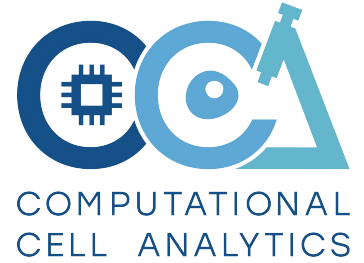




GEORG-AUGUST-UNIVERSITÄT  
GÖTTINGEN IN PUBLICA COMMODA  
SEIT 1737



# Vision Transformers & Vision Foundation Models

## Segment Anything & Applications in Microscopy

Constantin Pape

Institut für Informatik, Georg August Universität Göttingen



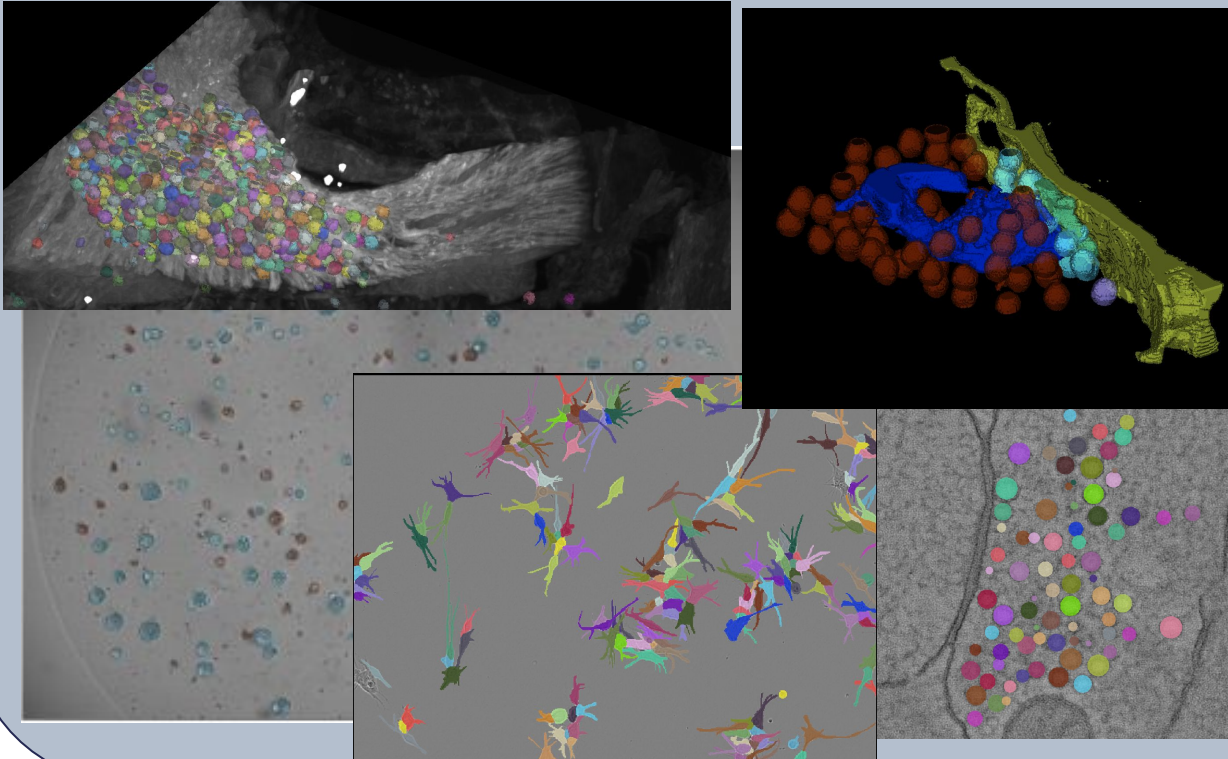
ccpape

<https://user.informatik.uni-goettingen.de/~pape41/>

# Group vision: from images to insight and clinical relevance in collaboration with life scientists

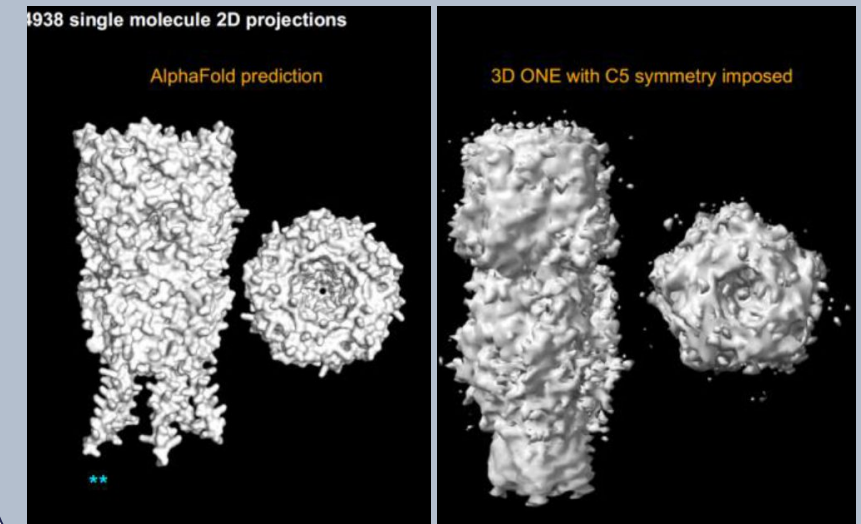
What we do ...

Segmentation and tracking lots of stuff in microscopy



Representation learning for microscopy and multi-modal data

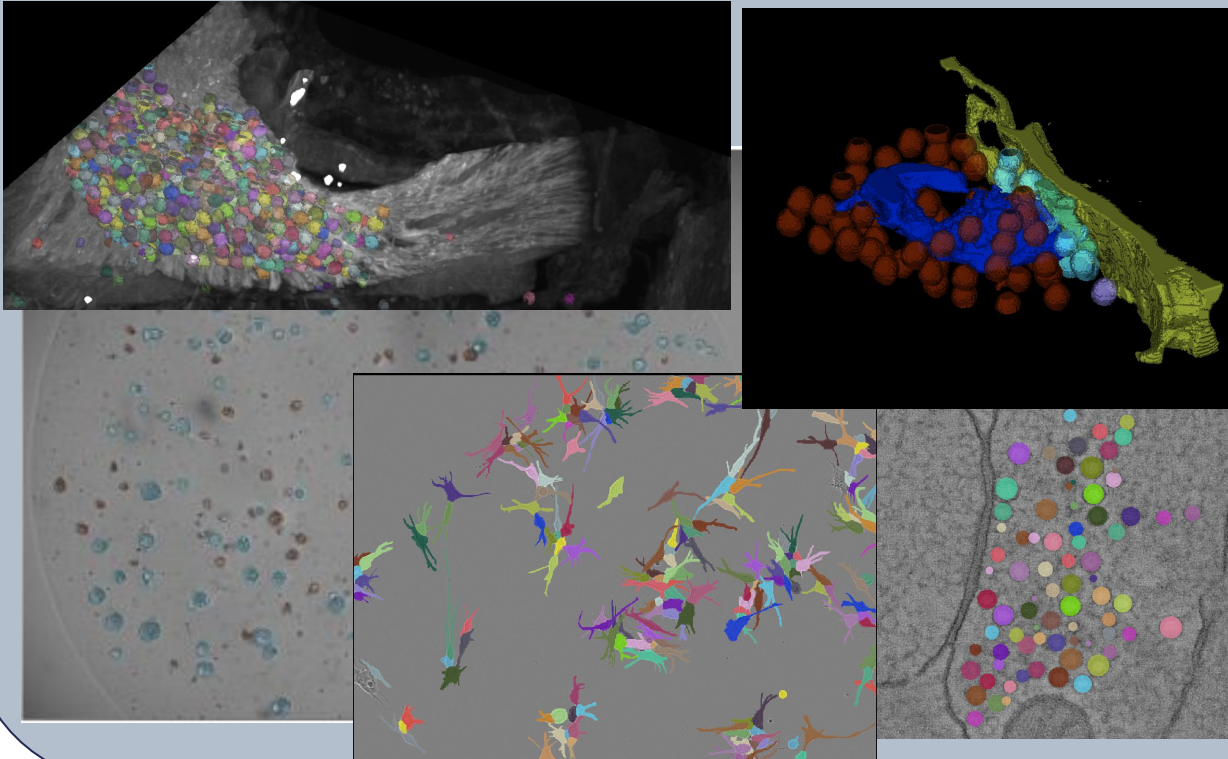
Protein structure analysis in cryo ET and **optical microscopy**



# Group vision: from images to insight and clinical relevance in collaboration with life scientists

---

## Segmentation and tracking lots of stuff in microscopy



Feasible with established DL methods  
(U-Net -> CellPose / StarDist, ...)  
**IF sufficient annotations / ground-truth!**

**Data annotation is a big bottleneck!**

**“Zero-shot” generalization is limited  
-> need retraining unless their training  
data is quite similar**

# Vision Transformers & Vision Foundation Models

---

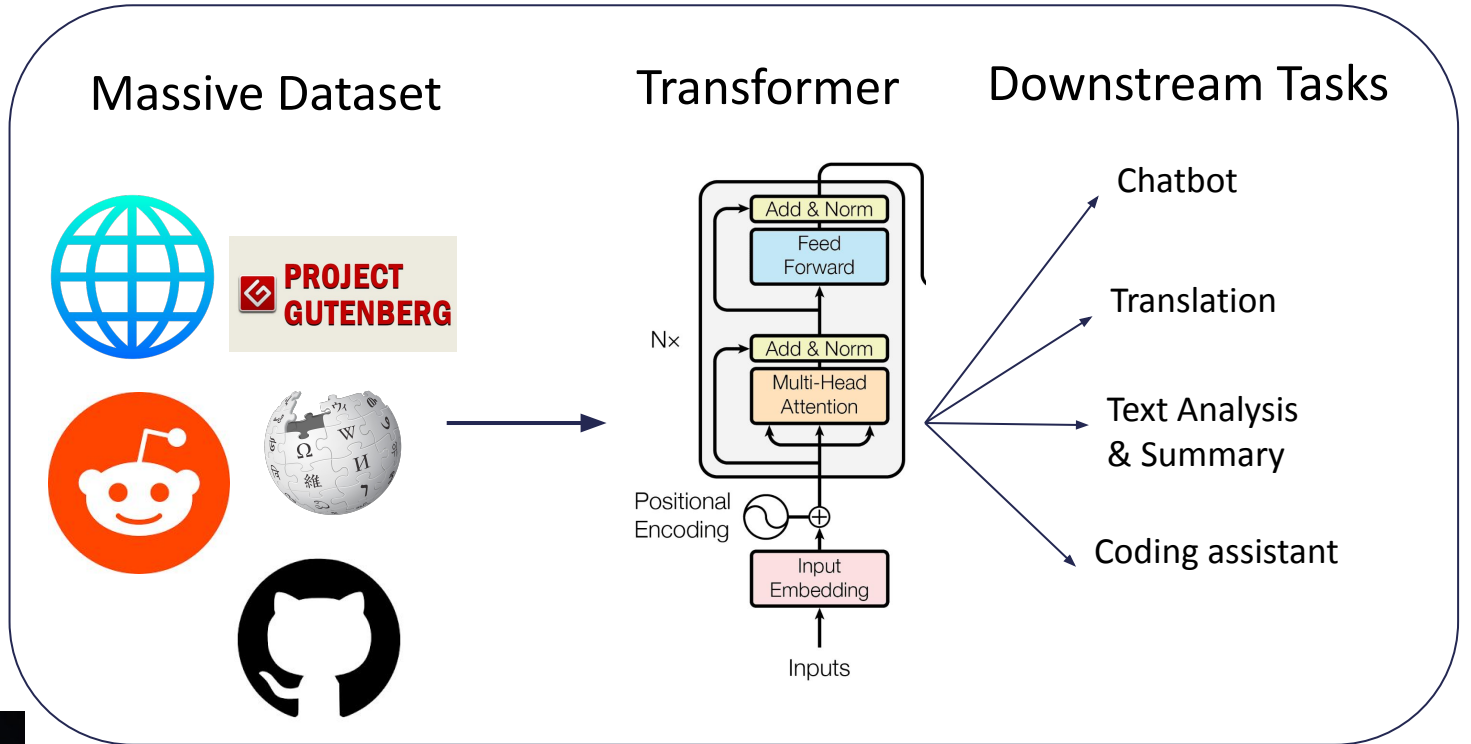


# Large language models

---



# Large language models



# Can we do the same for vision?

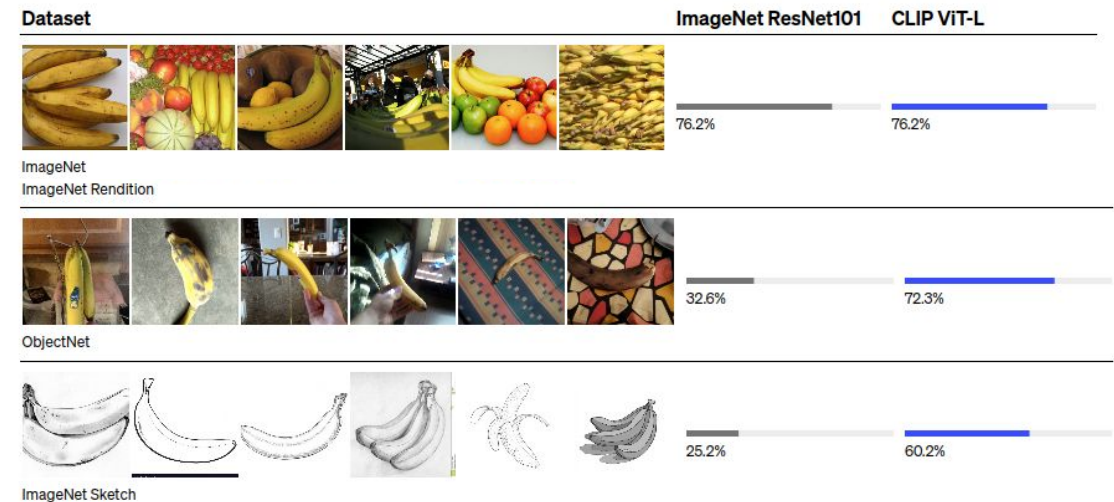
## Vision foundation models



Segment Anything

CLIP: Connecting text and images

An astronaut riding a horse in photorealistic style.



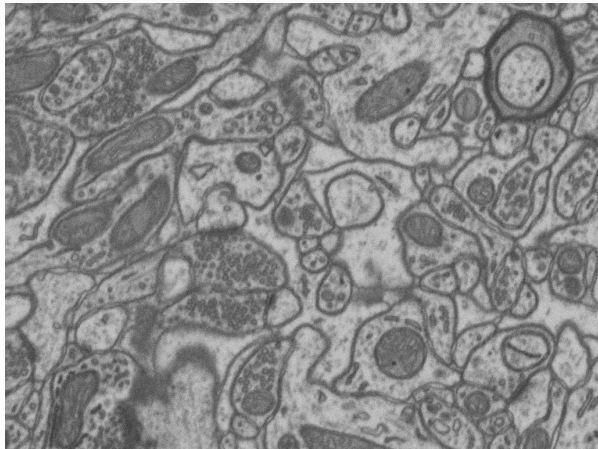


# Can we do the same for vision?

## Vision foundation models

### GPT4 Vision

You



What is the content of this image?

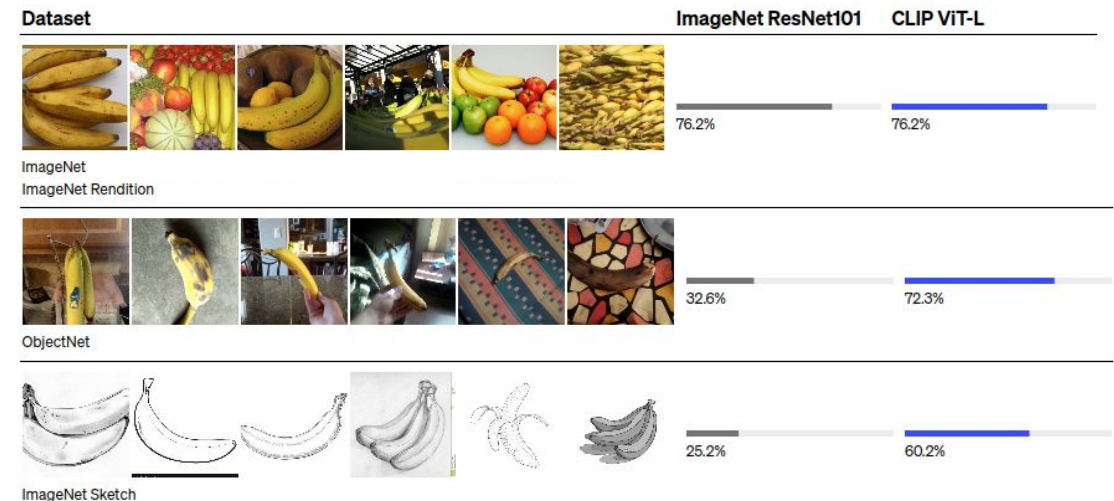
**ChatGPT**  
The image you've uploaded appears to be a microscopic view of biological cells or tissue, showing various cellular structures at high magnification. These could include organelles such as mitochondria, vesicles, and parts of the cell nucleus among others, commonly observed in electron microscopy images. The detailed textures and contrasts typically indicate different components within the cells, which are essential for various biological functions.



Segment Anything

CLIP: Connecting text and images

An astronaut riding a horse in photorealistic style.



# What is a transformer?

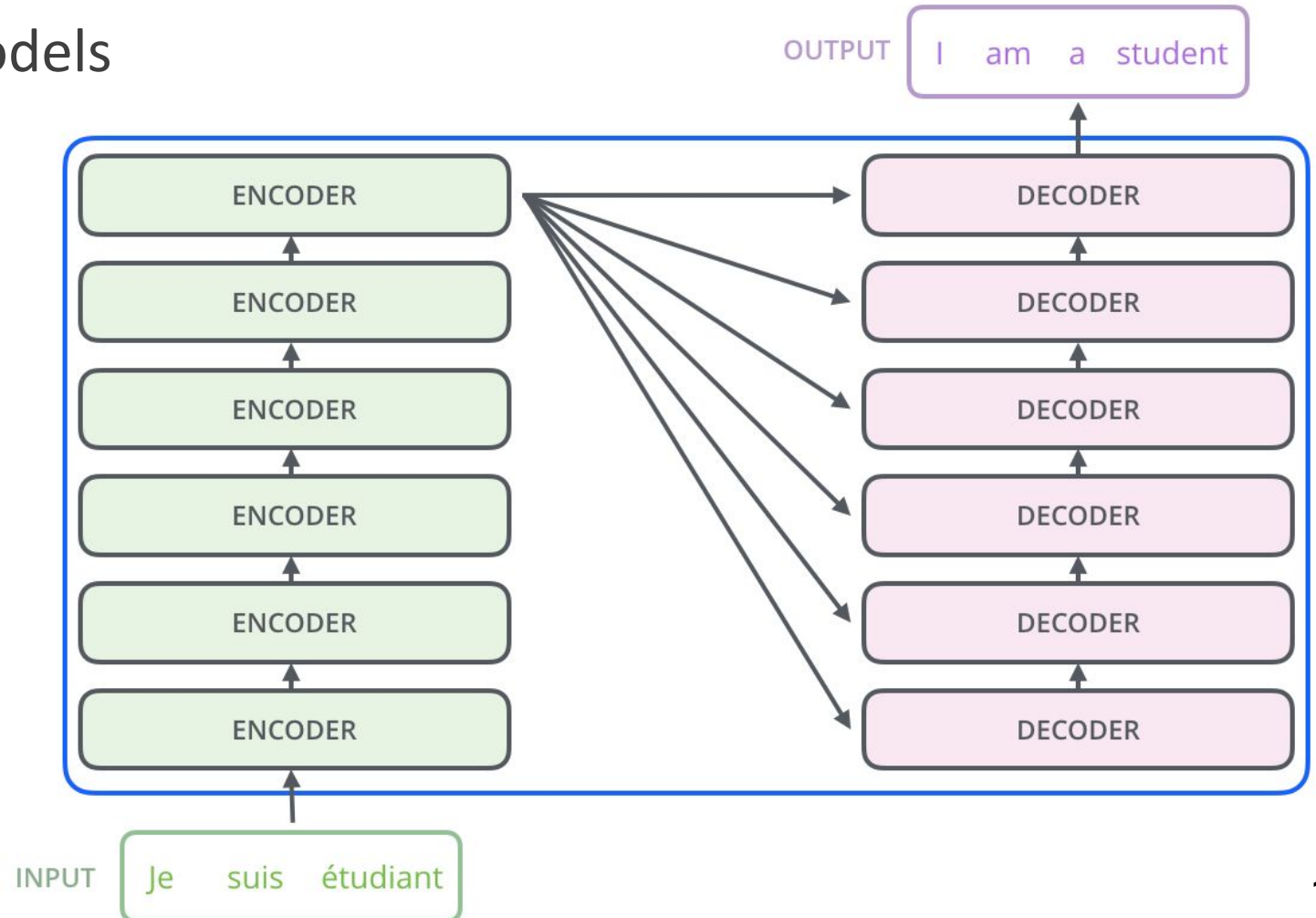
---



# Transformer

## Sequence to sequence models

- Translation
- Text generation
- ...

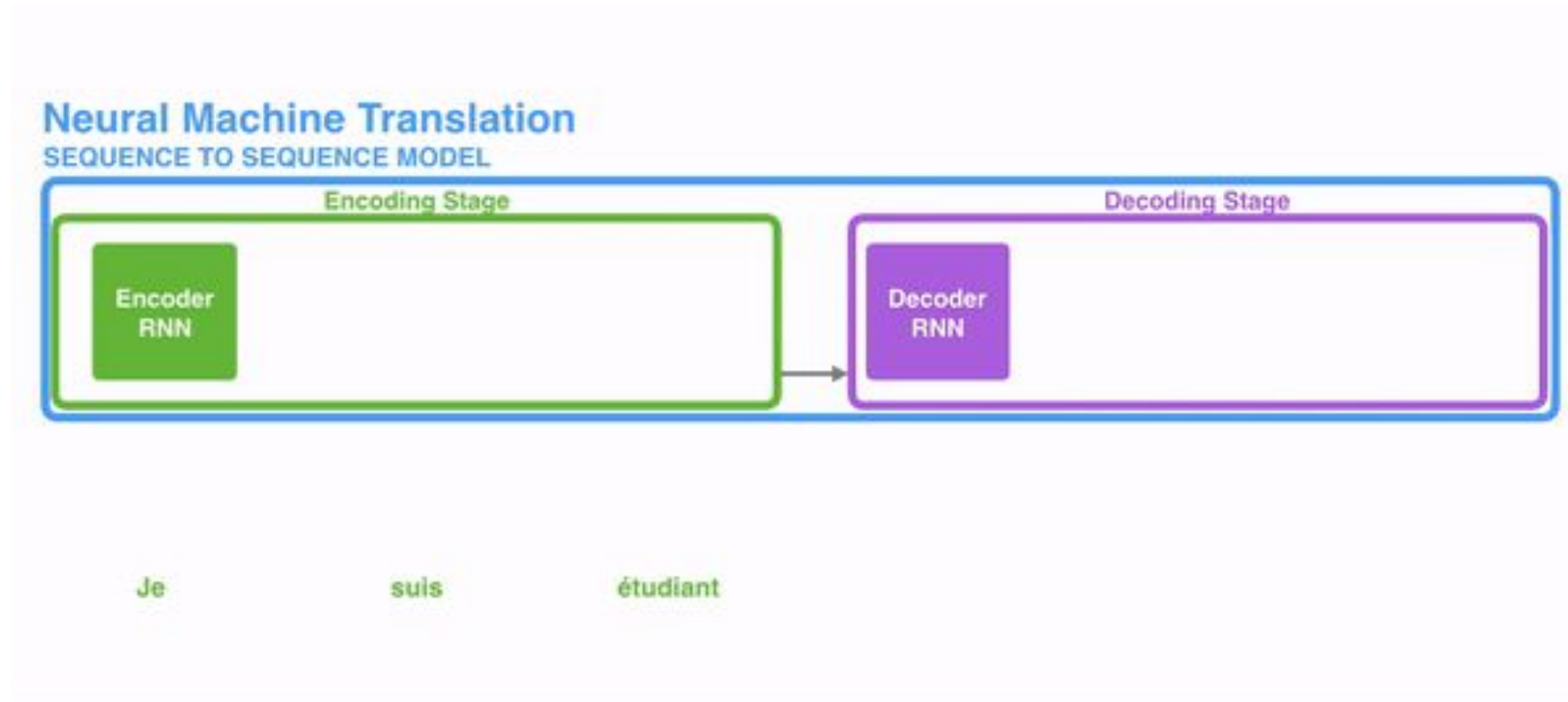


# Why transformers?

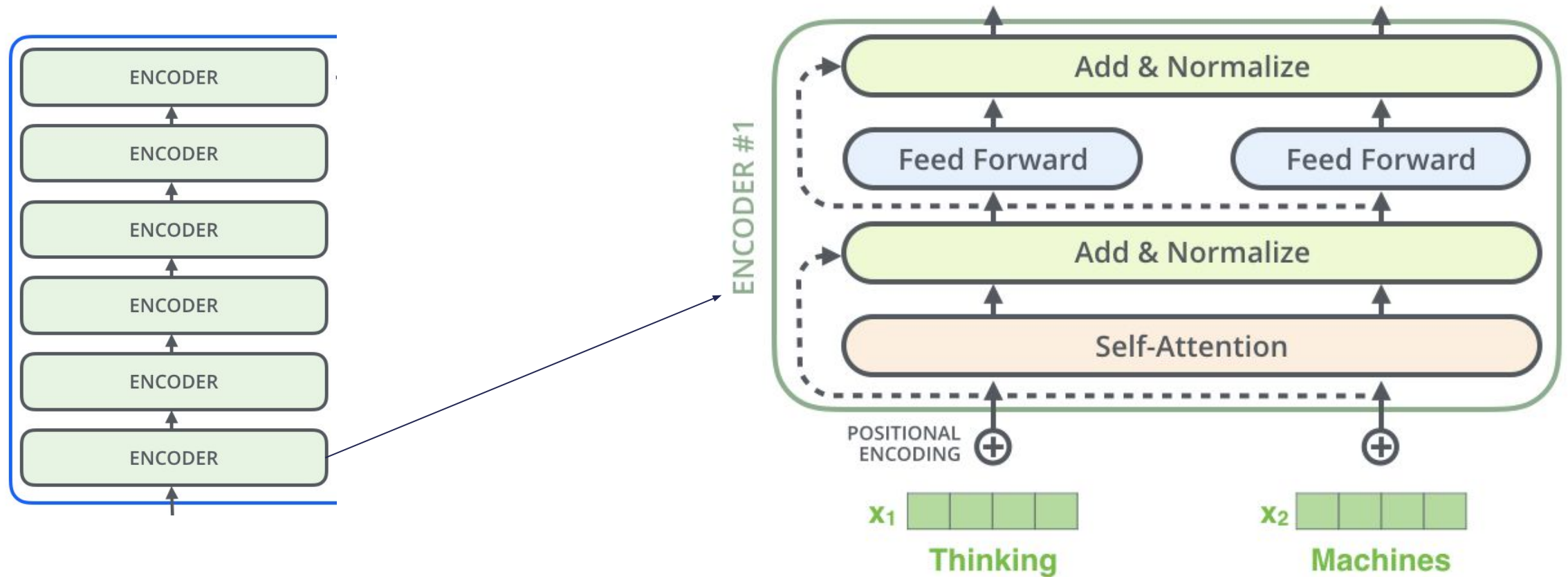
---

Previously: Recurrent Neural Networks (RNNs) like LSTM

- Problem: information decays over steps / “not enough memory”



