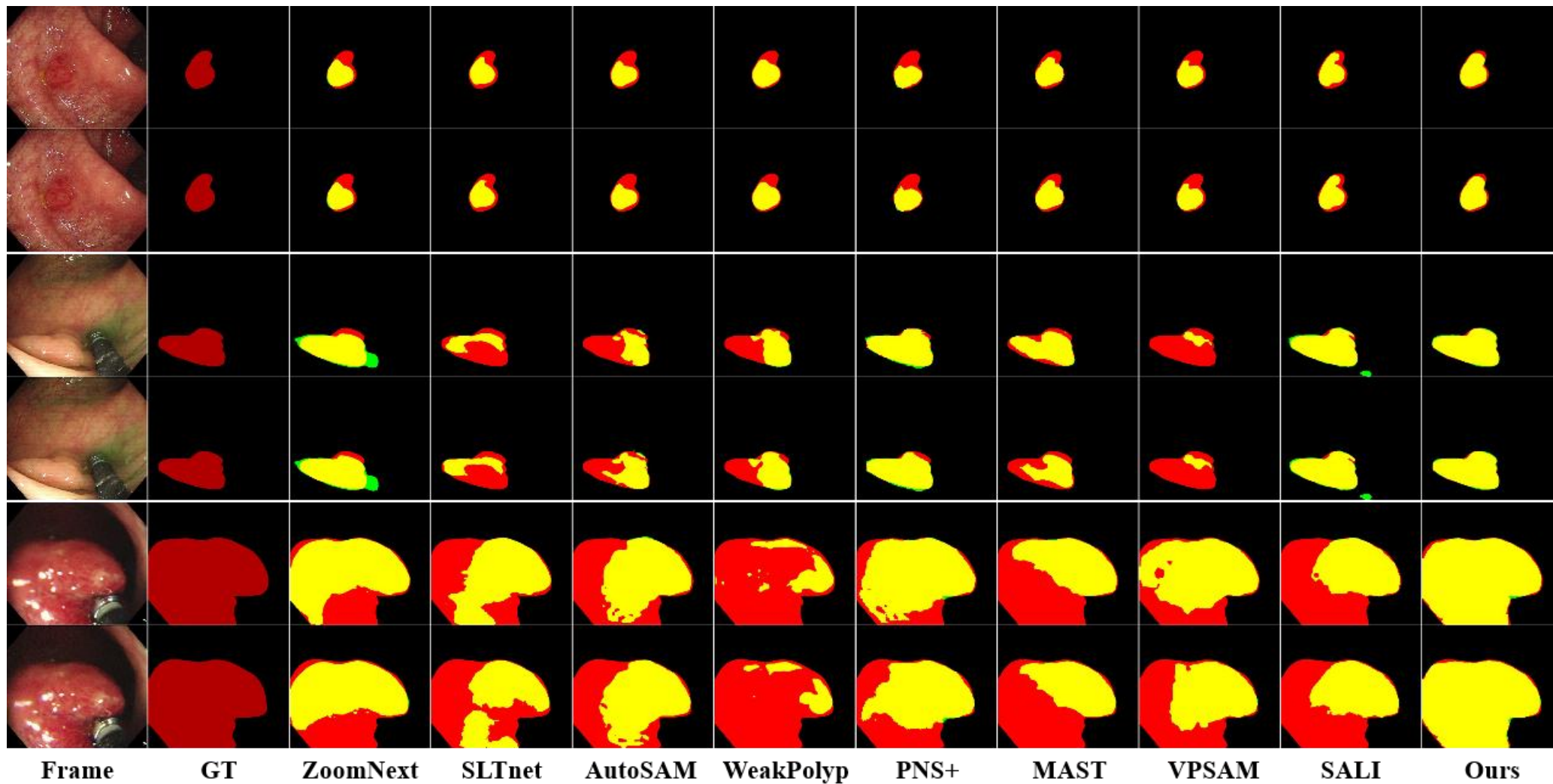
# Experiment

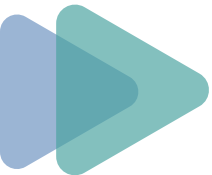
- Evaluation spans SUN-SEG and CVC-612 datasets
- HFSTI-Net achieves 86.27% Dice on SUN-SEG-Hard, outperforming all methods and leading across Dice, IoU, and boundary metrics.

