

Deep Learning for Image Analysis

Segmentation, Tracking & Transfer Learning

Constantin Pape

Institut für Informatik, Georg August Universität Göttingen



@cppape.bsky.social

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

Plan for today

- 9:00 - 10:30 - Lecture, Q&A & Intro to Exercise 2
- 10:30 - 12:30 - Work on Exercise 2
- 12:00 - 13:00 - Lunch
- 13:00 - 16:00 - Work on exercise 2
- 16:00 - Recap Exercise 2

Semantic Segmentation

and the U-Net architecture

Segmentation tasks: Scene understanding

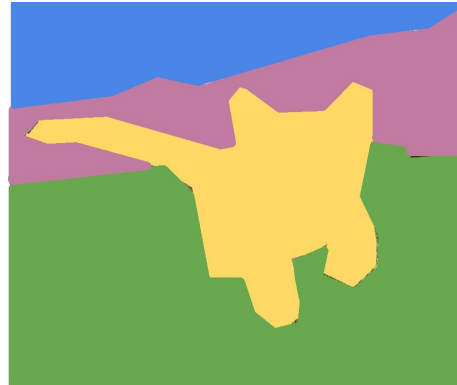
Image
Classification



Cat

No spatial extent

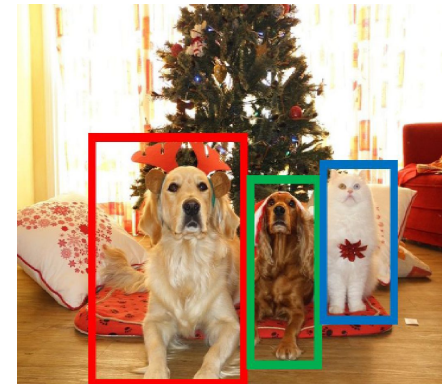
Semantic
Segmentation



Sky, Trees, Grass, Cat

No objects, just pixels

Object
Detection



Dog, Dog, Cat

Multiple objects

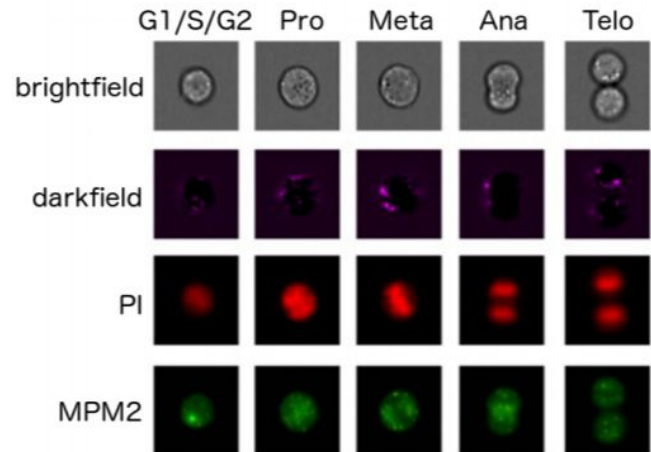
Instance
Segmentation



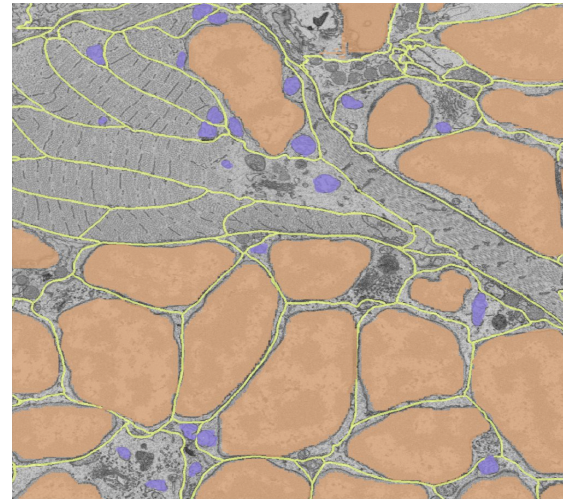
Dog, Dog, Cat

Segmentation for biomedical images

Image
Classification

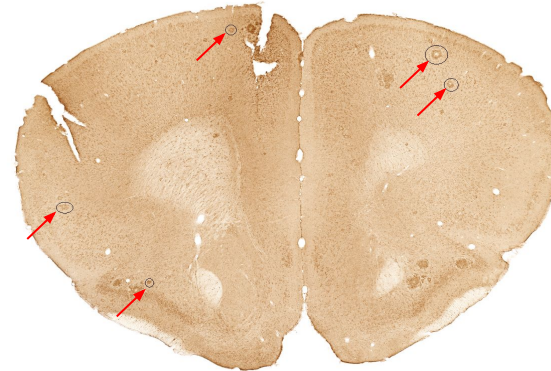


Semantic
Segmentation

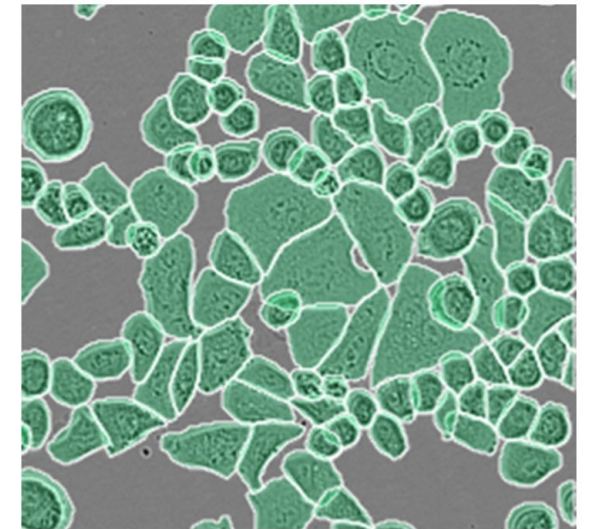


Cell Membrane
Nuclei
Mitochondria

Object
Detection

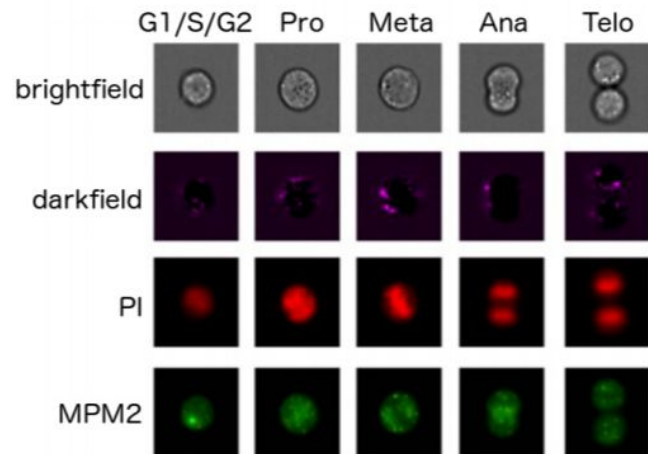