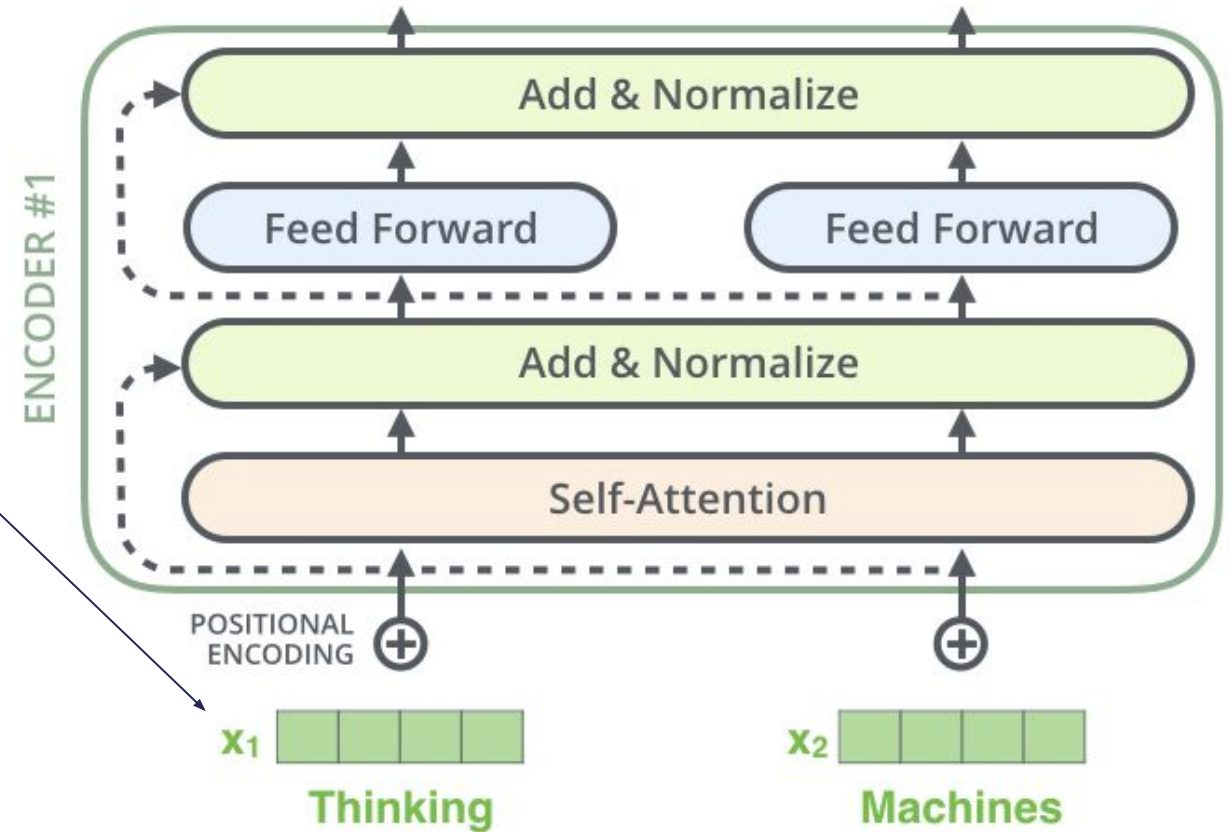
# Transformer Block



# Transformer Block

The input: “tokenize” elements of your input  
sequence: compute a vector for each  
element

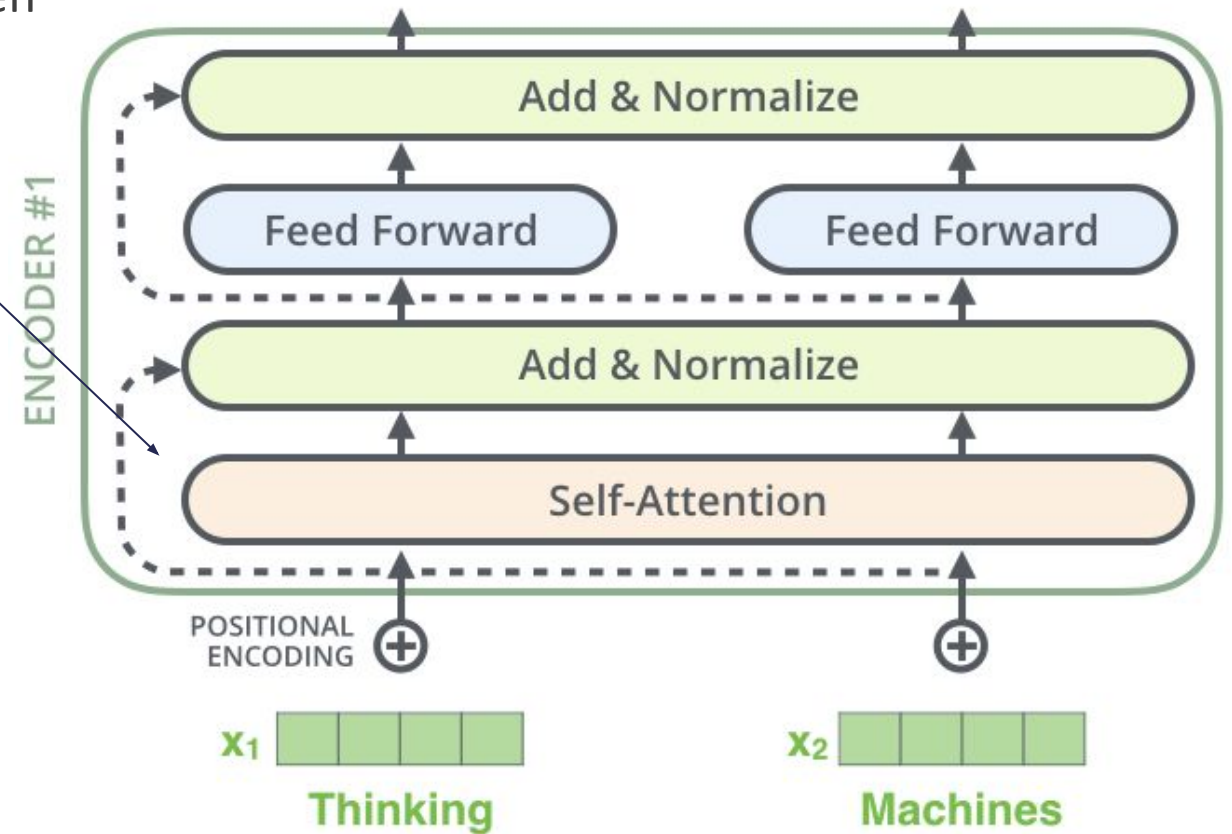
**Problem dependent!**



# Transformer Block

Self-attention: propagate information between tokens, all-to-all connectivity

Attention: learns “affinities” between pairs of elements

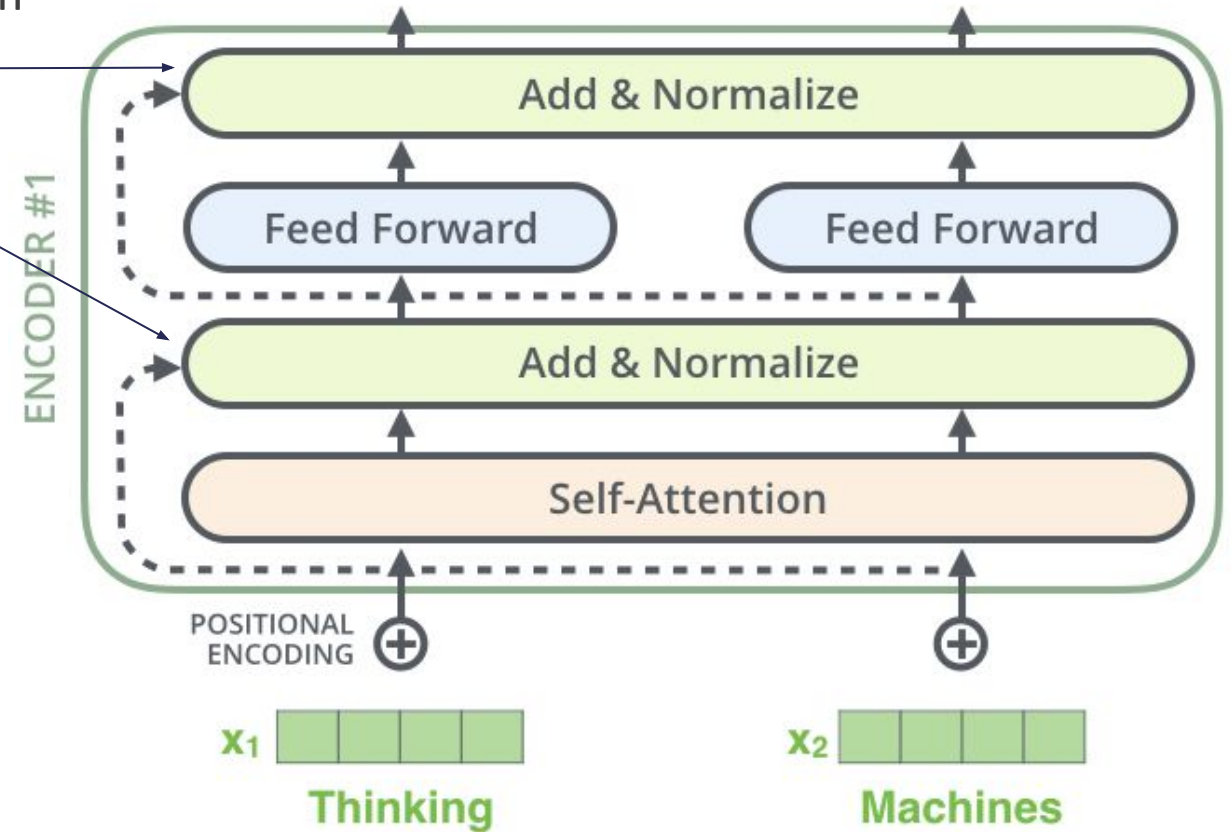




# Transformer Block

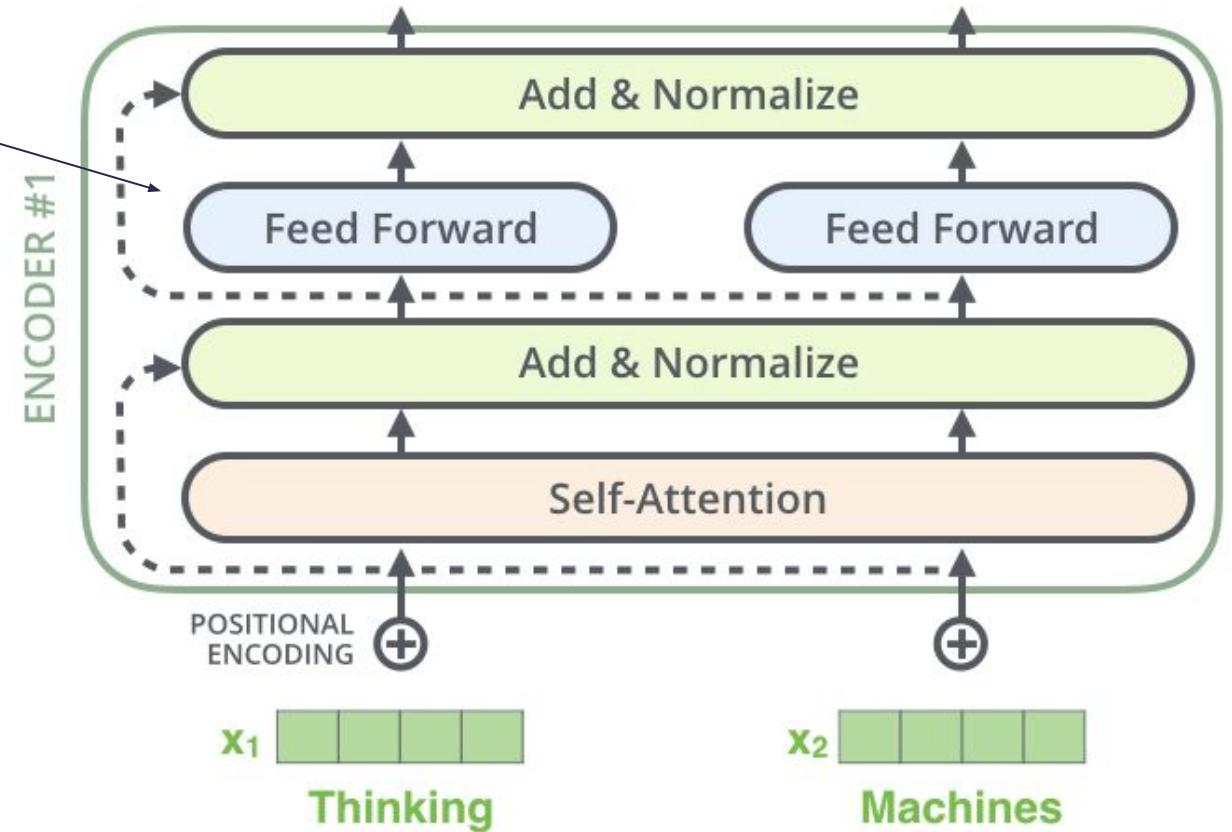
Add & Normalize: propagate information from before the layer, like ResNet

Learn residual transformations, helps gradient flow



# Transformer Block

Feed forward: process the individual tokens  
= individual elements



# Self-Attention

---

$$\text{softmax}\left(\frac{\begin{matrix} \text{Q} \\ \begin{array}{|c|c|c|} \hline \square & \square & \square \\ \hline \square & \square & \square \\ \hline \end{array} \end{matrix} \times \begin{matrix} \text{K}^T \\ \begin{array}{|c|c|} \hline \square & \square \\ \hline \square & \square \\ \hline \end{array} \end{matrix}}{\sqrt{d_k}}\right) \begin{matrix} \text{V} \\ \begin{array}{|c|c|c|} \hline \square & \square & \square \\ \hline \square & \square & \square \\ \hline \end{array} \end{matrix}$$
$$= \begin{matrix} \text{Z} \\ \begin{array}{|c|c|c|} \hline \square & \square & \square \\ \hline \square & \square & \square \\ \hline \end{array} \end{matrix}$$

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

# Self-Attention

$$\text{softmax}\left(\frac{\begin{matrix} \text{Q} \\ \begin{matrix} \square & \square & \square \\ \square & \square & \square \end{matrix} \end{matrix} \times \begin{matrix} \text{K}^T \\ \begin{matrix} \square & \square \\ \square & \square \\ \square & \square \end{matrix} \end{matrix}}{\sqrt{d_k}}\right) \begin{matrix} \text{V} \\ \begin{matrix} \square & \square & \square \\ \square & \square & \square \end{matrix} \end{matrix}$$

$$= \begin{matrix} \text{Z} \\ \begin{matrix} \square & \square & \square \\ \square & \square & \square \end{matrix} \end{matrix}$$

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Input

Embedding

Queries

Keys

Values

Score

Divide by 8 (  $\sqrt{d_k}$  )

Softmax

Thinking

Machines

$x_1$

$x_2$

$q_1$

$q_2$

$k_1$

$k_2$

$v_1$

$v_2$

$q_1 \cdot k_1 = 112$

$q_1 \cdot k_2 = 96$

14

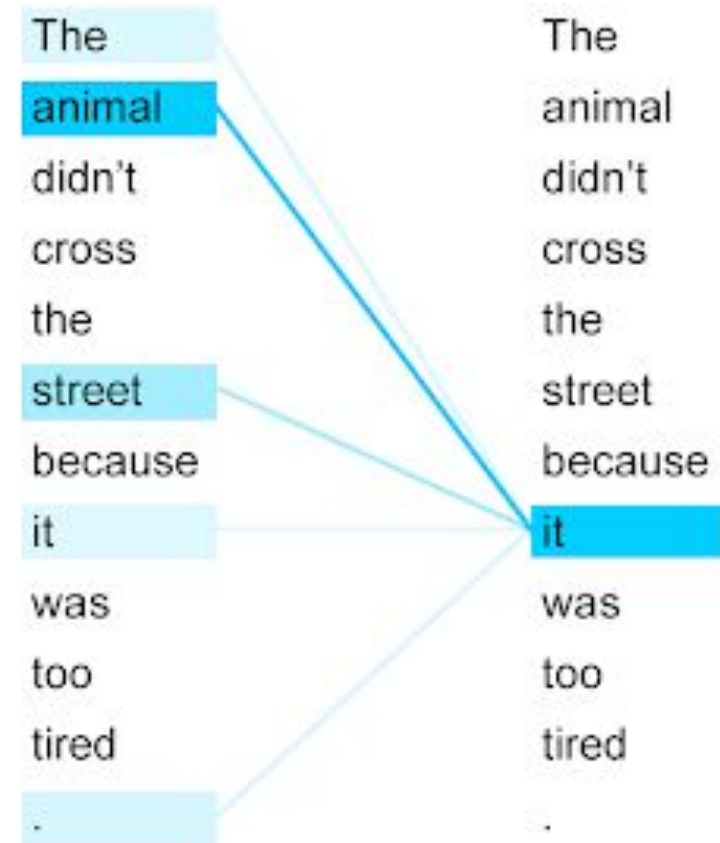
12

0.88

0.12

# Self-Attention

---





# Where is the learning happening???

---

# Where is the learning happening???

---

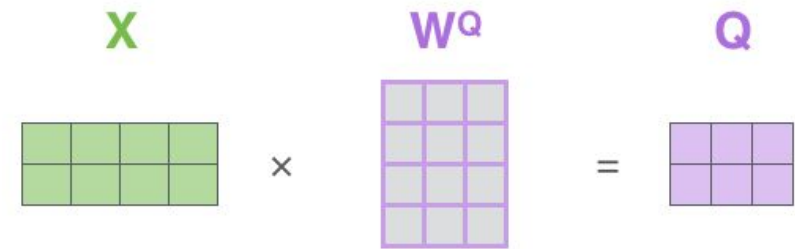
How do we obtain **Q**, **K**, **V**?

# Where is the learning happening???

How do we obtain **Q**, **K**, **V**?

Computed from **X** (input tokens) with *learned* weights.

A fully connected layer.



# Self-Attention vs. Feed-Forward

---

Self-Attention:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

$$Q = XW_q$$

$$K = XW_k$$

$$V = XW_v.$$

Fully connected layer:

$$\text{Out} = \text{activation}(WX)$$

# Self-Attention vs. Feed-Forward

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

$$Q = XW_q$$

$$K = XW_k$$

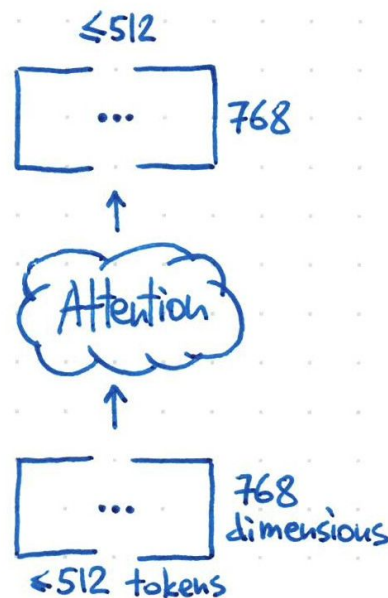
$$V = XW_v$$

$$\text{Out} = \text{activation}(WX)$$

Self-attention:

-  $3 \cdot 768^2 = 2\text{M}$  params

- works on sequence of any length



Fully-connected:

-  $(512 \cdot 768)^2 = 154\text{B}$  params

- sequence length (512) is baked in



# Self-Attention vs. Feed-Forward

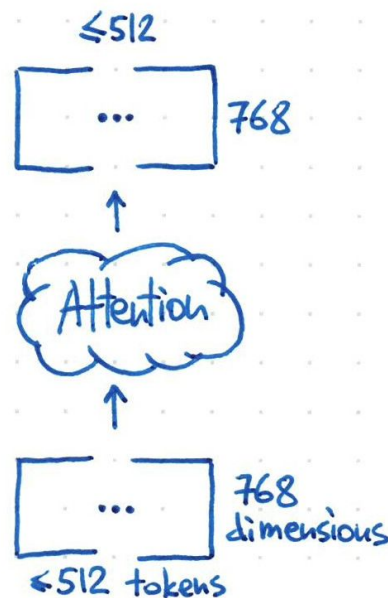
$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

$$Q = XW_q$$

$$K = XW_k$$

$$V = XW_v$$

**Learnable  
parameters**



$$\text{Out} = \text{activation}(WX)$$

**Learnable  
parameters**

Fully-connected:

$$-(512 \cdot 768)^2 = 154 \text{ B params}$$

- sequence length (512) is baked in

Self-attention:

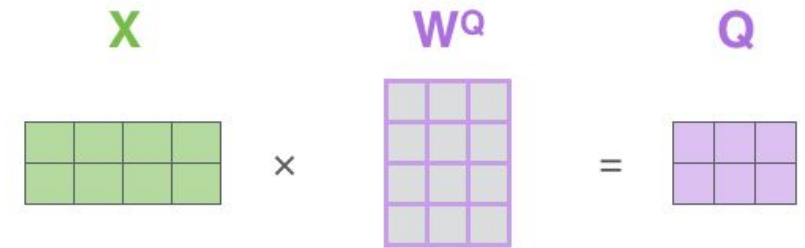
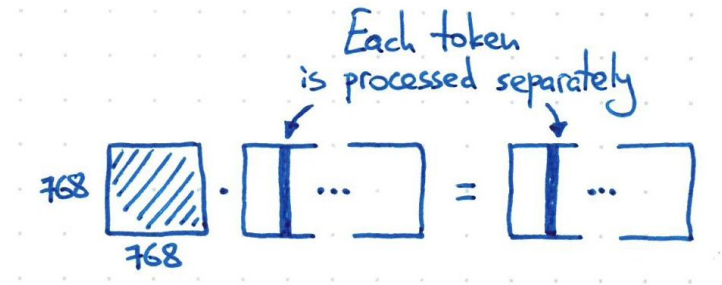
$$- 3 \cdot 768^2 = 2 \text{ M params}$$

- works on sequence of any length

# Self-Attention in a nutshell

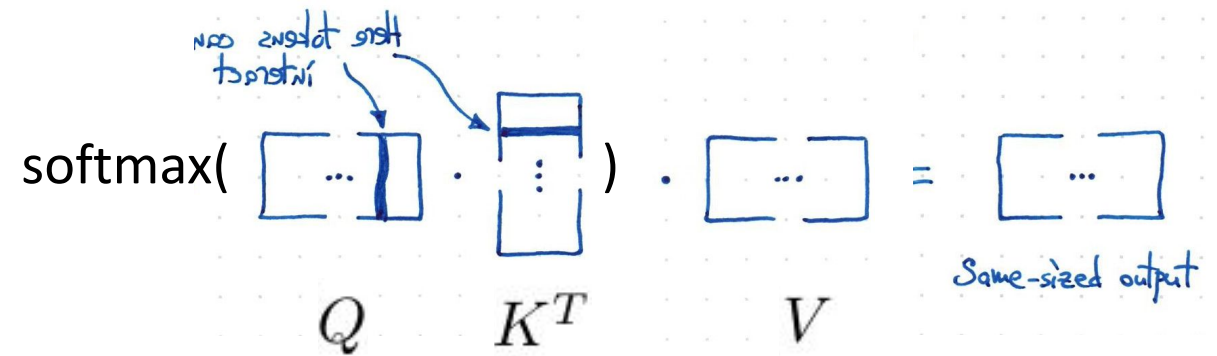
The weights:

$W_q \ W_k \ W_v$



## Self-attention:

- arbitrary sequence length
- fixed number of learnable parameters
- encodes pairwise affinity between sequence elements



# Self-Attention in a nutshell

The weights:

$W_q \quad W_k \quad W_v$



Arbitrary sequence length, fixed number of parameters!  
Did we get a free lunch???

