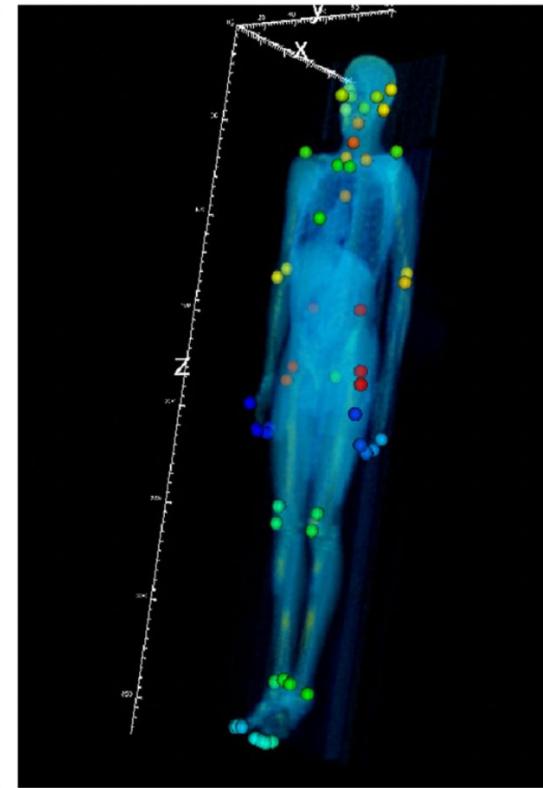
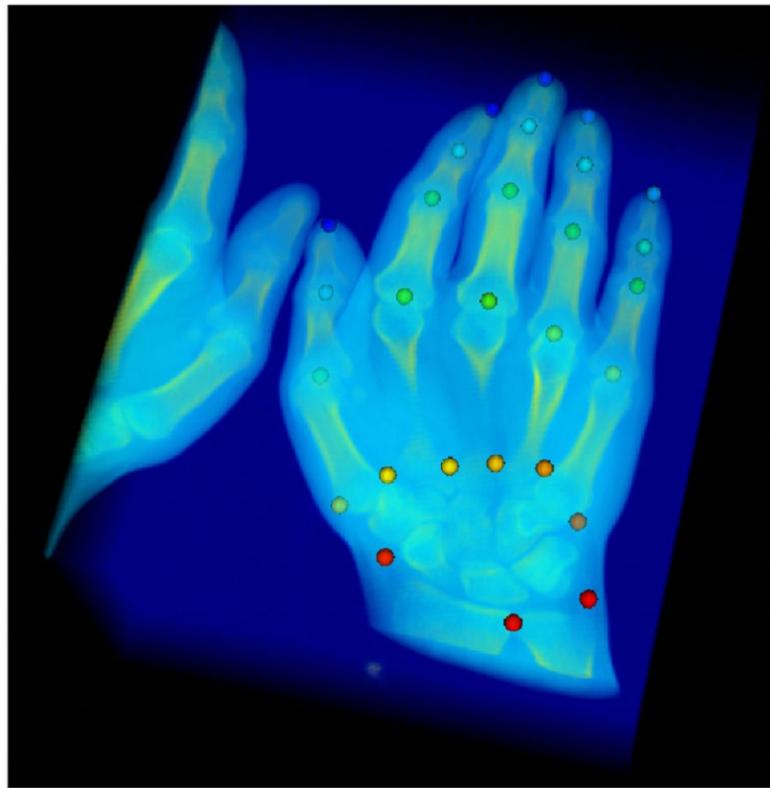


Segment Anything and the Rise of Foundation Models

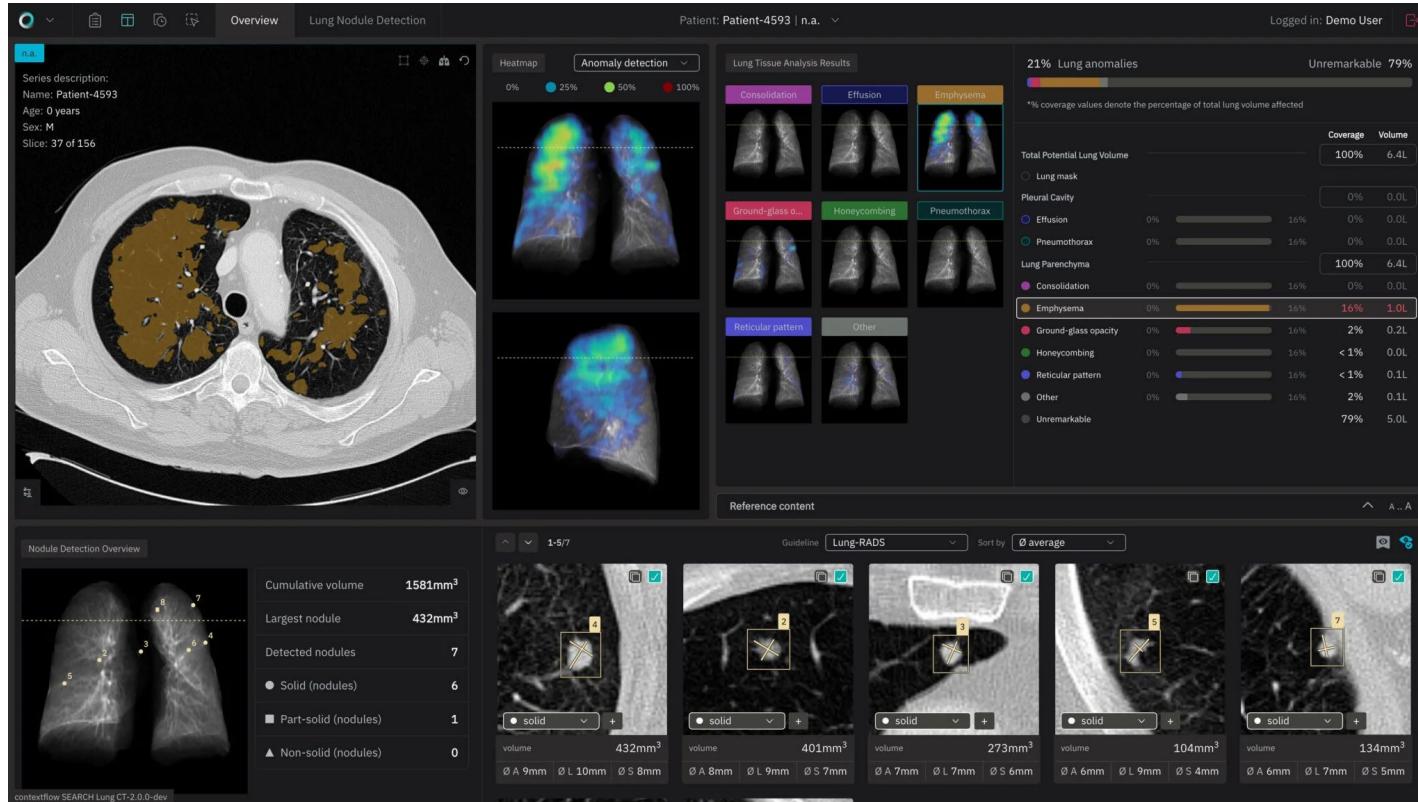
MEDICAL VOLUME ANNOTATOR



René Donner



contextflow



MEDICAL VOLUME ANNOTATOR



Medical Volume Annotator – On-Site Annotation Platform

Menu

- ▼ Datasets
- DevDataset
- ▼ Projects
- New project
- ▼ test
- To annotate (45)
- ▼ Administration
- Forms
- Labels
- Roles
- Suppliers
- Users