Instance
Segmentation

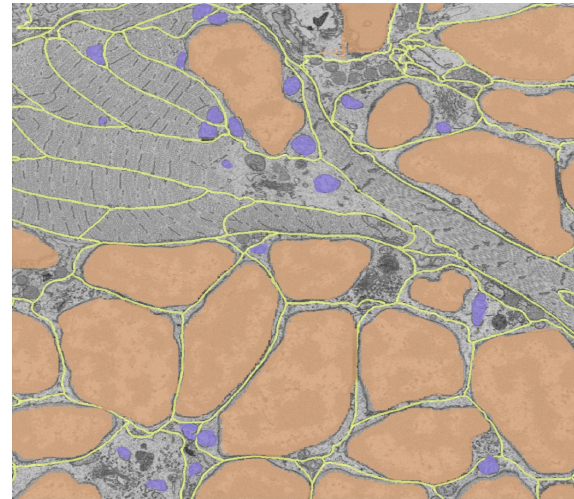


Segmentation for biomedical images

Image
Classification

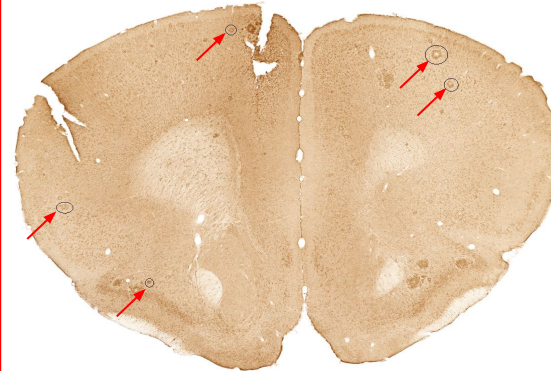


Semantic
Segmentation

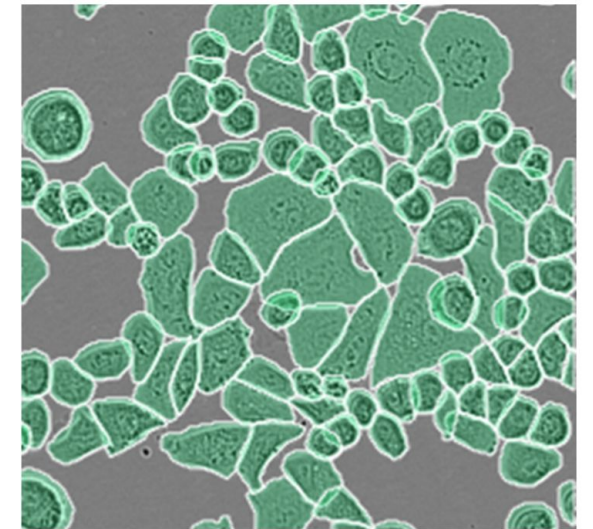


Cell Membrane
Nuclei
Mitochondria

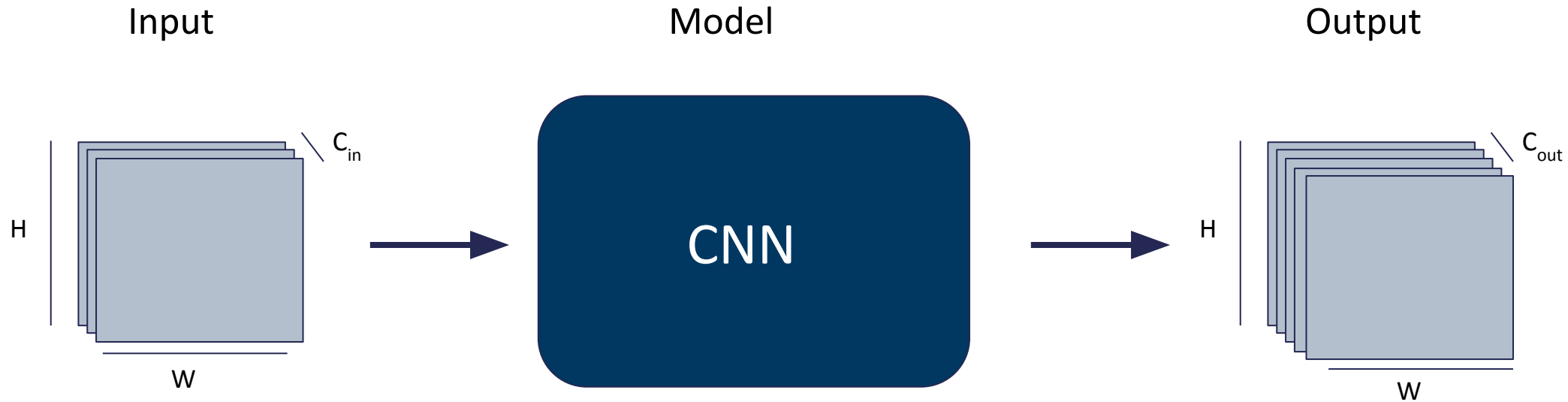
Object
Detection



Instance
Segmentation



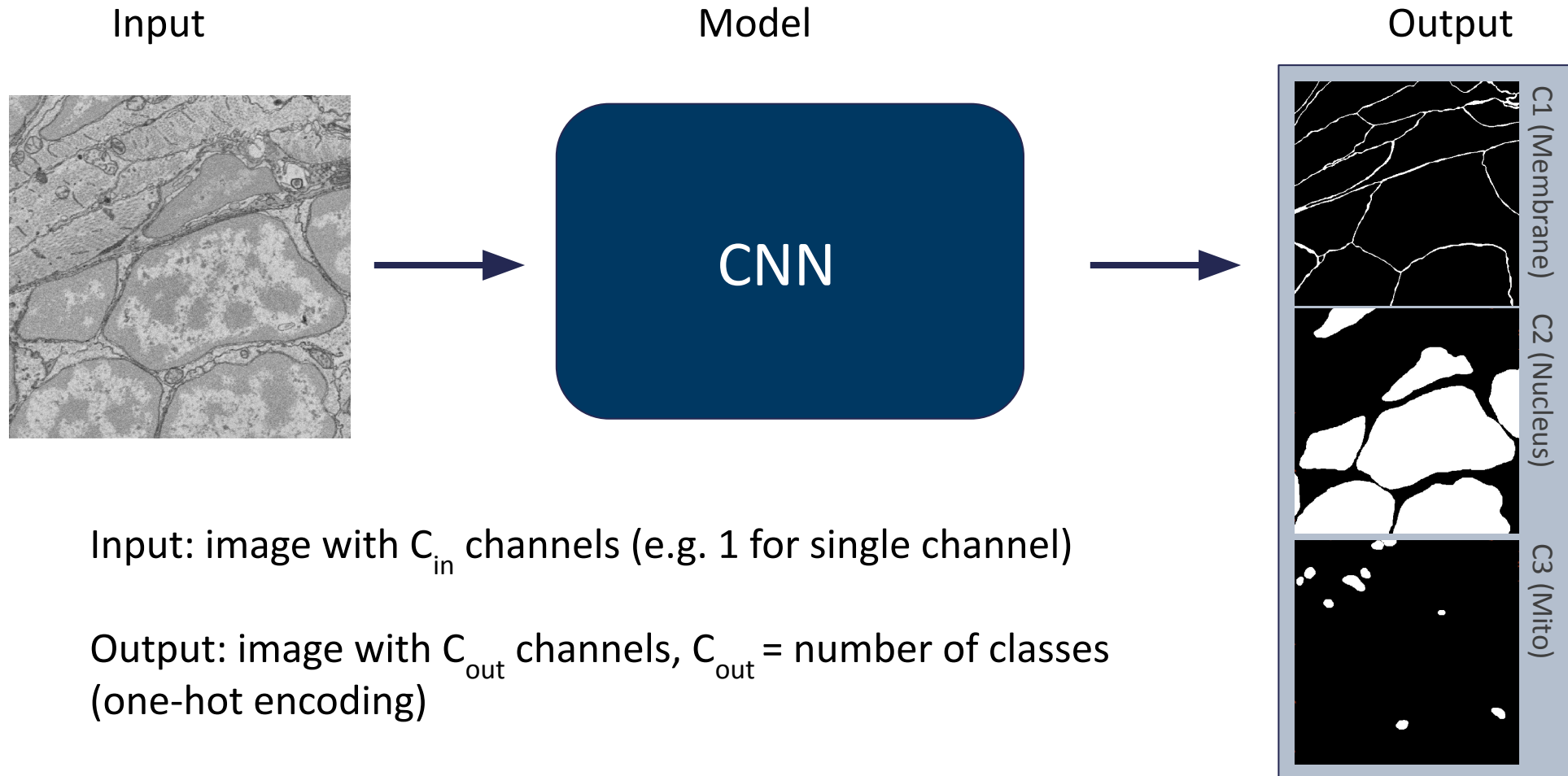
Architectures for semantic segmentation



Input: image with C_{in} channels (e.g. 1 for single channel)

Output: image with C_{out} channels, C_{out} = number of classes (one-hot encoding)