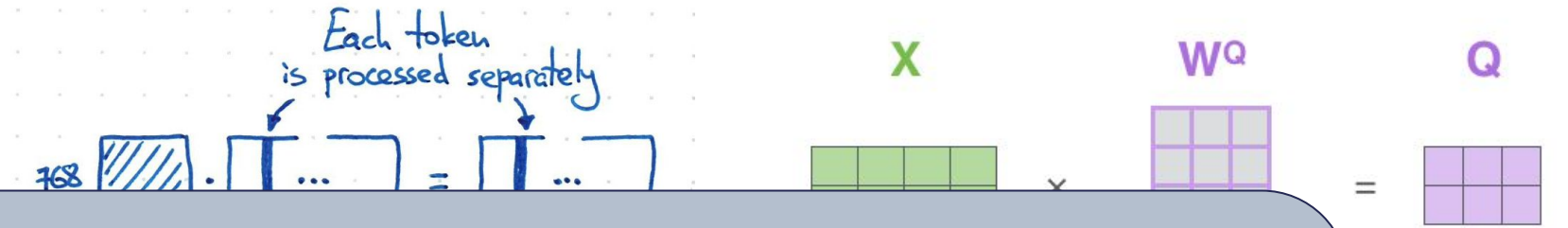
## Self-attention:

- arbitrary sequence length
- fixed number of parameters
- encodes pairwise affinity between sequence elements

# Self-Attention in a nutshell

The weights:

$W_q \ W_k \ W_v$



Arbitrary sequence length, fixed number of parameters!  
Did we get a free lunch???

## Self-attention:

- arbitrary sequence length
- fixed number of parameters
- encodes pairwise affinity between sequence elements

Quadratic complexity in sequence length (at runtime):  
Self-attention matrix is of size:  $n_{\text{elements}} \times n_{\text{elements}}$

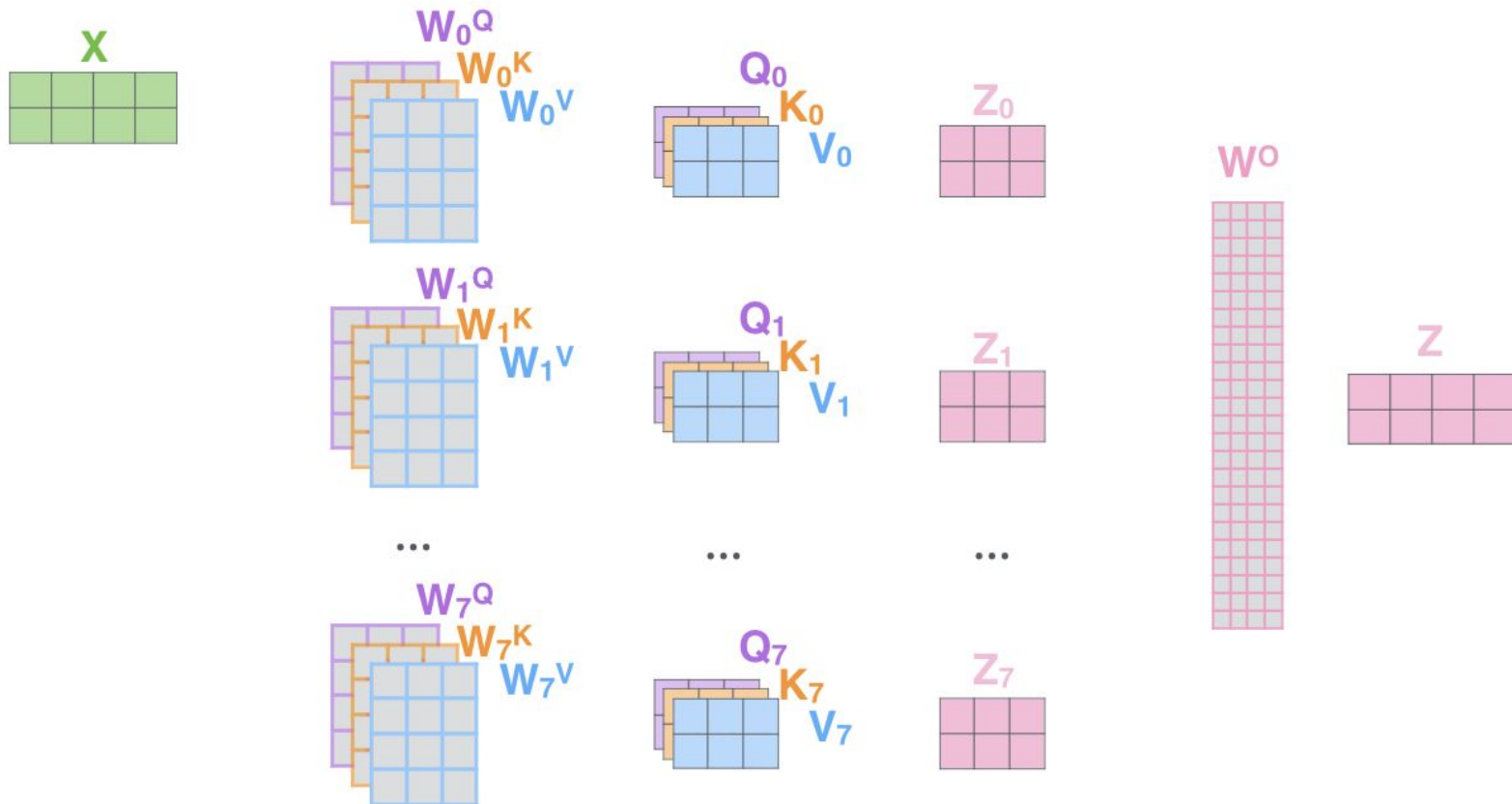
$Q \ K^T \ V$

Same-sized output

# Multi-head self-attention

One attention head is limited  
in learning interactions.

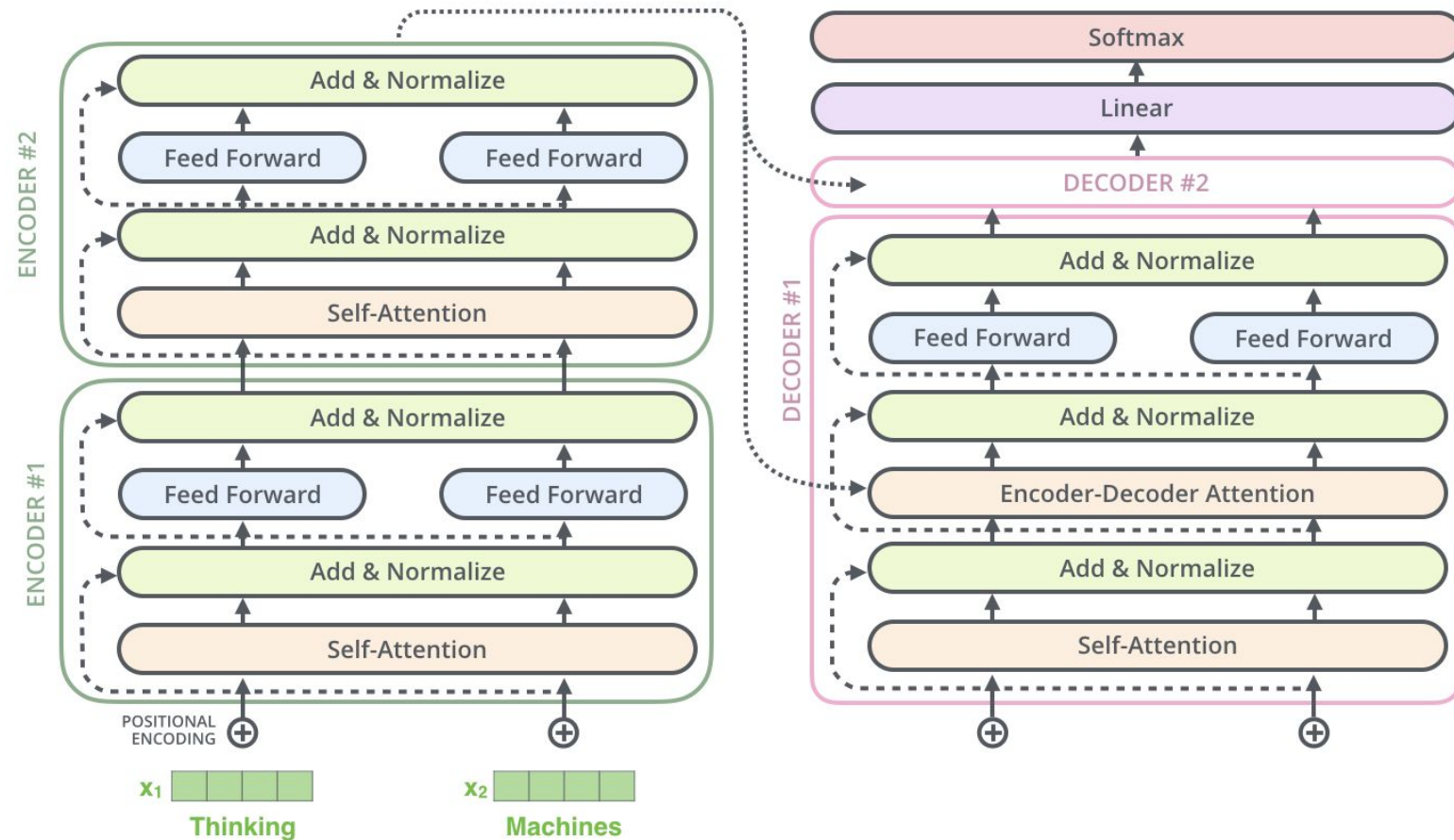
Combine multiple attention  
heads!



# The transformer architecture

Encoder:  
stack transformer blocks.  
Maps to internal sequence  
representation

Decoder:  
*Encoder-Decoder Attention*  
uses outputs from last  
layer as keys



BERT: <https://arxiv.org/abs/1810.04805>

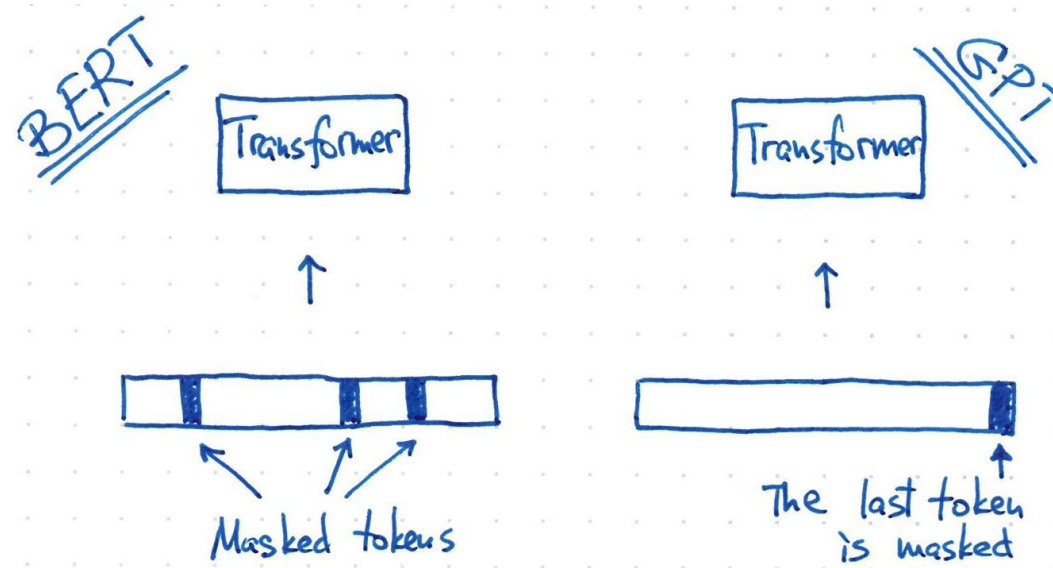
GPT3: <https://arxiv.org/abs/2005.14165>

Instruction Finetuning: <https://arxiv.org/abs/2203.02155>

RLHF: <https://arxiv.org/abs/1909.08593>

# Transformers in NLP: Large Language Models

- Autoregressive pretraining: predict part of sequence from the rest, trained on large corpus
  - BERT: Mask random tokens
  - GPT: Mask last token
  - ChatGPT, GPT4, LLAMA, ...



- Application to downstream tasks (Translation, Chatbot, ...)
  - Instruction fine-tuning (training on supervised samples)
  - Reinforcement learning from human feedback (RLHF) (rating of responses)



# Vision transformers

---

# Why and how?

---

Motivations:

- Success in NLP
- “Scaling laws”: transformer performance scales with dataset size
- Sequence representation is universal -> combine vision, text etc.
- Remove inductive biases?

# Why and how?

---

Motivations:

- Success in NLP
- “Scaling laws”: transformer performance scales with dataset size
- Sequence representation is universal -> combine vision, text etc.
- Remove inductive biases?

**How can we represent an image as a sequence?**

# Why and how?

---

Motivations:

- Success in NLP
- “Scaling laws”: transformer performance scales with dataset size
- Sequence representation is universal -> combine vision, text etc.
- Remove inductive biases?

**How can we represent an image as a sequence?**

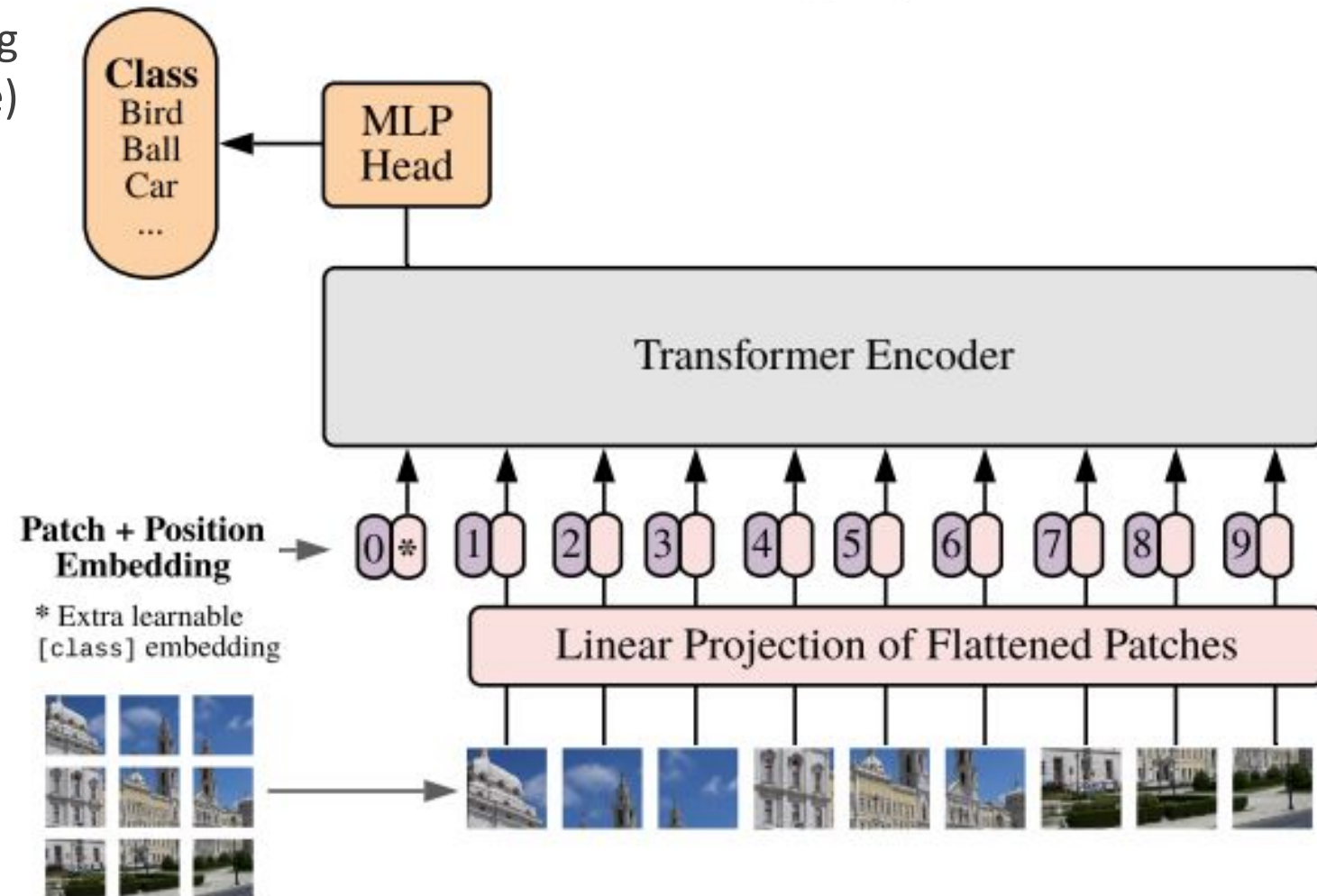
- Tokenize every pixel: too many elements! (Remember quadratic complexity)
- Tokenize small patches!

# The Vision Transformer (ViT)

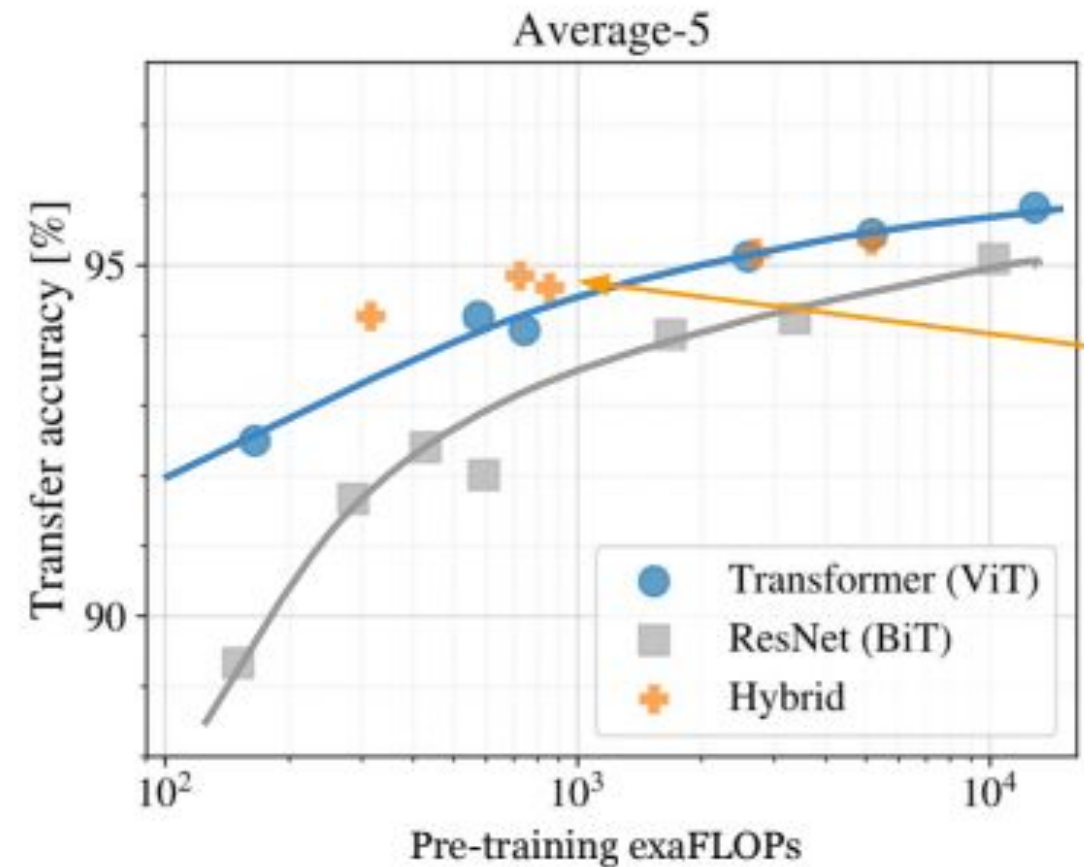
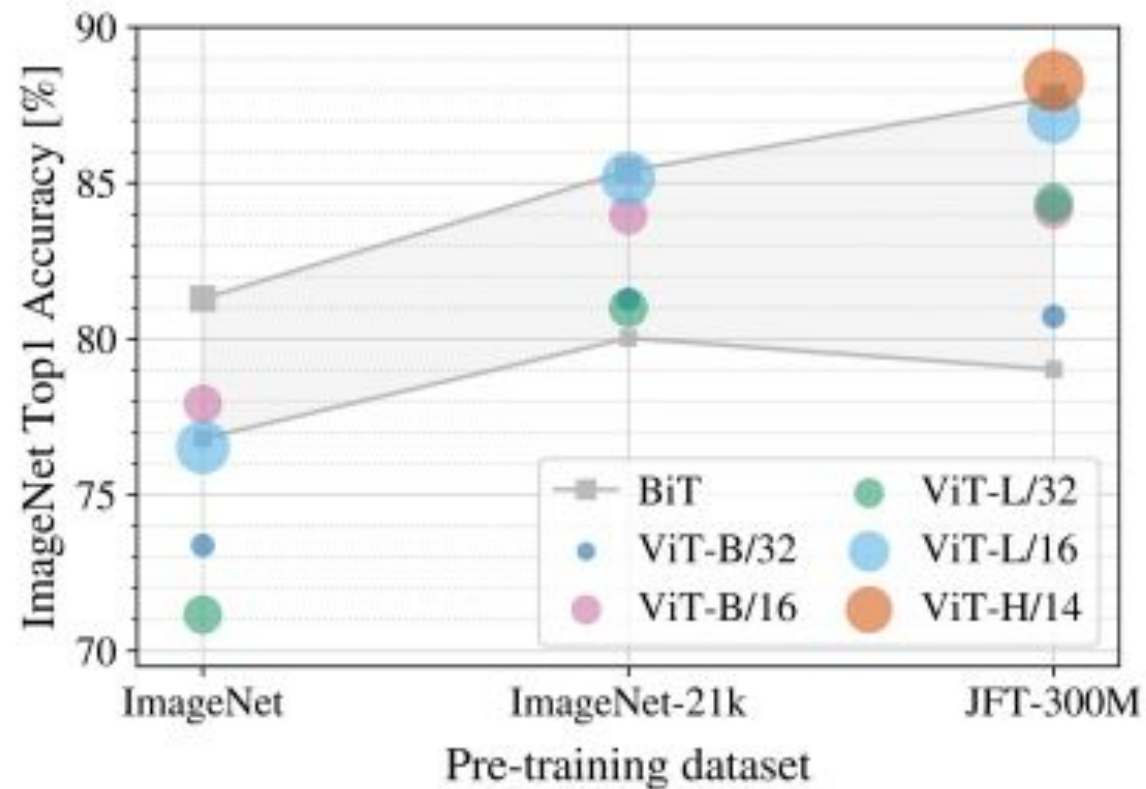
Tokenize patches + position embedding  
(encode position of patch in the image)

Encoder: same as in NLP

Classification based head on top of  
transformer output.



