







Foundation Models

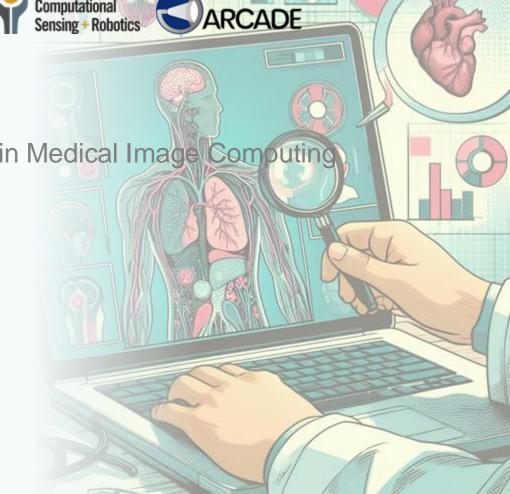
Preparing for the iPhone Moment in Medical Image Computing

MICCAI MedAGI, Marrakech, Morocco

October 6th, 2024

Mathias Unberath, PhD

John C. Malone Associate Professor Department of Computer Science Johns Hopkins University





2007: The iPhone moment

→ New technology resulting in profound transformation

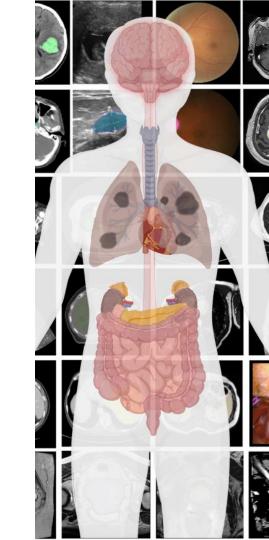




What makes foundation models different?

Formulated to solve "foundational tasks"

- No uber-specific task definition
- Building blocks rather than "complete" solutions



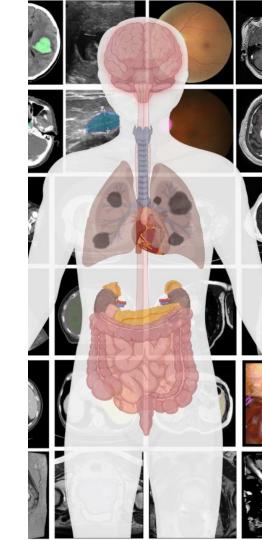
What makes foundation models different?

Formulated to solve "foundational tasks"

- No uber-specific task definition
- Building blocks rather than "complete" solutions

Flexibility enables scaling

- Dataset mixing (→ no "semantics")
- Self-supervised learning



What makes foundation models different?

Formulated to solve "foundational tasks"

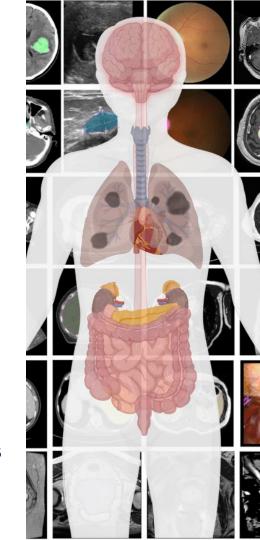
- No uber-specific task definition
- Building blocks rather than "complete" solutions

Flexibility enables scaling

- Dataset mixing (→ no "semantics")
- Self-supervised learning

Scaling results in (better) generalization

- Larger variety in training sets
- Often: Applicable "out of the box" to new problem domains



Foundation Models – iPhone Moment for MIC?

