

Segmentation Models

Tanveer Syeda-Mahmood

Segmentation models

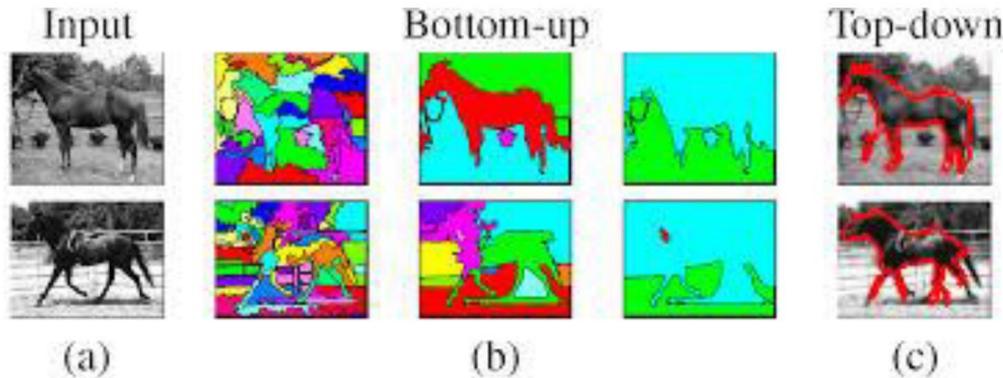
- Types of segmentation models
 - Semantic, instance, panoptic
- Major segmentation architectures
 - CNN, U-net
 - Transformer-based
 - VLM-based
 - Hybrid
- Segmentation foundational model (input-wise)
 - SAM
- Medical imaging-specific architectures
 - MedSAM, Instruction-tuned segmentation models
 - Custom architectures

The different notions of segmentation

- Computer vision
 - Image partitioning into objects and background
- Medical imaging
 - Anatomy segmentation
 - Anomaly segmentation

Bottom-up?
Or top-down?

Supervised?
Or unsupervised?



The different notions of segmentation

- Computer vision
 - Image partitioning into objects and background
- Medical imaging
 - Anatomy segmentation
 - Anomaly segmentation

Semantic segmentation

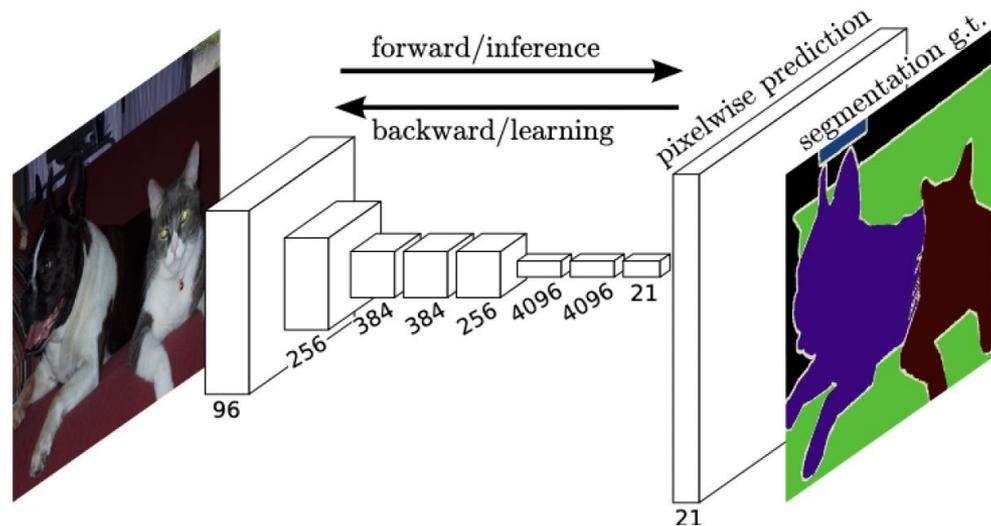


- Labels each pixel by class
- E.g. FCNet, DeepLab, PSPNet

Semantic segmentation example: FCN

- Fully-supervised method
- Used in many other networks for the segmentation head
- Using purely convolutional setup

1. Feature extraction (Encoder)
Edges, shapes, organs
2. Upsampling/decoding
transposed convolutions or interpolations
3. 1x1 convolution
4. Pixel-wise classification



<https://arxiv.org/abs/1411.4038>

**Fully Convolutional Networks for
Semantic Segmentation**

The different notions of segmentation

- Computer vision
 - Image partitioning into objects and background
- Medical imaging
 - Anatomy segmentation
 - Anomaly segmentation

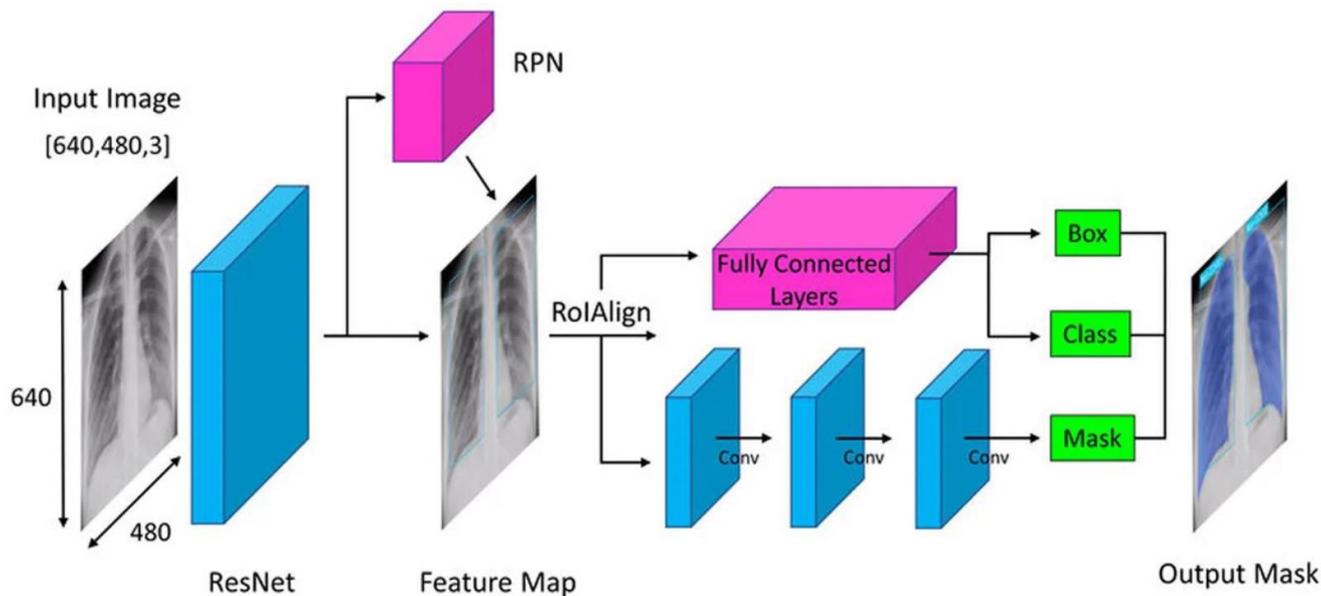
Instance segmentation



- Each region has a distinct label
- Combine object detection and semantic segmentation
- E.g. Mask R-CNN, Yolo-8

Instance segmentation example: Mask RCNN

- Up to 1500 class labels can be detected and labeled



<https://arxiv.org/abs/1703.06870>, K.He et al. Mask R-CNN

>48,000 citations

The different notions of segmentation

- Computer vision
 - Image partitioning into objects and background
- Medical imaging
 - Anatomy segmentation
 - Anomaly segmentation

panoptic segmentation



- Labels each pixel by class and also identifies different instances of the same class
- E.g. Efficient PS, Panoptic DeepLab, RS-DINO, HyperDETR, DinoV2, FCN+MaskRCN

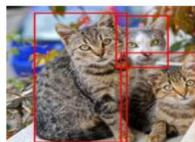
Panoptic segmentation(DETR)

Detection transformer

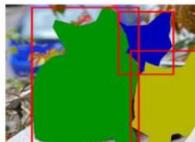
- Uses the transformer approach for capturing spatial relationship between regions
- Uses feature encoding
- Class and regional predictions
- Supervision needed
- Limited object support
- Panoptic segmentation



Classification



Detection



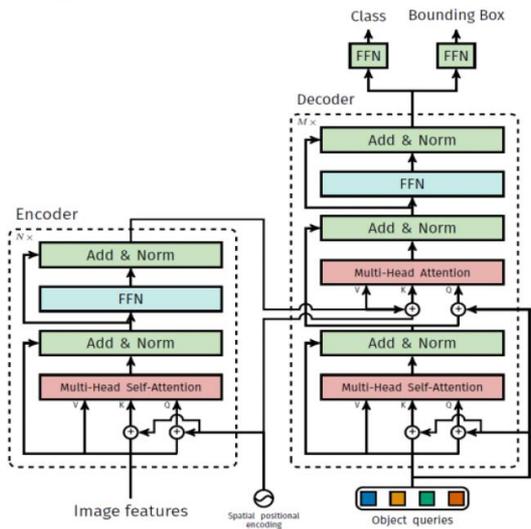
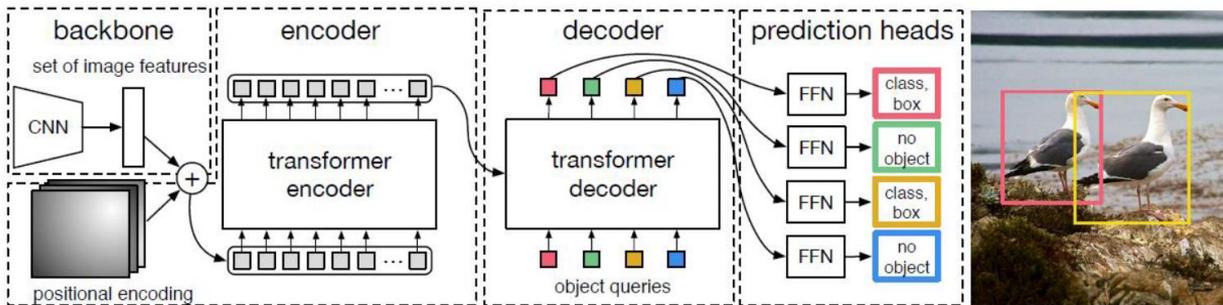
Instance segmentation