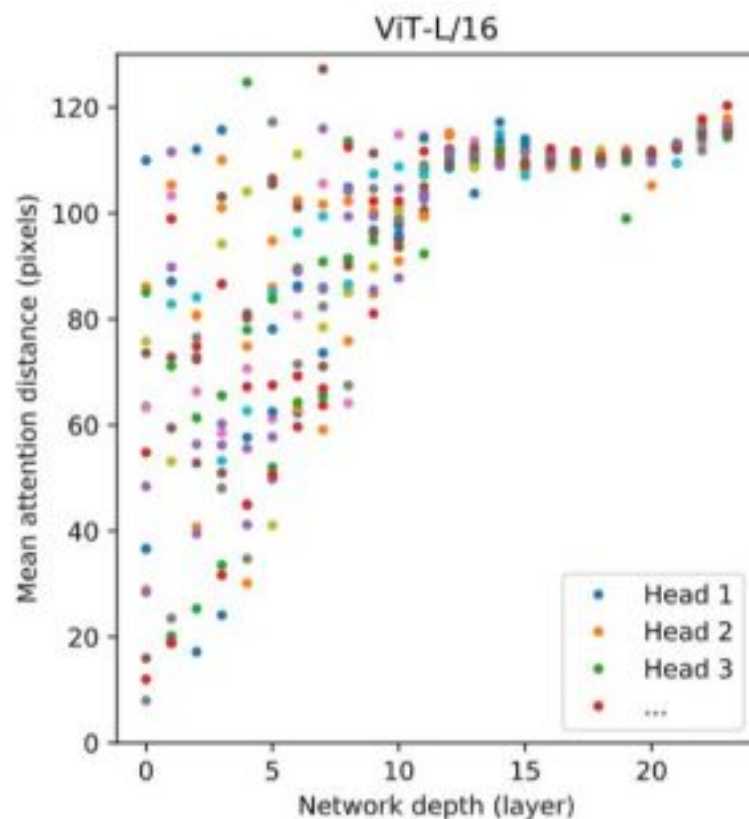
# The Vision Transformer (ViT)



# What does attention learn?

Spatially global features

Spatially local features



Inspection of learned attention heads

Input

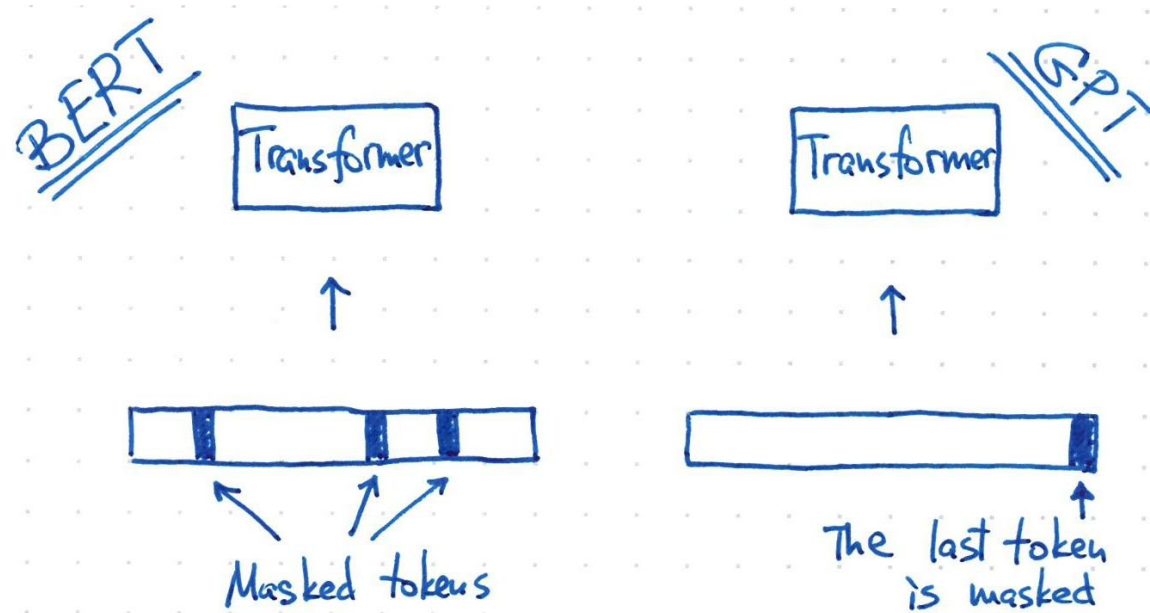
Attention





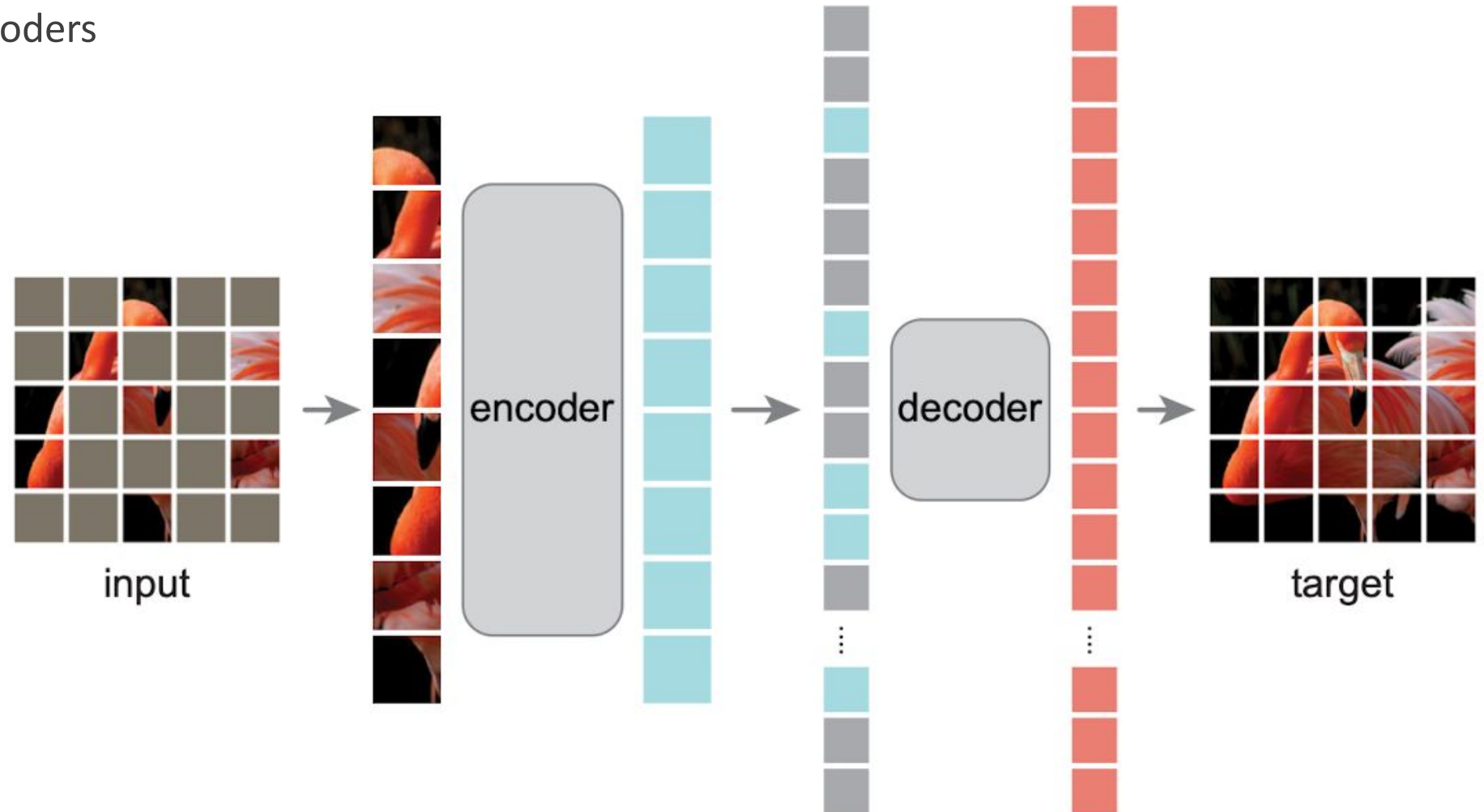
# Autoregressive training with ViT?

---



# Autoregressive training with ViT!

## Masked Auto Encoders



# Vision Foundation Models

---

- CLIP: Transformers that encode text and image to shared representation
  - Main ingredient of generative image methods: DALL-E, Stable Diffusion, ...
- DINO v2: Large vision transformer trained with different self-supervised losses
  - Powerful image features that enable many downstream tasks
- GPT4 Vision / Gemini Vision / ...
  - Presumably similar architecture to CLIP (no publication!)
- **Segment Anything Model**
  - Towards universal and promptable segmentation

# Segment Anything

---

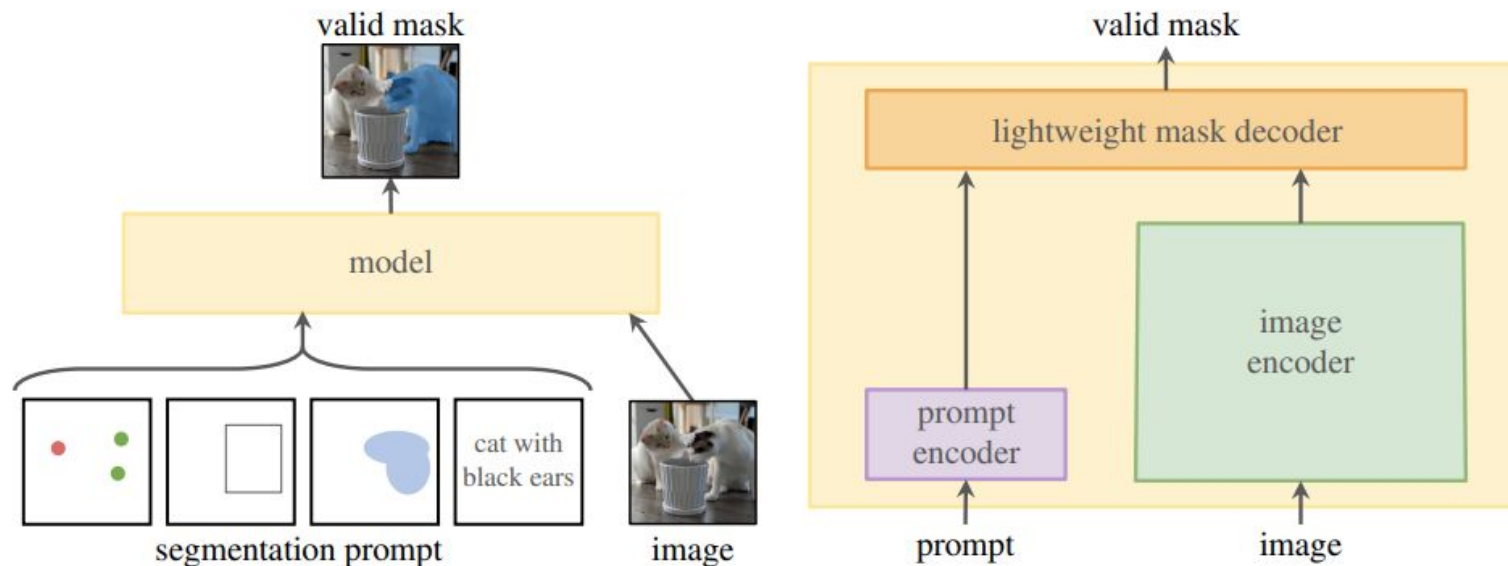
Vision foundation model for segmentation

# Segment Anything

<https://arxiv.org/abs/2304.02643>

Pretrained model for interactive segmentation from Meta.AI

SAM: Interactive segmentation



# Segment Anything

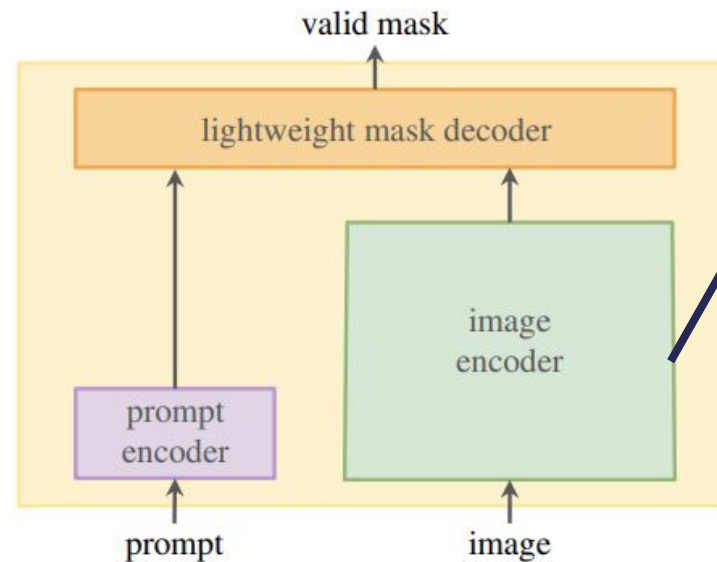
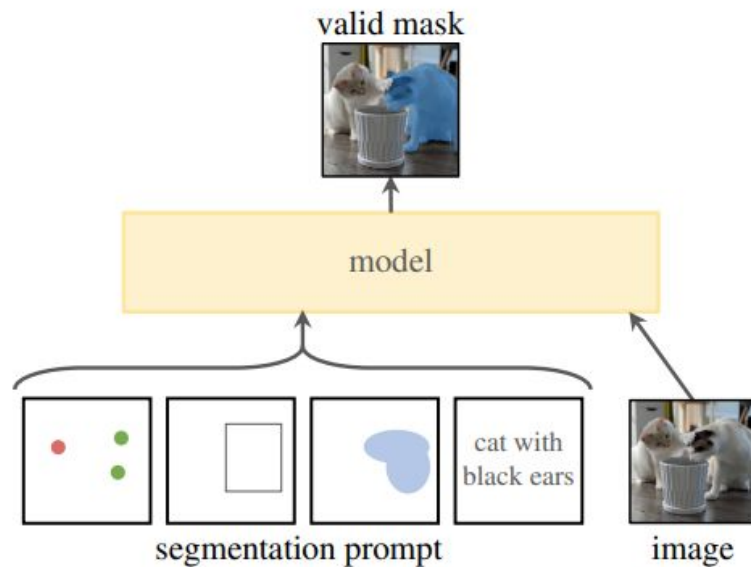
<https://arxiv.org/abs/2304.02643>

\* MobileSAM:

<https://arxiv.org/abs/2306.14289>

Pretrained model for interactive segmentation from Meta.AI

SAM: Interactive segmentation



4 different sizes:

- VIT-B (Base)
- VIT-L (Large)
- VIT-H (Huge)
- VIT-T (Tiny)\*

# Segment Anything: What's special?

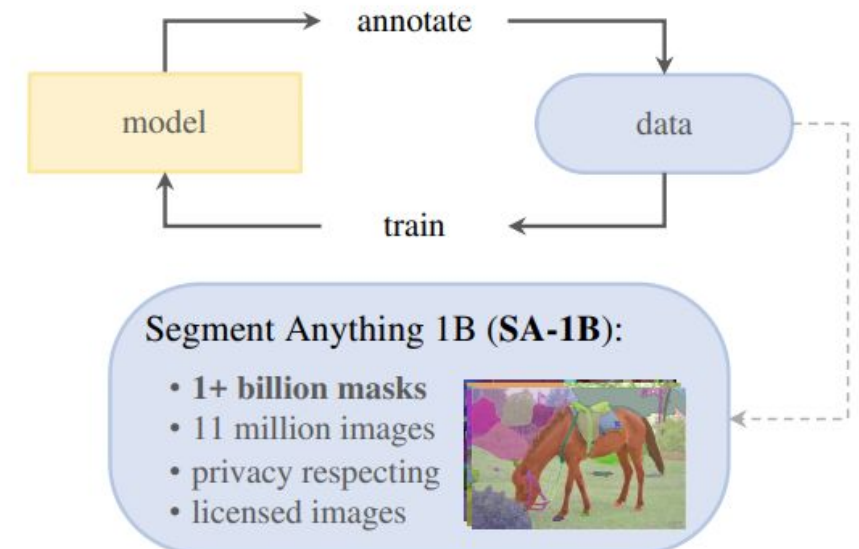
---

- Interactive segmentation: segment (almost) arbitrary objects from annotations
  - “prompts”: points and/or box and/or mask
  - more prompts improve the predictions
- Versatile: can be integrated within pipelines that provide prompts
  - From user inputs, object detectors, nucleus seeds, ...
  - Model is fully open-source!



# Segment Anything: What's special?

- Interactive segmentation: segment (almost) arbitrary objects from annotations
  - “prompts”: points and/or box and/or mask
  - more prompts improve the predictions
- Versatile: can be integrated within pipelines that provide prompts
  - From user inputs, object detectors, nucleus seeds, ...
  - Model is fully open-source!
- **How?**
  - **Large dataset with diverse images and objects**
  - Iterative training loop

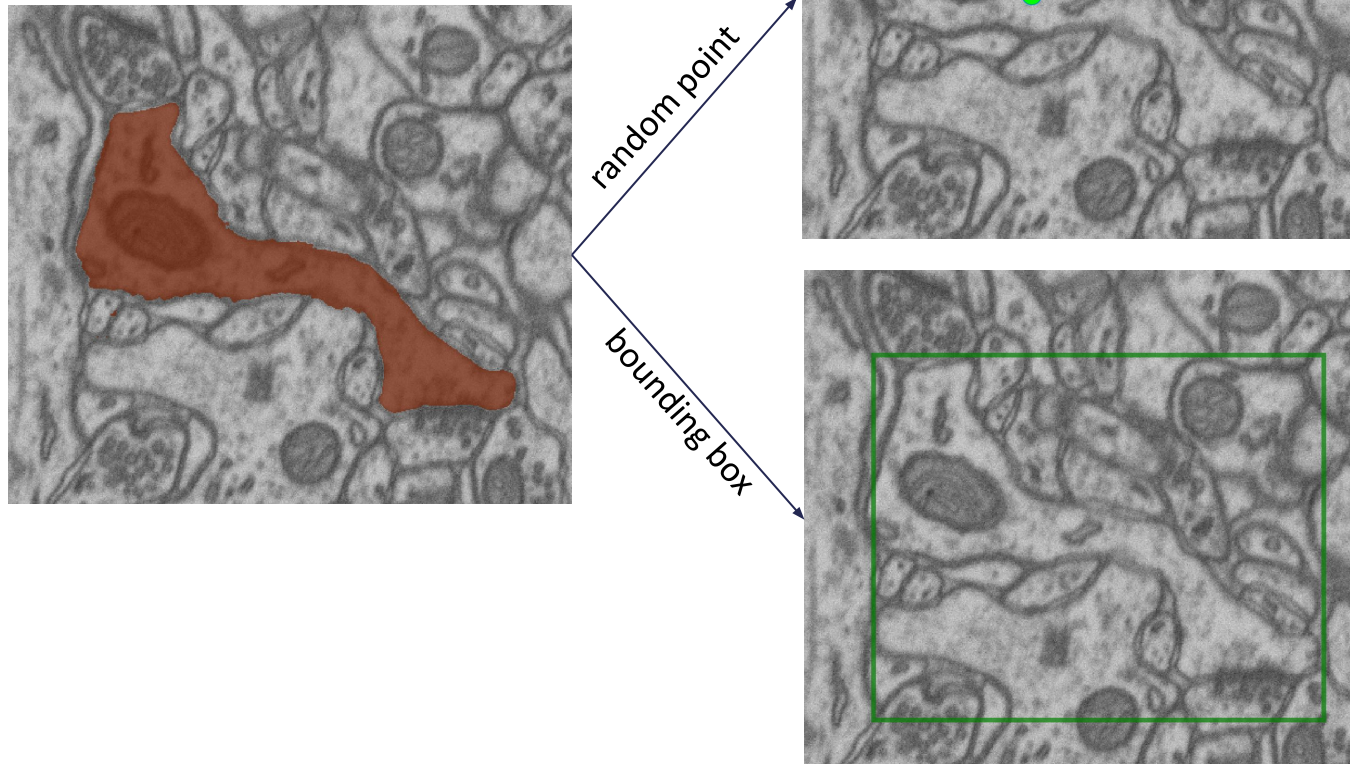


# Segment Anything: Training iteration

---

Given image and ground-truth mask

- Compute image embeddings, sample positive point or box

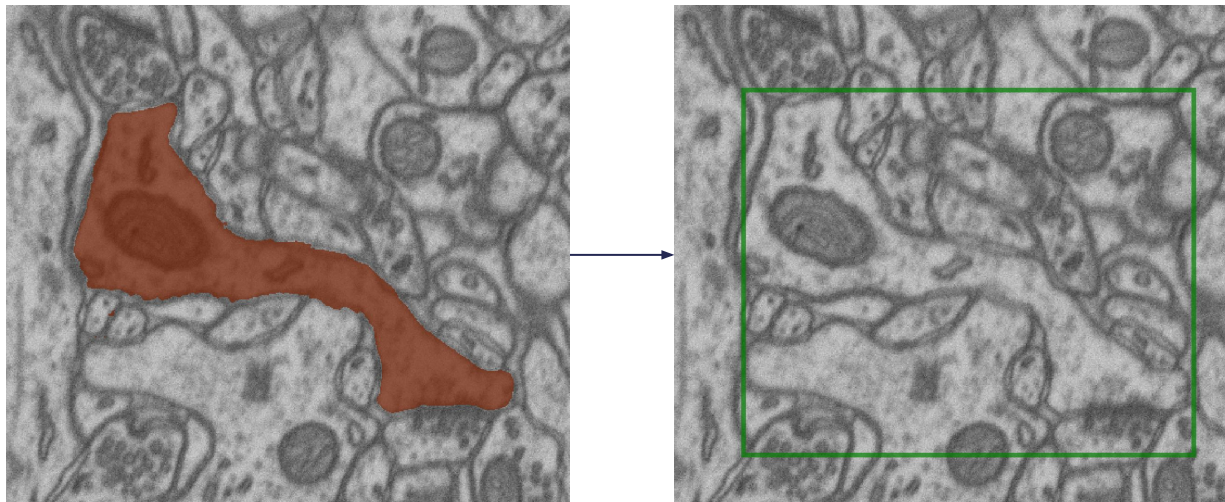


# Segment Anything: Training iteration

---

Given image and ground-truth mask

- Compute image embeddings, sample positive point or **box**

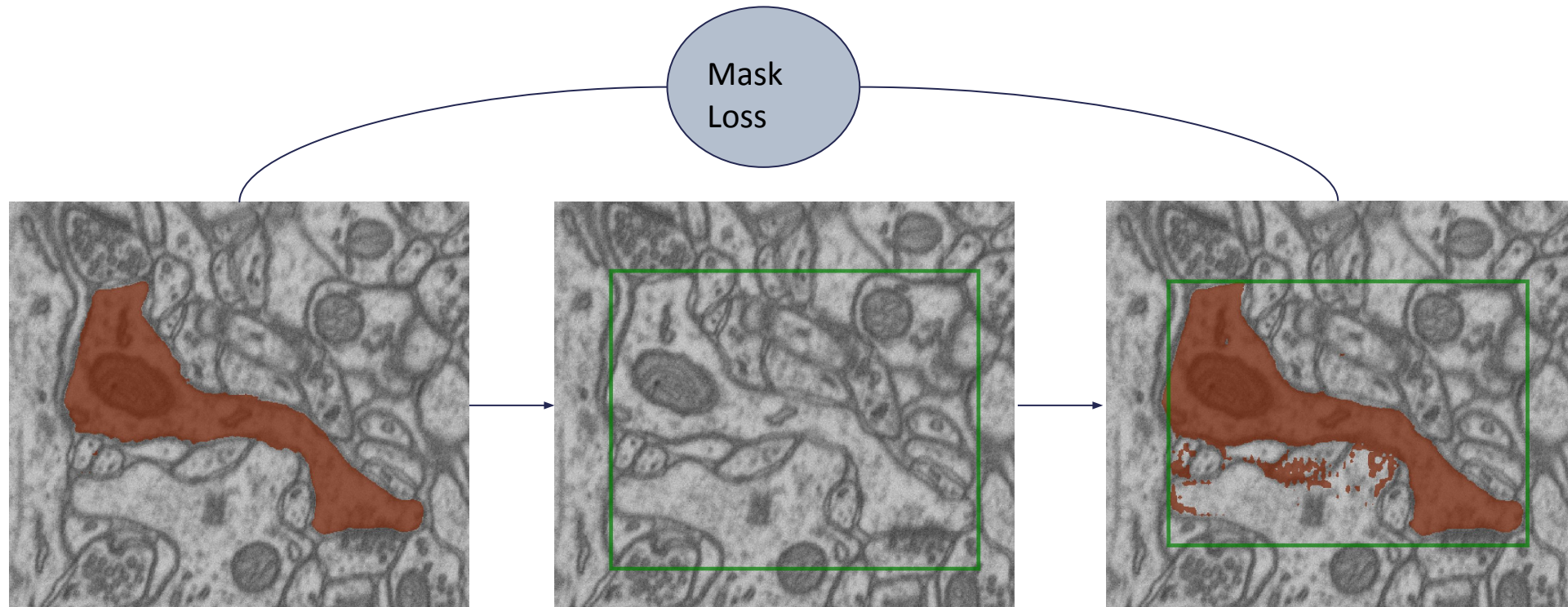


# Segment Anything: Training iteration

---

Given image and ground-truth mask

- Compute image embeddings, sample positive point or box
- Run prediction, compute loss for object and IOU estimate