► Dev Settings

Show Help / Shortcuts

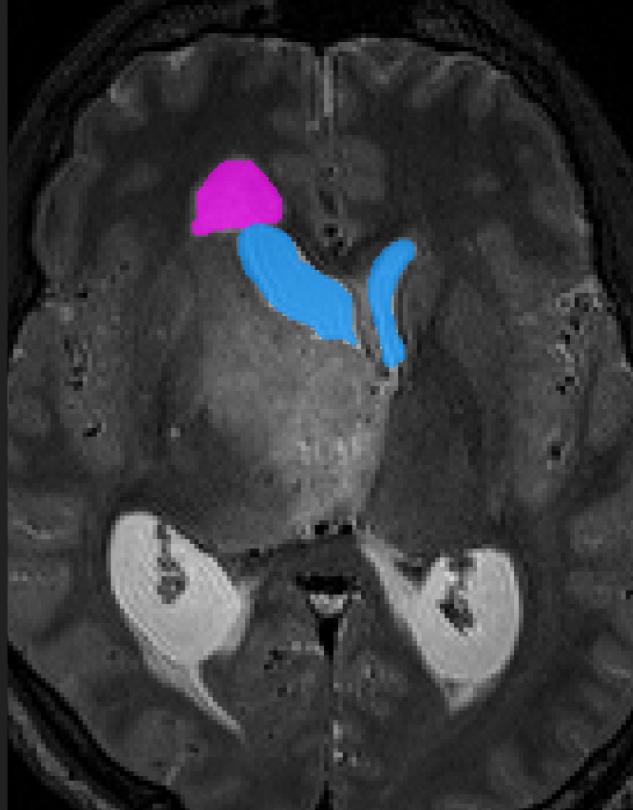
Toggle fullscreen

Logout

Volume loading took: 1 sec

test > Annotate > case2

<< Hide Menu Reset zoom Slice 161 of 320



Details

Viewer Windowing

- Pixels & Ann.
- Annotations
- No overlays

Annotation outlines only

0.70 Label opacity

Settings

64 Scroll speed

Linear interpolation

Labeling

- Brush
- Superpixels

1 Brush radius

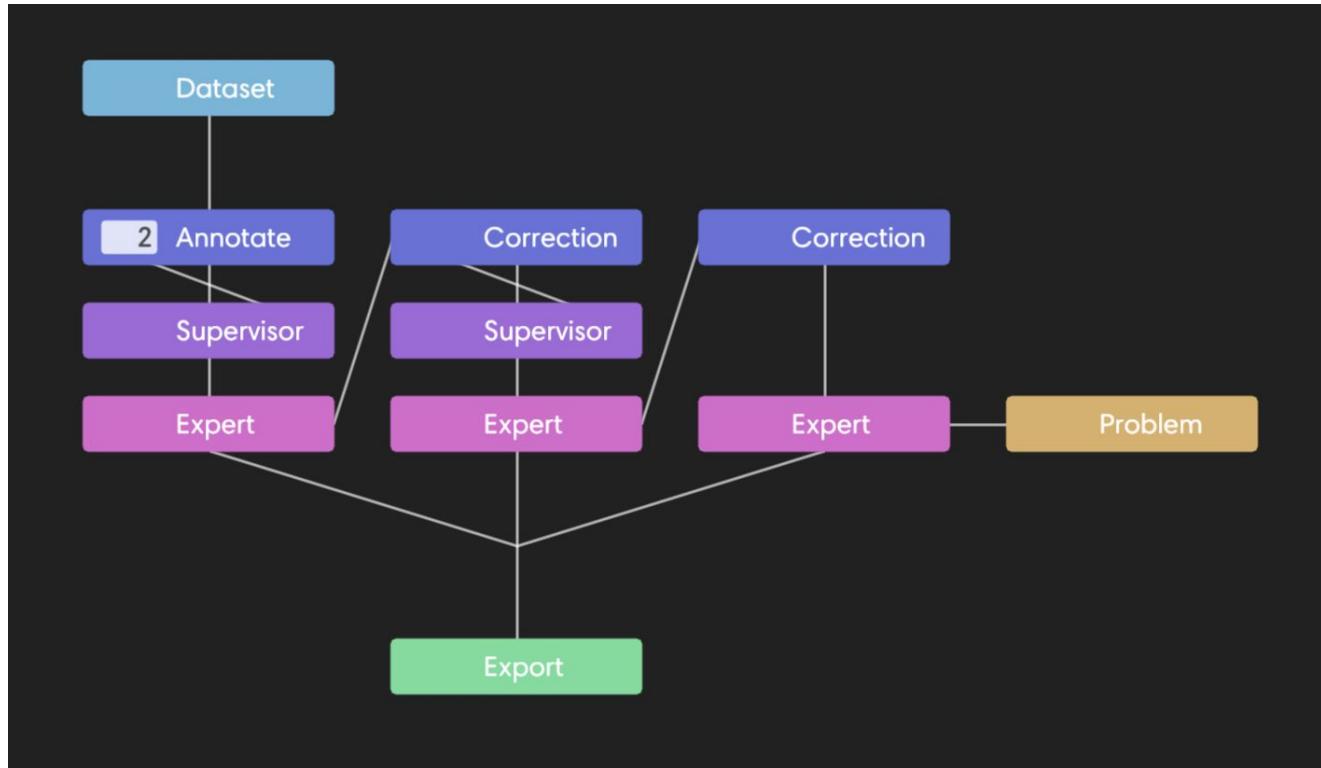
Existing labels:

- Cavity
- Consolidation
- Cyst
- Effusion

Submit case for review

Medical Volume
Annotator

Medical Volume Annotator – Workflow Management, Statistics ...



Smart Brushes – online learning of labeling models

Smart Propagation – distribute annotations spatially

Smart Regions – Super-pixels, over-segmentation, ...

Medical Volume Annotator – Efficient Annotation

Smart Brushes – online learning of labeling models

Smart Propagation – propagation annotations spatially

Smart Regions – Super-pixels, over-segmentation, ...

“Segment Anything Model (SAM): a new AI model from Meta AI that can ‘cut out’ any object, in any image, with a single click”

Segment Anything Model (SAM)

Segment Anything Model

Dataset – 11 Mio Images, 1.1 Billion masks // “Foundation Model”



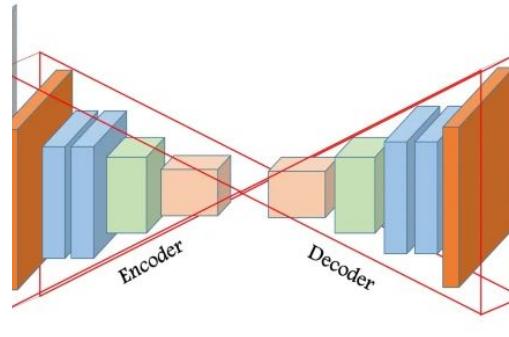
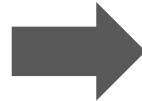
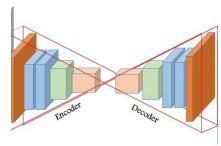
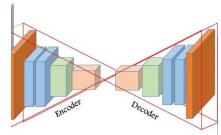
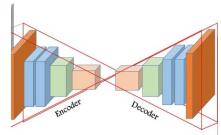
Foundation Models?

Large, pre-trained models

Strong generalization & zero-shot performance

Huge datasets, massive compute

Foundation Models?



+



Adapter

Segment Anything

From scratch

facebook DENOv2: Self-supervised Vision Transformer Model

Image-level visual tasks

Image classification

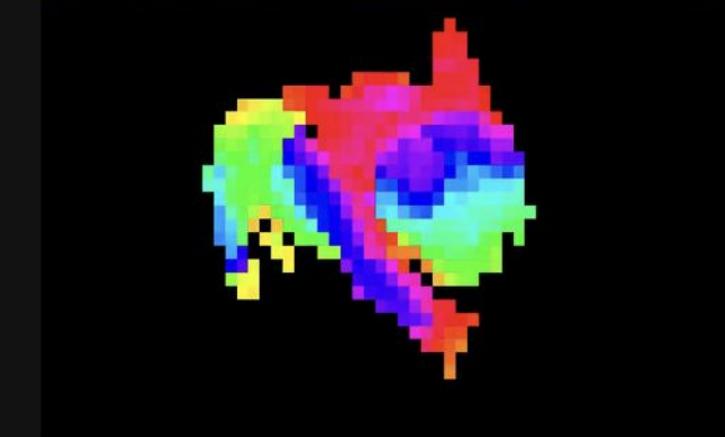
Instance retrieval

Video understanding

Pixel-level visual tasks

Depth estimation

Semantic segmentation



DENOv2: Self-supervised Vision Transformer Model



Depth Estimation

State-of-the-art results and strong generalization on estimating depth from a single image.

⊕ [Try the demo](#)

DENOv2: Self-supervised Vision Transformer Model



Instance Retrieval

Directly use frozen features to find art pieces similar to a given image from a large art collection.

↻ Try the demo

Microsoft's offerings

Language & Multilingual

- | **UniLM**: unified pre-training for language understanding and generation
- | **InfoXLM/XLM-E**: multilingual/cross-lingual pre-trained models for 100+ languages
- | **DeltaLM/mT6**: encoder-decoder pre-training for language generation and translation for 100+ languages
- | **MiniLM**: small and fast pre-trained models for language understanding and generation
- | **AdaLM**: domain, language, and task adaptation of pre-trained models
- | **EdgeLM** (NEW): small pre-trained models on edge/client devices
- | **SimLM** (NEW): large-scale pre-training for similarity matching
- | **E5** (NEW): text embeddings

Vision

- | **BEiT/BEiT-2**: generative self-supervised pre-training for vision / BERT Pre-Training of Image Transformers
- | **DiT**: self-supervised pre-training for Document Image Transformers
- | **TextDiffuser** (NEW): Diffusion Models as Text Painters

"Anything Models" – Recognize Anything



RAM

living room, dog, blanket, carpet, couch, desk, furniture, pillow, plant, sit, wood floor, lamp

Christmas market, Christmas tree, stall, market square, snow, people, stroll, town, building

Tag2Text

living room, dog, sit on, blanket, couch, plant, modern
Missing: lamp, carpet

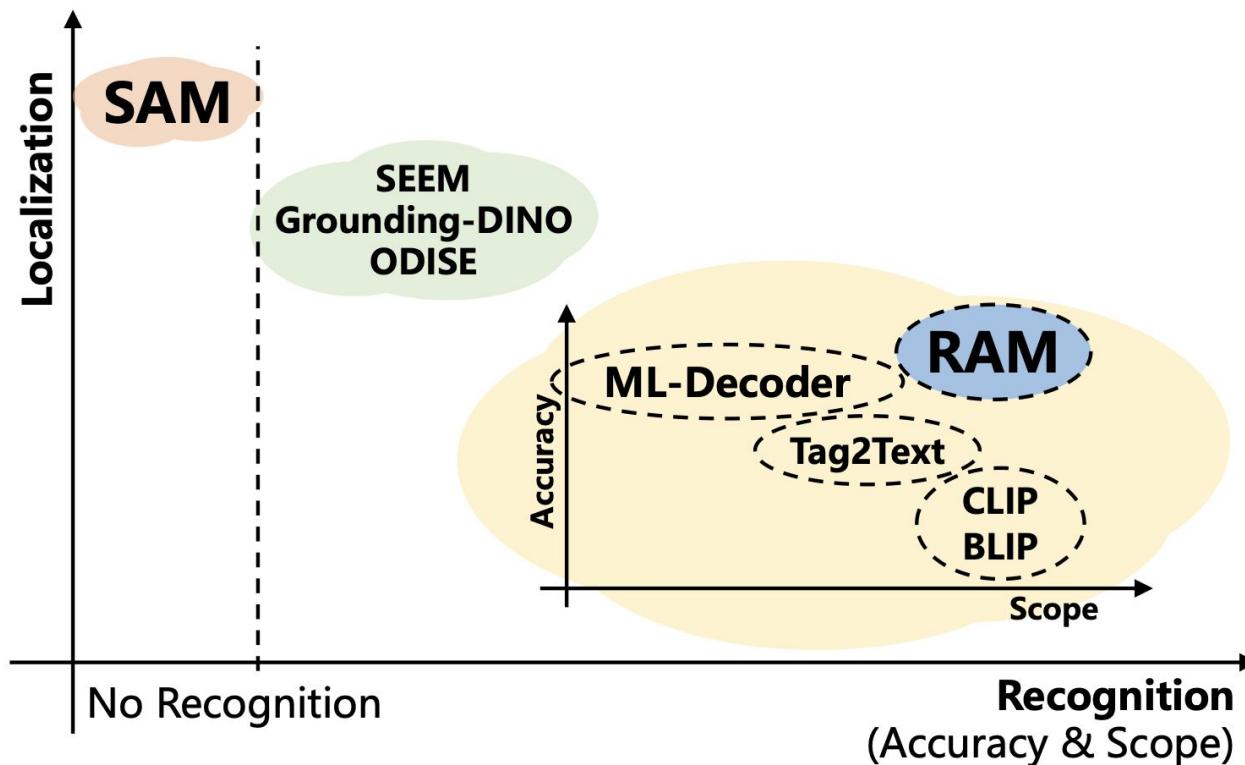
Christmas market, Christmas tree, snow, town, people
Missing: building

ML-Decoder

living room, lamp, houseplant, cushion, throw pillow, picture frame
Bad: property, design, throw
Missing: dog, couch, carpet, blanket

Christmas decoration, town square, market, snow, building
Bad: human hair, human head, mixed-use

"Anything Models" – Recognize Anything



Segment Anything Model

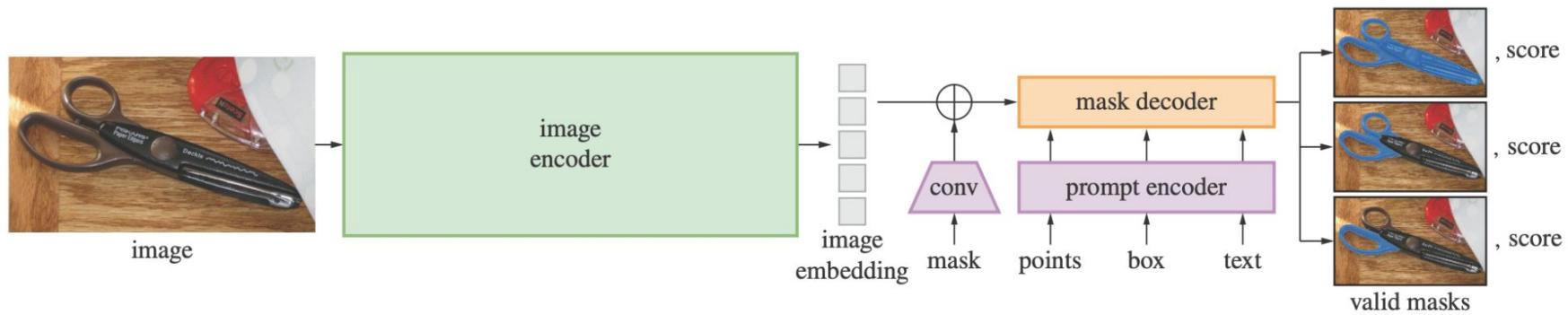
What can SAM do?