Semantic segmentation

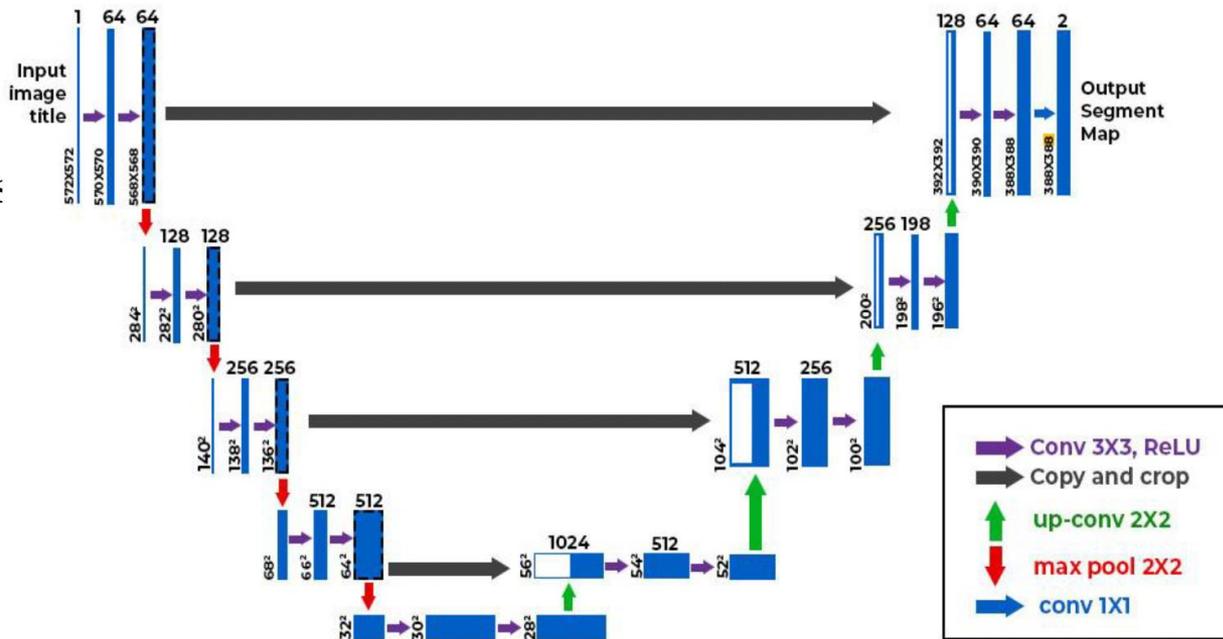


Panoptic segmentation



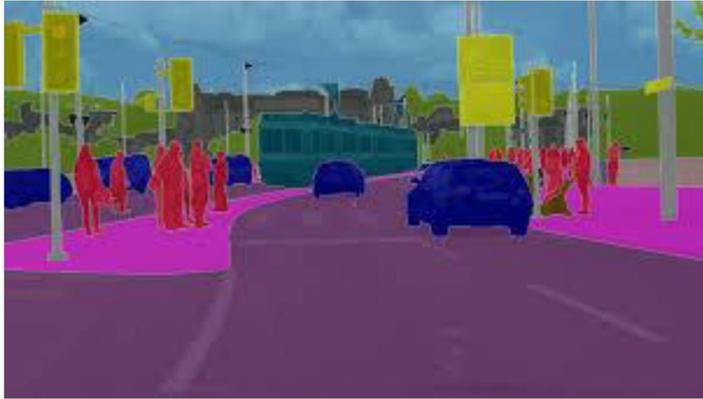
U-net - A classic CNN-style segmentation model

- Extensively used in medical imaging
- Works for both anatomy and anomaly detection
- Top-down segmentation works best for medical images since intensities are intermingled.
- Works well when the pattern to learn is in a fixed place
- Skip connections help remember features
- Fully CNN-based



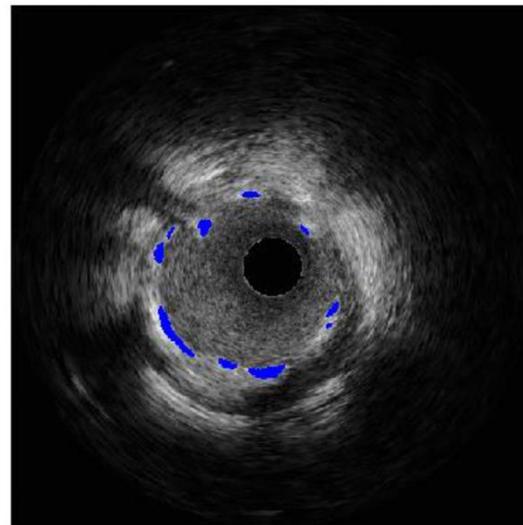
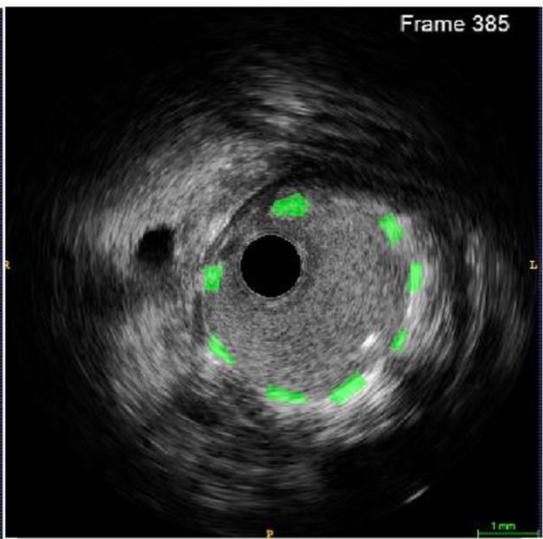
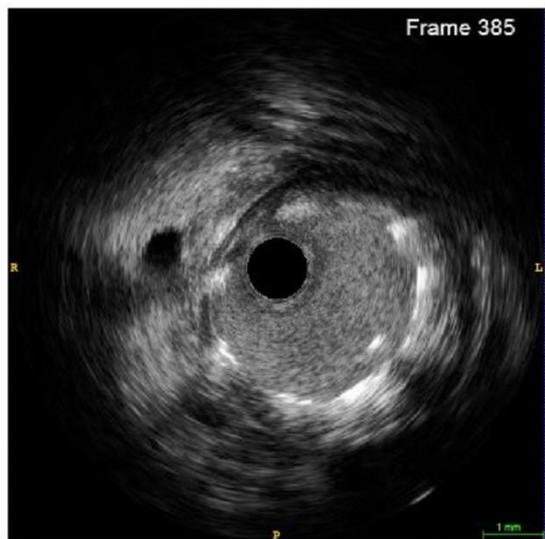
U-net applications

- Both medical and non-medical applications



U-Net issues

- Detections are not smooth
- Maintaining continuity is difficult
- Fragmentation can occur
- Small tumors and anatomical regions
- Stable contour detection



Basic architectures for segmentation FM

- Based on two major architectures:
 - CNN
 - Transformers
- Different combinations of encoders and decoders
 - U-net : Convolutional encoder and decoder
 - TransU-net: Transformer encoder and CNN decoder
 - SETR: Vision transformer and CNN decoder
 - Segformer: transformer encoder and decoder
 - CLIPSeg: CLIP encoder and transformer decoder
- Go from limited object labels to open vocabulary segmentation
- Mostly still supervised, but newer unsupervised open set models also coming back

Design choices for encoders and decoders for segmentation

1. Encoder choice depends on:

- Input modality (2D vs 3D)
- Need for global context (transformers excel)
- Pretraining availability

2. Decoder choice depends on:

- Desired output resolution
- Computational budget
- Multi-scale vs single-scale features

3. Skip connections are critical to:

- Recover fine-grained spatial details
- Preserve edges, organ boundaries, or object contours

Combinations of encoder and decoders

Model	Encoder	Decoder	Notes
U-Net	Custom CNN (Conv blocks)	Symmetric upsampling + skip connections	Widely used in medical imaging
UNet++	CNN	Nested skip connections	Better feature fusion
SegNet	VGG16	Upsampling with pooling indices	Lightweight, efficient
DeepLabV3+	ResNet / Xception	Atrous Spatial Pyramid Pooling (ASPP) + upsampling	Multi-scale context
FPN (Feature Pyramid Network)	ResNet / DenseNet	Top-down pyramid merging	Often used for detection + segmentation
SegFormer	MiT (Transformer encoder)	Lightweight MLP-based decoder	High performance with attention + low compute
Mask2Former	Swin Transformer / ResNet	Transformer decoder with queries	Supports panoptic, instance, semantic segmentation
Attention U-Net	CNN	Skip connections + attention gates	Focuses decoder on relevant encoder features
3D U-Net / V-Net	3D CNN	3D upsampling + skip connections	Volumetric medical image segmentation

Making segmentation model foundational

- Train with a variety of datasets
- Adapt transformer or CNN-based encoders and decoders
- An early example of a foundational model for medical image segmentation – based on V-net

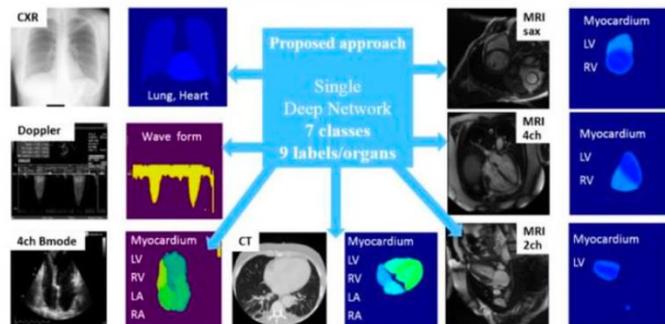
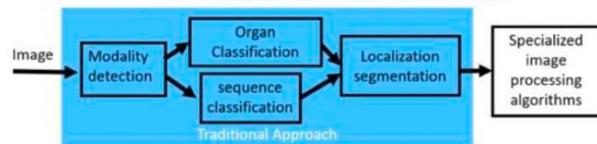
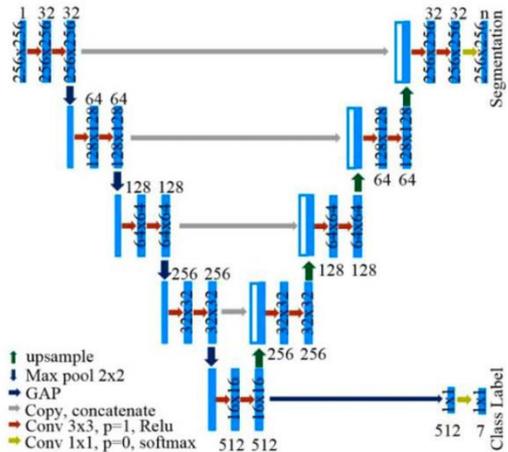