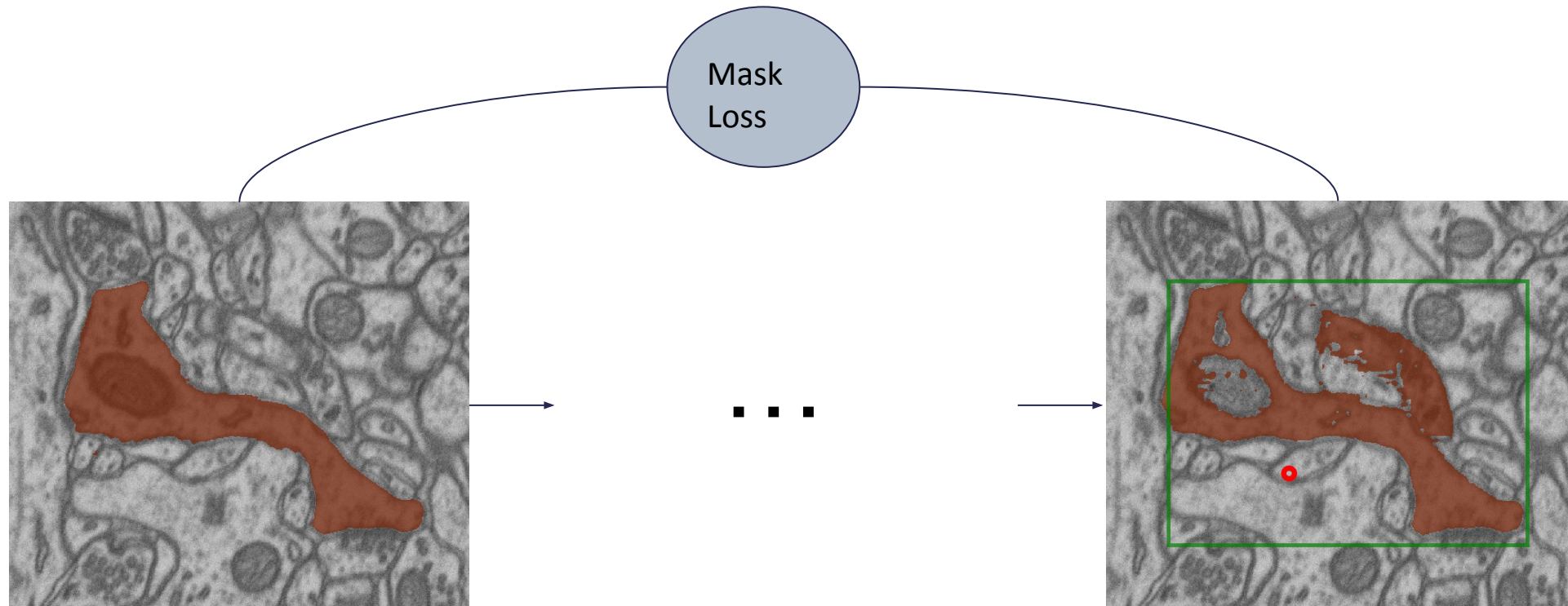




# Segment Anything: Training iteration

Given image and ground-truth mask

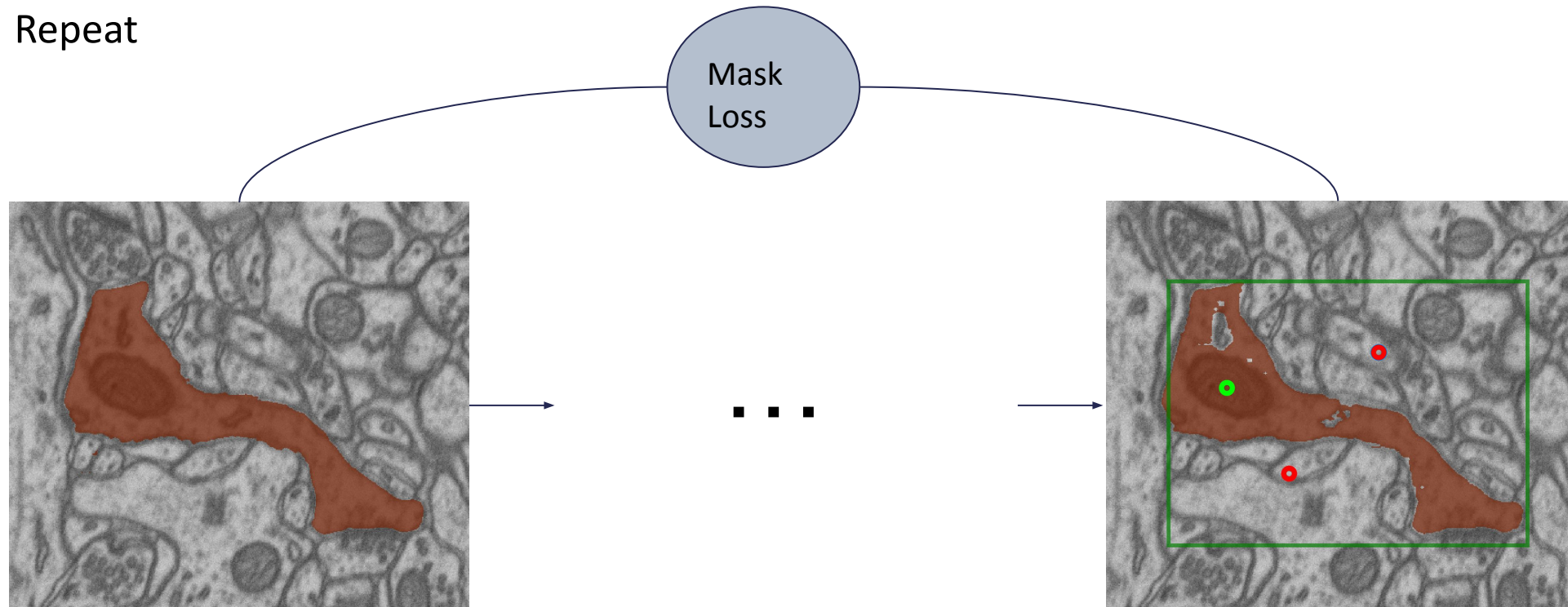
- Compute image embeddings, sample positive point or box
- Run prediction, compute loss for object and IOU estimate
- Sample point prompts where prediction is wrong, rerun prediction with all prompts + mask



# Segment Anything: Training iteration

Given image and ground-truth mask

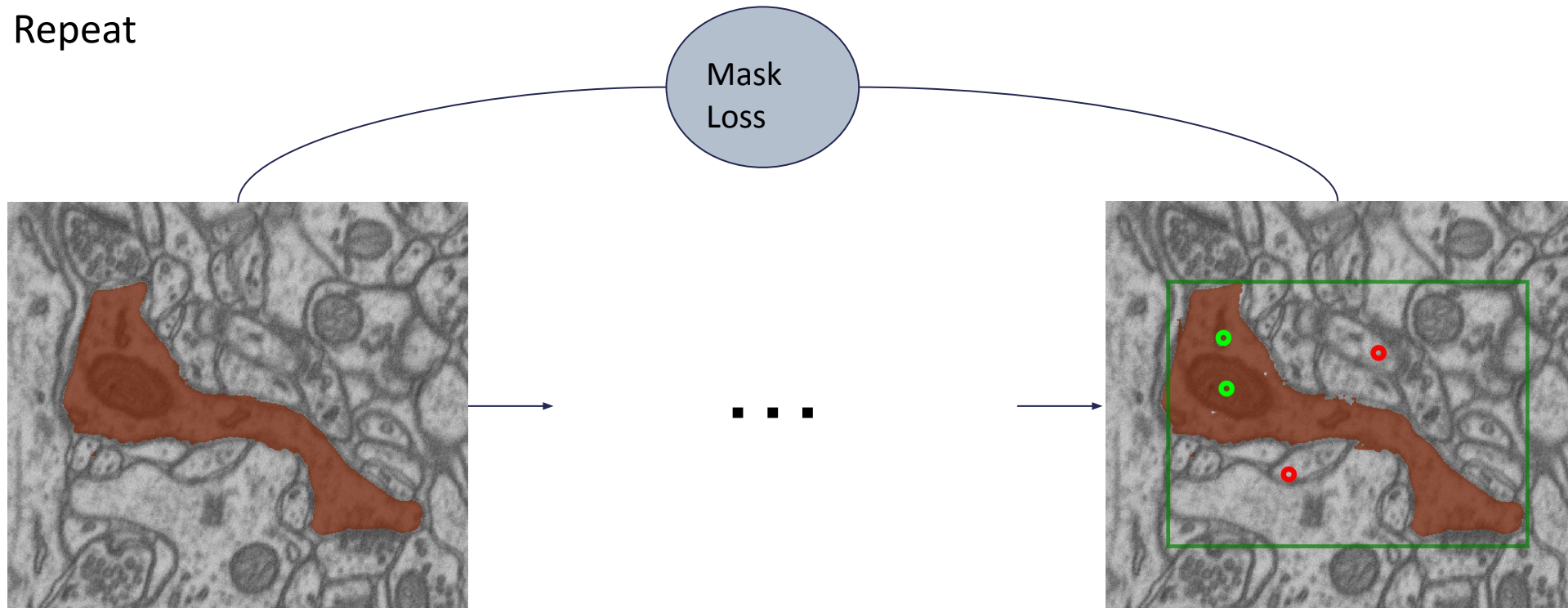
- Compute image embeddings, sample positive point or box
- Run prediction, compute loss for object and IOU estimate
- Sample point prompts where prediction is wrong, rerun prediction with all prompts + mask
- Repeat



# Segment Anything: Training iteration

Given image and ground-truth mask

- Compute image embeddings, sample positive point or box
- Run prediction, compute loss for object and IOU estimate
- Sample point prompts where prediction is wrong, rerun prediction with all prompts + mask
- Repeat





# Segment Anything: Training iteration

---

Given image and ground-truth mask

- Compute image embeddings, sample positive point or box
- Run prediction, compute loss for object and IOU estimate
- Sample point prompts where prediction is wrong, rerun prediction with all prompts + mask
- Repeat
- Average losses, update weights



# Segment Anything: Capabilities

<https://segment-anything.com/>

Segmentation from user inputs (prompts)





# Segment Anything: Capabilities

<https://segment-anything.com/>

Segmentation from user inputs (prompts)



Automatic Mask Generation (AMG)



# Segment Anything for Microscopy

---



Anwai  
Archit

# Our aims & contributions

Archit, ..., **Pape**, *bioRxiv* (2023)  
<https://doi.org/10.1101/2023.08.21.554208>

---

- How well does SAM work for microscopy data? Which model size is best?
- Can we improve it (by finetuning) on microscopy data?
- Build a napari-based tool for interactive and automatic segmentation and tracking.

Collaboration between my group and DFKI; + several open source contributions.



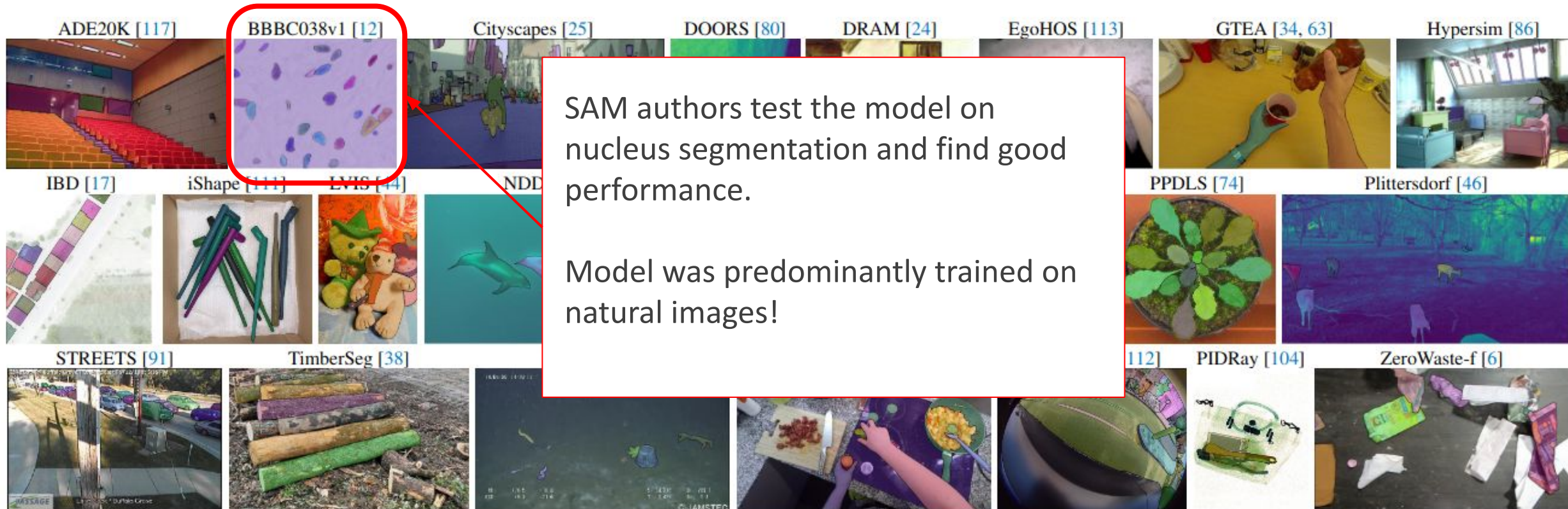


Anwai  
Archit

# Our aims & contributions

Archit, ..., Pape, *bioRxiv* (2023)  
<https://doi.org/10.1101/2023.08.21.554208>

- How well does SAM work for microscopy data? Which model size is best?





Anwai  
Archit

# Our aims & contributions

Archit, ..., **Pape**, *bioRxiv* (2023)  
<https://doi.org/10.1101/2023.08.21.554208>

- How well does SAM work for microscopy data? Which model size is best?
- Can we improve it by finetuning on microscopy data?
- Build a napari-based tool for interactive and automatic segmentation and tracking.

**We just released our stable v1.0 version!**

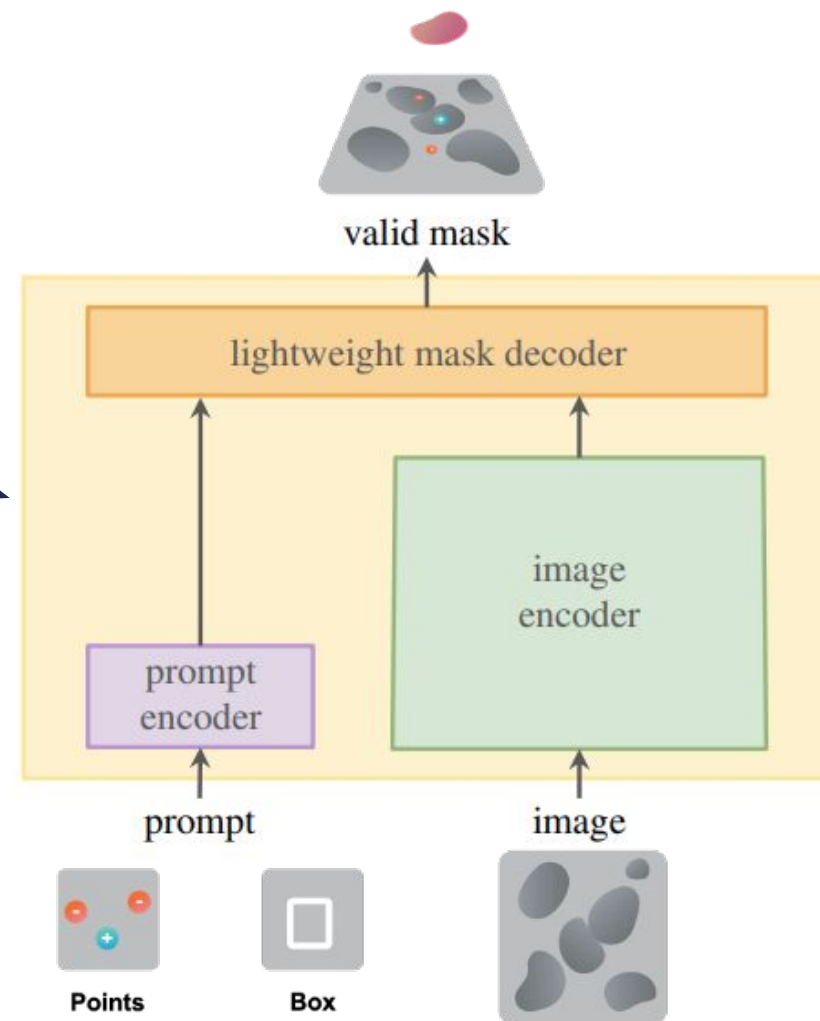
And submitted revision -> experiments are not in our preprint yet.



# Finetuning SAM

Our contributions:

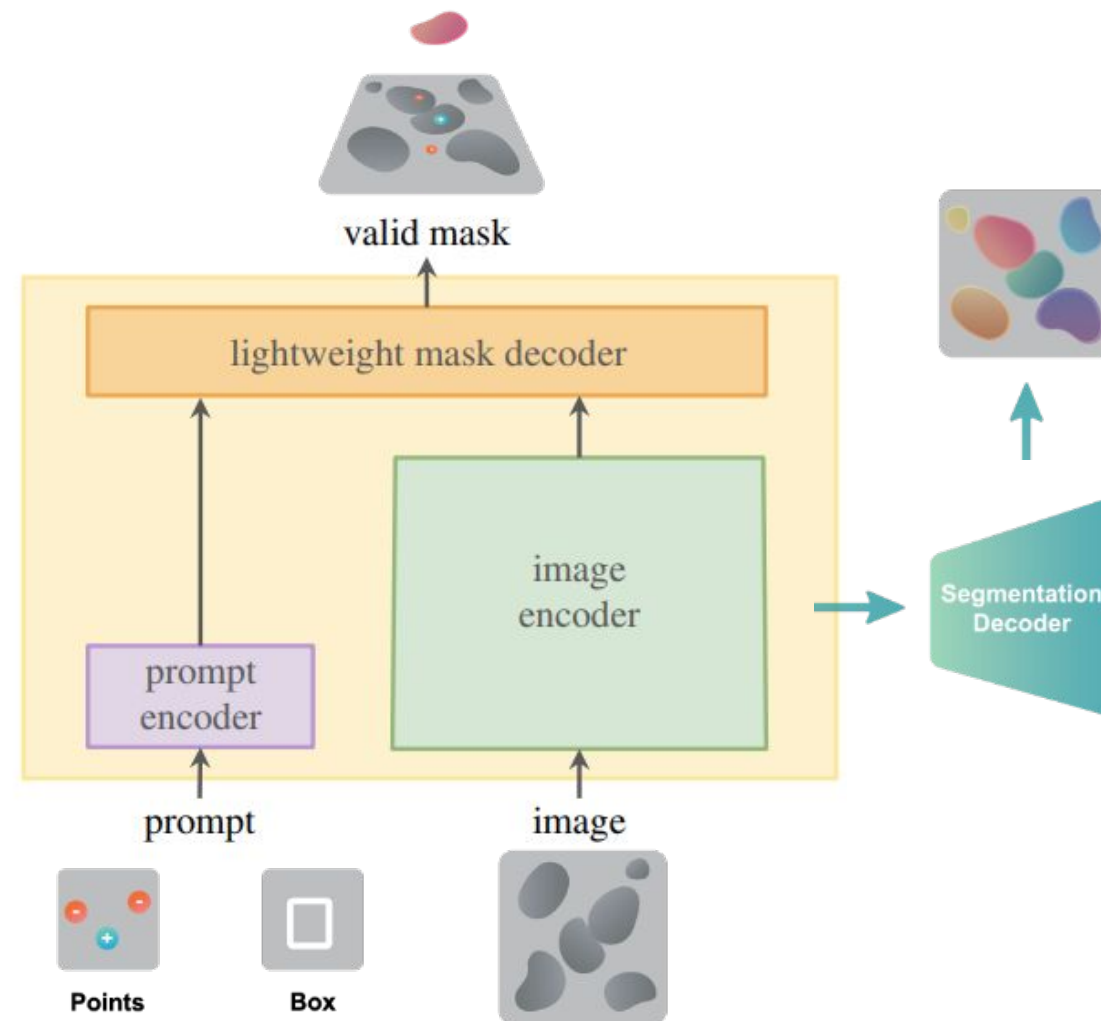
- Re-implement iterative training
  - Original code not published
  - Complex procedure
  - Use to finetune **SAM components**



# Finetuning SAM + improve instance seg

Our contributions:

- Re-implement iterative training
  - Original code not published
  - Complex procedure
  - Use to finetune SAM components
- Add decoder for instance segmentation (AIS)
  - Predicts foreground
  - Regresses distances to boundary + centroid
  - Input for watershed

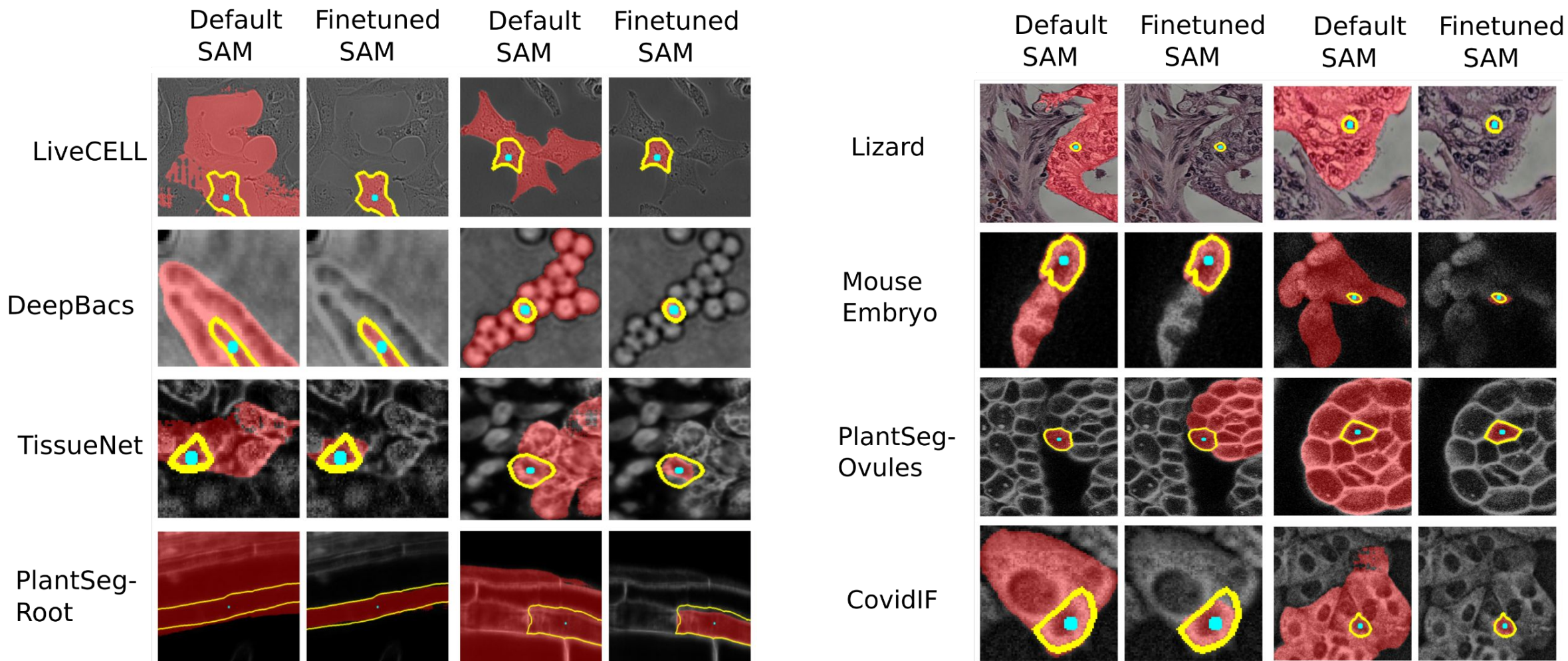


# Finetuning for light microscopy

---

- Training data: cell and nucleus segmentation (published datasets)
  - Cells in Phase-contrast (LiveCELL)
  - Cells in Tissue (TissueNet)
  - Cells and Nuclei in Fluorescence (Neurips Cell Seg, DSB)
  - Cells in LightSheet (PlantSeg-Roots)
  - Bacteria in labelfree imaging (DeepBacs)
- Evaluate on test-split of training datasets (“in domain”) and unseen datasets (“out of domain”):
  - Nuclei and cells in confocal, cells in immunofluorescence, nuclei in histopathology, ...
- Compare interactive and automatic instance segmentation
  - Compare to CellPose baseline for automatic segmentation

# Interactive Segmentation: In domain & Out-of-domain



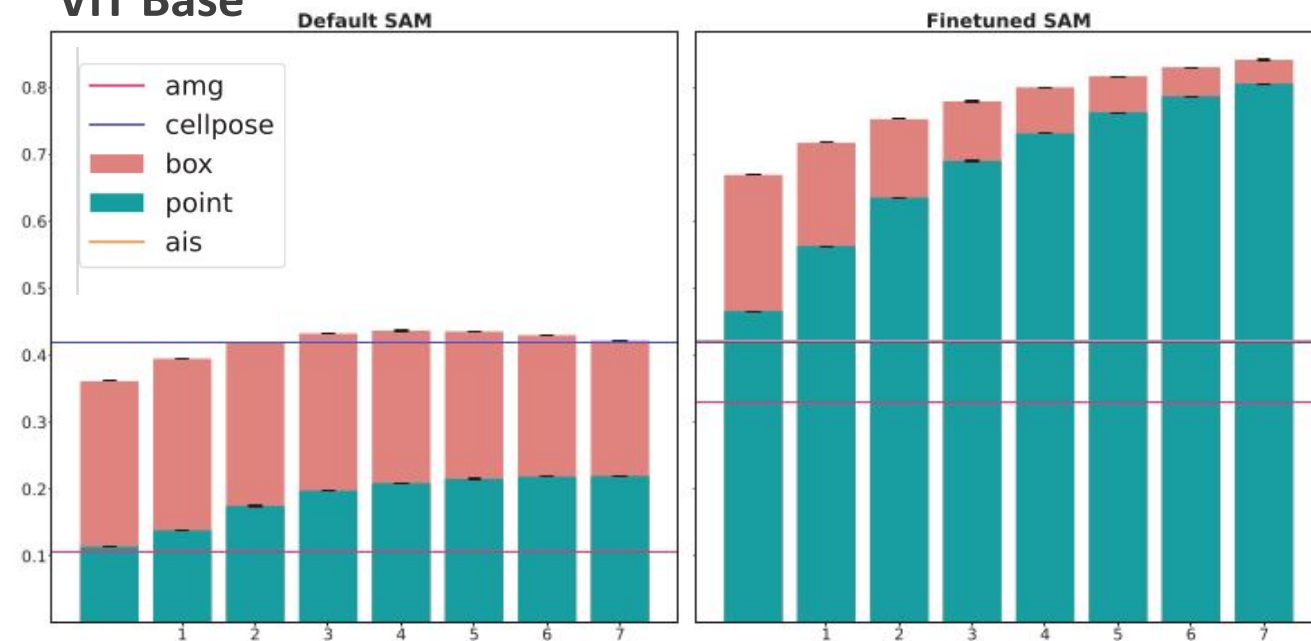
# Results: In Domain

Results for LIVECell Dataset  
(In Domain; Test Split)