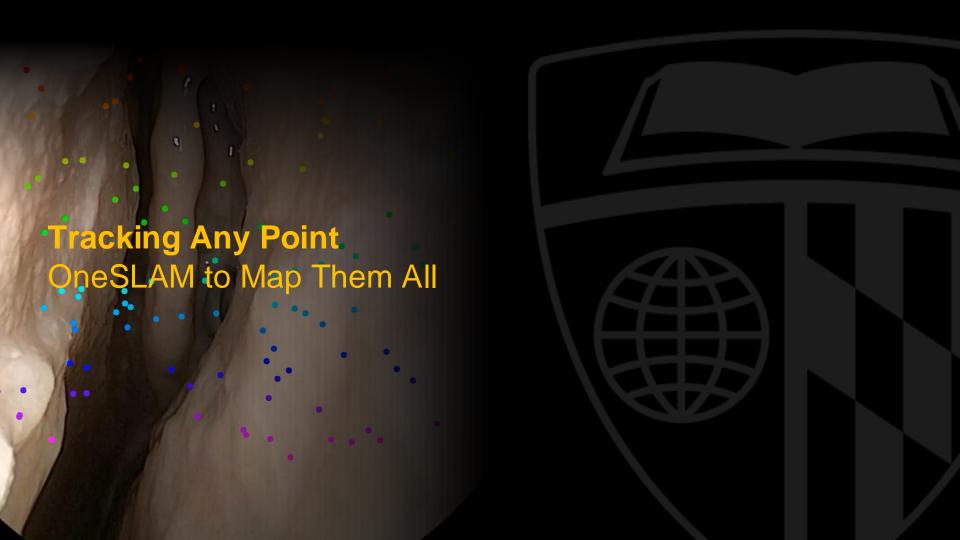
Potential for catalytic effect in medical image computing and beyond

Flexible task formulations

→ Strong, generalizable building blocks

"Out of the box generalization"

- → Dramatically reduced dev. Costs
- → Potential for ubiquitous adoption of quantitative tools



Quantitative Endoscopy

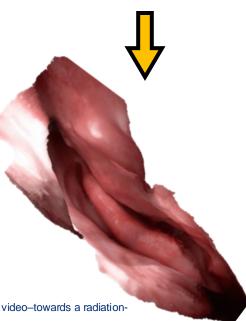
Endoscopy (and related modalities) is qualitative

- Want: Endoscope tracking and reconstruction (SLAM)
- Enables quantitative analysis

Persistent challenges for computer vision in endoscopy

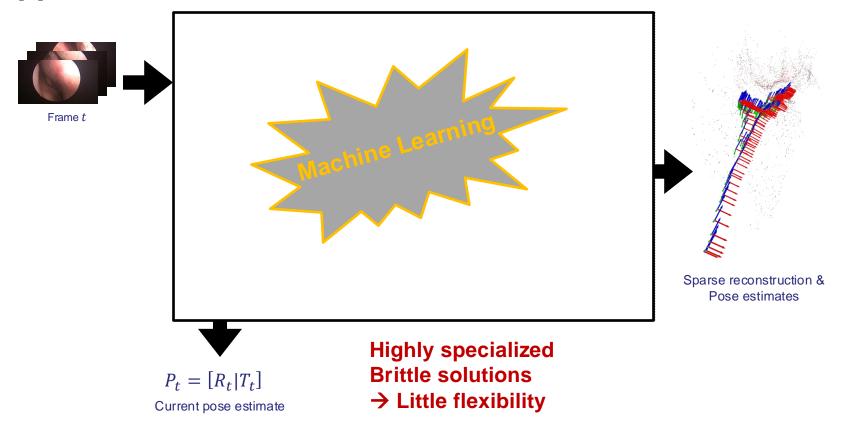
- Lack of photometric constancy (same looks different)
- Scarce features, repetitive texture