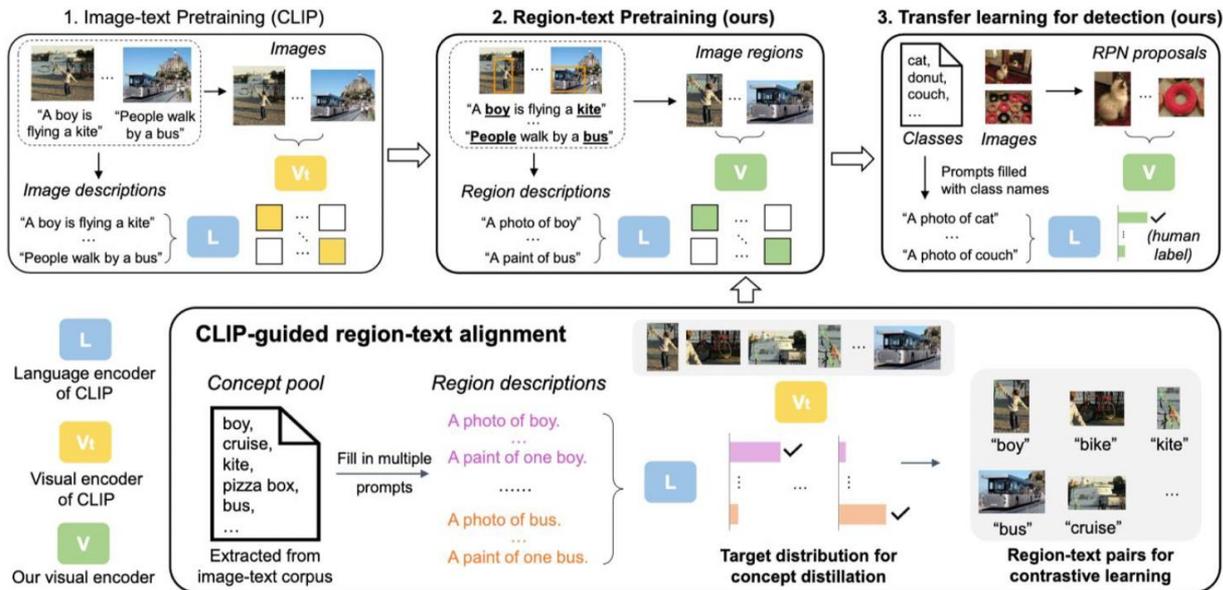
Figure 1: Top: Traditional architecture tackling one problem at a time. Bottom: our proposed network to both classify different modalities with different viewpoints (X-ray, CT, Ultra sound, 2 chamber MRI, 4 chamber MRI, short axis MRI) as well as segment different structures as Lung, heart, Doppler wave form, Myocardium (Myo), left ventricular (LV), right ventricular (RV), left atrium (LA), and right atrium (RA).

Harouni et al., "Universal multimodal deep network for classification and segmentation of medical images," in ISBI 2018.

Segmentation models derived from VLM models

- Use joint visual-textual information to aid segmentation
- Use CLIP underneath for the encoding
- Built separate decoders to aid in segmentation
 - CNN-based
 - Transformed-based decoders
- Segmentation at the level of bounding boxes
 - RegionCLIP
 - VILD
 - GLIP
- Segmentation as full regions
 - CLIPSeg

VLM – Visual Grounding



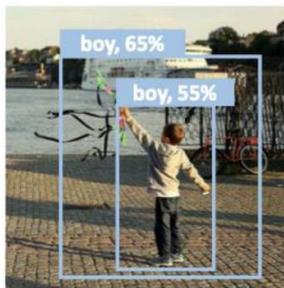
RegionCLIP – Contrastively learned region-text alignment

Using VLM for generalized instance segmentation

- Use region crops from RPN
- Initial CLIP image-to-text to label
- What if we do region detection and CLIP on regions fails to recognize objects

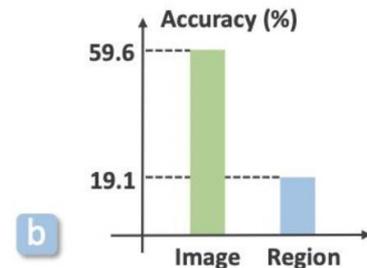
The labels should be known at inference time

Cropped image regions recognized by CLIP

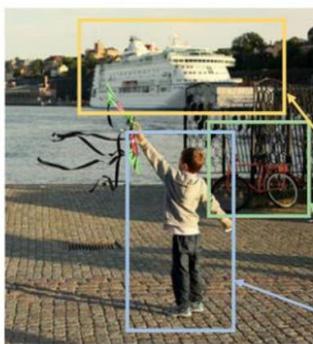


a

Image classification (ImageNet)
Region classification (LVIS)



b



c



"A boy is flying a kite."

Image-text matching (CLIP)

"A photo of one cruise."

"A bad photo of a bike."

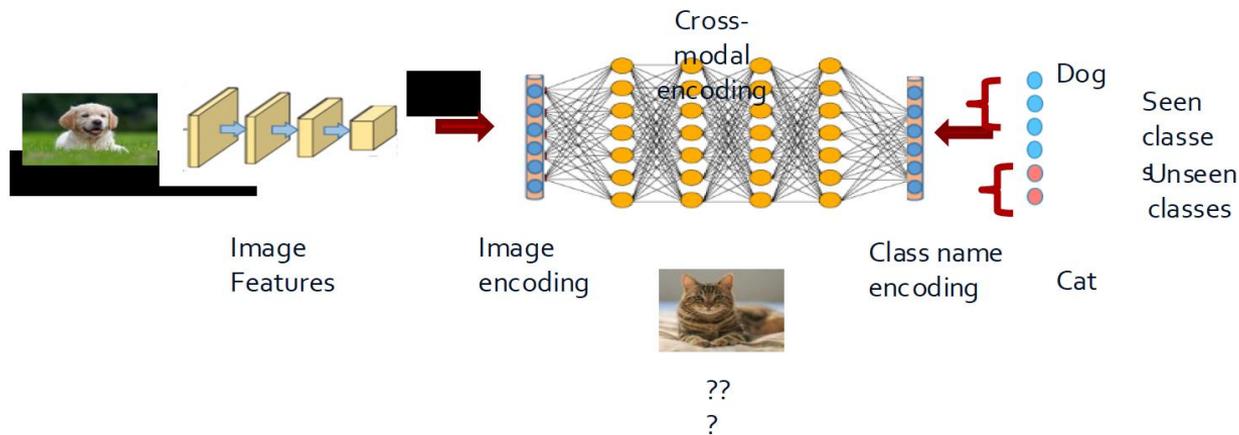
"A photo of a boy."

Region-text matching (Ours)

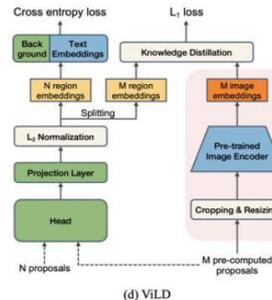
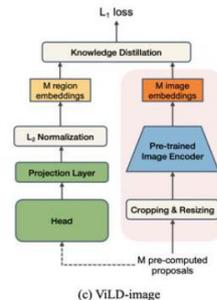
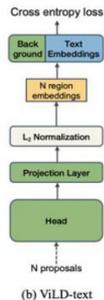
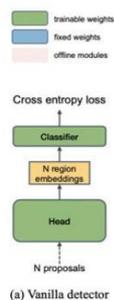
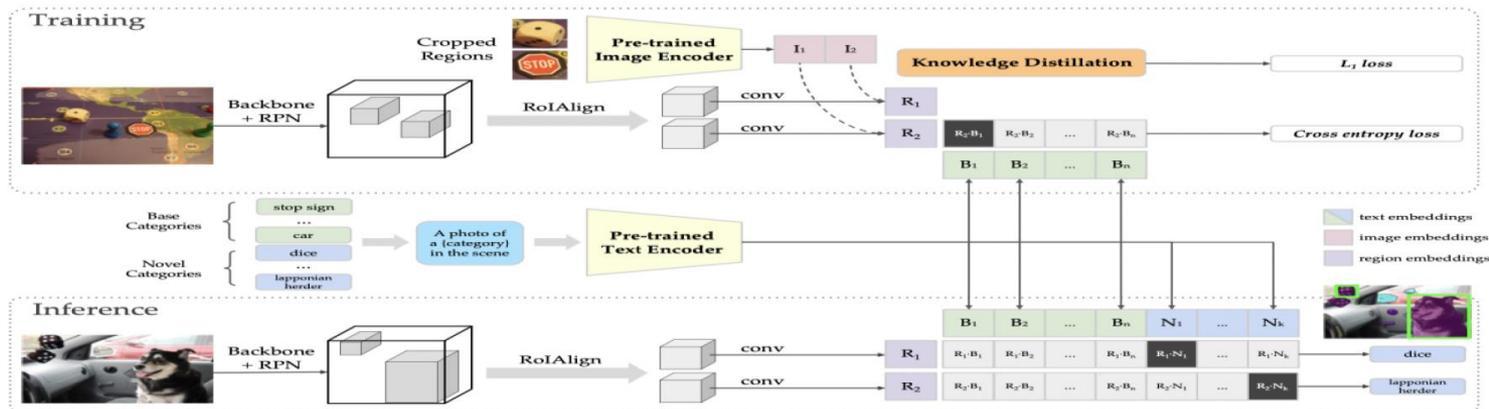
RegionCLIP – Contrastively learned region-text alignment

VILD architecture

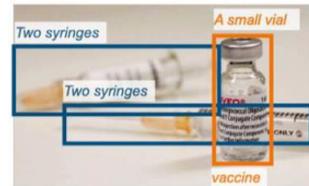
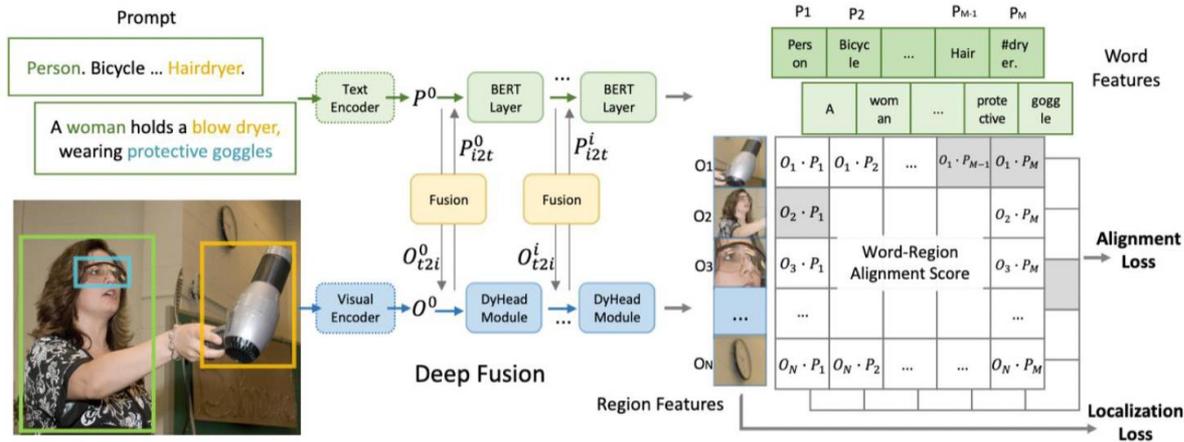
- Combines pre-trained FM (CLIP) with MaskRCNN-style RPN and zero-shot learning methods
- Distill the knowledge from the alignment of region embedding to image embeddings of cropped regions during training.