Architectures for semantic segmentation

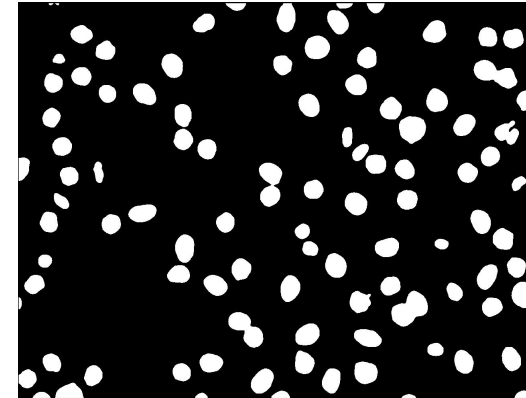
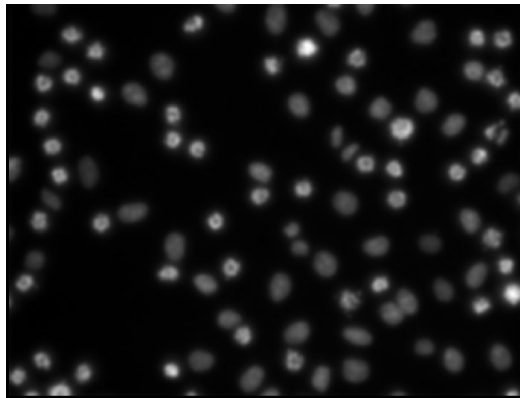


Architectures for semantic segmentation

Input

Model

Output

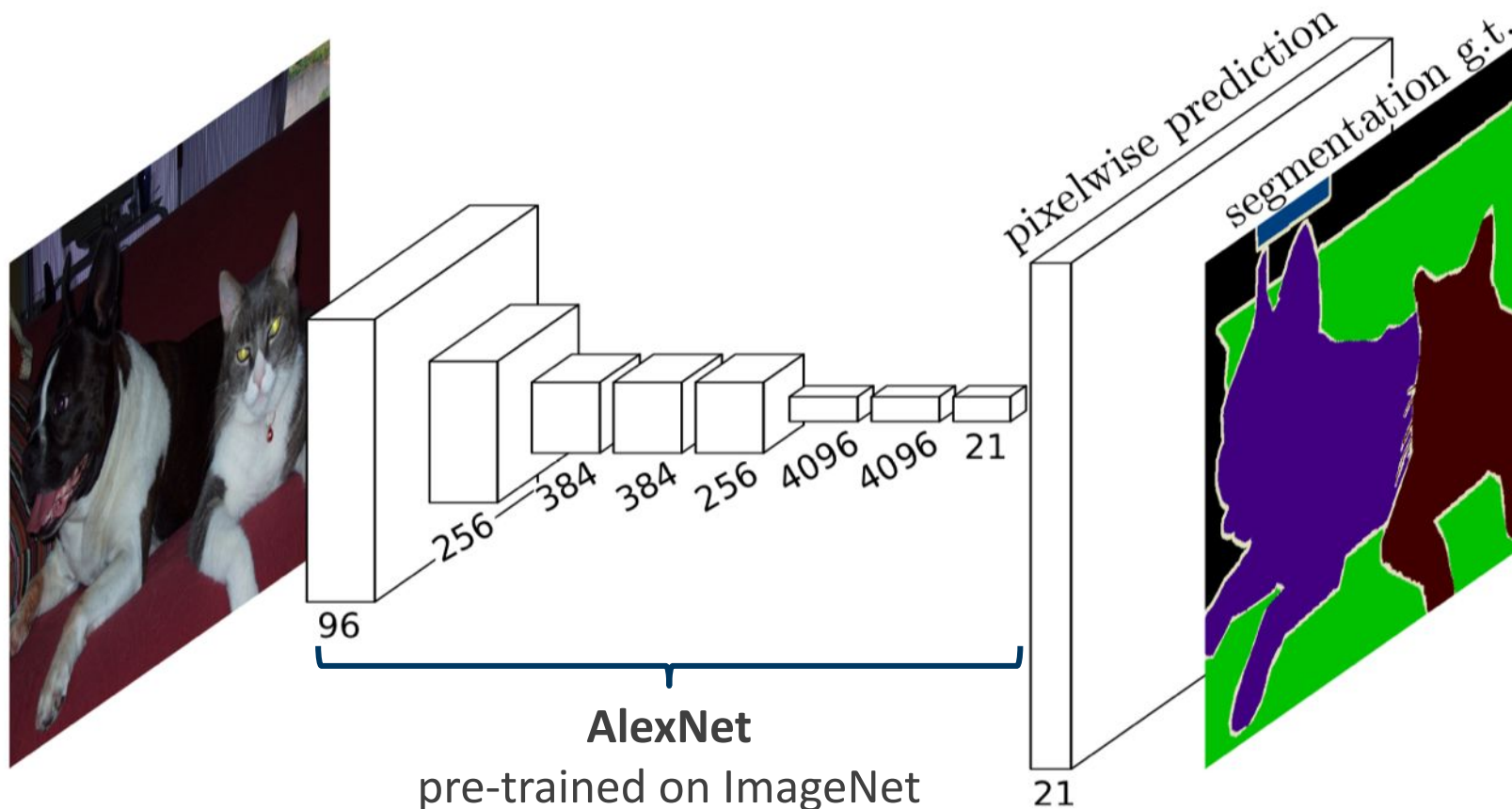


Binary nucleus segmentation

Special case: foreground vs. background segmentation:

$C_{\text{out}} = 1$ (binary segmentation)

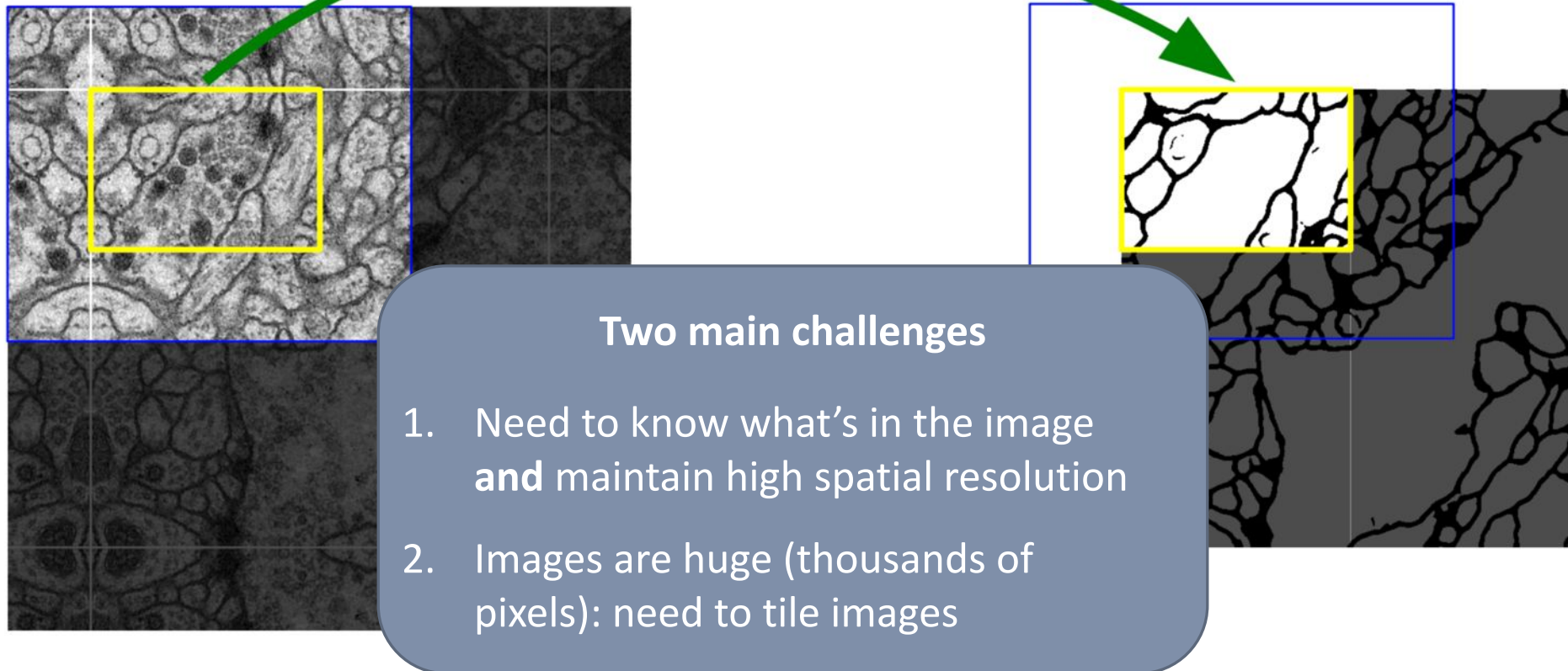
Fully convolutional network



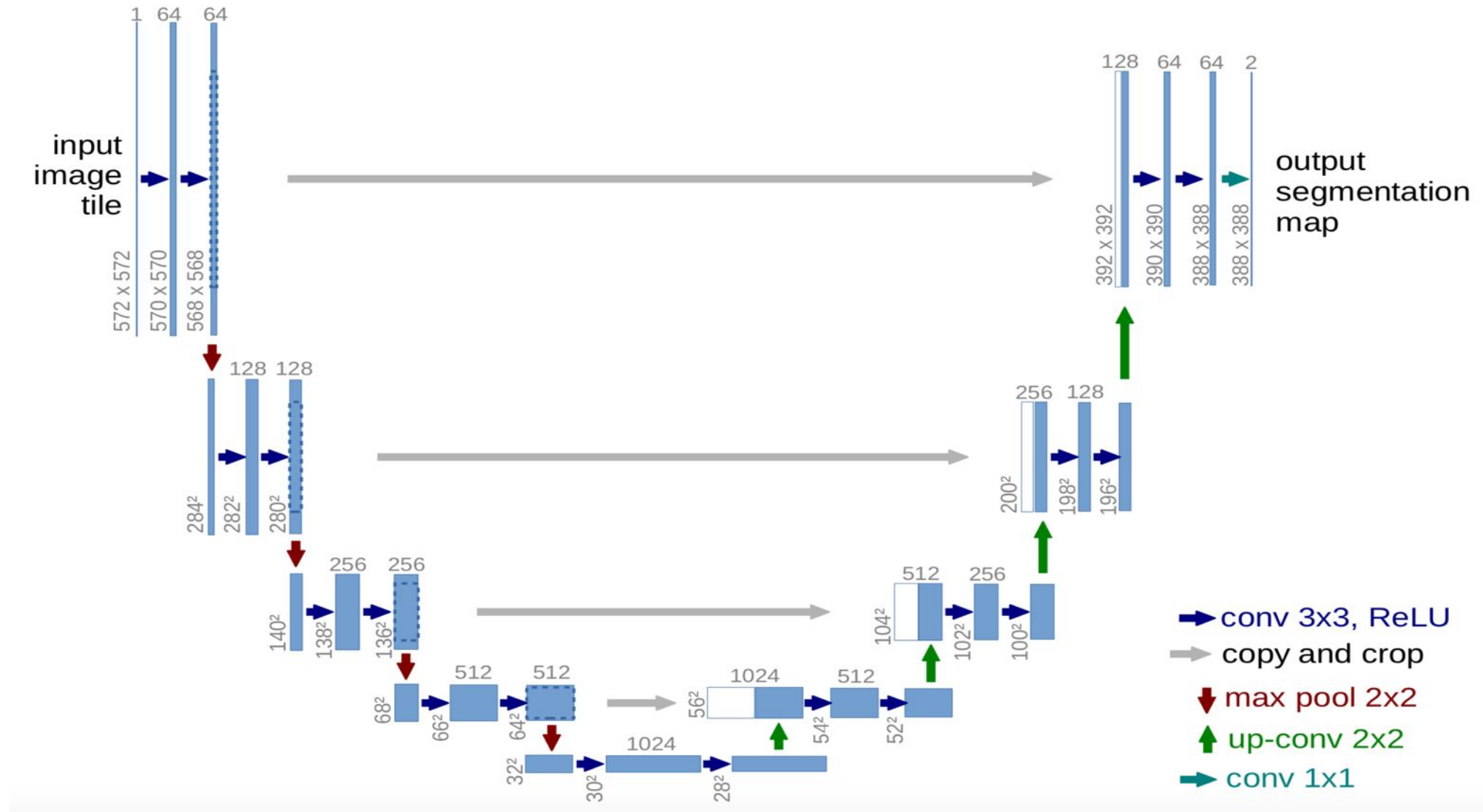
- Pretrained classification network, adapted for segmentation
 - “Transfer learning” (tomorrow)