"A new AI model from Meta AI that can 'cut out' any object, in any image, with a single click"

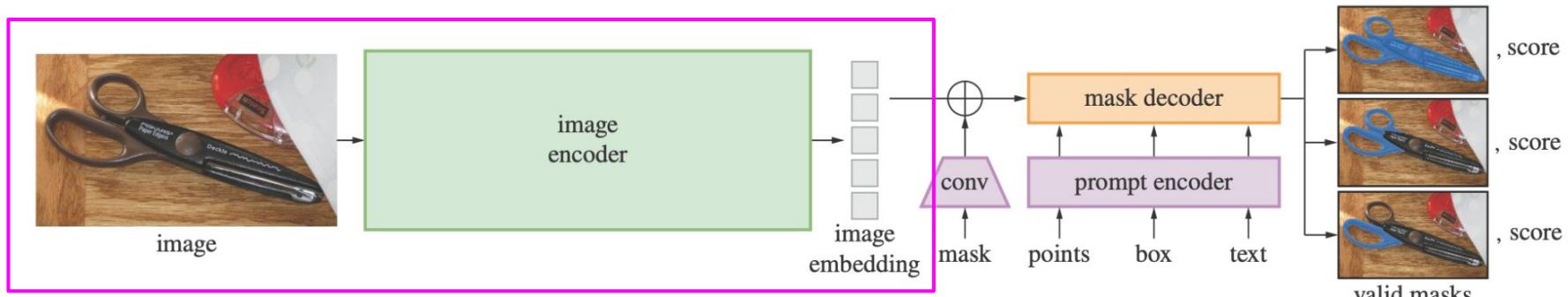


Figure 3: Each column shows 3 valid masks generated by SAM from a single ambiguous point prompt (green circle).

SAM Architecture



SAM Architecture



SAM Architecture

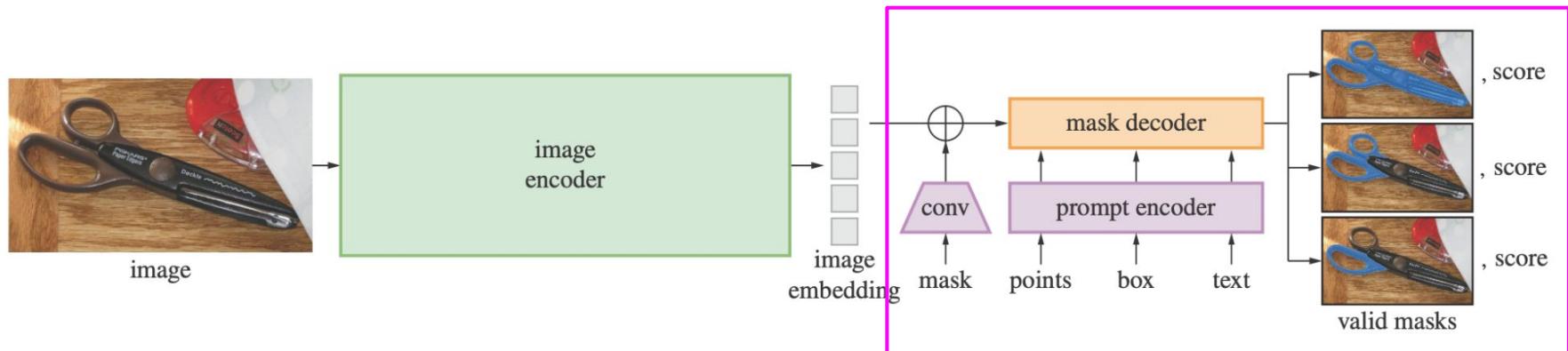
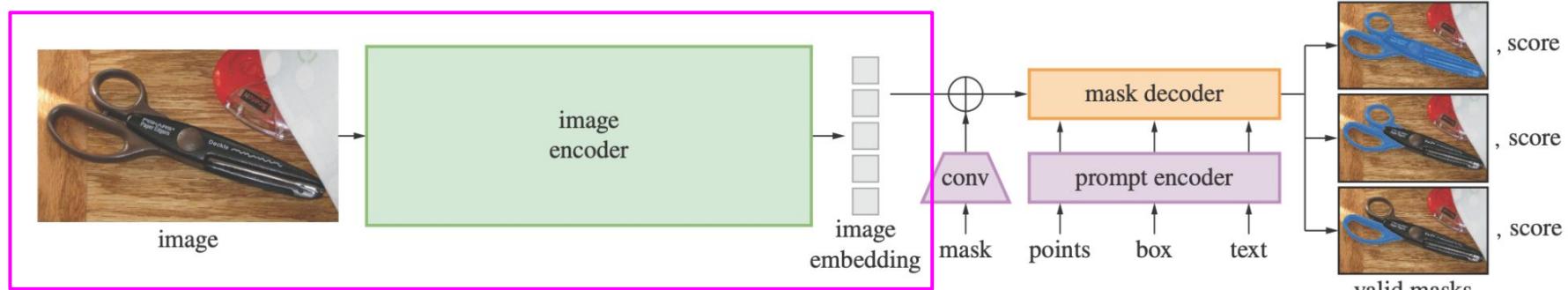


Image Encoder

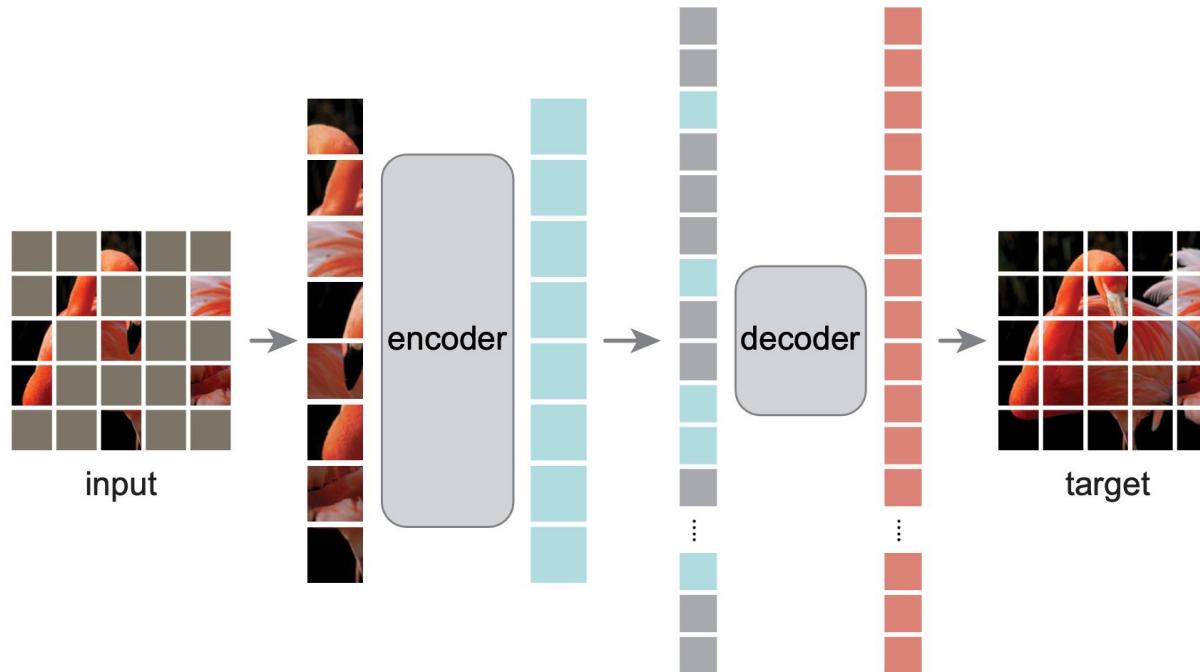


Encoding performed off-line, before user interaction

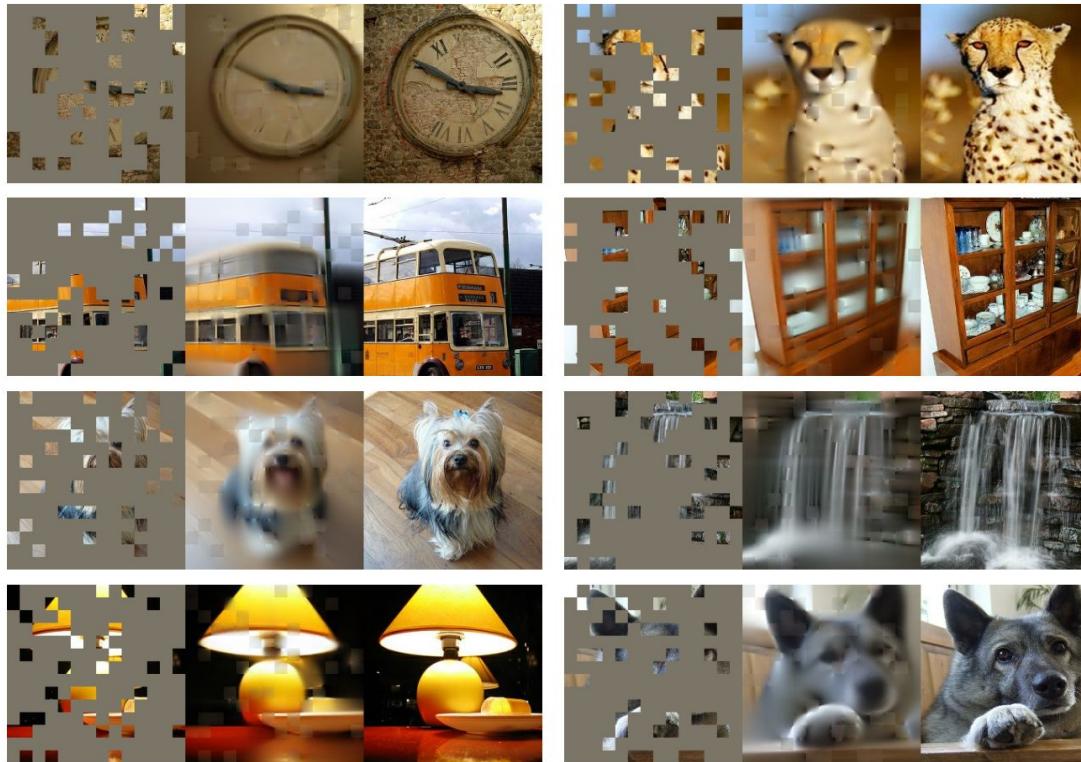
1024 x 1024 RGB image \Rightarrow 64 x 64 x 256 channels