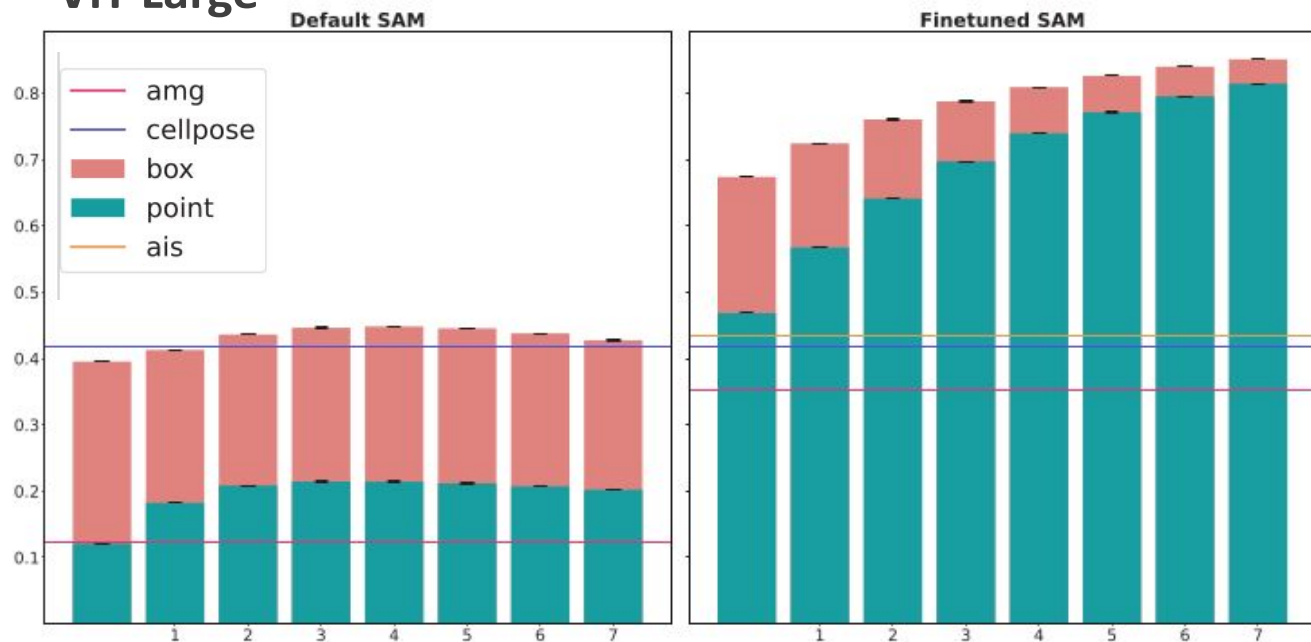
Evaluation:

- Interactive Segmentation:
  - Derive prompts from ground-truth, improve iteratively
- Instance segmentation:
  - Compare with CellPose
- Both: compute segmentation accuracy (compared to ground-truth)

## ViT Base



## ViT Large



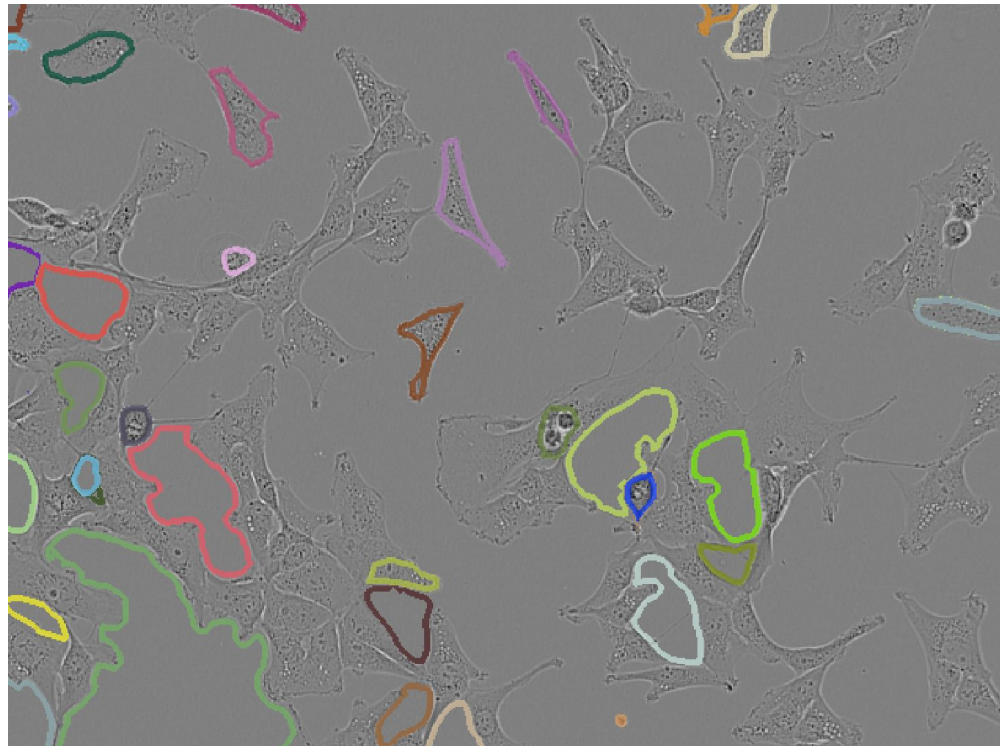


# Automatic Segmentation

---

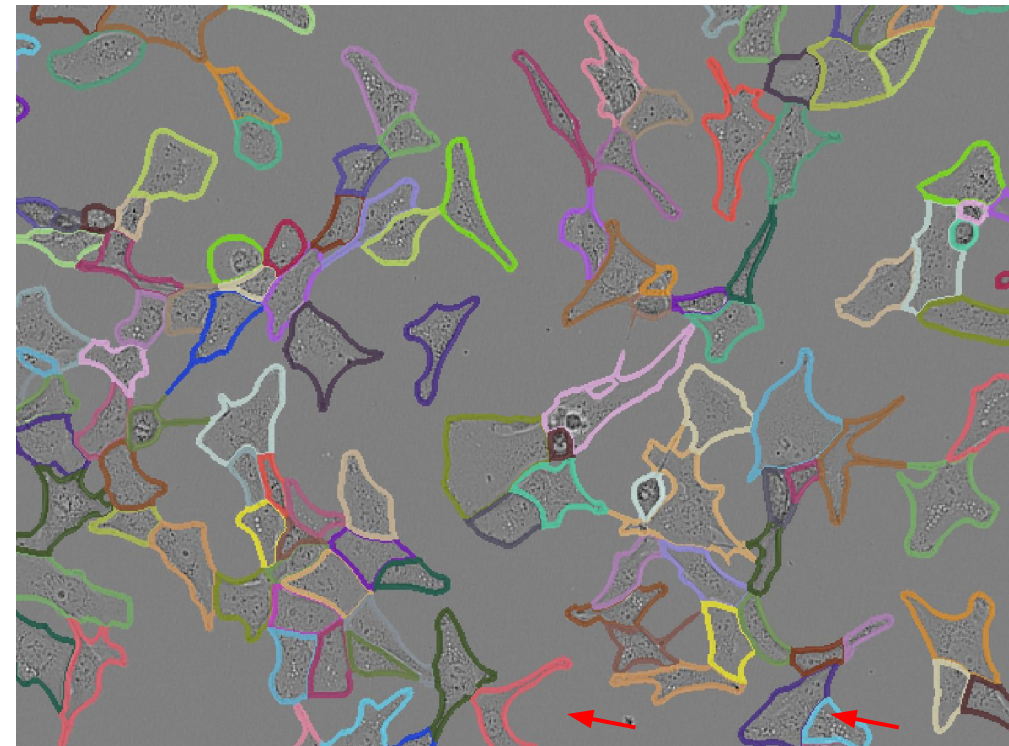
Instance segmentation on LIVECell Dataset

Runtimes on laptop (CPU);  
including embedding computation (dominates for AIS)



VIT-B  
AMG: 75 sec

VIT-B-LM  
AIS: 9 sec



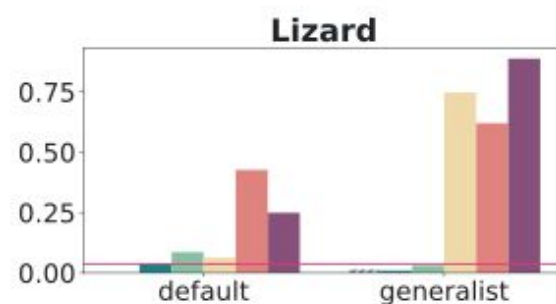
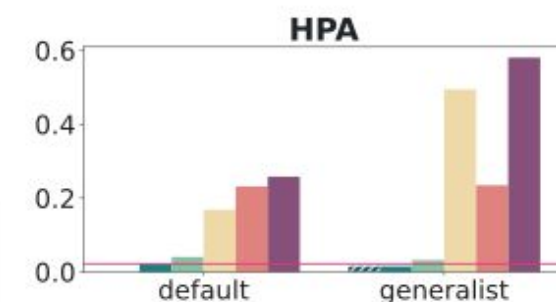
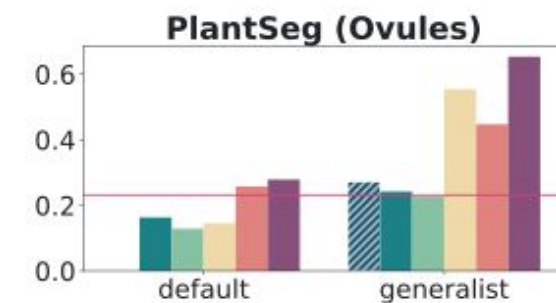
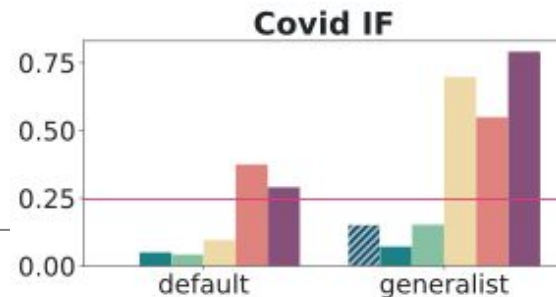
# Results: Out of domain

Results for out of-domain datasets.

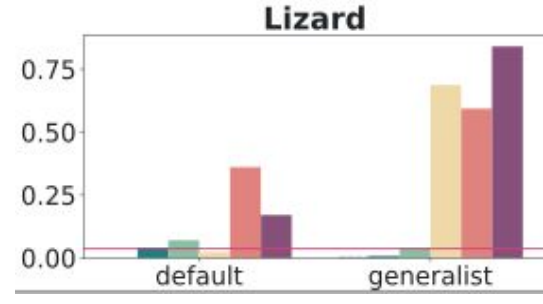
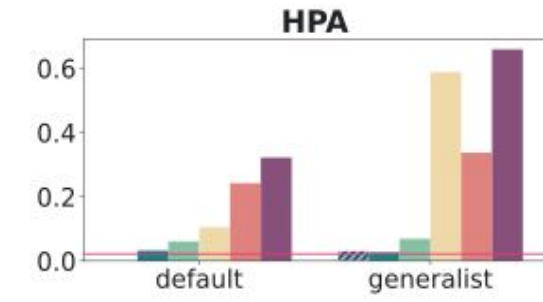
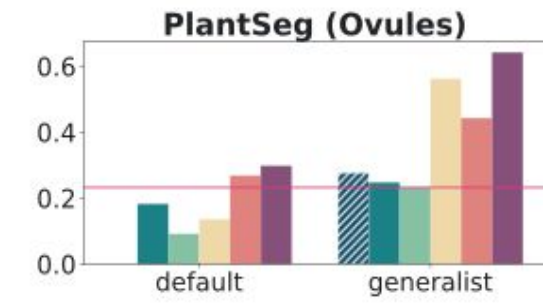
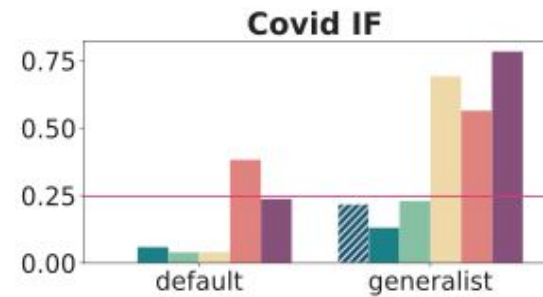
Same evaluation procedure as before.



ViT Base



ViT Large





# Results: Out of domain

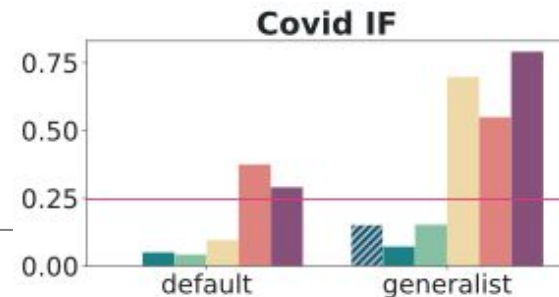
Results for out of-domain datasets.

Same evaluation procedure as before.

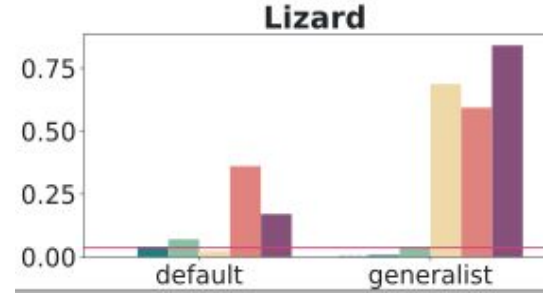
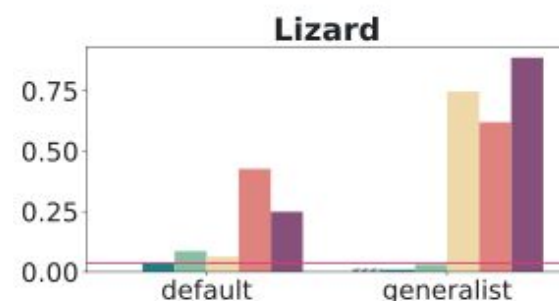
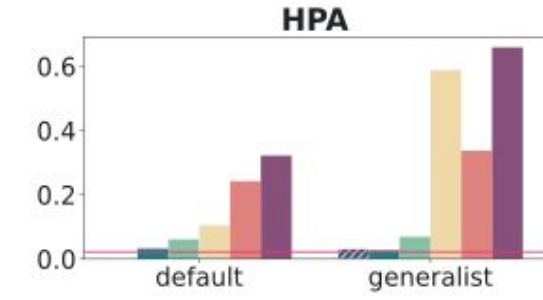
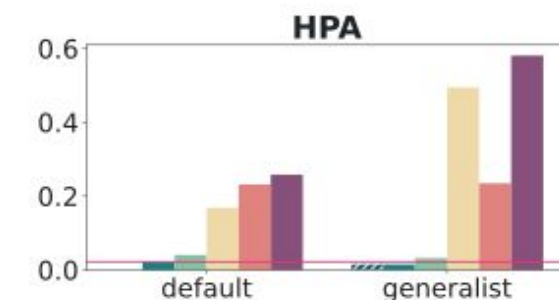
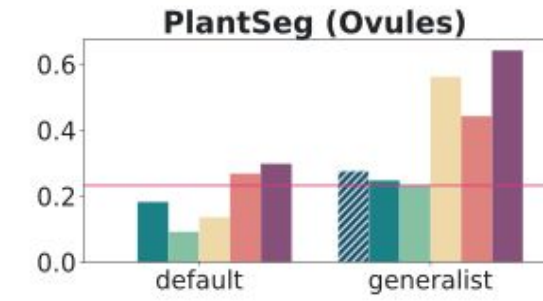
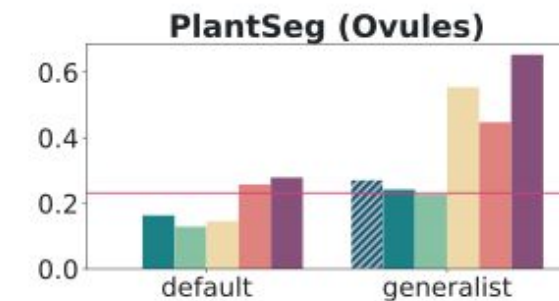
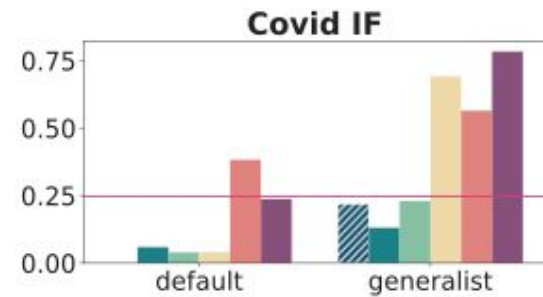
## Conclusions:

- Finetuning improves models!
- Best model: vit\_l
  - If runtime matters: vit\_b / vit\_t
- Comparison to CellPose (automatic seg.):
  - Similar performance on most out of domain datasets (cyto2 model)

ViT Base



ViT Large



# Finetuning for electron microscopy

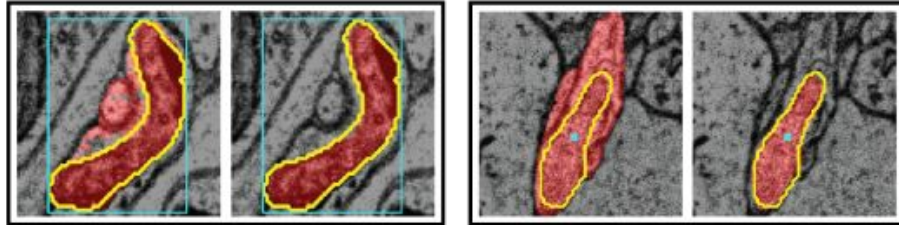
---

- Training data: Mitochondria and nucleus segmentation in electron microscopy
  - Most training data from MitoNet (<https://doi.org/10.1016/j.cels.2022.12.006>).
- Compare default and finetuned model.
  - Compare automated segmentation with MitoNet.
- Evaluate on test-split of training datasets (“in domain”) and unseen datasets (“out of domain”)
  - Application to EM mitochondria from non-training data.

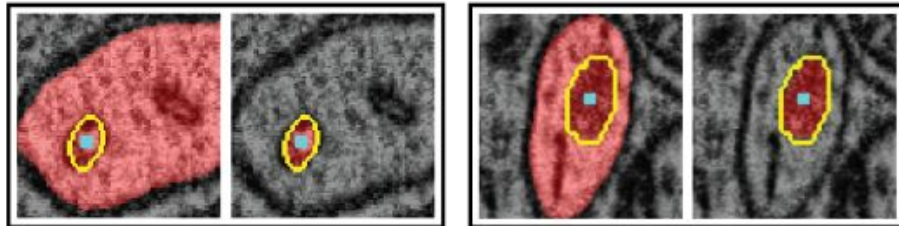
# Interactive Segmentation: In domain & Out-of-domain

Default SAM    Finetuned SAM    Default SAM    Finetuned SAM

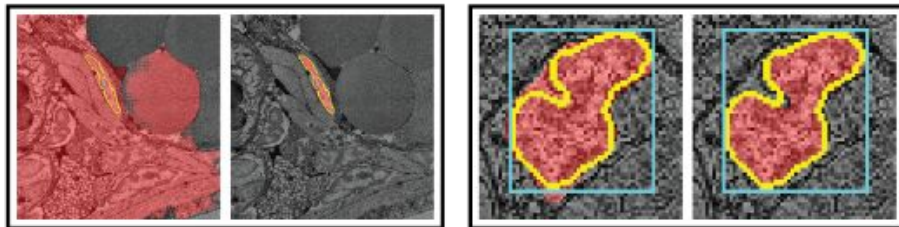
*MitoEM  
(Human)*



*MitoEM  
(Rat)*

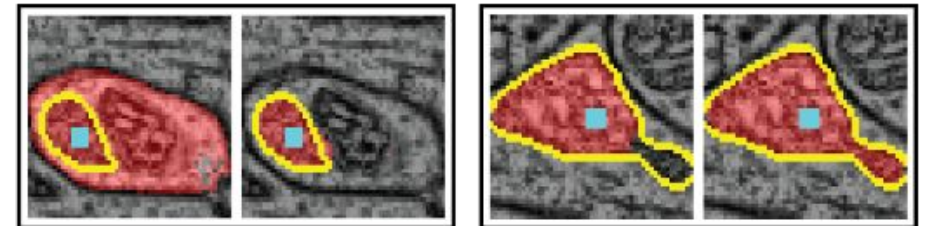


*Platynereis  
(Nuclei)*

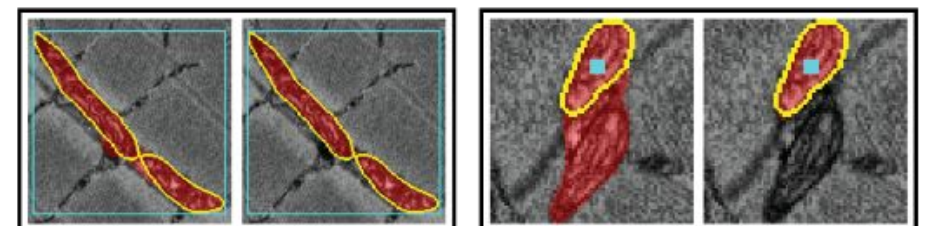


Default SAM    Finetuned SAM    Default SAM    Finetuned SAM

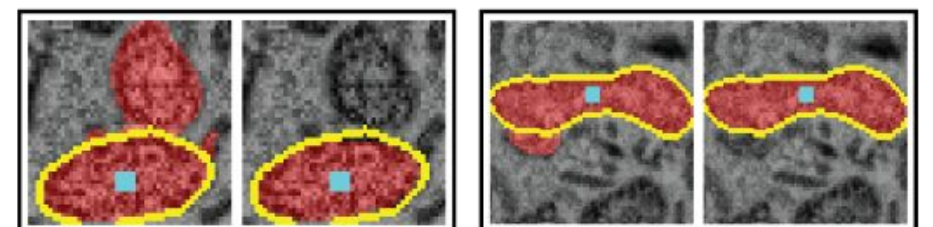
*MitoLab  
(Fly Brain)*



*MitoLab  
(Glycotic  
Muscle)*



*MitoLab  
(HeLa Cell)*

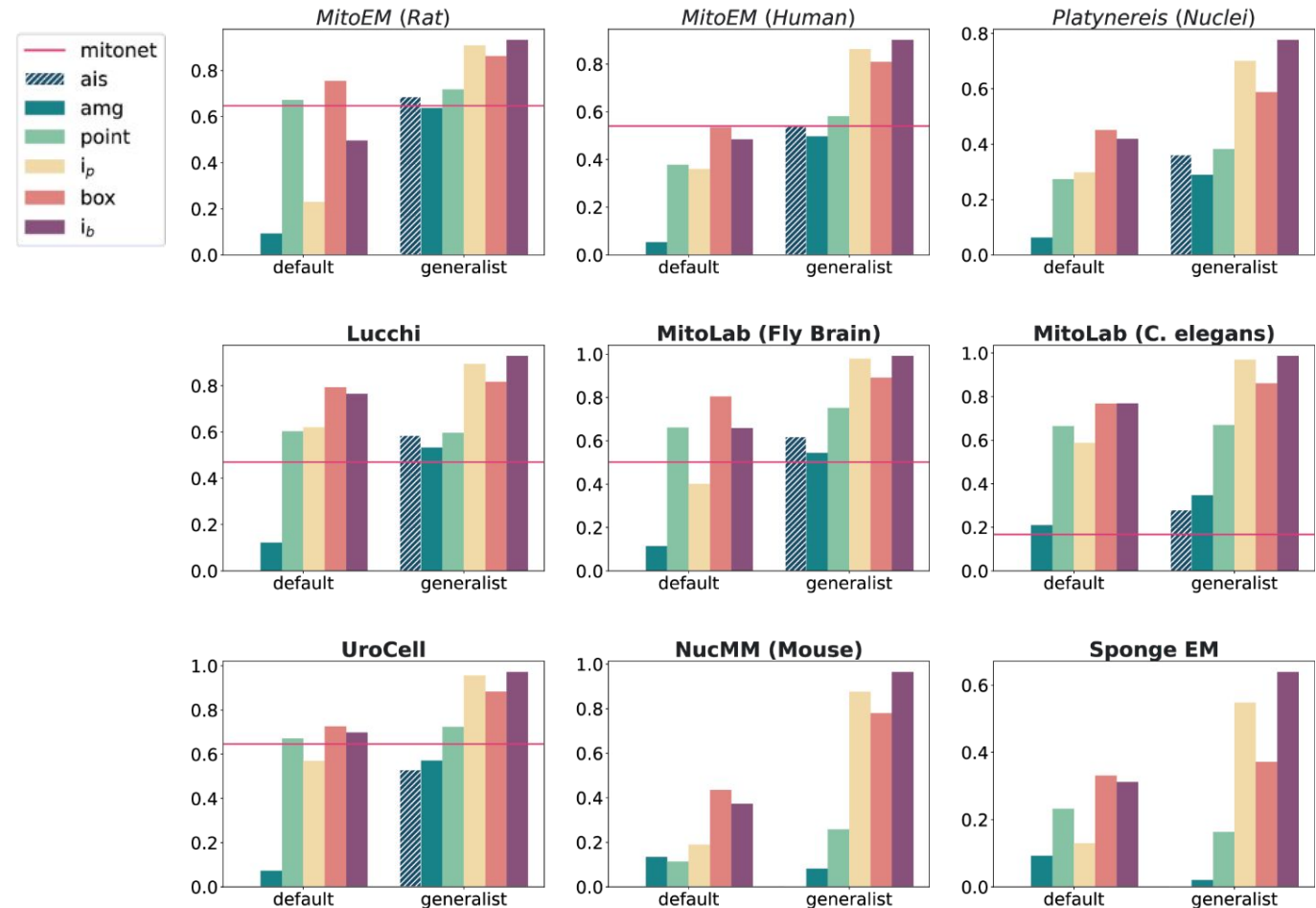


# Results: In & out-of domain

ViT Large

Evaluation: Same approach as for LM

- In domain (top row)
- Out of domain (rest)