- Outputs from convolutional layers is upsampled to predict segmentation map
- **Not high-resolution enough for biomedical images!**

Biomedical image segmentation

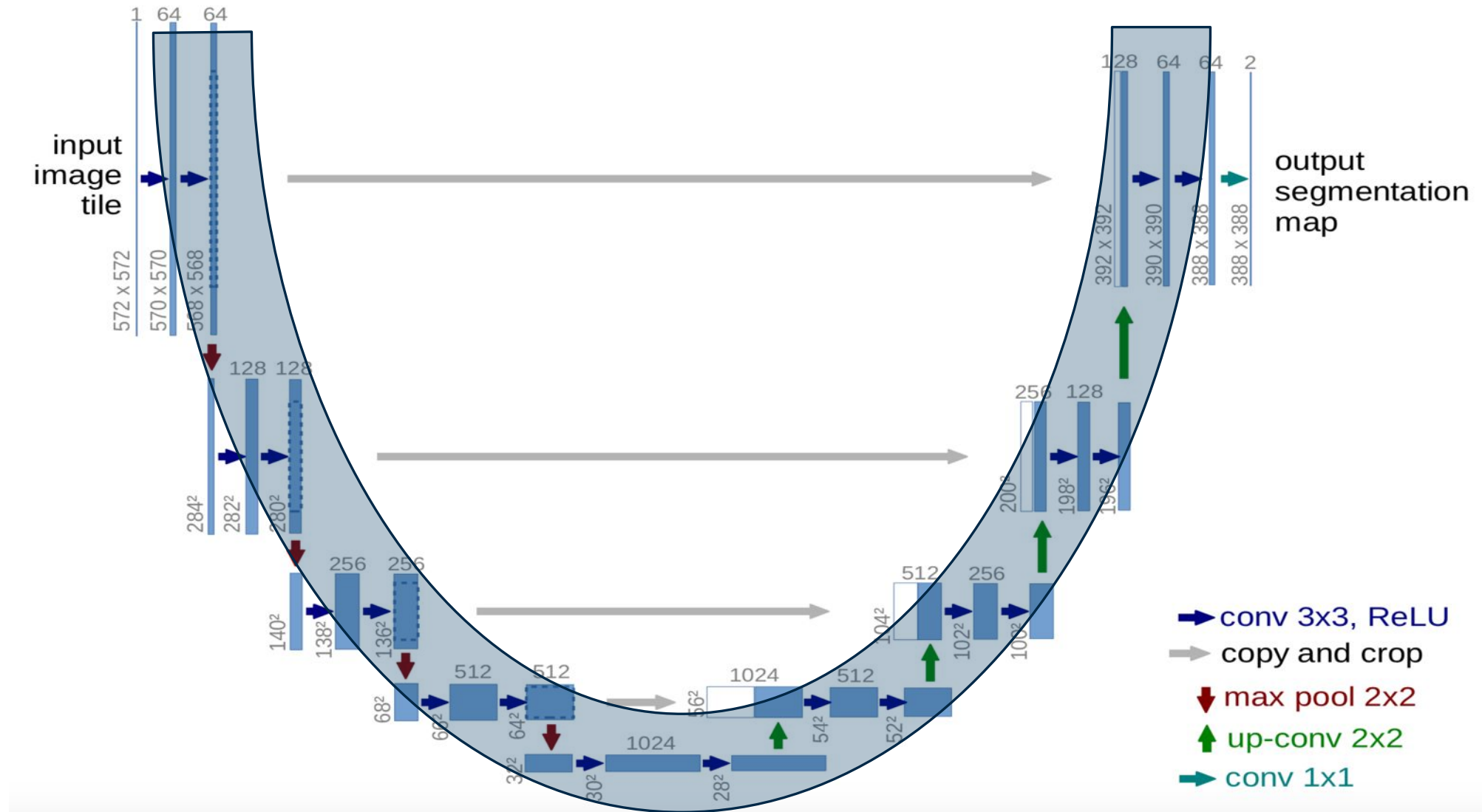
Task: segment neuronal structures from electron microscopy images