Initialized with a pre-trained Masked Autoencoder, ViT-L backend

Masked Autoencoder



Masked Autoencoder



Masked Autoencoder

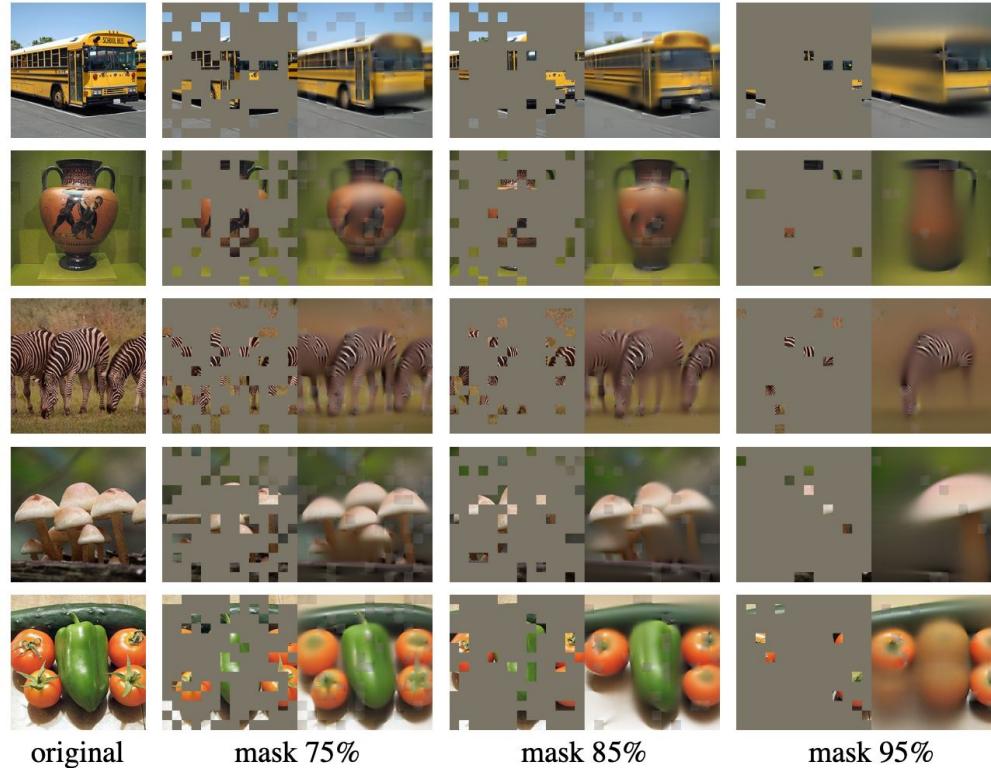
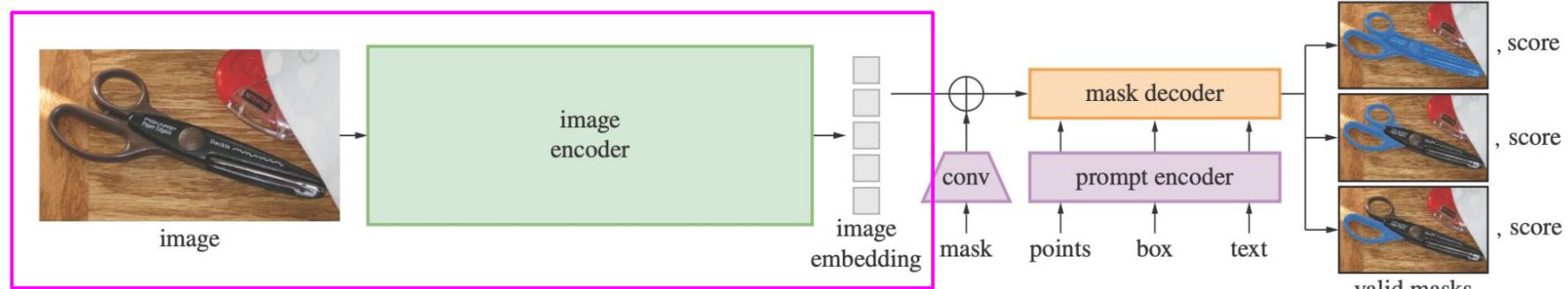


Image Encoder

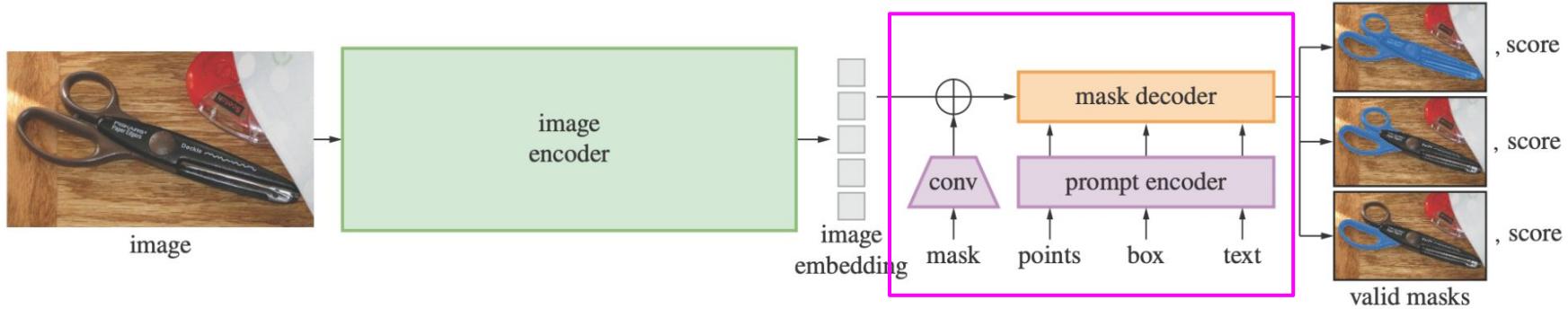


Encoding performed off-line, before user interaction

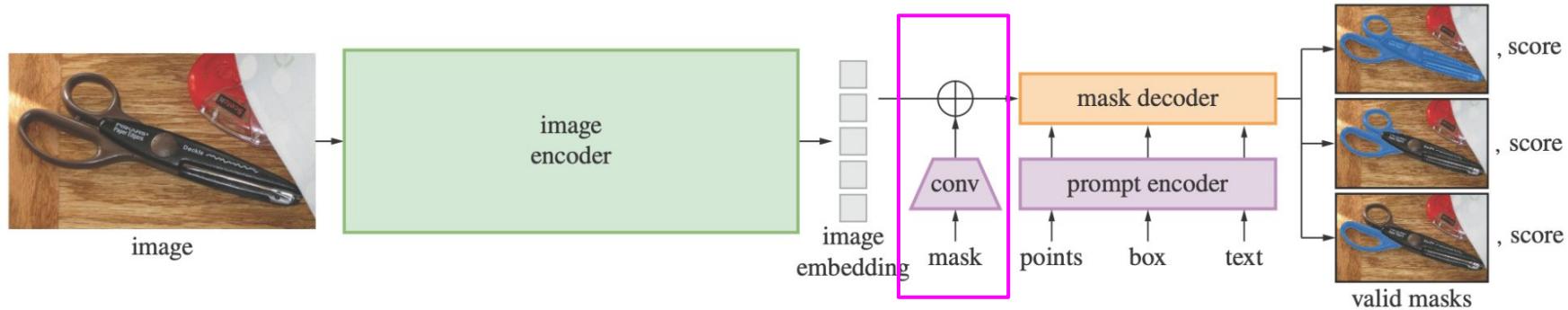
1024×1024 RGB image \Rightarrow $64 \times 64 \times 256$ channels

Initialized with a pre-trained Masked Autoencoder, ViT-L backend

Prompt Encoder



Mask Encoder



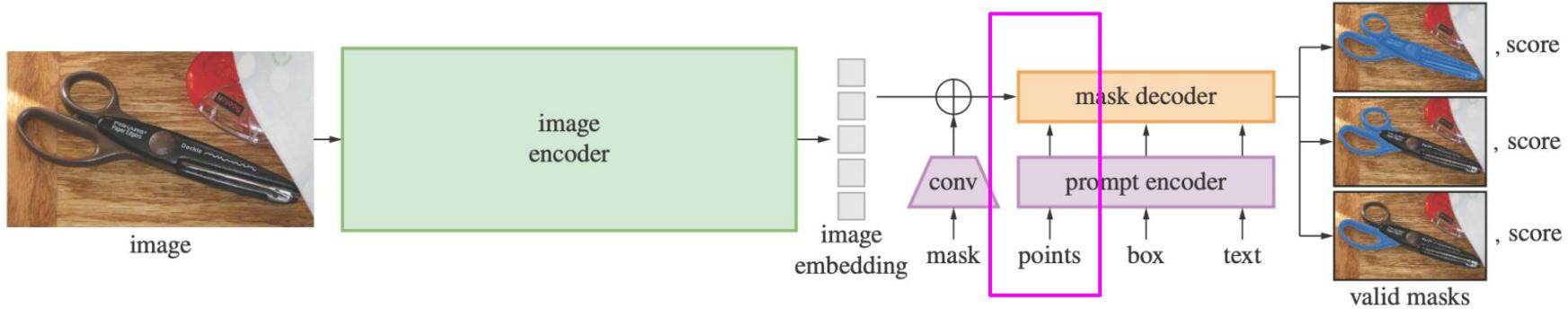
256 x 256 probability mask

Two 2×2 , stride-2 convolutions, then 1×1 conv \Rightarrow $64 \times 64 \times 256$

Element-wise addition

Learned (!) embedding for “no mask”

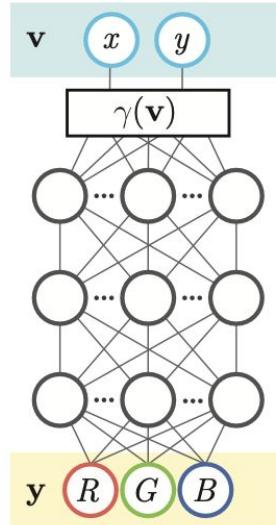
Prompt Encoder



Positional Encoding using Fourier Features

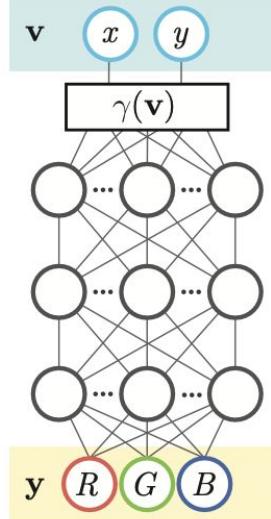
Summed with one of two learned embeddings for foreground/background