ViLD – Visual Grounding



VLM – Phrase Grounding



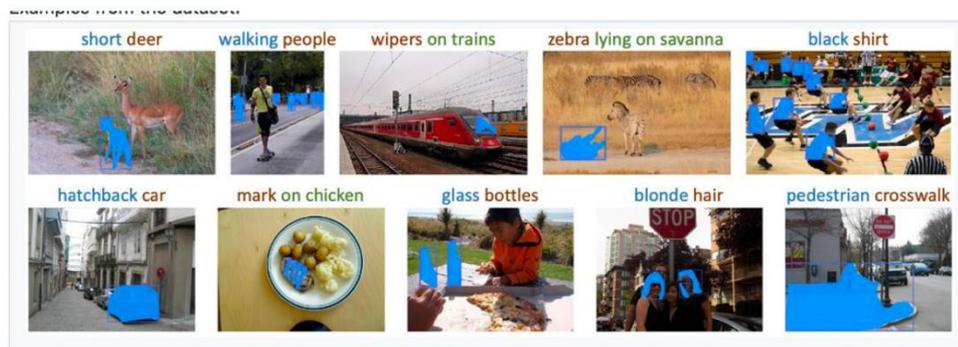
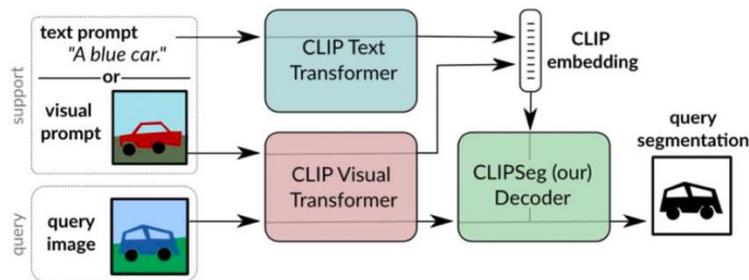
Two syringes and a small vial of vaccine.



playa esmeralda in holguin, cuba. the view from the top of the beach. beautiful caribbean sea turquoise

- Given an image and a corresponding caption, the **Phrase Grounding** task aims to ground each entity mentioned by a noun phrase in the caption to a region in the image.
- GLIP – grounded language pre-training

CLIPSeg

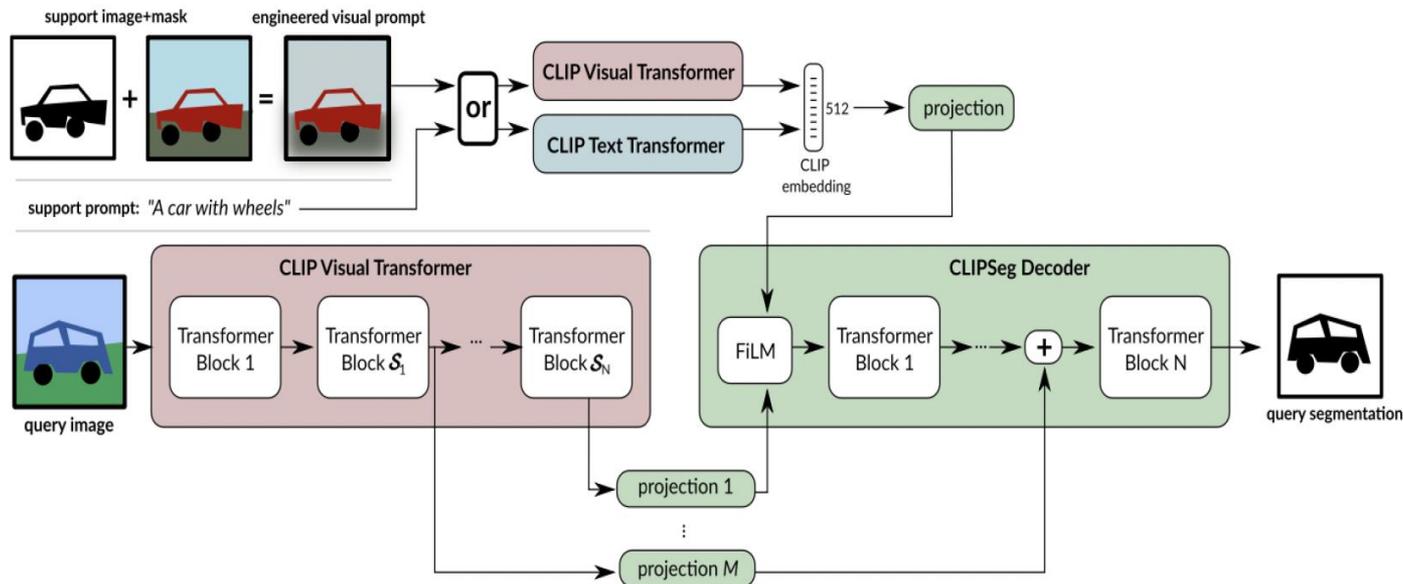


- Uses the PhrasCut dataset of regional masks derived from VisualGenome
- Combined CLIP with a lightweight decoder

CLIPSeg: Image segmentation using text and image prompts

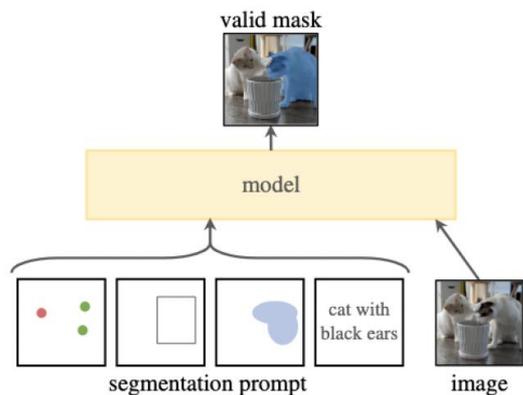
ClipSeg

- Skip connections to preserve the visual context from 3 of the layers of the ViT to the decoder as in U-net
- The decoder incoming dimension is 64

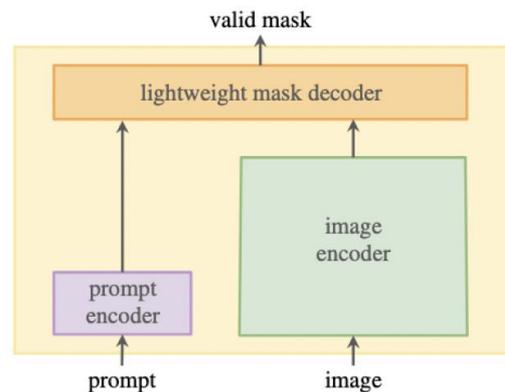


Segment Anything (SAM)

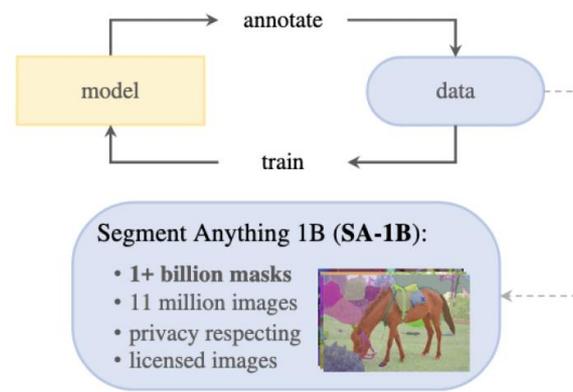
- A foundation model for image segmentation.
- A promptable model and pre-train it on a broad dataset using a task that enables powerful generalization



(a) **Task:** promptable segmentation



(b) **Model:** Segment Anything Model (SAM)



(c) **Data:** data engine (top) & dataset (bottom)

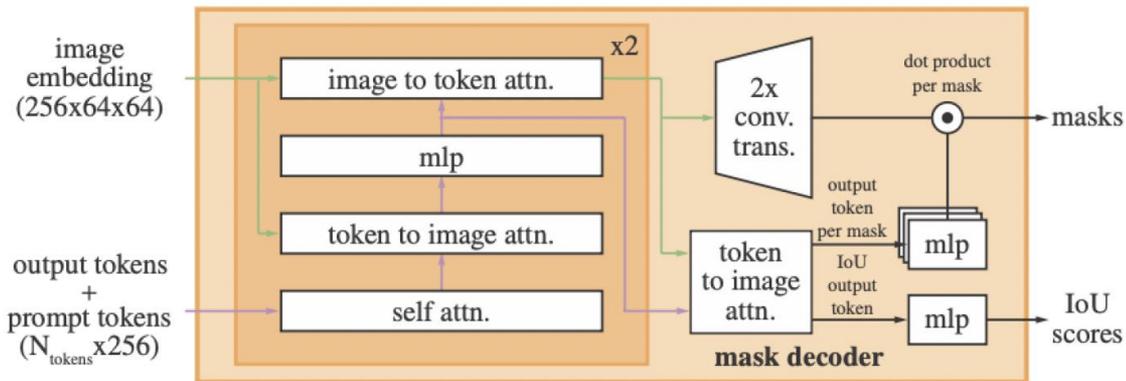
Different kinds of prompts to aid in segmentation