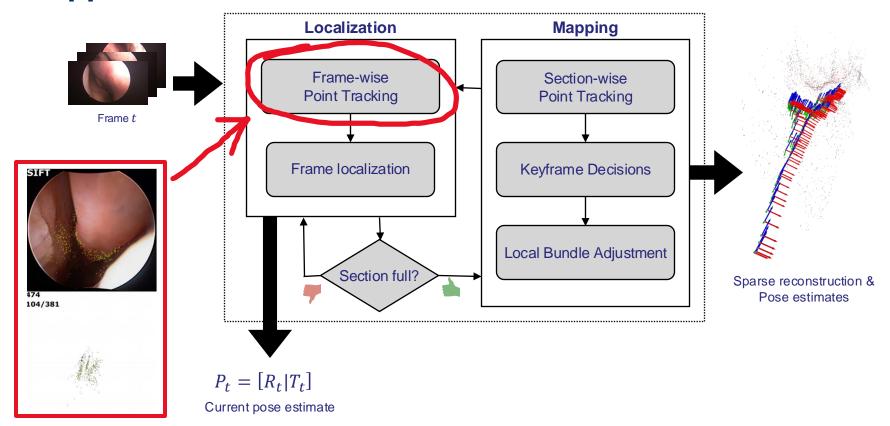
Liu, X., Stiber, M., Huang, J., Ishii, M., Hager, G. D., Taylor, R. H., & Unberath, M. (2020). Reconstructing sinus anatomy from endoscopic video—towards a radiation-free approach for quantitative longitudinal assessment. In MICCAI 2020.

Opportunities for Foundation Models



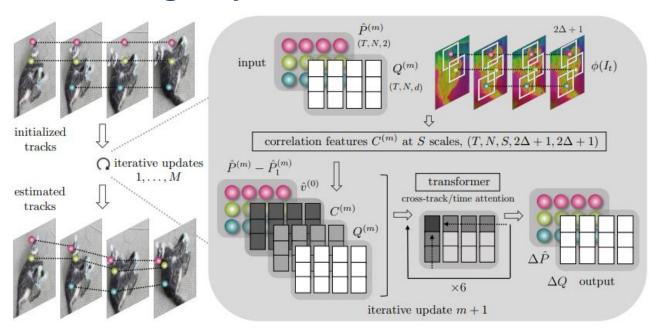
Liu, X., Zheng, Y., Killeen, B., Ishii, M., Hager, G. D., Taylor, R. H., & Unberath, M. (2020). Extremely dense point correspondences using a learned feature descriptor. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition (pp. 4847-4856).

Opportunities for Foundation Models



Liu, X., Zheng, Y., Killeen, B., Ishii, M., Hager, G. D., Taylor, R. H., & Unberath, M. (2020). Extremely dense point correspondences using a learned feature descriptor. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition (pp. 4847-4856).

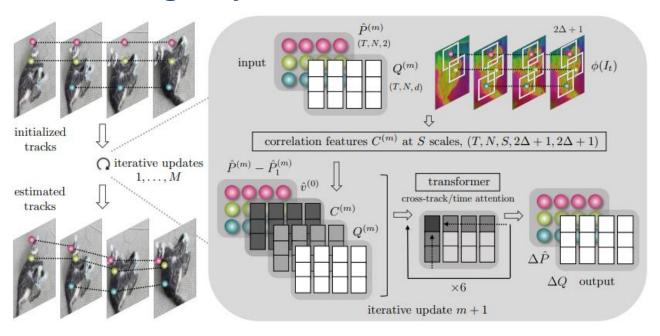
Tracking Any Point Foundation Models



Iteratively update point tracks of arbitrary points

- Cross-track and cross-time attention
- Modern, transformer-based architecture

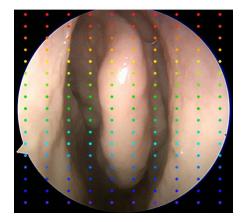
Tracking Any Point Foundation Models



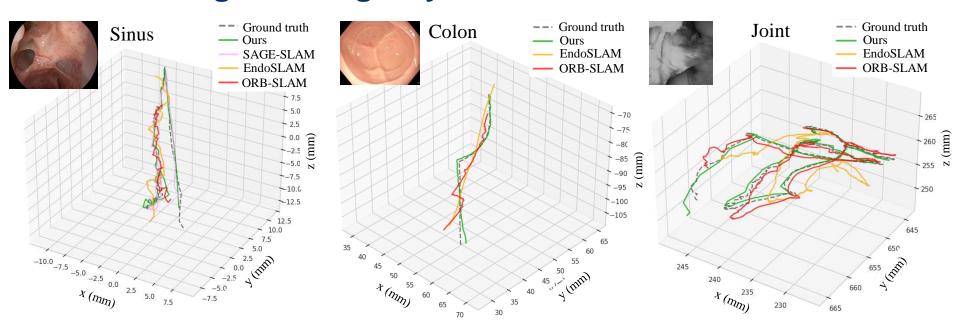


- Cross-track and cross-time attention
- Modern, transformer-based architecture





SLAM using Tracking Any Point Foundation Models



OneSLAM based on Tracking Any Point Foundation Model

- First FM-based SLAM (at least for endoscopy)
- Outperforms both: Conventional approaches & domain-specific ML methods