Positional Encoding – Fourier Features



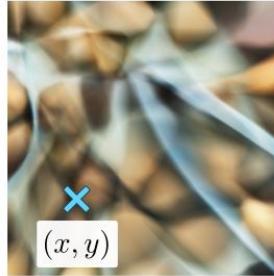
(a) Coordinate-based MLP

Positional Encoding – Fourier Features



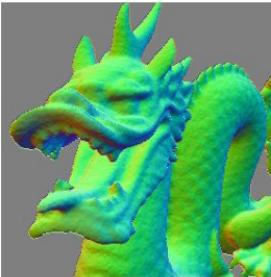
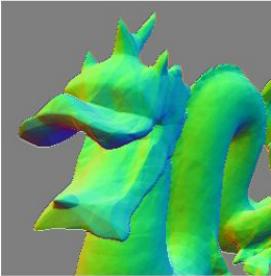
(a) Coordinate-based MLP

No Fourier features
 $\gamma(\mathbf{v}) = \mathbf{v}$

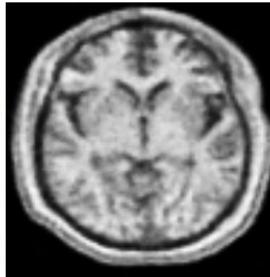
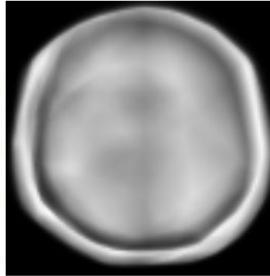


(b) Image regression
 $(x,y) \rightarrow \text{RGB}$

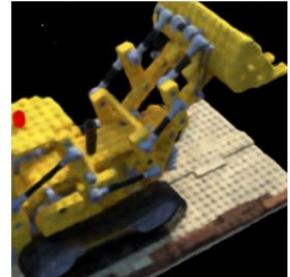
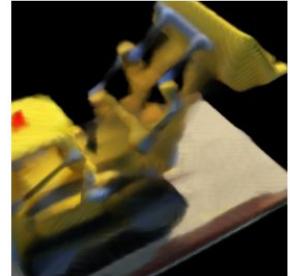
With Fourier features
 $\gamma(\mathbf{v}) = \text{FF}(\mathbf{v})$



(c) 3D shape regression
 $(x,y,z) \rightarrow \text{occupancy}$

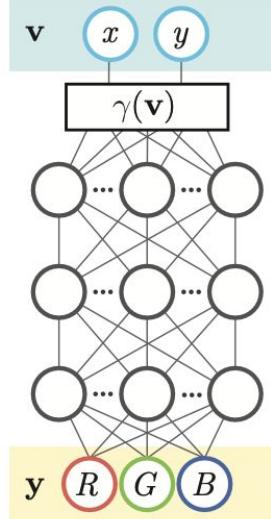


(d) MRI reconstruction
 $(x,y,z) \rightarrow \text{density}$



(e) Inverse rendering
 $(x,y,z) \rightarrow \text{RGB, density}$

Positional Encoding – Fourier Features

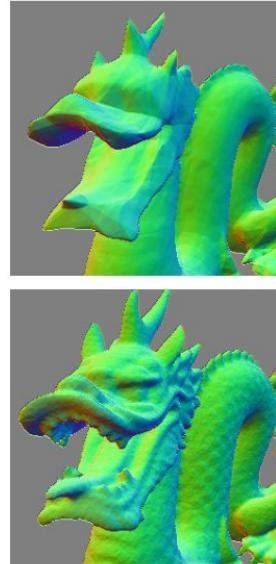


(a) Coordinate-based MLP

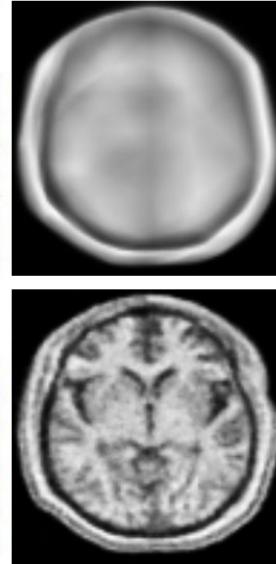
$$\text{No Fourier features} \quad \gamma(\mathbf{v}) = \mathbf{v}$$
$$\text{With Fourier features} \quad \gamma(\mathbf{v}) = \text{FF}(\mathbf{v})$$



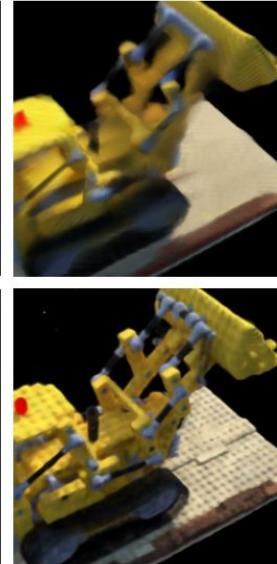
(b) Image regression
 $(x,y) \rightarrow \text{RGB}$



(c) 3D shape regression
 $(x,y,z) \rightarrow \text{occupancy}$



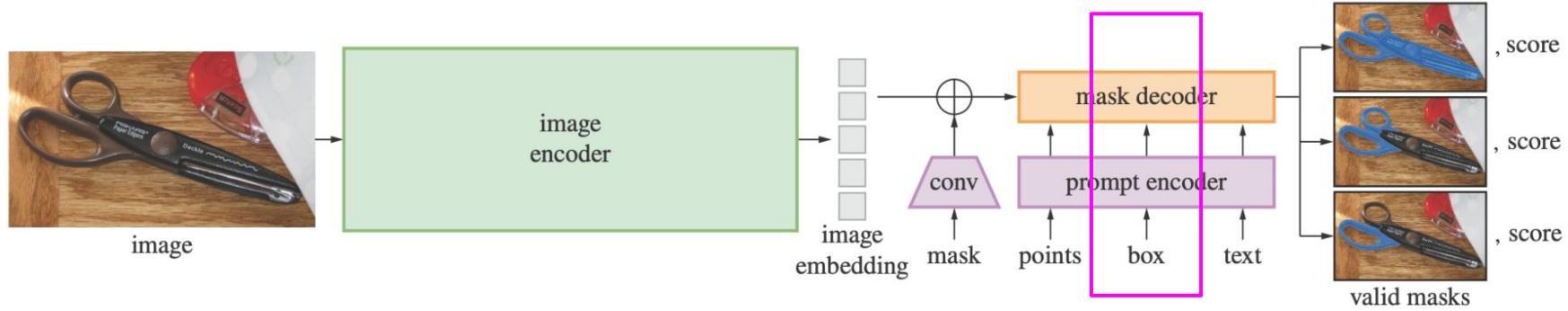
(d) MRI reconstruction
 $(x,y,z) \rightarrow \text{density}$



(e) Inverse rendering
 $(x,y,z) \rightarrow \text{RGB}, \text{density}$

$$\gamma(\mathbf{v}) = [a_1 \cos(2\pi \mathbf{b}_1^T \mathbf{v}), a_1 \sin(2\pi \mathbf{b}_1^T \mathbf{v}), \dots, a_m \cos(2\pi \mathbf{b}_m^T \mathbf{v}), a_m \sin(2\pi \mathbf{b}_m^T \mathbf{v})]^T$$

Prompt Encoder

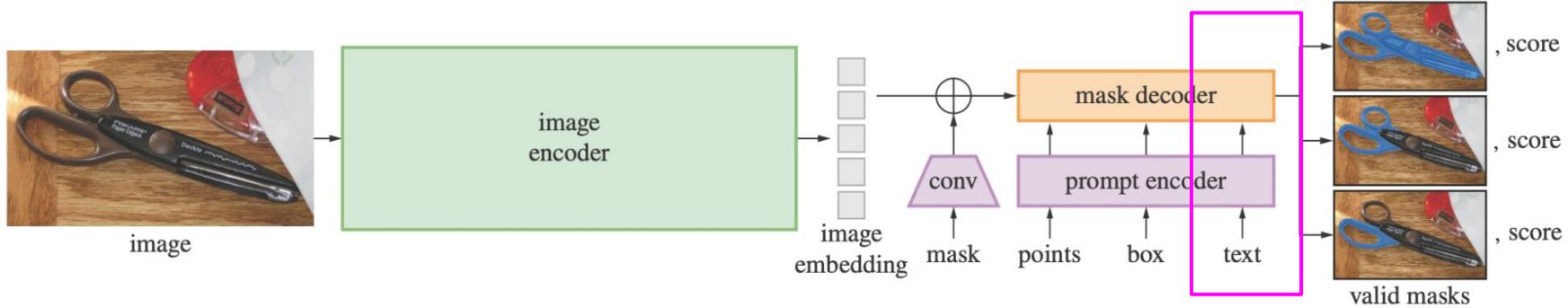


Embedding pair:

Positional encoding of top-left corner,
summed with a learned embedding “top-left corner”

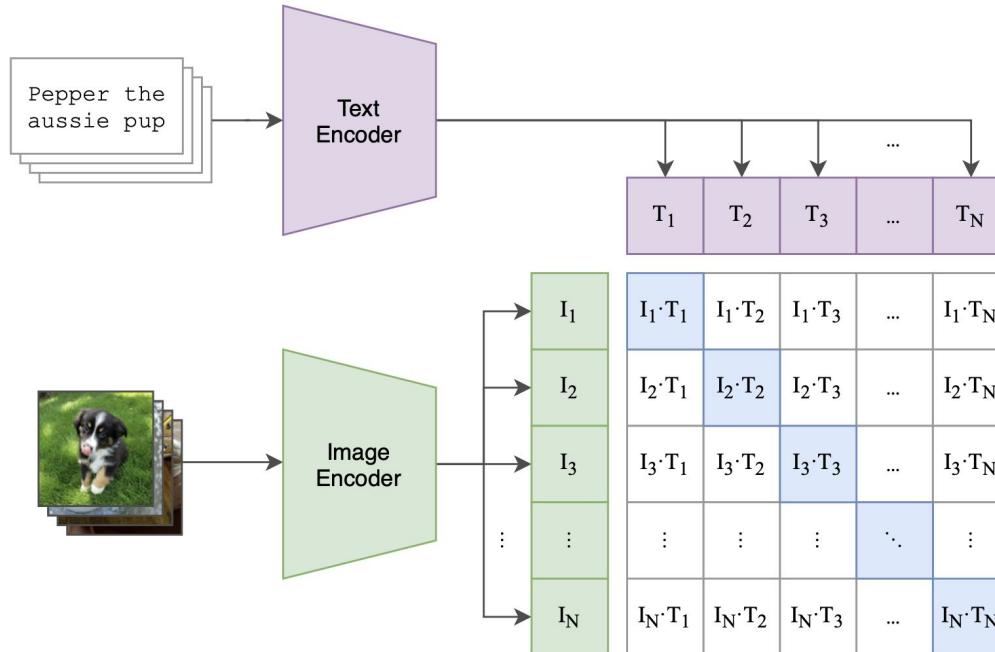
Same for “bottom-right corner”