SAM components

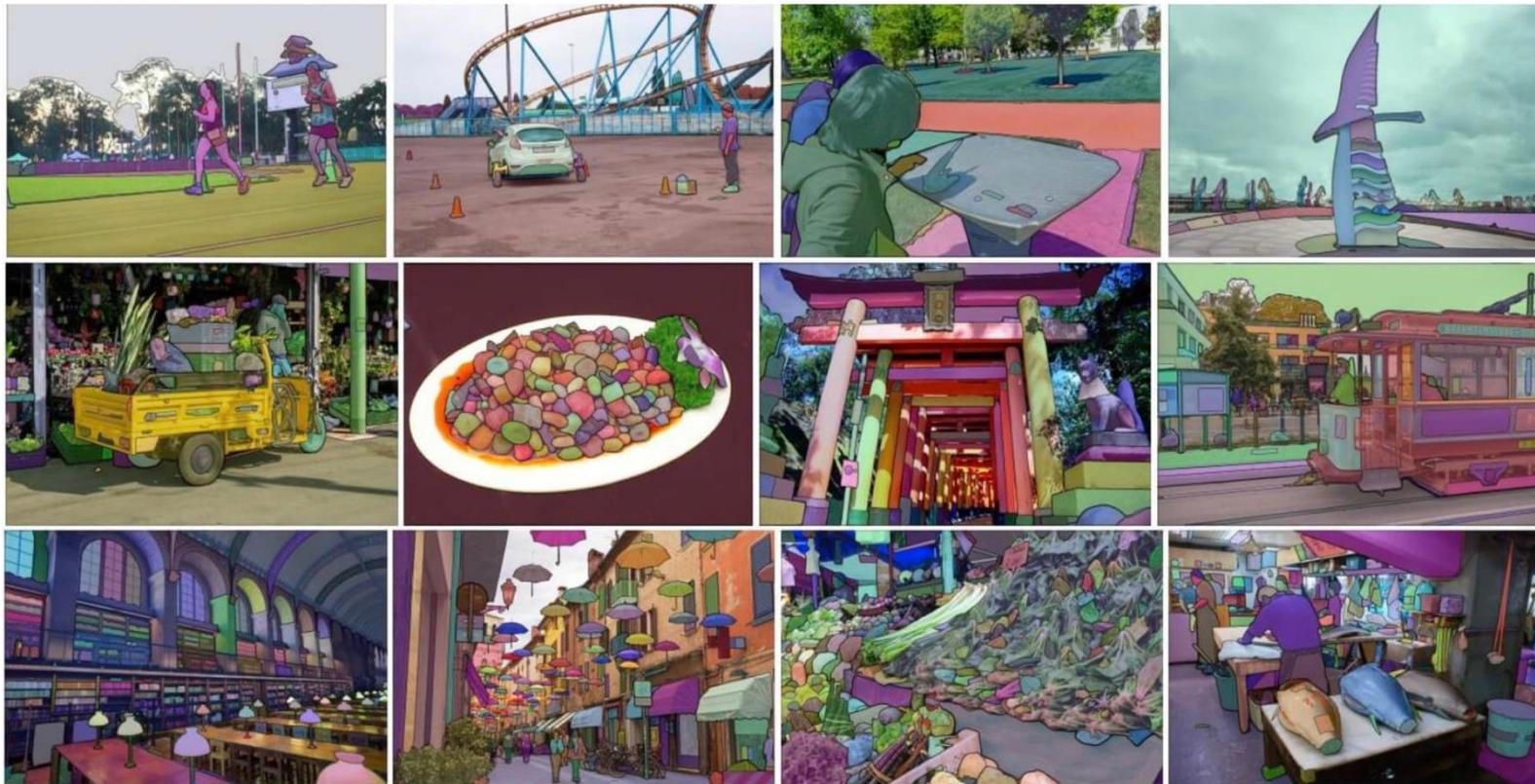
- Image Encoder:
 - Encoder portion of the Masked Auto-encoder which is a pre-trained Vision Transformer (ViT) with adaptations to process high resolution inputs,
 - input resolution of 1024×1024 obtained by rescaling the image and padding the shorter side. The image embedding is therefore 64×64×256.
- Prompt encoder:
 - Produces a 256-dimensional vector
 - Points
 - Point is represented as a sum of positional encoding of the points's location and one of the two learned embeddings to indicate either a **foreground point** or **background point**.
 - Bounding boxes:
 - Boxes are represented by an **embedding pair of corners**.
 - Text:
 - CLIP text-encoding
 - Masks:
 - Downscaled versions of masked images CNN-style and flattened to 256-dimensional vectors

Segment Anything – Mask Decoder

- Image embeddings and prompt embeddings are mapped to the final mask
- uses prompt self-attention and cross-attention in two directions (prompt-to-image embedding and vice-versa) to update *all* embeddings.
- MLP maps the output token to a dynamic linear classifier,
- which then computes the mask foreground probability at each image location.

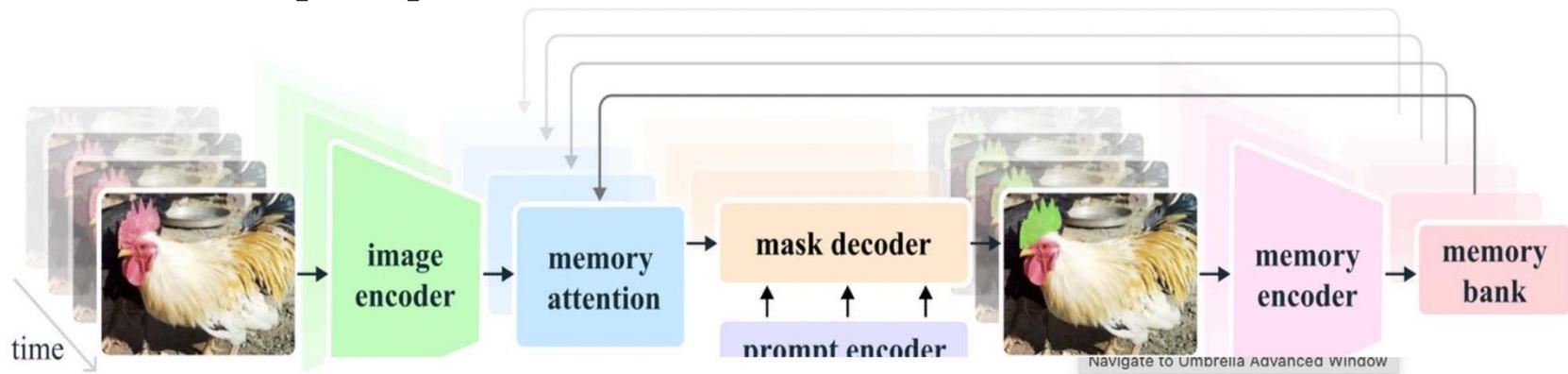


SAM results

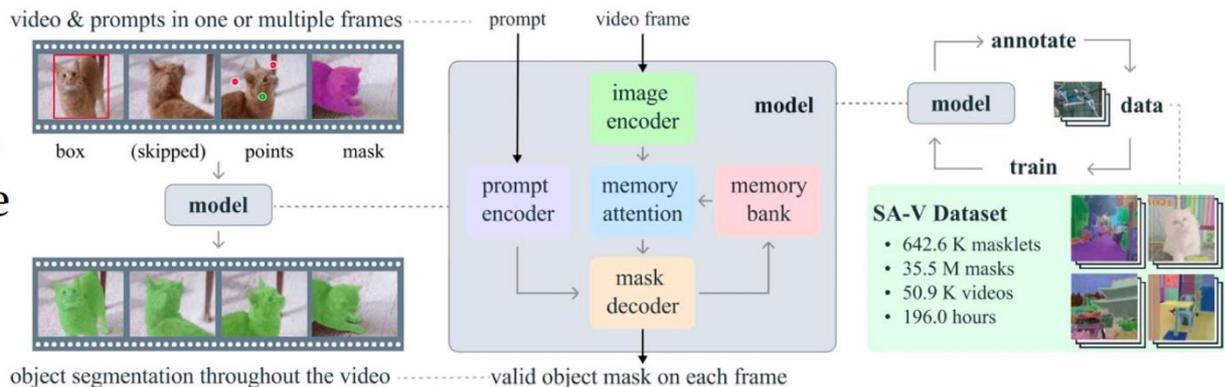


SAM-2 and SAM-3

SAM-3 can take text prompts



Extends to video sequences by retaining previous frame information



object segmentation throughout the video ----- valid object mask on each frame
(a) Text prompt and object segmentation (b) Model Segment Anything: Model 2.0 (c) Dataset for training and dataset

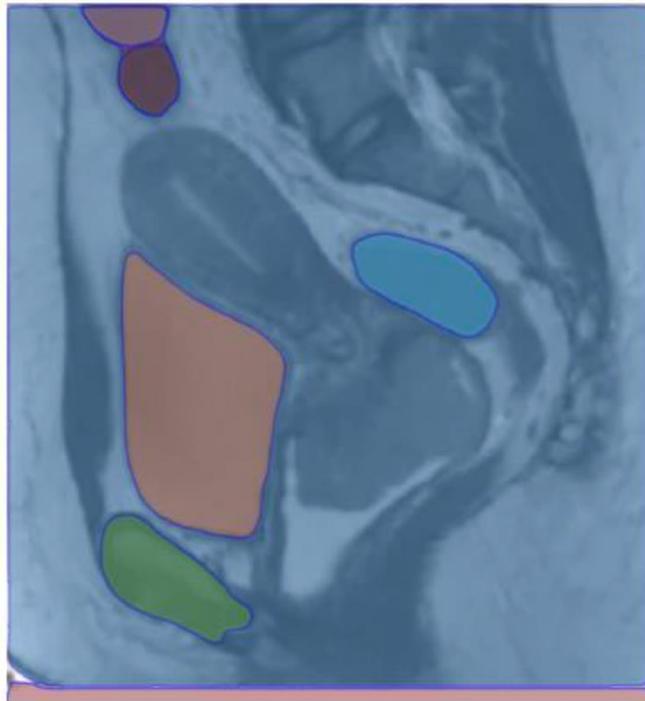
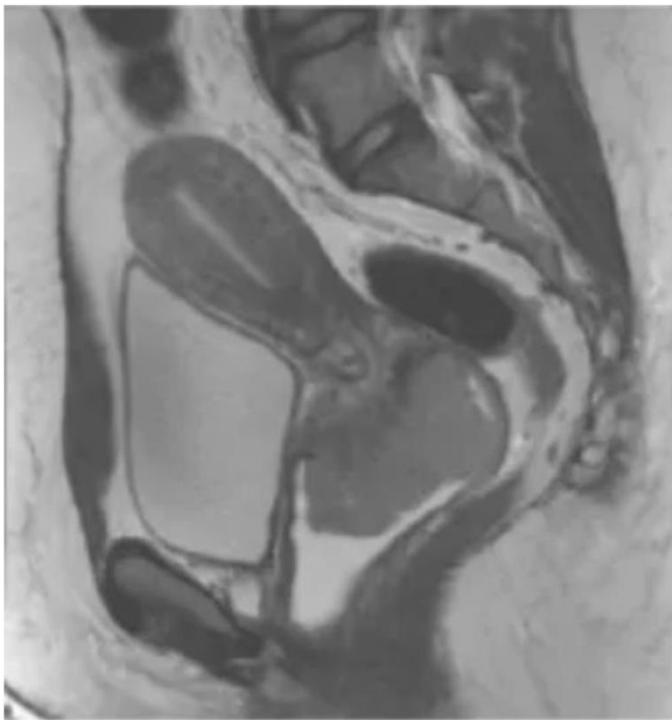
State of segmentation

- Automatic class-agnostic object segmentation
 - SAM in automatic mask generation
 - No recognition of objects though
- Fully automatic semantic segmentation
 - output pixel-wise class labels automatically (e.g. Mask2Former, DeepLabV3)
- Open vocabulary automatic segmentation
 - Automatically segment and label arbitrary objects without prompts. – an open research problem (still needs to know the vocabulary even if not seen an object)

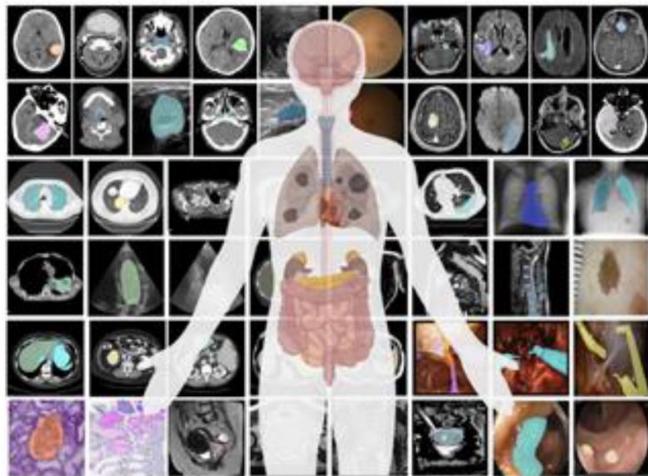
There are three levels of "automatic":

Type	Needs Prompt?	Needs Class List?	Open Vocabulary?
SAM automatic mode	✗	✗	✗ (no labels)
DeepLab-style models	✗	✓	✗
CLIPSeg-style	✓ (text)	✗	✓
GLIP-style	✓ (text)	✗	✓
Research hybrids	Sometimes	Sometimes	Partially

How well does SAM do on medical images?