# Results: In & out-of domain

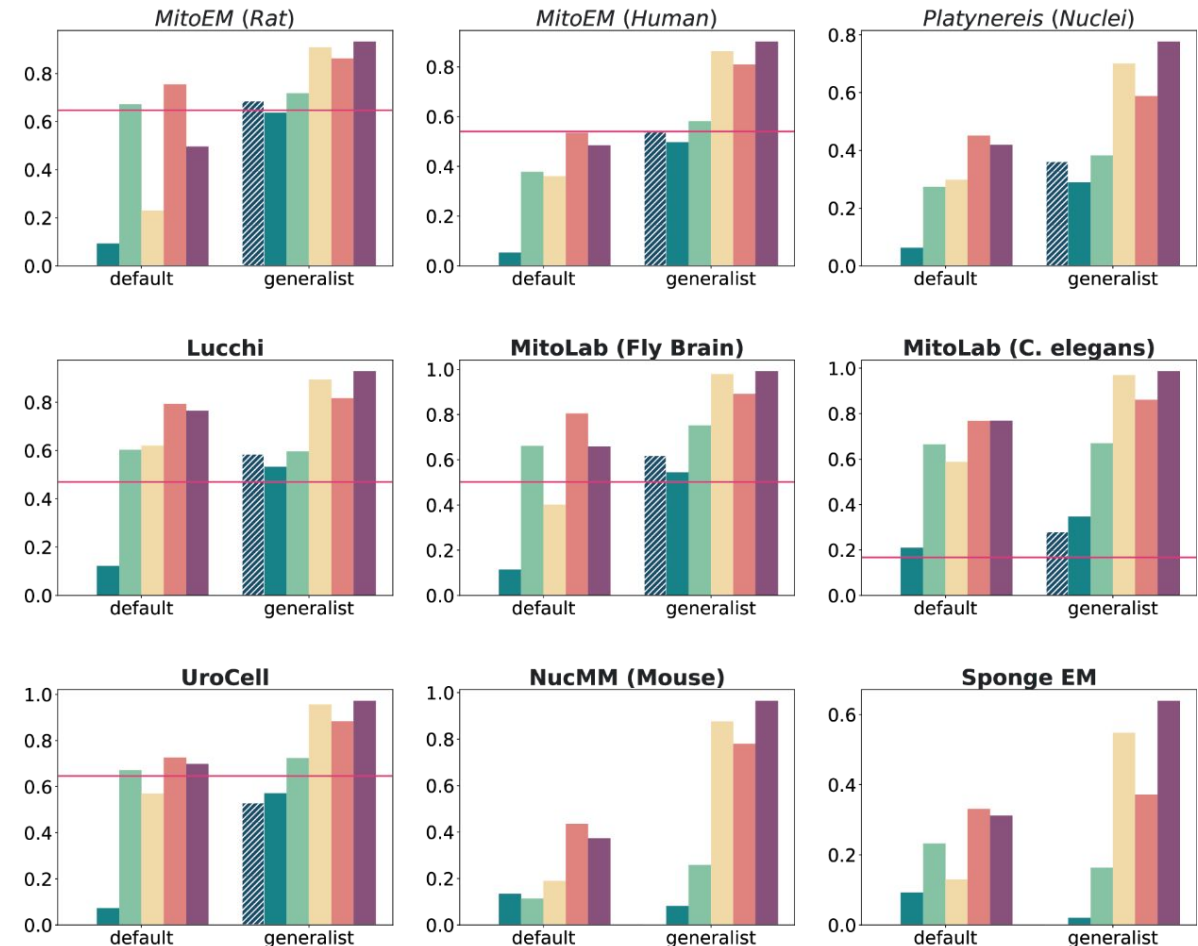
ViT Large

Evaluation: Same approach as for LM

- In domain (top row)
- Out of domain (rest)

## Conclusions:

- Finetuning improves, best model is vit\_l
- Similar performance to MitoNet on most datasets (AIS/AMG)
- Improves segmentation for some other organelles (cilia, microvilli), but worsens it for cellular compartments
  - Bigger diversity in EM!



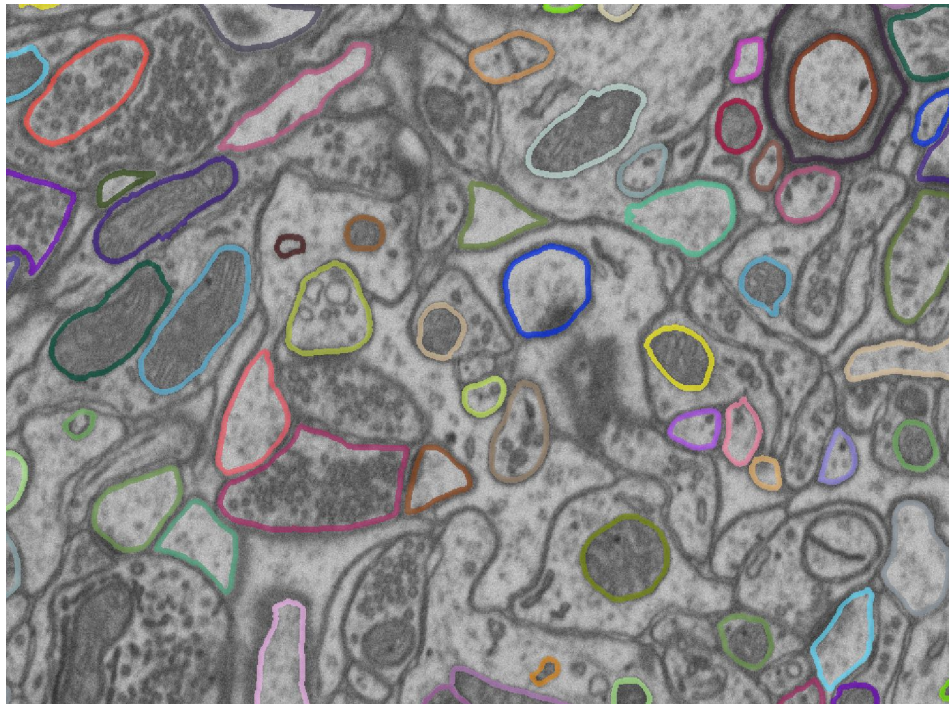


# Mitochondria

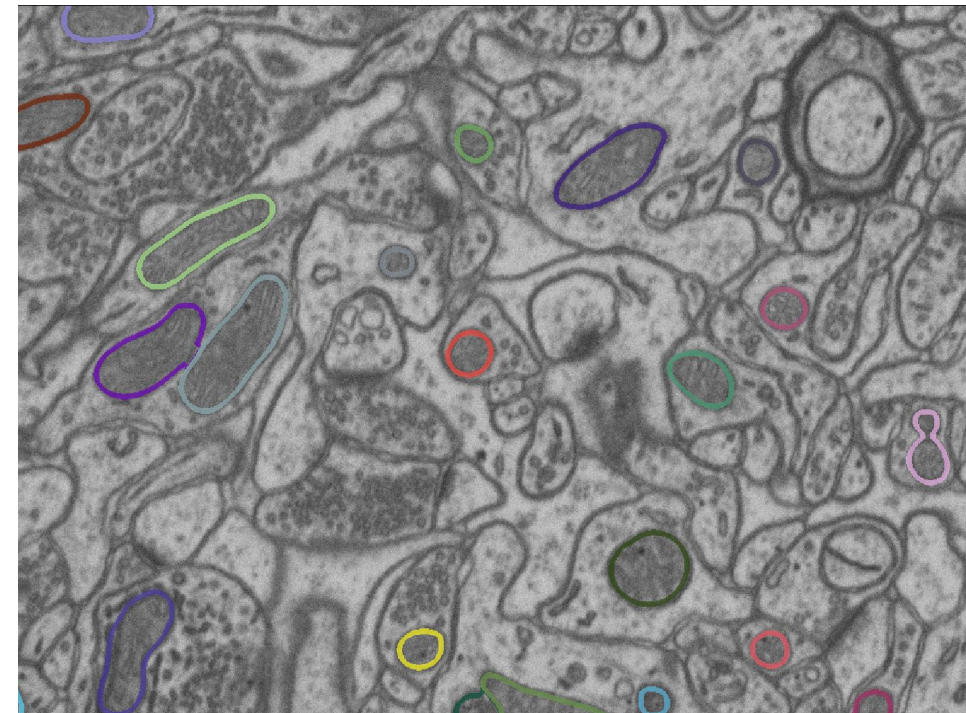
---

Instance segmentation on Lucchi Dataset

Runtimes on laptop (CPU);  
including embedding computation (dominates for AIS)



VIT-B-EM  
AIS: 10 sec

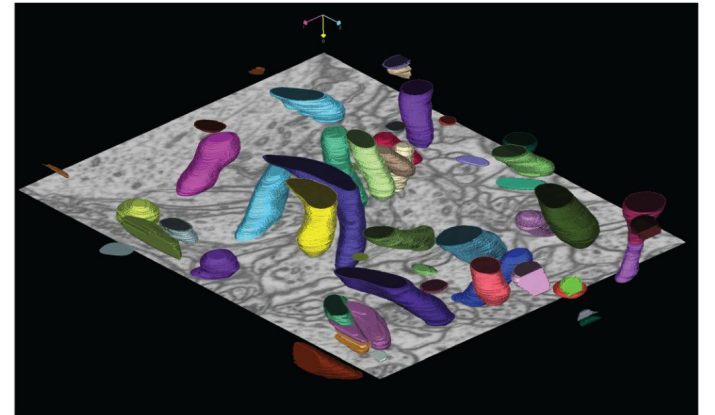
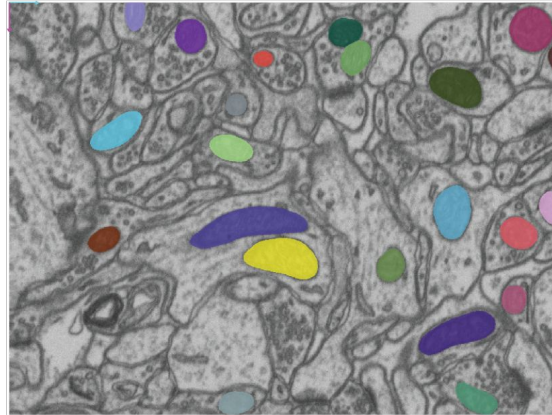


VIT-B  
AMG: 80 sec

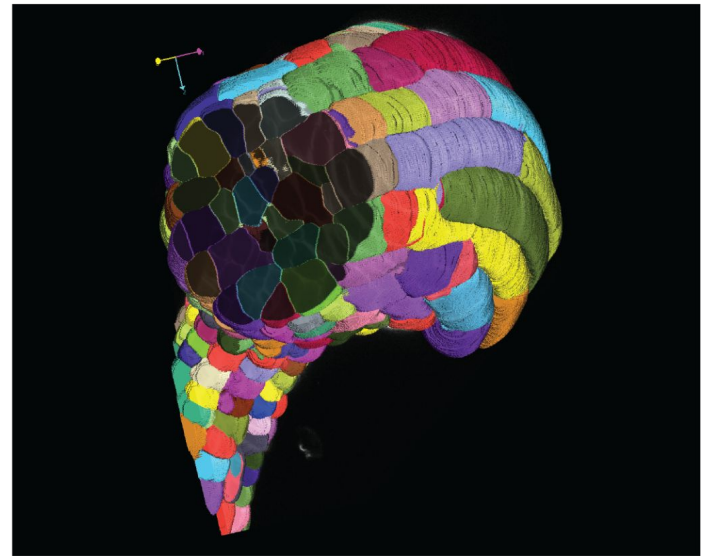
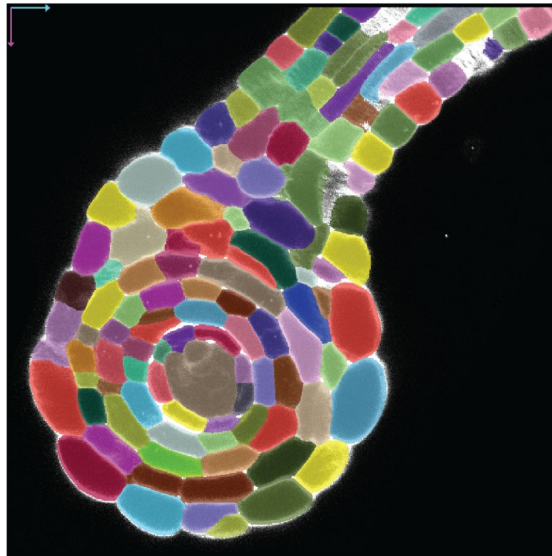
# 3D Segmentation with AIS

Segment objects slice by slice and  
merge across 3D

Lucchi



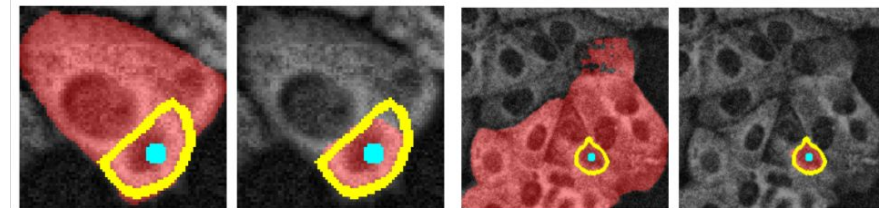
PlantSeg (Ovules)





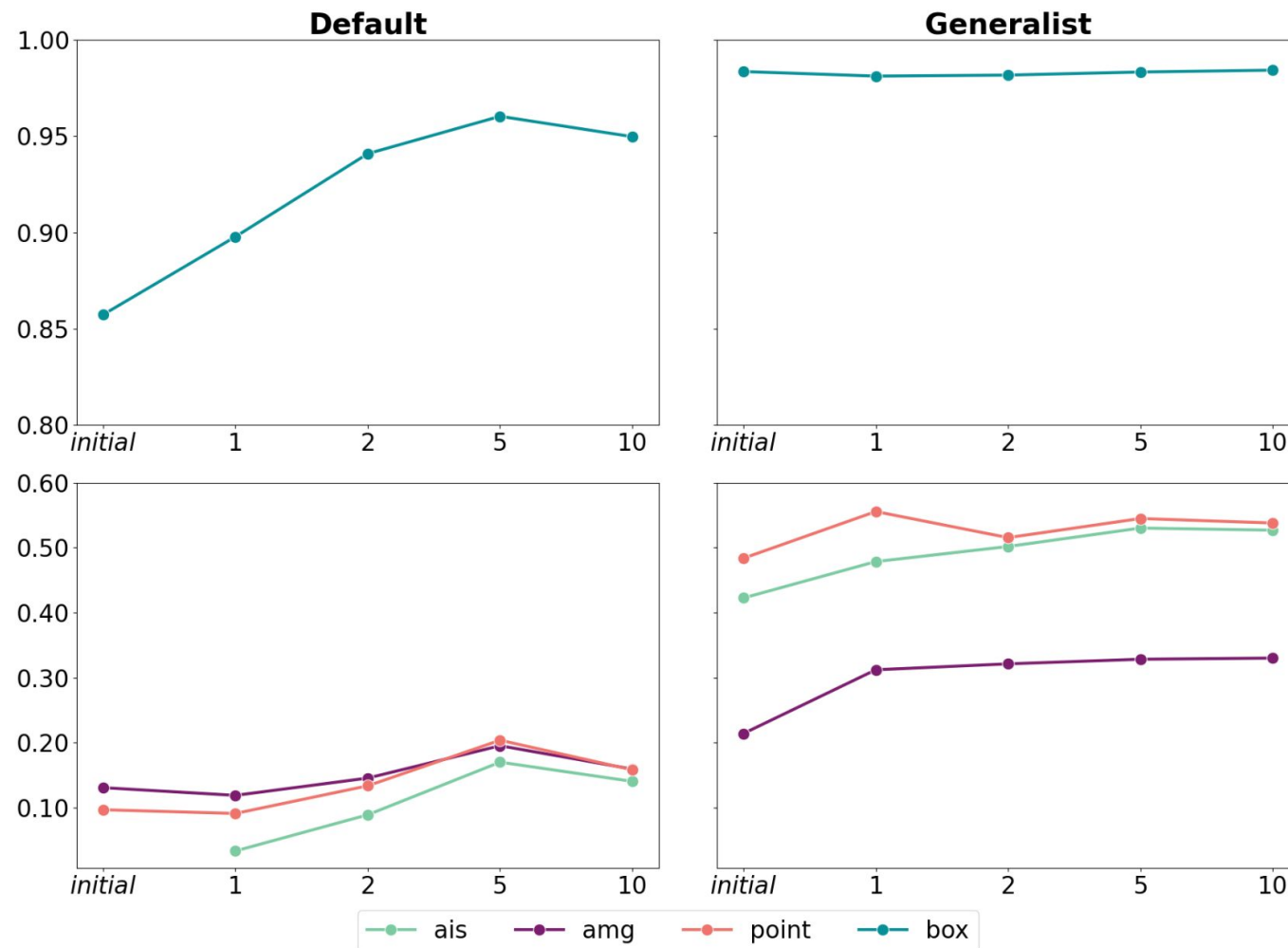
# Finetuning as a user

CovidIF



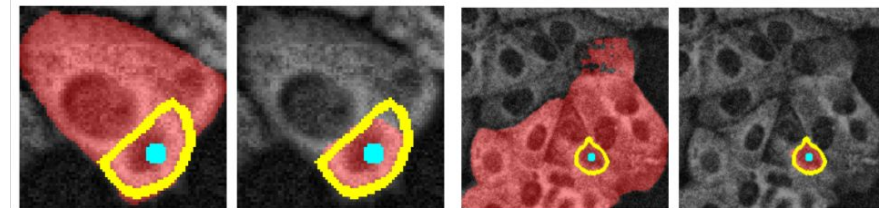
Improve models further for your data?

- How much data is needed?
- Which computational resources are required?



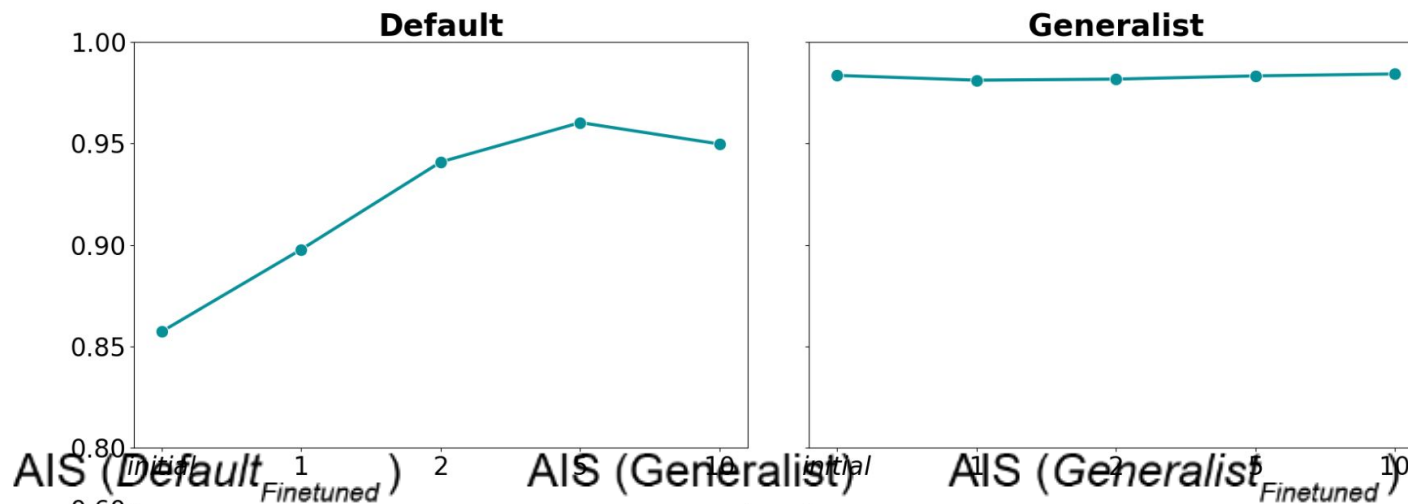
# Finetuning as a user

CovidIF



Improve models further for your data?

- How much data is needed?
- Which computational resources are required?



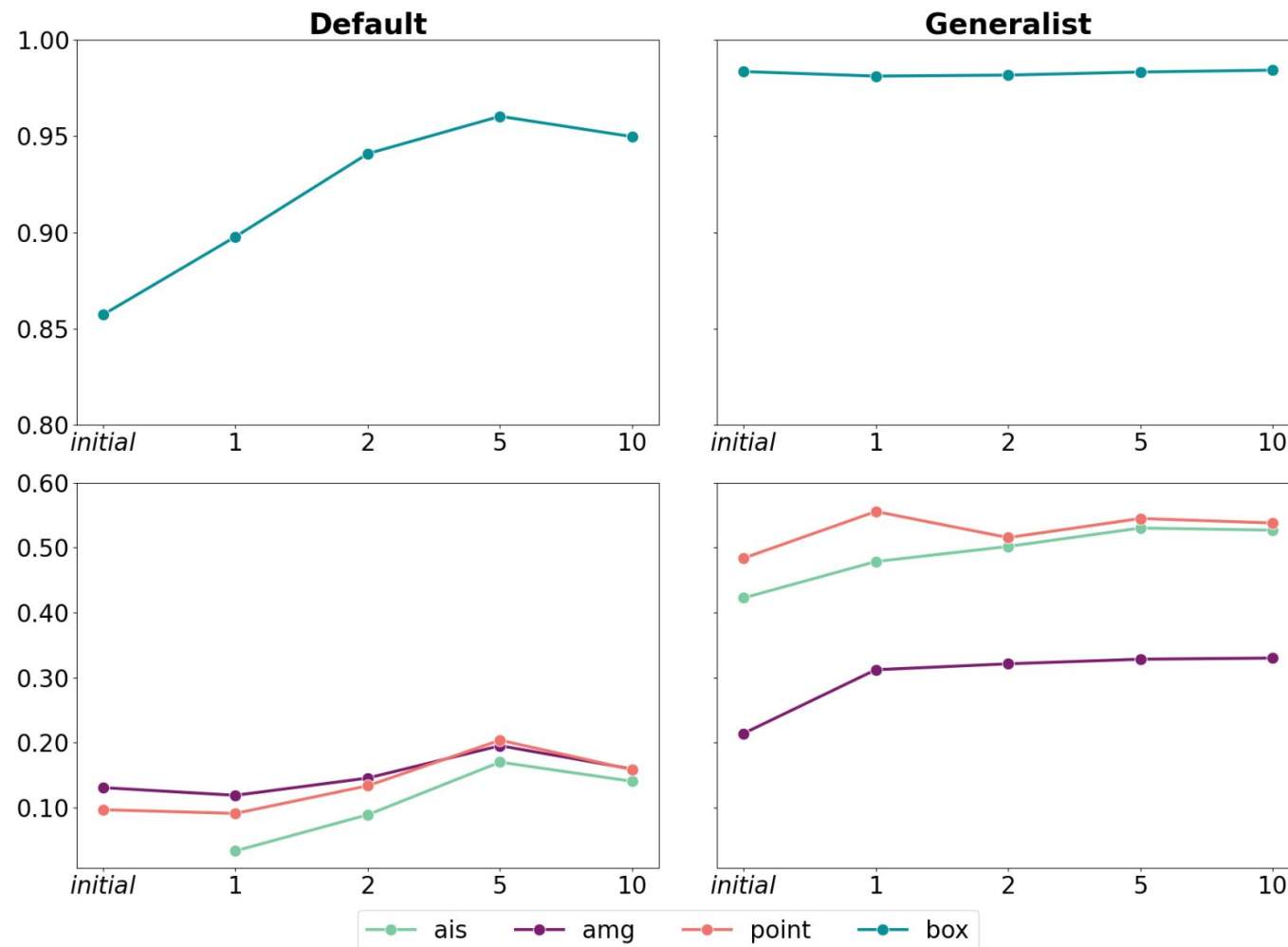
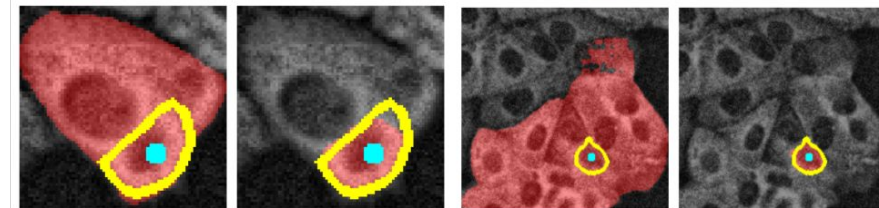
# Finetuning as a user

Improve models further for your data?