Prompt Encoder

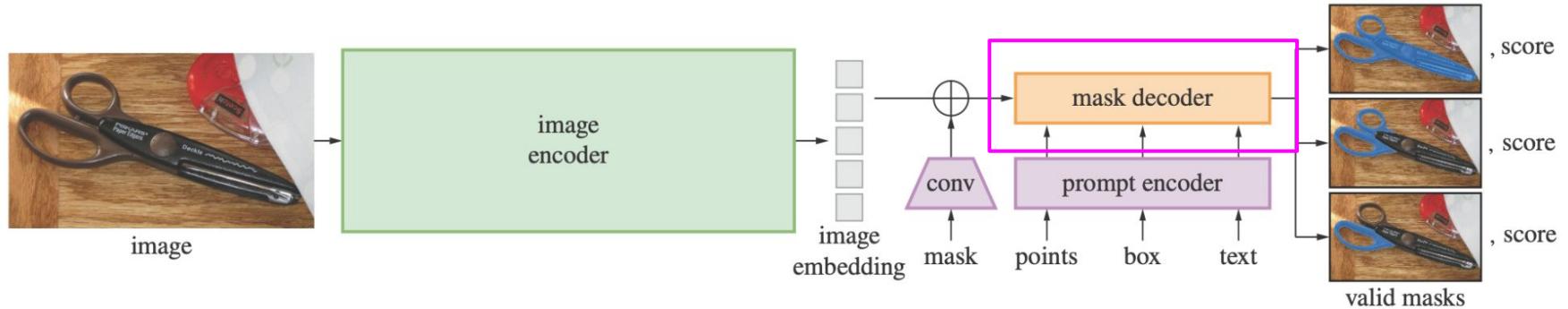


CLIP (Contrastive Language-Image Pre-training) text encoder

CLIP (Contrastive Language-Image Pre-training) text encoder



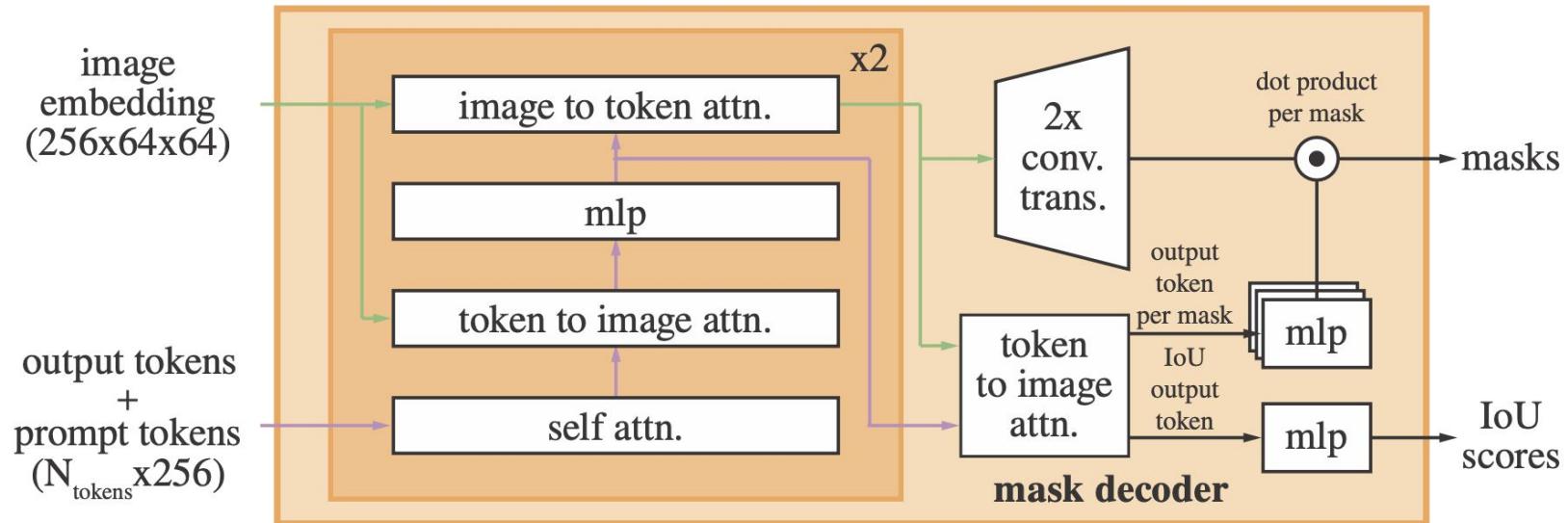
Mask Decoder



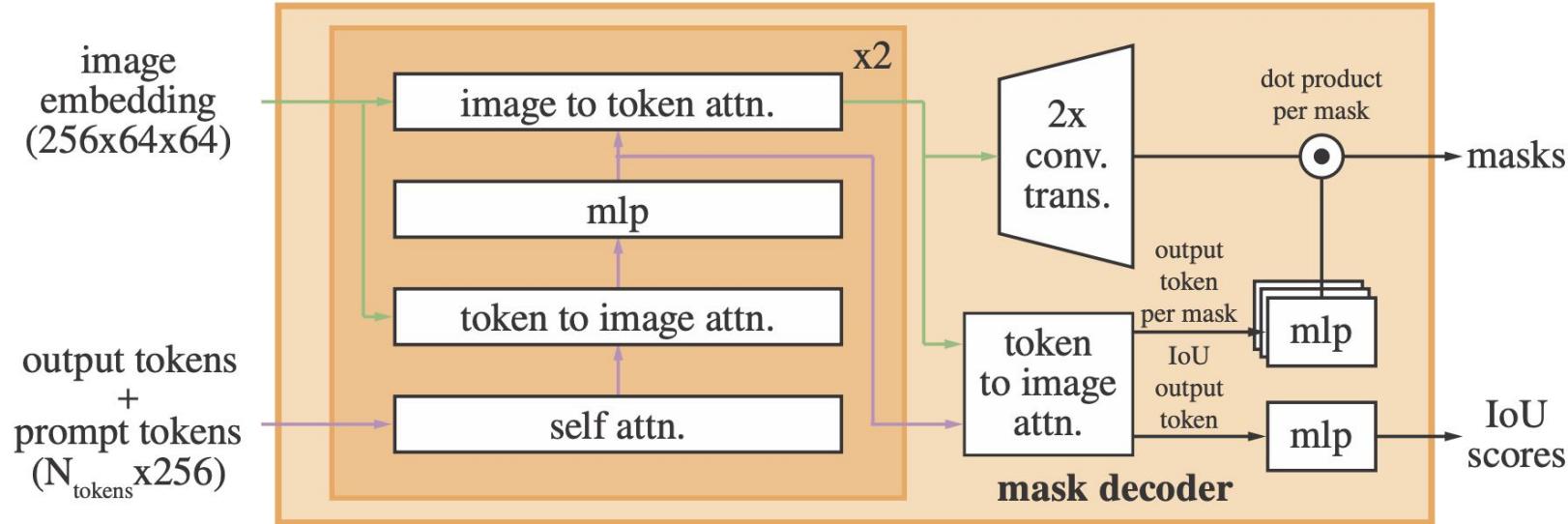
Based on Transformer segmentation models

Produces 3 masks and predicted IoU (Intersection over Union) scores

Mask Decoder



Mask Decoder



To ensure the decoder has access to critical geometric information the positional encodings are added to the image embedding whenever they participate in an attention layer.

Additionally, the entire original prompt tokens (including their positional encodings) are re-added to the updated tokens whenever they participate in an attention layer.

Simulated User interactions

Three predicted masks: Whole, part, sub-part

Focal loss + dice loss + mean-square-error for IoU predictions

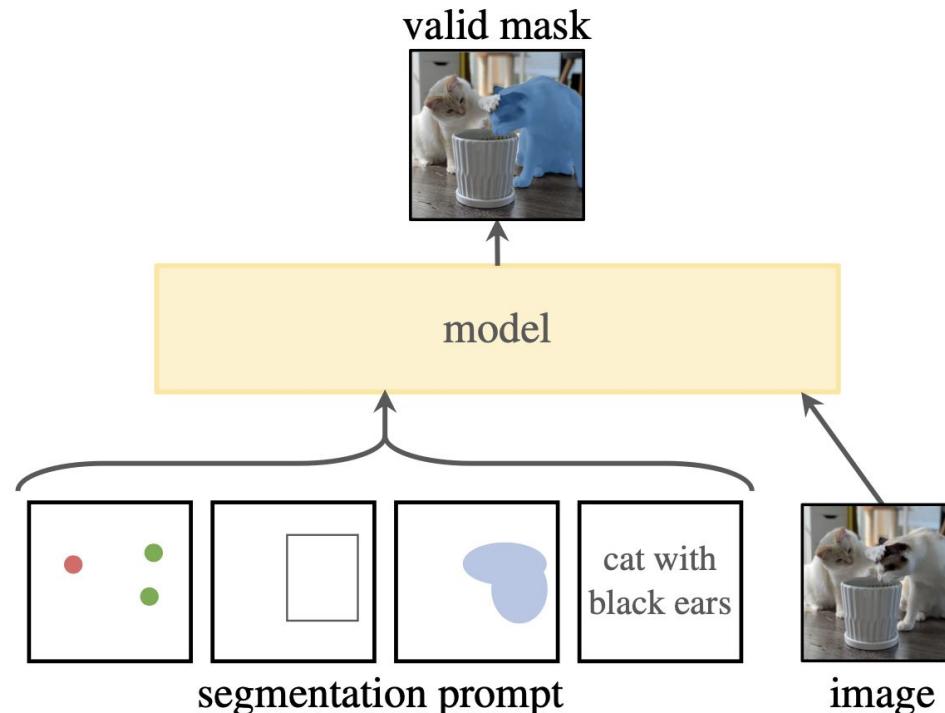
Compute 3 losses, but only backpropagate the smallest one