Model	Backbone	Class	SUN-SEG-Easy				SUN-SEG-Hard				CVC-612			
			$S_\alpha$	$E_\phi^{mn}$	$F_\beta^w$	Dice	$S_\alpha$	$E_\phi^{mn}$	$F_\beta^w$	Dice	$S_\alpha$	$E_\phi^{mn}$	$F_\beta^w$	Dice
ZoomNext	PVT-B2	NVS	88.33	90.48	80.66	85.49	87.64	90.84	80.25	83.51	94.66	<u>97.83</u>	92.45	93.17
SLTnet	PVT-B2	NVS	88.13	91.75	83.09	85.91	87.04	90.89	80.98	83.36	<u>94.84</u>	<u>97.37</u>	<u>92.73</u>	<u>93.62</u>
AutoSAM	VIT-B	IPS	86.28	91.67	78.36	81.28	83.57	89.93	73.59	77.37	91.52	95.38	87.47	88.73
WeakPolyp	PVT-B2	IPS	89.04	92.77	<u>83.83</u>	85.27	88.41	<u>92.57</u>	82.93	84.59	91.44	95.78	88.54	88.79
PNS+	Res-50	VPS	86.20	86.17	76.28	82.23	84.29	86.13	72.98	79.60	94.81	96.75	89.63	93.06
MAST	PVT-B2	VPS	84.53	89.81	77.04	78.43	86.17	91.42	77.76	80.32	92.03	95.38	87.47	90.84
VPSAM	VIT-B	VPS	89.31	92.34	82.86	85.62	<u>88.93</u>	92.13	<u>82.98</u>	<u>85.28</u>	93.26	95.75	89.63	92.33
SALI	PVT-B2	VPS	<u>89.54</u>	<u>93.07</u>	83.68	<u>86.17</u>	87.58	91.93	80.56	83.87	91.73	95.21	86.54	88.77
<b>Ours</b>	PVT-B2	VPS	<b>90.73</b>	<b>94.86</b>	<b>85.82</b>	<b>88.03</b>	<b>89.63</b>	<b>93.92</b>	<b>83.26</b>	<b>86.27</b>	<b>95.02</b>	<b>98.46</b>	<b>93.58</b>	<b>94.31</b>



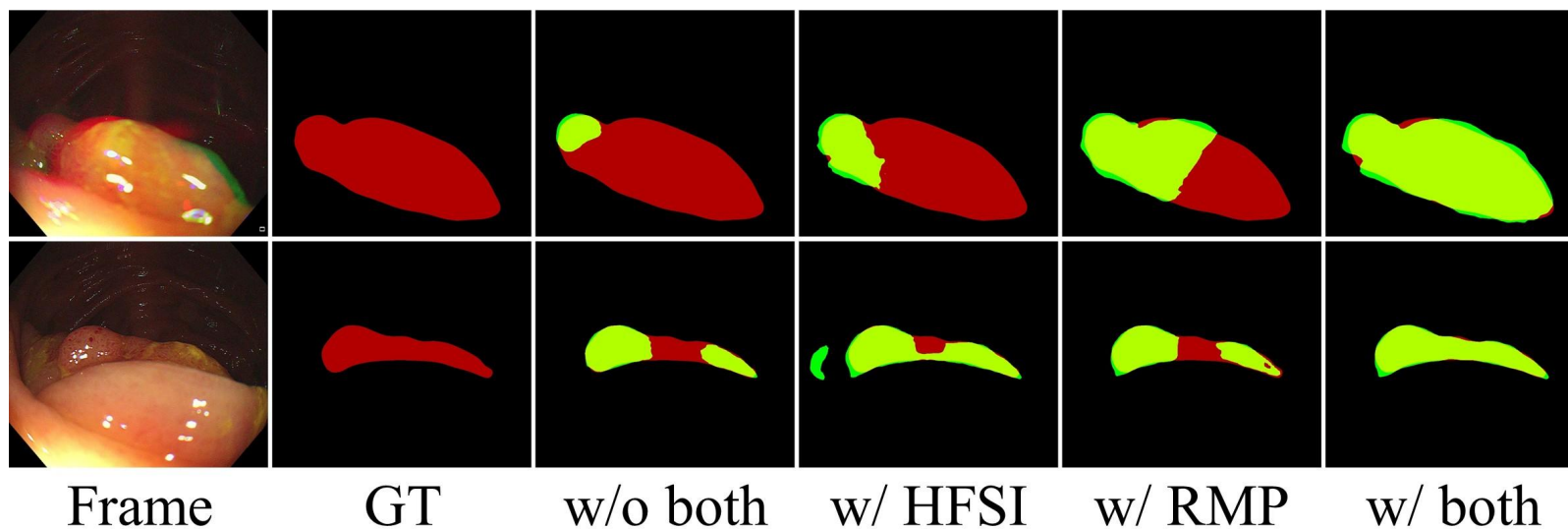
# Experiment

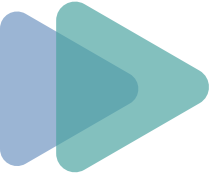




# Ablation Study

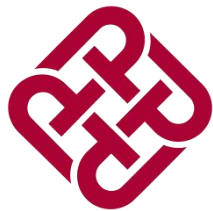
HFSI	RMP	SUN-SEG-Easy				SUN-SEG-Hard			
		$S_a$	$E_\varphi^{mn}$	$F_\beta^\omega$	Dice	$S_a$	$E_\varphi^{mn}$	$F_\beta^\omega$	Dice
		89.51	92.93	83.72	86.20	88.03	92.29	80.62	83.67
	✓	89.71	93.68	84.15	87.04	88.59	92.91	81.24	84.03
✓		90.51	93.87	84.53	87.24	88.97	93.02	81.86	85.27
✓	✓	<b>90.73</b>	<b>94.86</b>	<b>85.82</b>	<b>88.03</b>	<b>89.63</b>	<b>93.92</b>	<b>83.26</b>	<b>86.27</b>





## Conclusion

- HFSTI-Net pioneers joint frequency-spatial-temporal modeling for VPS.
- HFSI eliminates shape collapse via spectral fusion;
- RMP cures episodic amnesia through explicit memory tracking, achieving breakthrough performance.
- State-of-the-art accuracy combined with 31.27 FPS real-time speed overcomes historical efficiency-precision trade-offs.



**ICLR**

**THANK YOU**