Segment Anything Models (SAM)

Segmentation is a key enabling task

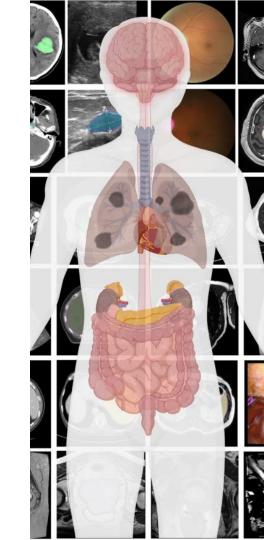
- "Making sense" of medical images
- Enables quantitative analysis

Foundation models for segmentation

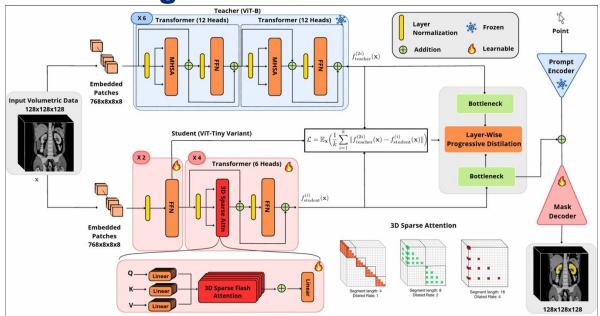
- Replace semantics with user query
- Interactivity ←→ Generalization

Two challenges

- 1. Scaling SAM to volumetric images
- 2. Clawing back semantic segmentation benefits



Scaling SAM to High-resolution Volumes



Iteratively update point tracks of arbitrary points

- Layer-wise progressive distillation to a VIT-Tiny
- 3D sparse flash attention

FastSAM3D

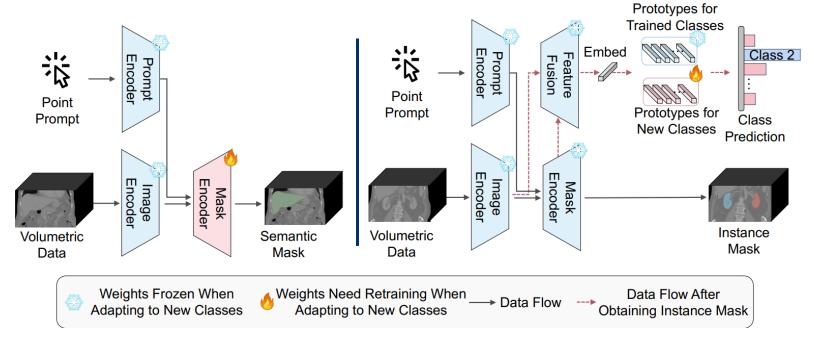
Dim	Method			Encode	er		Decode	er	Acceleration		
		Resolution	Time	FLOPs	Memory			_	To 2D ↑	To 3D ↑	
			$(ms)\downarrow$	$(G)\downarrow$	$(Gb)\downarrow$	$(ms)\downarrow$	$(G)\downarrow$	$(\mathrm{Gb})\!\!\downarrow$	10 2D		
2D	SAM [16]	1024×1024	3980	369.0	7.87	239	3.0	5.57	1.00×	/	
	MobileSAM [32]	1024×1024	584	36.7	5.48	233	3.0	5.27	$5.16 \times$	/	
	TinySAM [24]	1024×1024	609	36.7	5.48	246	3.0	5.27	4.93×	/	
	MedSAM [18]	1024×1024	3983	369.0	7.87	241	2.9	5.57	1.00×	/	
	SAM-Med2D [4]	256×256	1063	32.0	6.32	216	0.21	5.55	$3.30 \times$	/	
3D	SAM-Med3D [27]	$128 \times 128 \times 128$	70	89.5	6.58	20	2.8	5.53	$60.27 \times$	1.00×	
	FastSAM3D	$128 \times 128 \times 128$	3	21.9	0.78	5	2.8	0.71	$527.38 \times$	$8.75 \times$	