Simulated User interactions

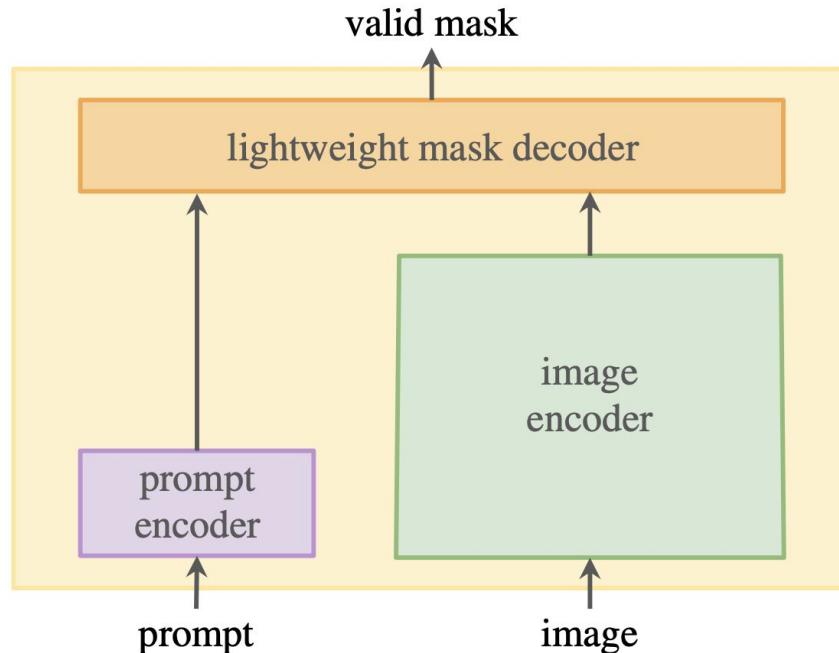
Three predicted masks: Whole, part, sub-part

Focal loss + dice loss + mean-square-error for IoU predictions

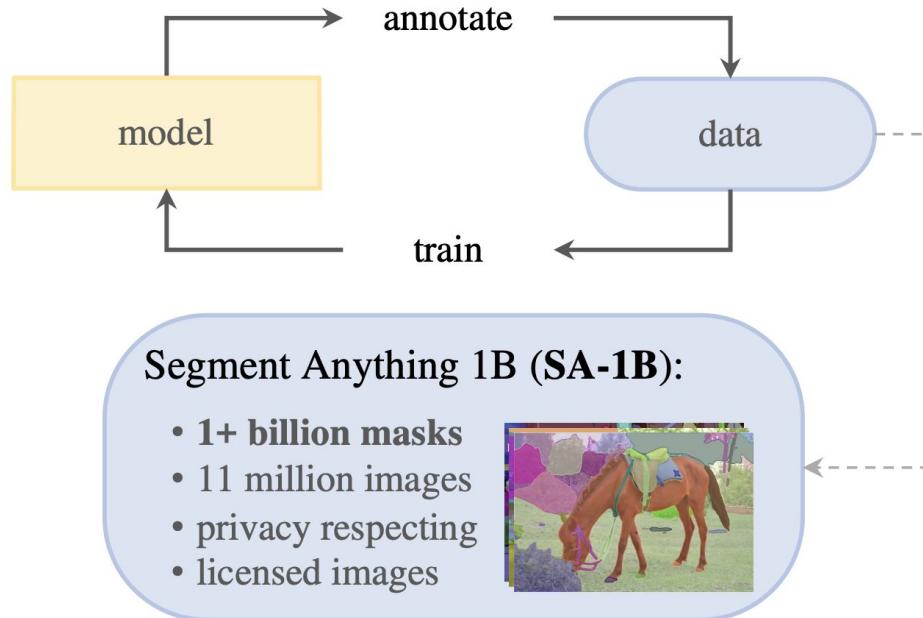
Compute 3 losses, but only backpropagate the smallest one



(a) **Task:** promptable segmentation



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



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

Manual

4.3M masks / 120k images

Semi-automatic

5.9M masks / 180k images

Fully automatic

1.1B masks / 11M images

SAM Dataset – Responsible AI

Data set card

Annotation card

Model card

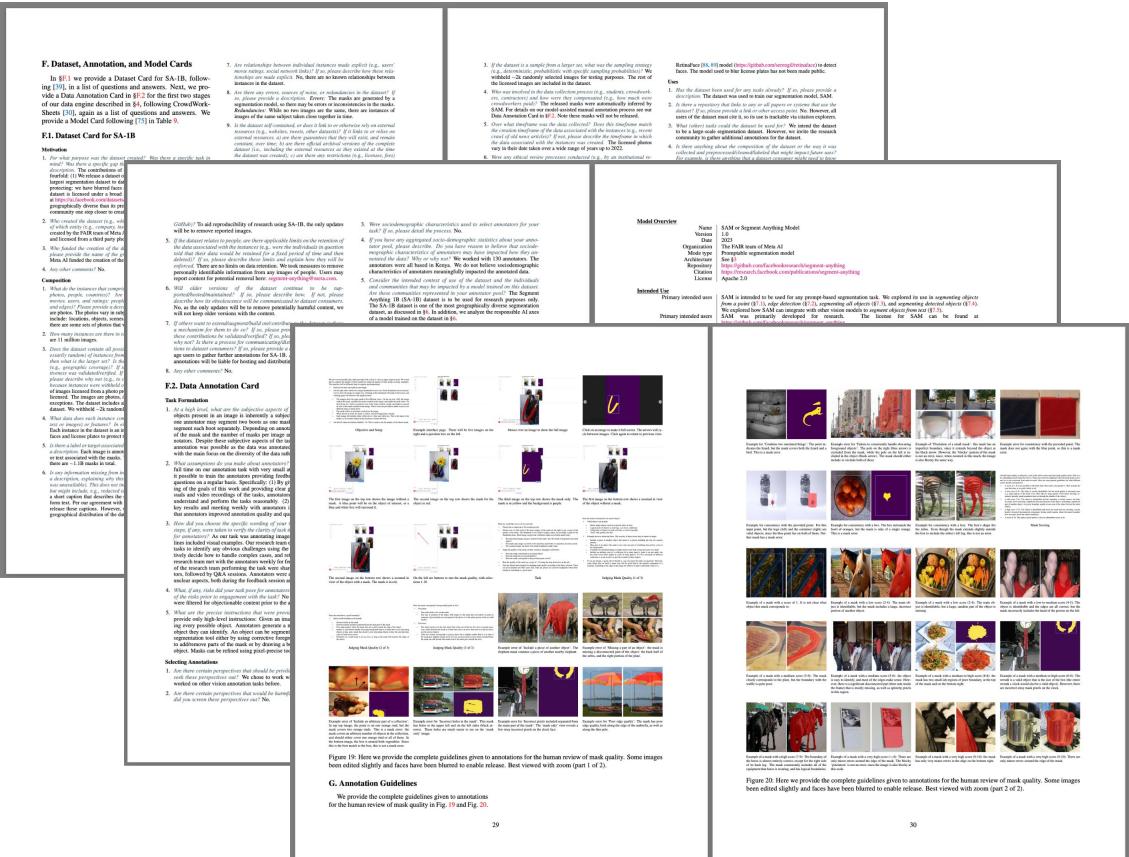


Figure 19: Here we provide the complete guidelines given to annotations for the human review of mask quality. Some images have been edited slightly and faces have been blurred to enable release. Best viewed with zoom (part 1 of 2).

G. Annotation Guidelines

We provide the complete guidelines given to annotations for the human review of mask quality in Fig. 19 and Fig. 20.

SAM Dataset – Responsible AI

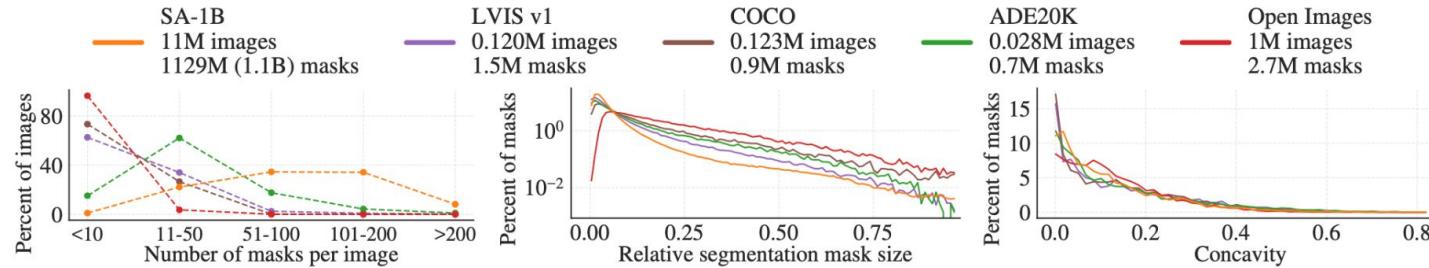


Figure 6: Dataset mask properties. The legend references the number of images and masks in each dataset. Note, that SA-1B has $11\times$ more images and $400\times$ more masks than the largest existing segmentation dataset Open Images [60].

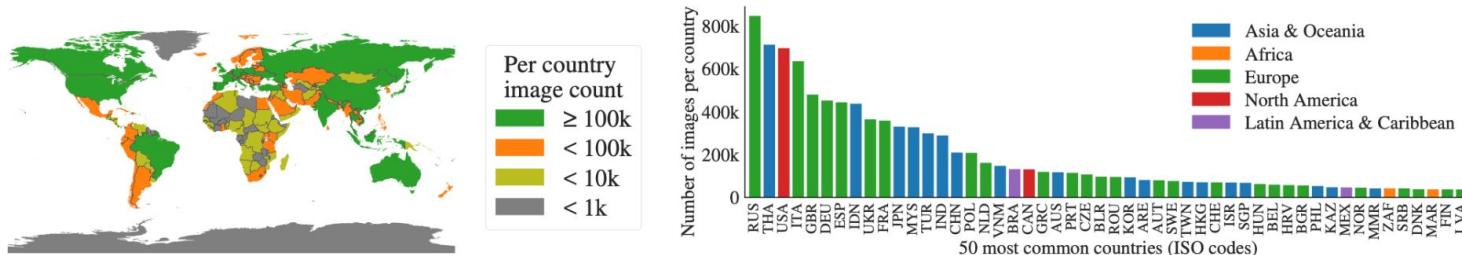
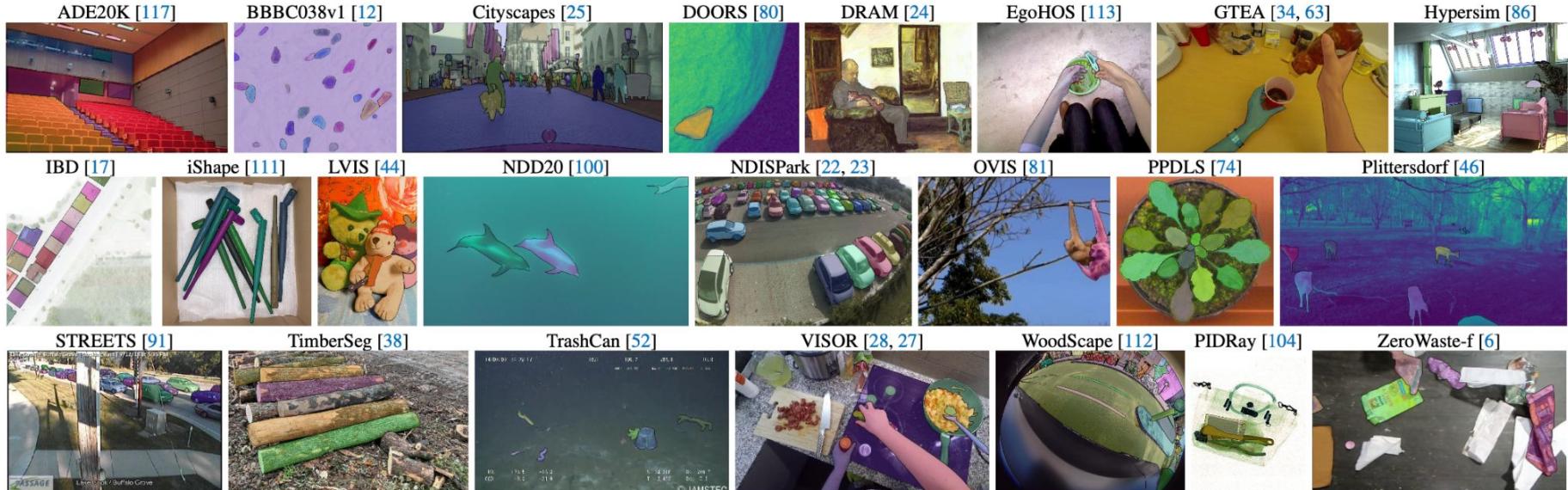