- How much data is needed?
- Which computational resources are required?

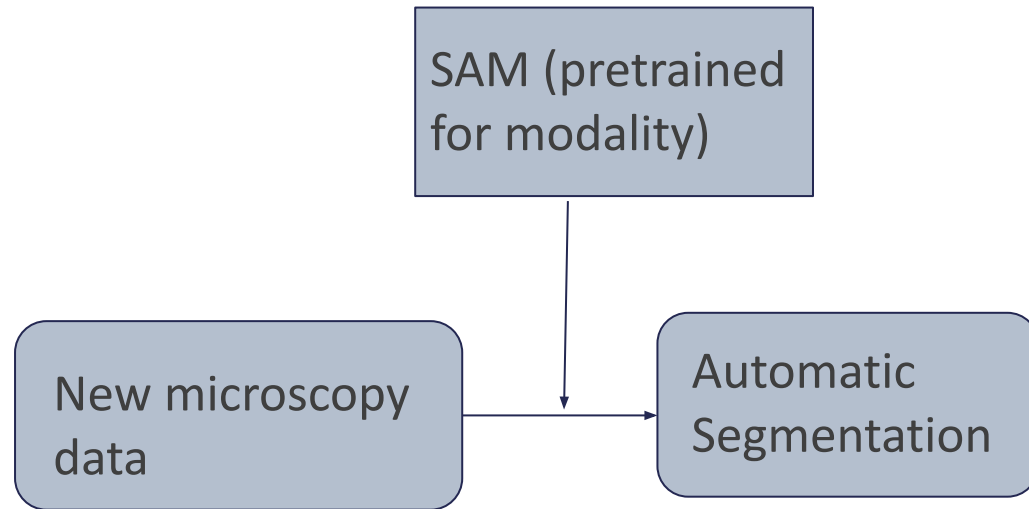
- Few images with annotations are sufficient!
- Possible on CPU
  - here: ca. 6 hours
- Faster on GPU
  - here: ca. 30 minutes

CovidIF



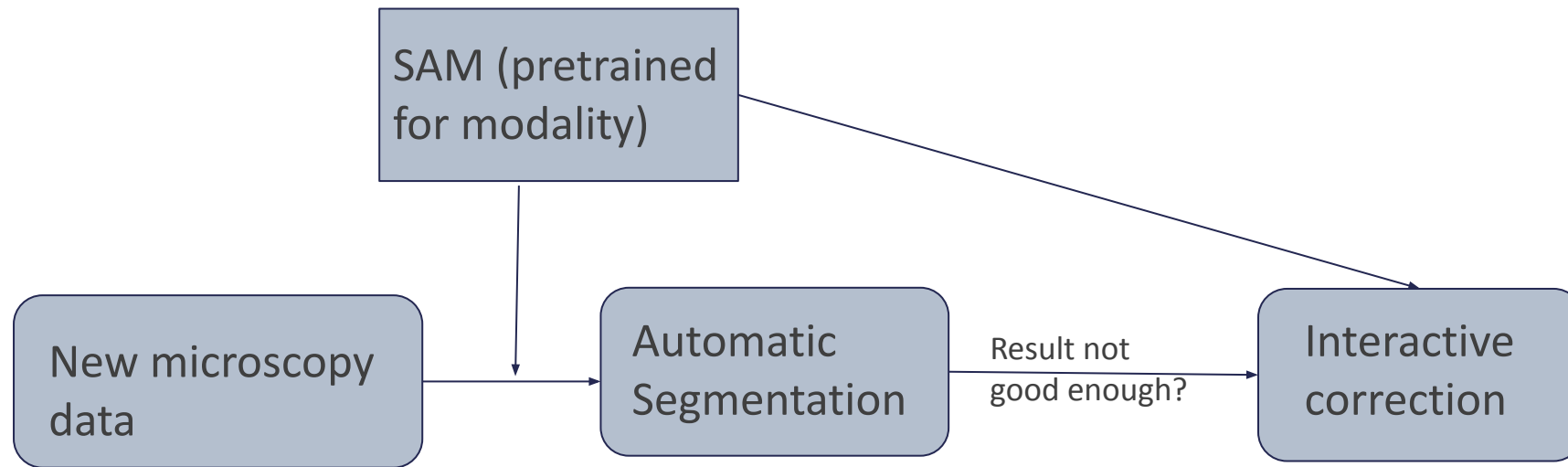
# Application in practice

---



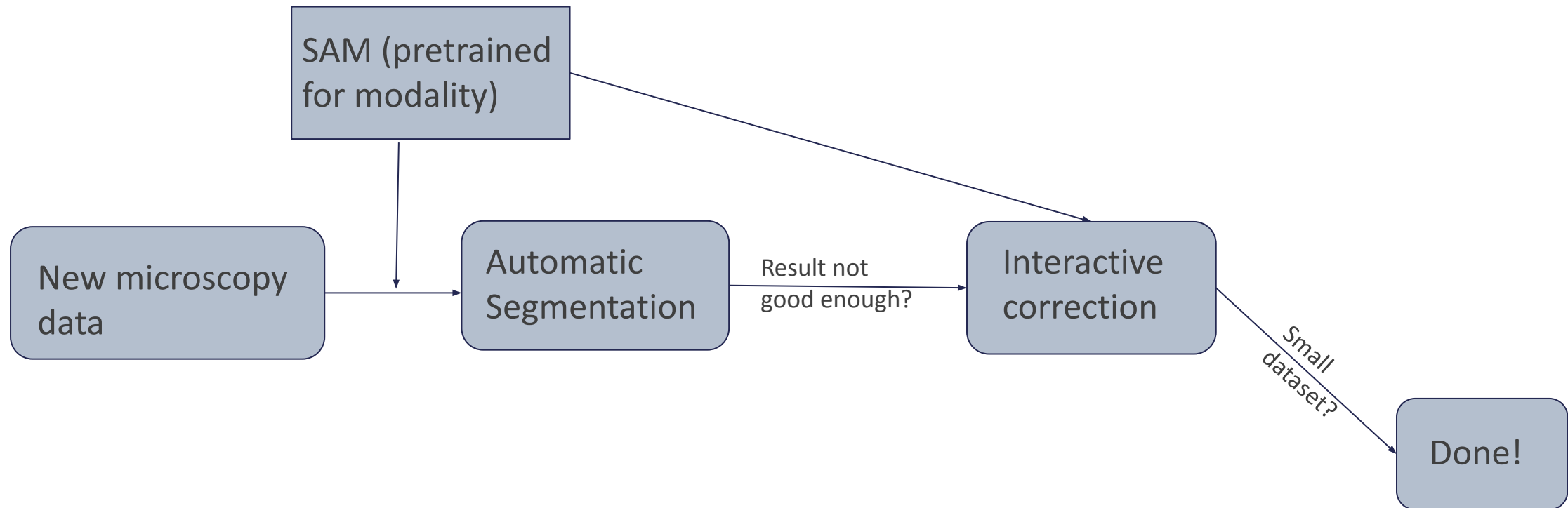
# Application in practice

---



# Application in practice

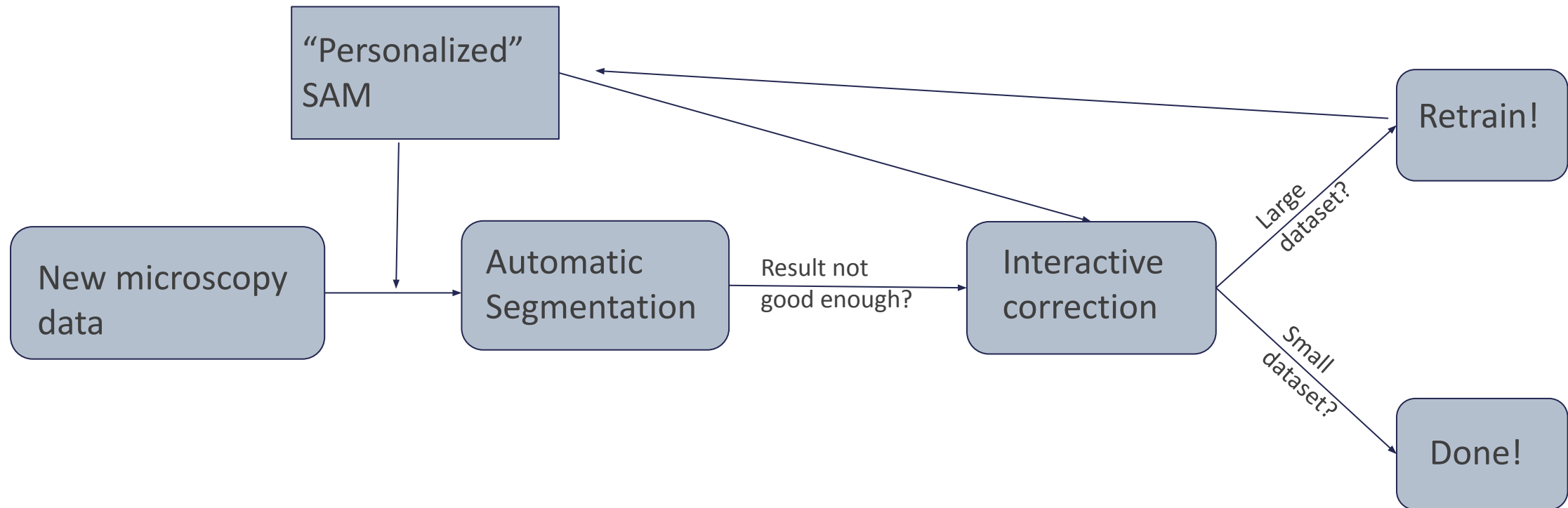
---





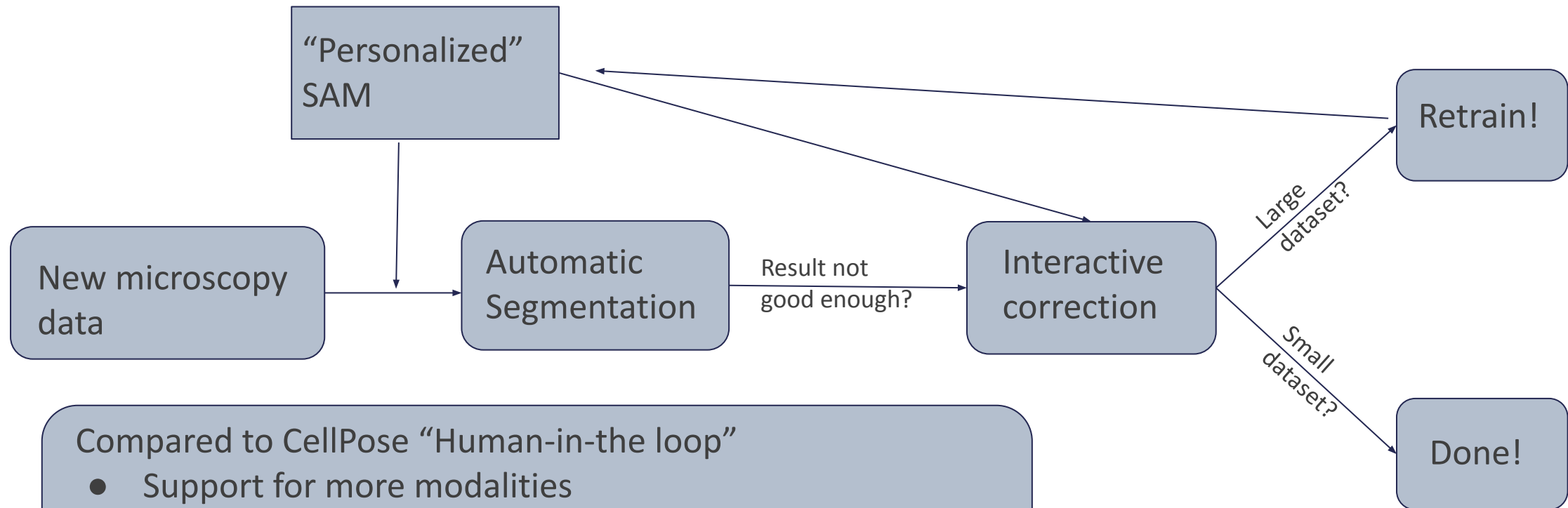
# Application in practice

---



# Application in practice

---



Compared to CellPose “Human-in-the loop”

- Support for more modalities
- Interactive correction speeds up annotation significantly!
- **BUT:** Training model takes longer (esp. on CPU)

# Application in practice

Segmenting mitochondria in electron tomography with Wiebke Möbius and Leonie Schadt (MPI Göttingen)



Segmenting organelles in phytoplankton with Karel Mocaer (EMBL Heidelberg)

# microSAM: Napari Integration

---

# microSAM: SAM for napari

---

- napari plugins that enable interactive and automatic:
  - 2D Segmentation
  - 3D Segmentation
  - Tracking (2D + time)
  - Finetuning on own data
- Core functionality:
  - Default + generalist models
  - Multidimensional segmentation / tracking (interactive and automatic)
  - Tiled prediction for large images

# microSAM: SAM for napari

---

- napari plugins that enable interactive and automatic:
  - 2D Segmentation
  - 3D Segmentation
  - Tracking (2D + time)
  - Finetuning on own data
- Core functionality:
  - Default + generalist models
  - Multidimensional segmentation / tracking (interactive and automatic)
  - Tiled prediction for large images

Code and documentation available at:

<https://github.com/computational-cell-analytics/micro-sam>

**New release (v1.0):**

- Latest microscopy models, compatible with BioImage.IO modelzoo.
- Updated and extended UI, napari plugin integration.





# Parallel Developments & Next Steps & Outlook

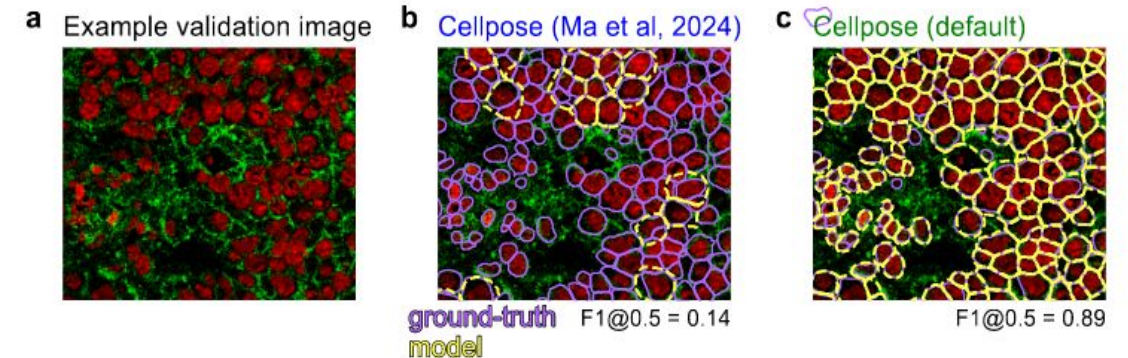
---

# Are vision transformers better than CNNs?

nnU-Net Revisited:  
A Call for Rigorous Validation  
in 3D Medical Image Segmentation

## Transformers do not outperform Cellpose

Carsen Stringer<sup>†</sup>, Marius Pachitariu<sup>†</sup>



# Are vision transformers better than CNNs?

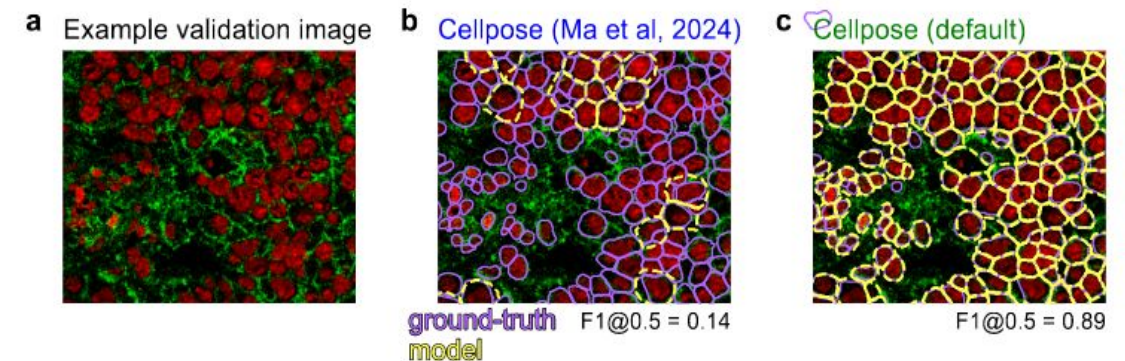
nnU-Net Revisited:  
A Call for Rigorous Validation  
in 3D Medical Image Segmentation

## Transformers do not outperform Cellpose

Carsen Stringer<sup>†</sup>, Marius Pachitariu<sup>†</sup>

Transformers are not inherently better than CNNs!

- Disadvantages for small data and runtimes
- Possible advantage for large data
- Advantage of same architecture as NLP
  - Can profit from improvements and **easier to build multi-modal models.**



Feedback and contributions on the tool are very welcome!

## Next steps

---

- Integration of efficient training procedures for finetuning (LoRA, **QLoRA**)
  - To enable better training on CPU and small GPUs
- Provide better and more models:
  - EM Organelle Generalist Model
    - Training on OpenOrganelle and other organelle segmentation datasets.
  - Histopathology Model

Check out our repository for all the details:

<https://github.com/computational-cell-analytics/micro-sam>

# Outlook:

## Universal microscopy segmentation and tracking

---

- Incorporate 3D (2D + time) segmentation in SAM-like model
  - Advantage ransformer: same model for 2d and 3d
- Vision Mamba: Investigate newer (more efficient) architectures
  - Our recent (**preliminary!**) work shows promise: <https://arxiv.org/abs/2404.07705>
- Semantic awareness (e.g. differentiate organelles in EM, one model for microscopy)
- Zero-shot adaptation (improve segmentation from examples)



# Acknowledgments

## EMBL Heidelberg

**Anna Kreshuk & her group**  
et al.

## Uni Göttingen & Campus

**Alexander Ecker**  
**Tobias Moser**  
Silvio Rizzoli  
et al.

## My group

**Anwai Archit**  
**Luca Freckmann**  
Sushmita Nair  
Marei Freitag  
Sagnik Gupta  
et al.



Göttingen  
Campus



From Molecular Machines to Networks of Excitable Cells

**DFG**



**SFB 1286**

**CIDAS**  
Campus-Institut Data Science





# Hands-on Session

---

# Before the break: Access VM, Install micro-sam

---

Access your VM and run:

```
$ mamba create -c pytorch -c nvidia -c conda-forge -n micro-sam pytorch micro_sam  
pytorch-cuda=11.6 -y
```

# Plan Hands-On-Session

---

- Starting the tool, explain plugins
  - See <https://youtu.be/gcv0fa84mCc> for a rough idea
- 2D Segmentation on LiveCELL
  - Compare default and finetuned model (vit\_b, show auto segmentation for vit\_b\_lm)
  - See [https://youtu.be/9xjJBg\\_Bfuc](https://youtu.be/9xjJBg_Bfuc) for a rough idea
- 3D Segmentation on Lucchi
  - Segment mitos in 3d with the em model
  - See <https://youtu.be/nqpyNQSyu74> for a rough idea
- If there is still time students can choose to:
  - Try annotation on their own data
  - Start finetuning a model, either on their own data or on our sample data. Using [https://github.com/computational-cell-analytics/micro-sam/blob/master/notebooks/sam\\_finnetuning.ipynb](https://github.com/computational-cell-analytics/micro-sam/blob/master/notebooks/sam_finnetuning.ipynb)

# Application on your data

---

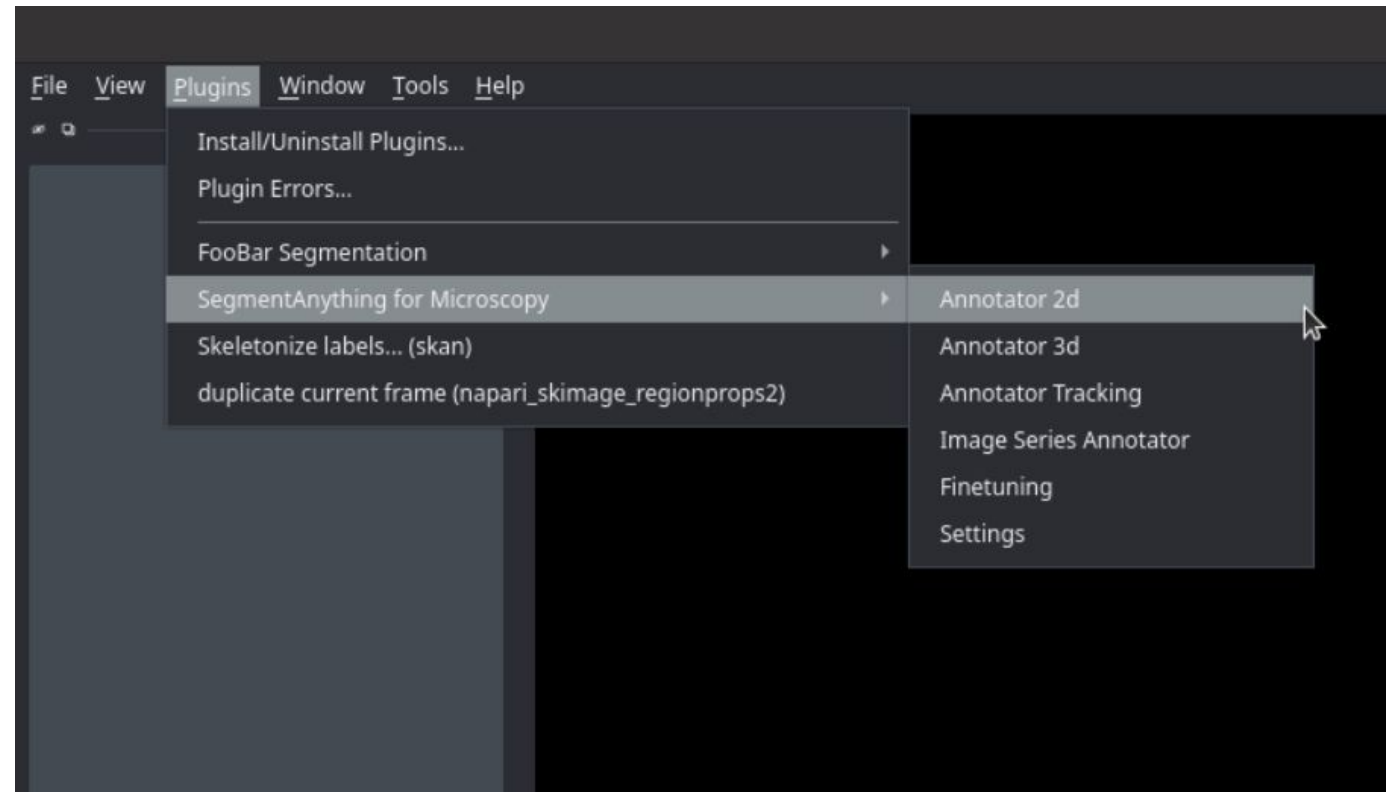
- 2D segmentation -> “Annotator 2d” or “Image Series Annotator”
- 3D segmentation -> “Annotator 3d”
- Tracking -> “Annotator Tracking”

Which model should I use?

LM -> vit\_X\_lm

EM (potato-shaped structures)  
-> vit\_X\_em\_organelles

Other -> vit\_X



# Finetuning Notebook

[https://github.com/dl4mia/03\\_learned\\_representations/blob/main/sam\\_fineting\\_HT.ipynb](https://github.com/dl4mia/03_learned_representations/blob/main/sam_fineting_HT.ipynb)

To get it on your VM:

```
$ git pull origin main
```

metrics	the first commit	3 weeks ago
.gitignore	the first commit	3 weeks ago
LICENSE	Initial commit	3 weeks ago
README.md	applied some changes based on reviews	last week
generate_exercise.py	set default output file to None	3 weeks ago
sam_fineting_HT.ipynb	Add notebook for finetuning SAM on custom data	yesterday
setup.sh	add jupyterlab installation to base env	3 days ago
transformer_representation_part_1_exercis...	update environment name to 03	last week
transformer_representation_part_1_solutio...	update environment name to 03	last week
transformer_representation_part_1_solutio...	update environment name to 03	last week
transformer_representation_part_2_exercis...	update environment name to 03	last week
transformer_representation_part_2_solutio...	update environment name to 03	last week
transformer_representation_part_2_solutio...	update environment name to 03	last week
utils.py	applied some changes based on reviews	last week
README BSD-3-Clause license		