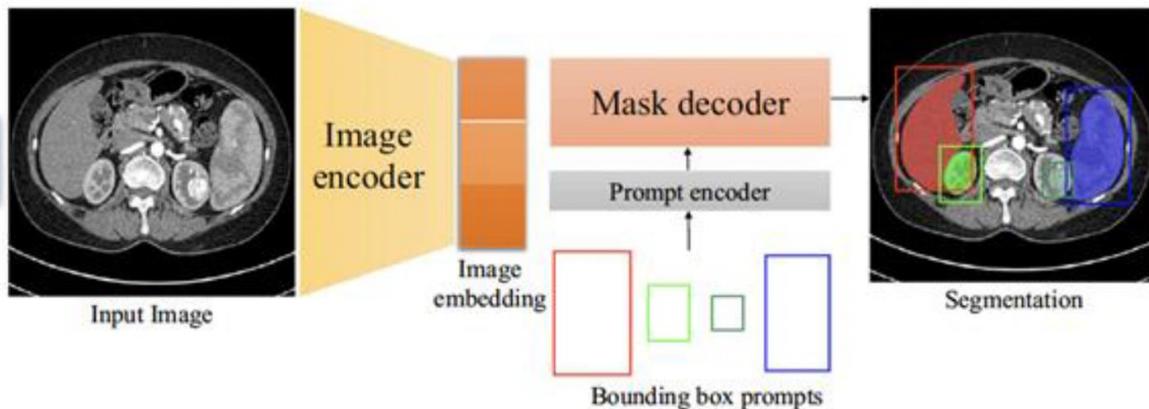
Foundational models for segmentation



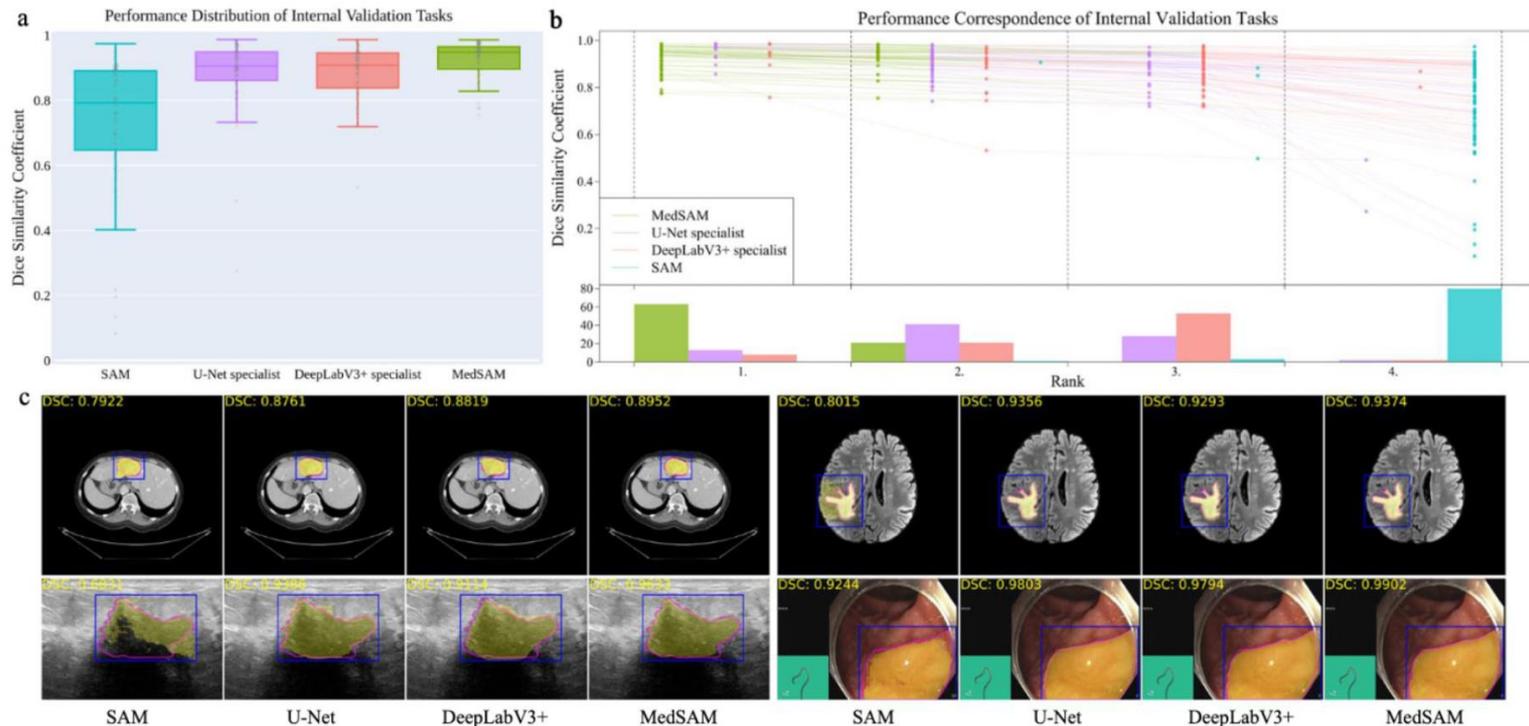
Trained on an huge amount of data

- 10 imaging modalities
- 30 cancer types
- Over 1.5M mask-image pairs

Pipeline

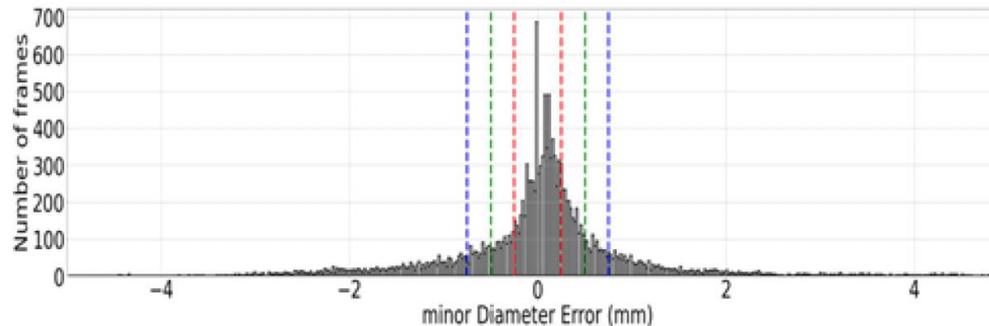
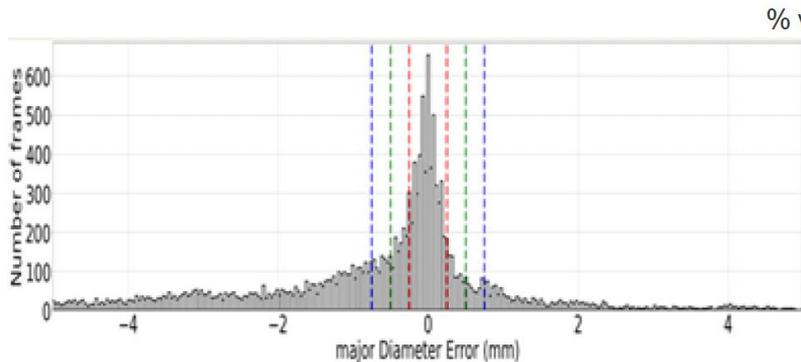
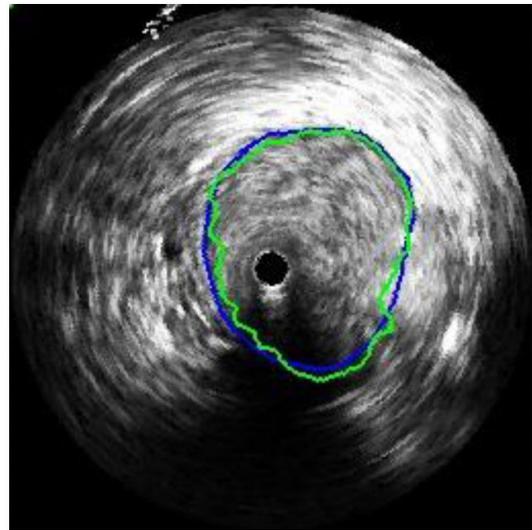


MedSAM generalization



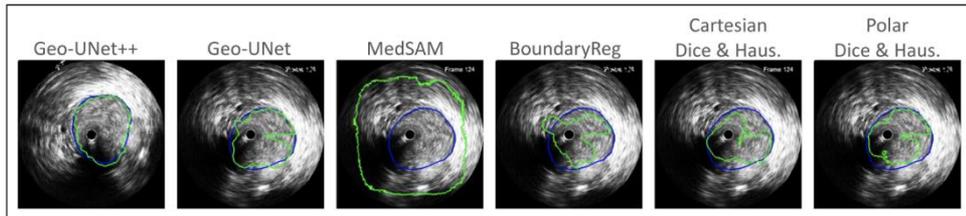
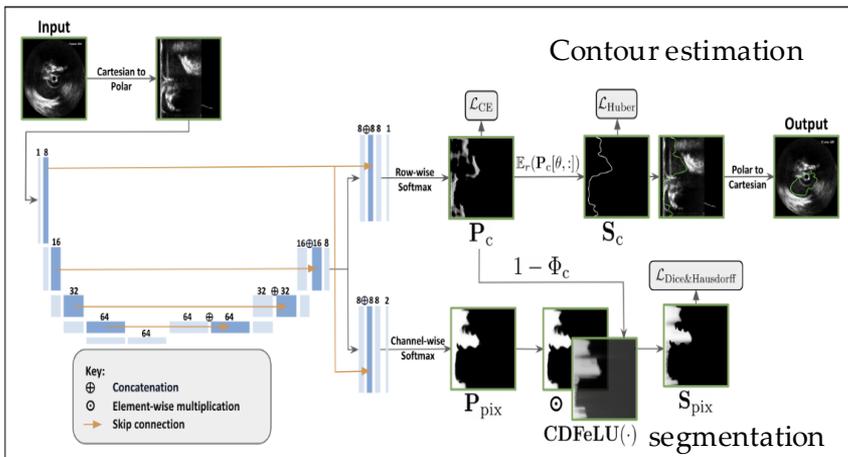
Are FMs ready for high-precision segmentation?