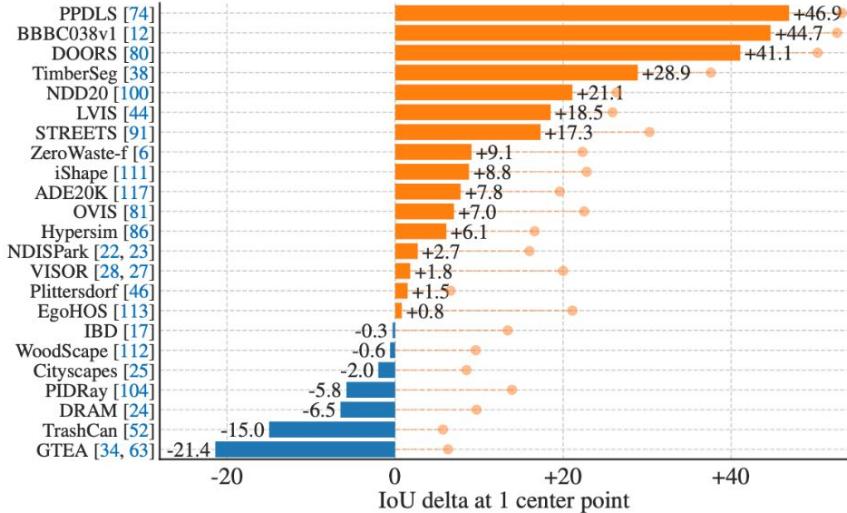
Figure 7: Estimated geographic distribution of SA-1B images. Most of the world's countries have more than 1000 images in SA-1B, and the three countries with the most images are from different parts of the world.

Zero-shot Generalization

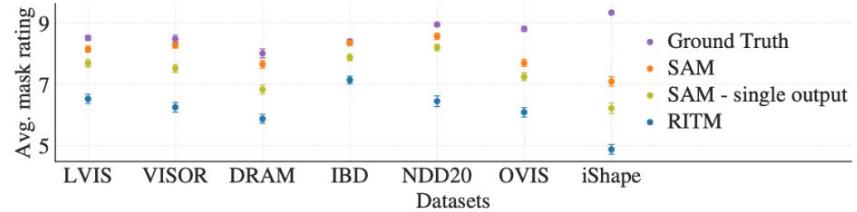
Zero-Shot Single Point Valid Mask Evaluation



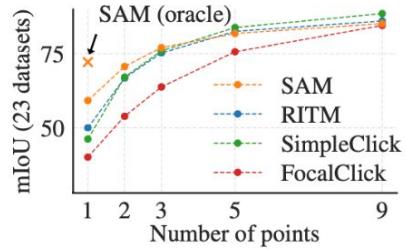
Zero-Shot Single Point Valid Mask Evaluation



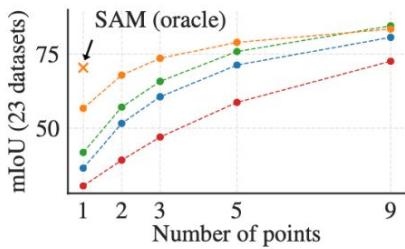
(a) SAM vs. RITM [92] on 23 datasets



(b) Mask quality ratings by human annotators



(c) Center points (default)



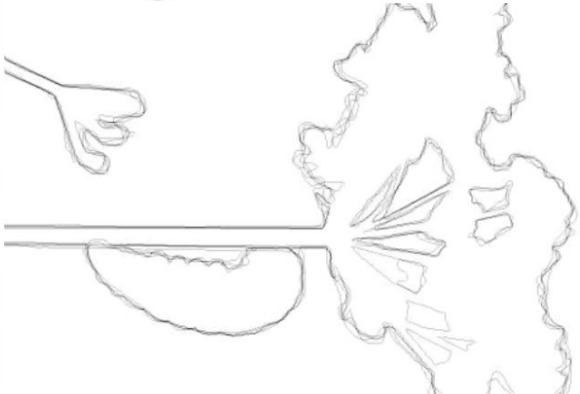
(d) Random points

Zero-Shot Edge Detection

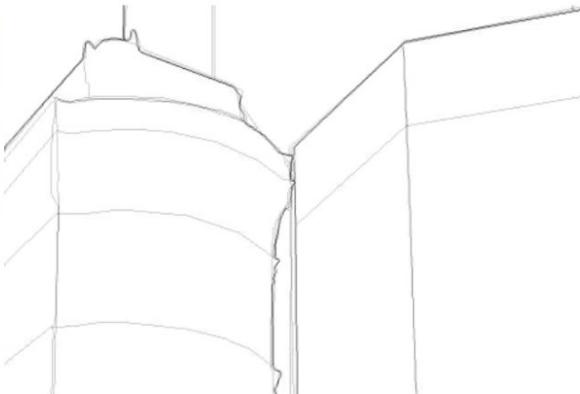
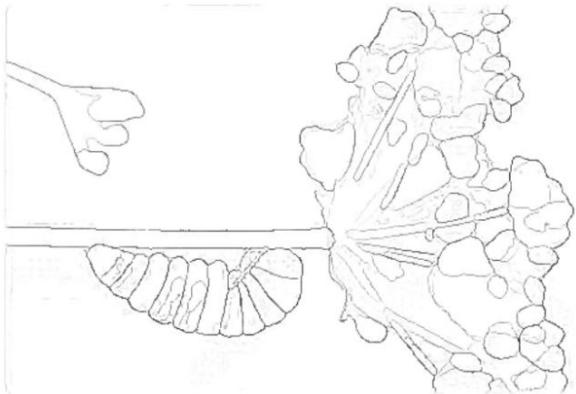
image



ground truth

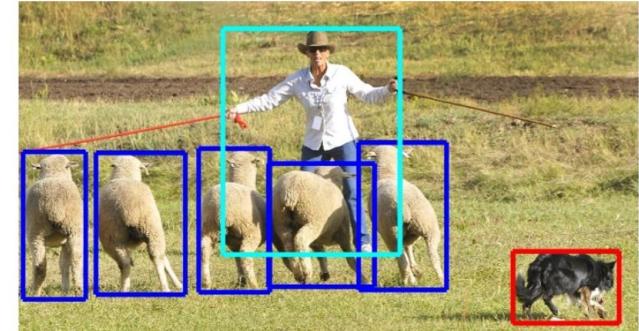


SAM



Generalization

Zero-Shot Object Proposals (Boxes)



Zero-Shot Instance Segmentation

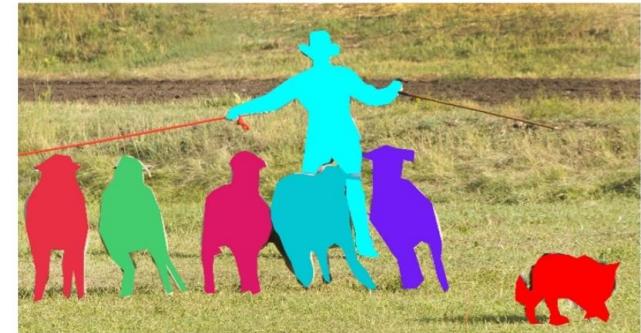
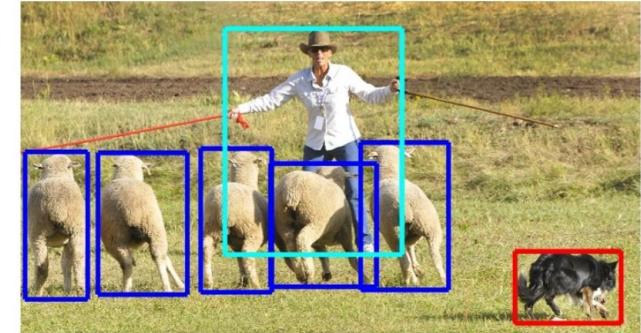


Generalization

Zero-Shot Object Proposals (Boxes)

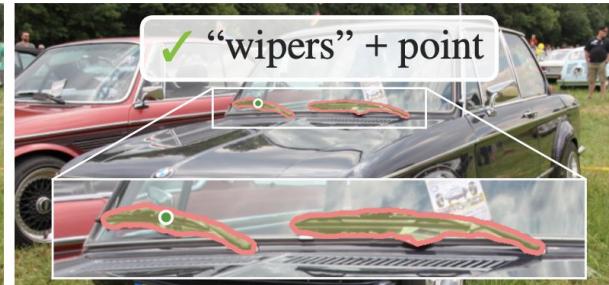
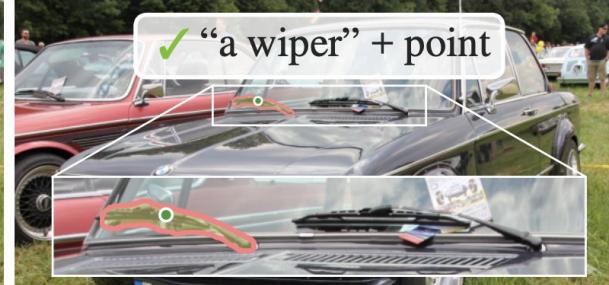


Zero-Shot Instance Segmentation



Microsoft COCO: Common Objects in Context
<https://arxiv.org/abs/1405.0312>

Zero-shot text-to-mask





Segment Anything and the Rise of Foundation Models

MEDICAL VOLUME ANNOTATOR



René Donner