Dim	Method	AMOS [13]			TotalSegmentator [[28] BraTS [20]					
		1pt	3pt	5pt	10pt	1pt	3pt	5pt	10pt	1pt	3pt	5pt	10pt
2D	SAM [16]	0.049	0.093	0.114	0.145	0.202	0.279	0.311	0.348	0.108	0.192	0.217	0.237
	MobileSAM [32]	0.041	0.056	0.063	0.070	0.149	0.170	0.182	0.212	0.079	0.132	0.156	0.186
	TinySAM [24]	0.049	0.077	0.089	0.101	0.171	0.225	0.243	0.262	0.103	0.165	0.187	0.211
	MedSAM [18]	0.004	0.051	0.060	0.074	0.006	0.069	0.090	0.111	0.008	0.059	0.064	0.071
	SAM-Med2D [4]	0.097	0.127	0.129	0.132	0.008	0.081	0.100	0.128	0.013	0.076	0.082	0.084
3D	SAM-Med3D [27]	0.289	0.386	0.418	0.448	0.252	0.400	0.463	0.522	0.328	0.395	0.418	0.446
	FastSAM3D	0.273	0.368	0.402	0.437	0.250	0.378	0.445	0.519	0.333	0.401	0.421	0.445

Convincing performance (for SAM models)

- Best or second-best performance across several benchmarks
- 3ms execution! >500x / >8x speed-up compared to 2D / 3D methods

Leveraging the Best of Both Worlds – ProtoSAM3D



Automated semantic segmentation while retaining SAM's flexibility

- Two-stages: 1) Coarse prototype identification; 2) Fine segmentation
- Prompting still possible! New prototypes can be added anytime!

Leveraging the Best of Both Worlds – ProtoSAM3D

AMOS-8 (Gen)

IoU

70.53

77.32

77.32

70.57

79.54

83.07

84.27

26.07

30.17

34.74

36.52

35.23

78.32

68.97

85.21

85.95

Acc

81.25

82.15

87.32

79,99

85.32

88.21

87.35

86.34

78.45

76.89

89,32

89.64

IoU

75.23

84.13

78.54

72.34

84,78

Acc

81.45

85.12

88.57

80.34

87.23

89.54

90.59

86.76

82.11

80.56

92.12

93.00

DSC

79.75

81.34

85.23

79,46

82.09

84.24

86.34

35.21

39.64

40.12

45.42

42.56

84.23

75.23

87.34

87.98

AMOS-7 (Zero)