- Needed for accurate sizing of the stents
- Based on normal frames identification
- Estimates the maximum and minimum diameter
- Major and minor diameter errors within:
 - 0.25/0.5/0.75mm for 50/90/95% of all N1 frames.
 - 0.5/0.75mm for 50/70% of frames for N2 frames
- N2 normal are mainly used for vessel compression detection and not stent sizing
- Simple U-net will not suffice as single regions cannot be ensured
- Contours may not closely follow the lumen boundaries



Geo-UNet(MLMI'24)

Clinical-grade lumen segmentation in intravascular ultrasound



Methodology	Test Dice (avg/std)	% Frames w. Maj. Dia. err. within 0.25/0.50/0.75mm	% Frames w. Min. Dia. err. within 0.25/0.50/0.75mm
Against Baselines (N1 frames)			
Geo-UNet++	0.95/0.045	66/84/90	73/89/94
Geo-UNet	0.95/0.034	69/84/90	69/85/91
MedSAM [13]	0.31/0.087	0/0/0	0/0/0
BoundaryReg [7]	0.94/0.043	60/78/86	70/86/91
Cart. Dice & Haus.	0.93/0.051	61/77/83	62/79/87
Polar Dice & Haus.	0.94/0.038	<u>66/80/87</u>	67/84/90

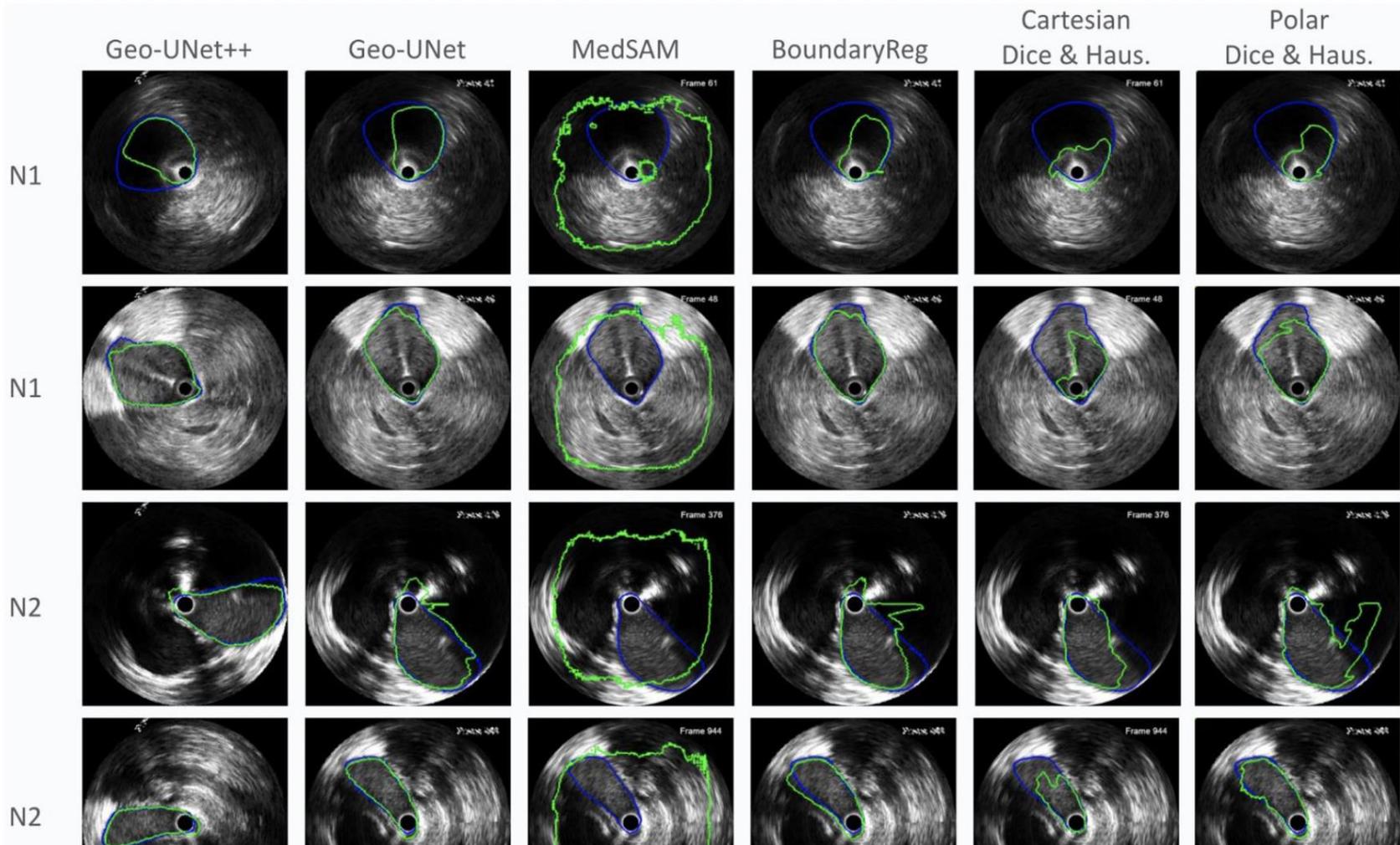
- Polar space input representation
- Anatomically constrained self-informing network
 - standard pixel-level segmentation
 - single lumen contour prediction

Boundary mismatch in segmentation space

$$\mathcal{L}_{Dice\&Hausdorff}(\cdot) = \lambda * \mathcal{L}_{Dice}(S_{pix}, Y_{pix}) + (1 - \lambda) * \mathcal{L}_{Haus.}(S_{pix}, Y_{pix})$$

smoothness of contour

$$\mathcal{L}_{Huber}(\cdot) = \sum_{\theta=0}^{R-1} \frac{d_{\theta}^2}{2} \mathbb{1}(|d_{\theta}| < 1) + (|d_{\theta}| - 0.5) \mathbb{1}(|d_{\theta}| \geq 1),$$



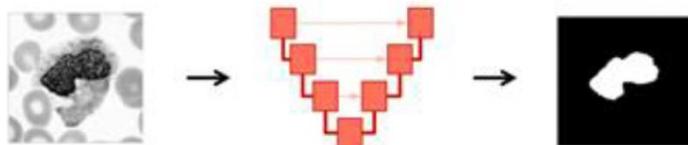
Results

Methodology	Test Dice (avg/std)	% Frames w. Maj. Dia. err. within 0.25/0.50/0.75mm	% Frames w. Min. Dia. err. within 0.25/0.50/0.75mm
Against Baselines (N1 frames)			
Geo-UNet++	<u>0.95/0.045</u>	<u>66/84/90</u>	73/89/94
Geo-UNet	0.95/0.034	69/84/90	69/85/ <u>91</u>
MedSAM [10]	0.31/0.087	0/0/0	0/0/0
BoundaryReg [4]	0.94/0.043	60/78/86	<u>70/86/91</u>
Cart. Dice & Haus.	0.93/0.051	61/77/83	62/79/87
Polar Dice & Haus.	0.94/0.038	<u>66/80/87</u>	67/84/90
Against Baselines (N2 frames)			
Geo-UNet++	0.88/0.094	<u>41/59/69</u>	60/80/87
Geo-UNet	<u>0.87/0.10</u>	47/64/73	<u>57/76/85</u>
MedSAM [10]	0.23/0.085	0/0/0	0/0/0
BoundaryReg [4]	<u>0.87/0.093</u>	36/54/65	55/74/84
Cart. Dice & Haus.	<u>0.83/0.12</u>	32/44/52	44/63/74
Polar Dice & Haus.	<u>0.86/0.12</u>	<u>40/58/69</u>	55/74/83
Against Ablations (N1 frames)			
Geo-UNet	0.95/0.034	69/84/90	69/85/91
w/o CDFeLU	0.94/0.035	<u>69/82/88</u>	<u>65/83/90</u>
w/o pixel-wise pred.	<u>0.95/0.039</u>	67/81/87	69/85/91
Against Ablations (N2 frames)			
Geo-UNet	<u>0.87/0.10</u>	47/64/73	57/76/85
w/o CDFeLU	<u>0.86/0.10</u>	<u>45/63/72</u>	<u>53/71/81</u>
w/o pixel-wise pred.	0.88/0.092	<u>46/62/71</u>	<u>57/76/85</u>

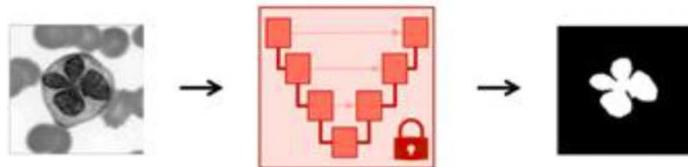
Meta learning in UniverSeg

Traditional Approach

1. Design and train a task-specific model.

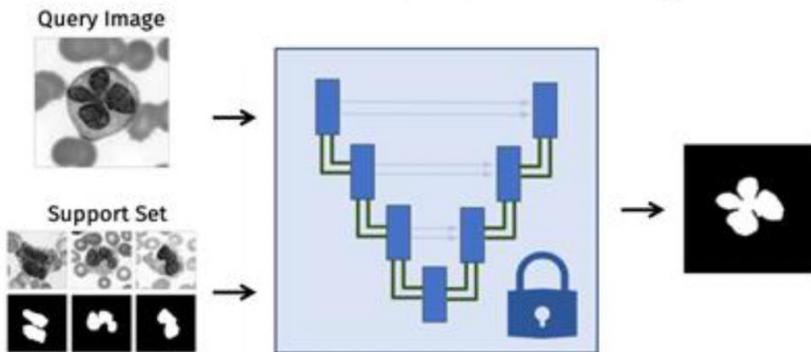


2. Predict new images with the trained model.

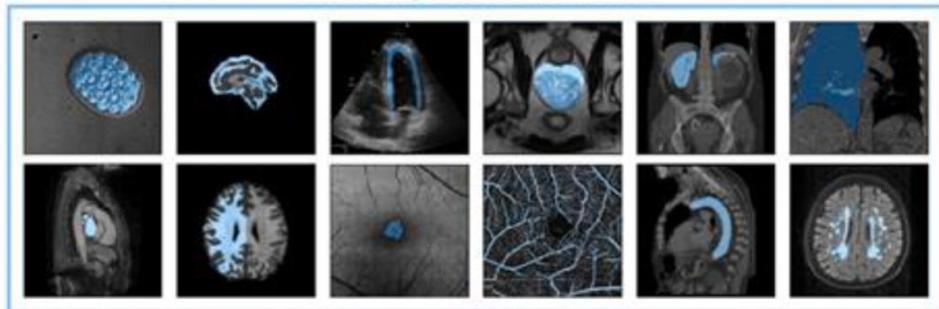


UniverSeg Approach

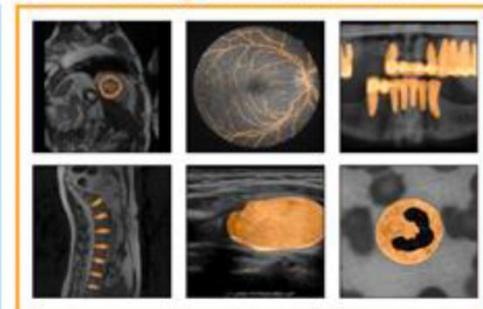
With a trained UniverSeg model, predict new images for the new task from a few labeled pairs without retraining.