IoU

23.76

33.56

33.07 31.56

28.45

31.74

32.65

DSC

36.43

40.21

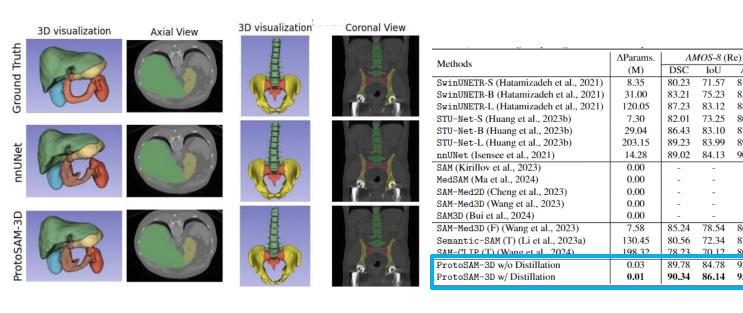
41.35

44.76

39.87

45.76

46.32

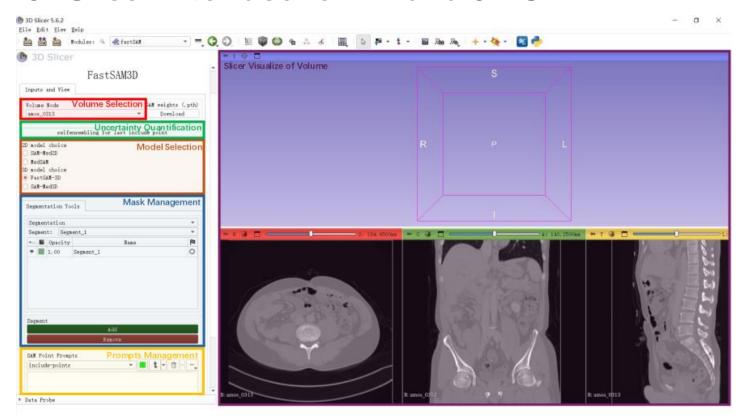


Automated semantic segmentation while retaining SAM's flexibility

- Convincing qualitative and quantitative performance
- Among first models for semantic and zero-shot instance segmentation

Shen, Y., Dreizin, D., Inigo Romillo, B., & Unberath, M. (2024). ProtoSAM3D: Interactive Semantic Segmentation in Volumetric Medical Imaging via a Segment Anything Model and Mask-Level Prototypes, Under review.

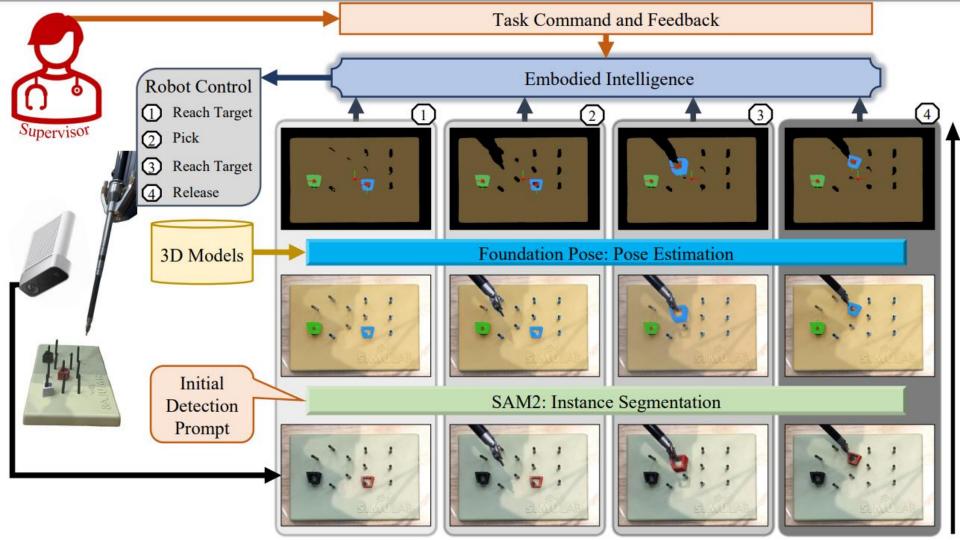
3DSlicer Interface for 2D and 3D SAM





Shen, Y., Shao, X., Romillo, B. I., Dreizin, D., & Unberath, M. (2024). FastSAM3DSlicer: A 3D-Slicer Extension for 3D Volumetric Segment Anything Model with Uncertainty Quantification. MICCAI MedAGI Workshop.











Towards Robust Automation of Surgical Systems via Digital Twin-based Scene Representations from Foundation Models

Hao Ding et al.

Department of Computer Science, Johns Hopkins University



Foundation Models for Medical Image Computing

This is the beginning – not the end!

- Movement towards foundation models is young
- Many open challenges

Opportunities and open problems

- Data Flexibility helps but not always sufficient
- Flexible problem formulations and task design to enable scaling

Many reasons for optimism and excitement!

- High accessibility
- Transformational impact on MIC and precision healthcare