Train Segmentation Tasks



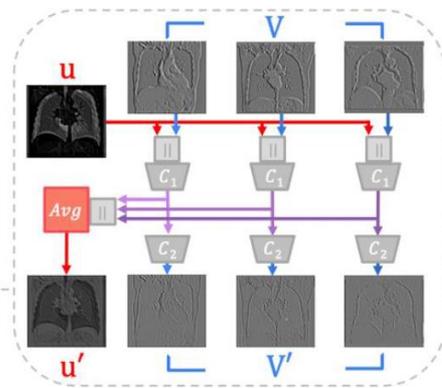
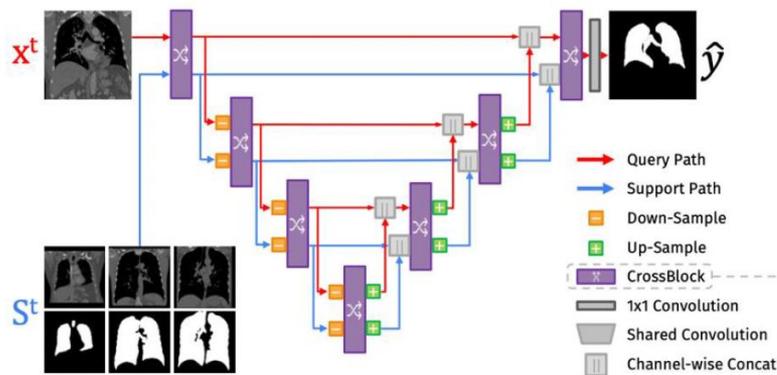
Test Segmentation Tasks



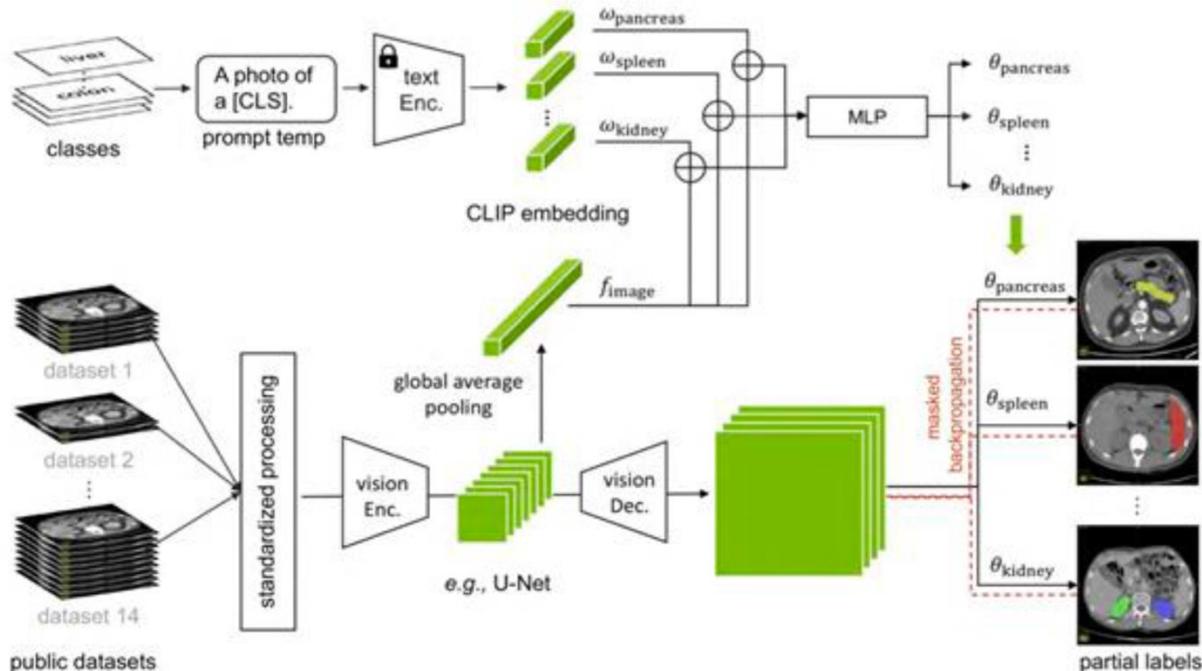
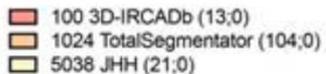
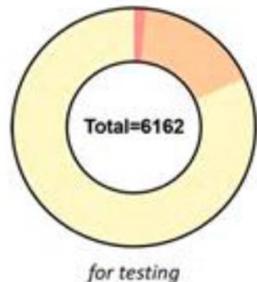
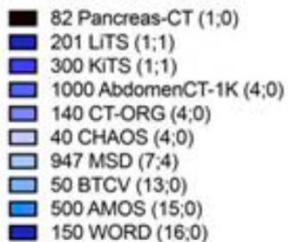
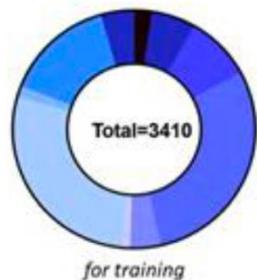
Instructional tuning of segmentation models

- Based on U-net style convolutional architecture
- Uses multi-scale cross-block features between instructional labeling sets and given image
- No retraining needed or fine-tuning needed!
- Trained on Megamedical dataset: 53 datasets, 23 medical domains, 16 modalities

- Cross-convlayer
- Convolve the query with each support
 - Reduce them together to incorporate information about what to segment based on the support labels
 - The result is updated representations for both query and support



Combining U-net and CLIP for anatomy and tumor segmentation



Slides from MICCAI 2024 tutorial

Foundational models for segmentation

CLIP-Driven

Main idea

Text branch

(generates text embedding for class k) \mathbf{w}_k

Visual branch-encoder

(generates visual embedding for image x) \mathbf{f}

Text-based controller MLP
(generates class parameters)

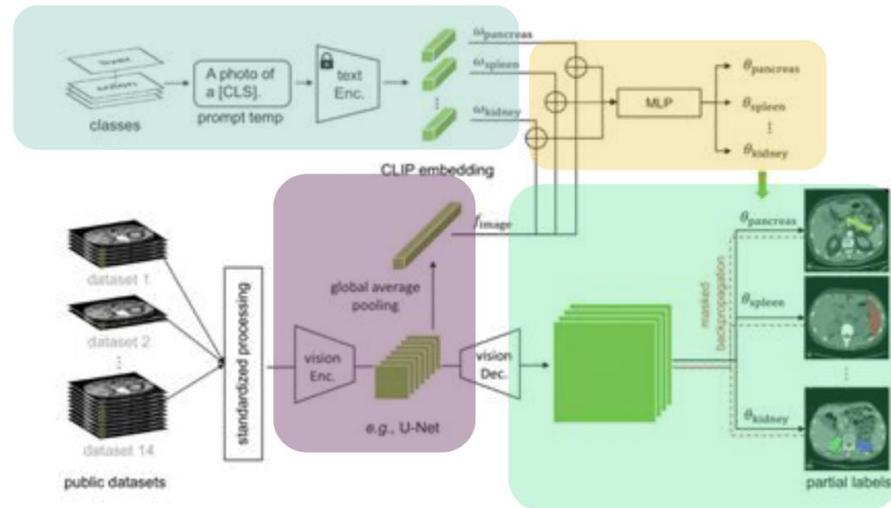
$$\theta_k = MLP(\mathbf{w}_k \oplus \mathbf{f})$$
$$\theta_k = \{\theta_{k_1}, \theta_{k_2}, \theta_{k_3}\}$$

Visual branch-decoder

(generates visual embedding for image x)

$$\mathbf{P}_k = \text{sigmoid}(((\mathbf{F} * \theta_{k_1}) * \theta_{k_2}) * \theta_{k_3})$$

It represents foreground class k vs background



Training loss

Binary cross-entropy per class (and terms masked for those classes not present)

$$\mathcal{L} = \sum_{k=1}^K \mathbf{1}_{\{k \in y\}} \cdot \text{BCE}_k$$

Segmentation models

Model / Framework	Type	Strength
SAM 3	General Vision Foundation Model	Open-vocabulary & multi-instance segmentation (images + video)
SAM 2	Vision Foundation Model	Interactive prompt + video temporal consistency
SAM 1	Vision Foundation Model	Zero-shot general segmentation
ENSAM	Medical/3D segmentation variant	Efficient 3D segmentation under limited data
SAMPO	Enhanced prompt adaptation research	Intent-aware segmentation
TopoLoRA-SAM	Adaptation framework	Efficient fine-tuning for specific tasks
SAM2-UNeXT	High-resolution foundation adaptation	Stronger features across domains

Summary

- Early approaches to segmentation were unsupervised and didn't scale
- Supervised approaches had limited labels issues
- Foundational models generalized across datasets and open vocabularies
- Key architectures are still based on CNN or transformers
- Active field of research in developed generalized foundation models for segmentation with extensions to video sequences
- Gap seen with applications to medical imaging leading to some rich innovations for medical imaging adaptations
- Unsupervised segmentation approaches may still be relevant
 - STEGO: Unsupervised Semantic Segmentation by Distilling Feature Correspondences, ICLR'2022 -> Learns with no labels