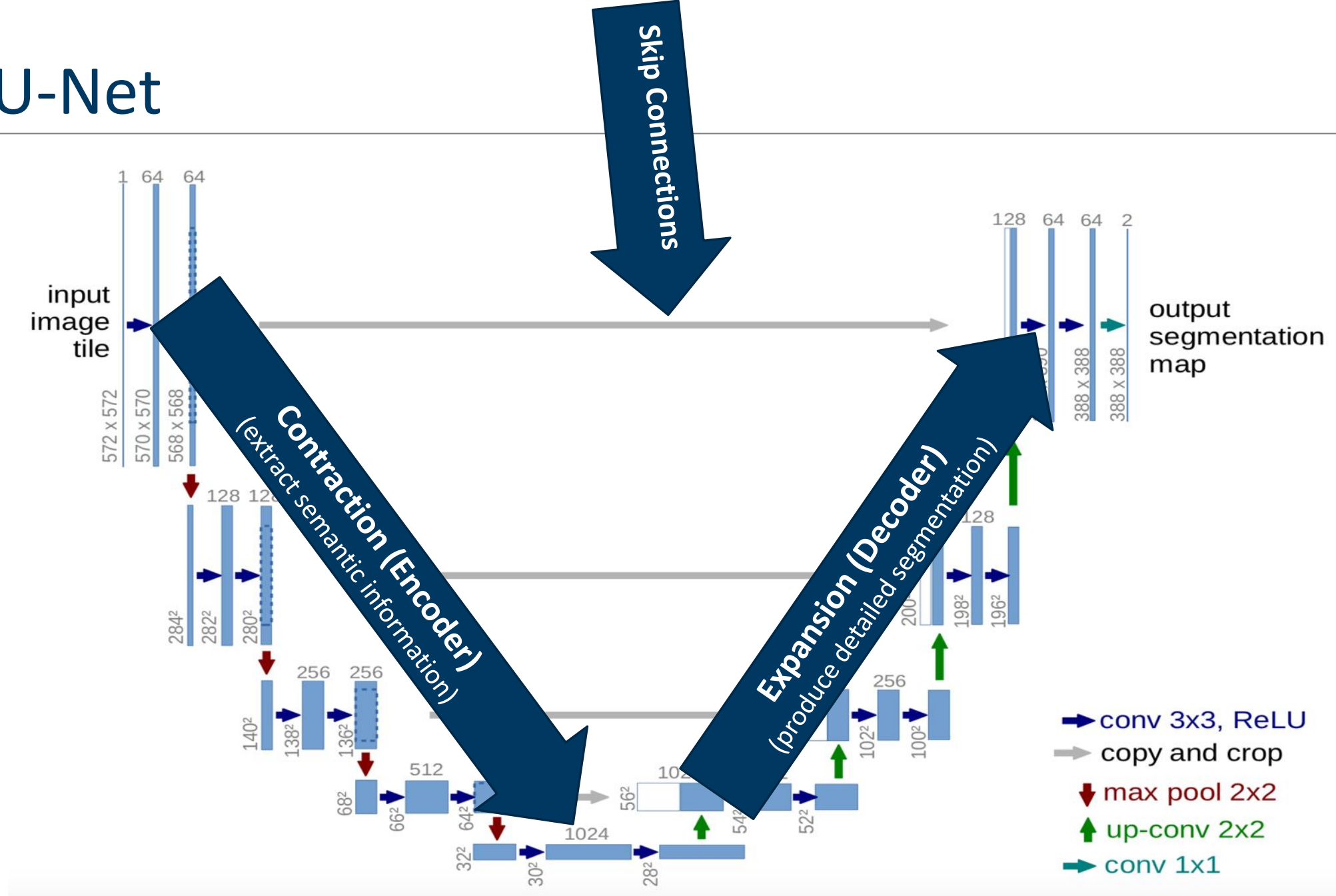
U-Net



U-Net



U-Net

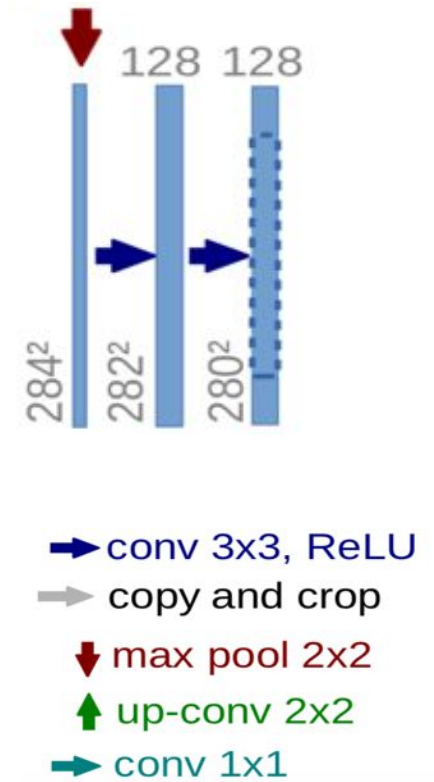


U-Net: Encoder

One level in the encoder:

- 2 convolutional layers with 3x3 kernel
 - Don't use padding -> lose 2 pixels per convolution (we will do this differently in the exercises)
- ReLU activation (the standard!)
- 2x2 max pooling to group features and reduce spatial dimension

We have seen all this yesterday!



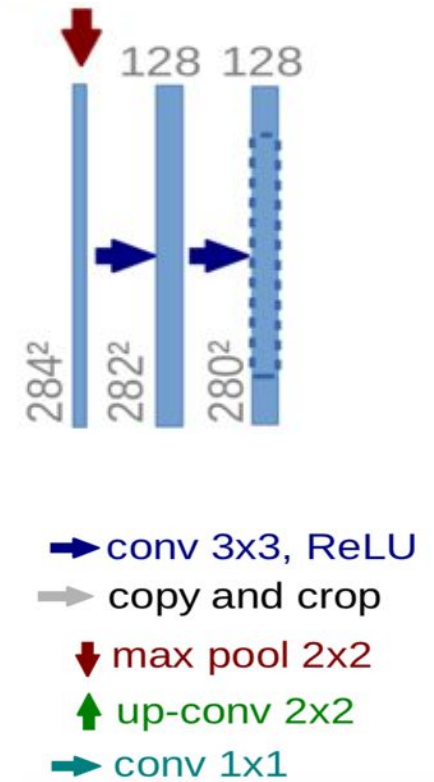
U-Net: Encoder

One level in the encoder:

- 2 convolutional layers with 3x3 kernel
 - Don't use padding -> lose 2 pixels per convolution (we will do this differently in the exercises)
- ReLU activation (the standard!)
- 2x2 max pooling to group features and reduce spatial dimension

We have seen all this yesterday!

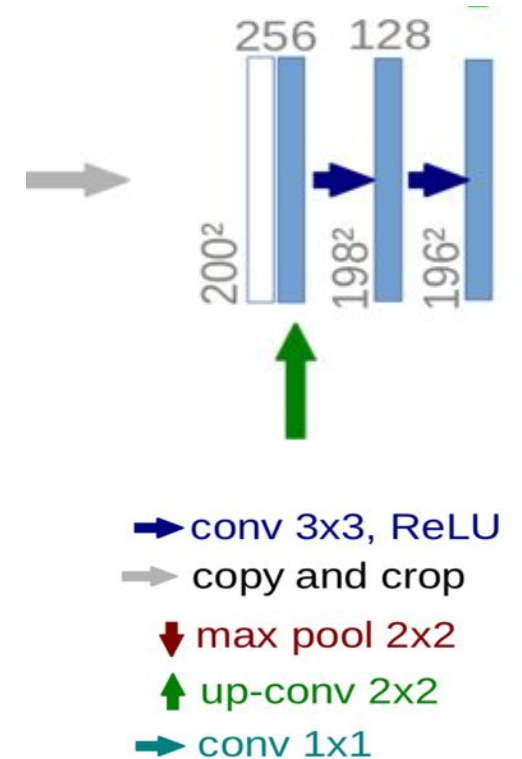
**In the exercise we will also add a normalization layer.
Keep the network activation in well-defined range and
make training more stable.**



U-Net: Decoder

One level in the decoder:

- 2 convolutional layers + ReLU (as before)
- 2x2 up-conv to (spatially upsample the features)

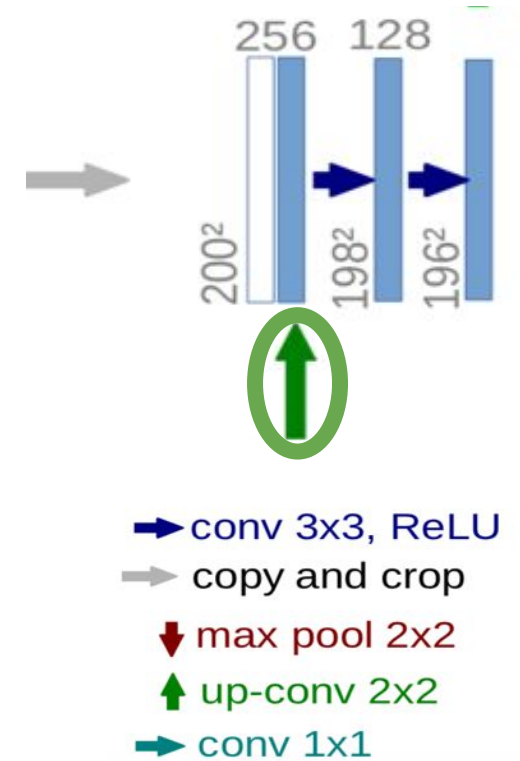


U-Net: Decoder

One level in the decoder:

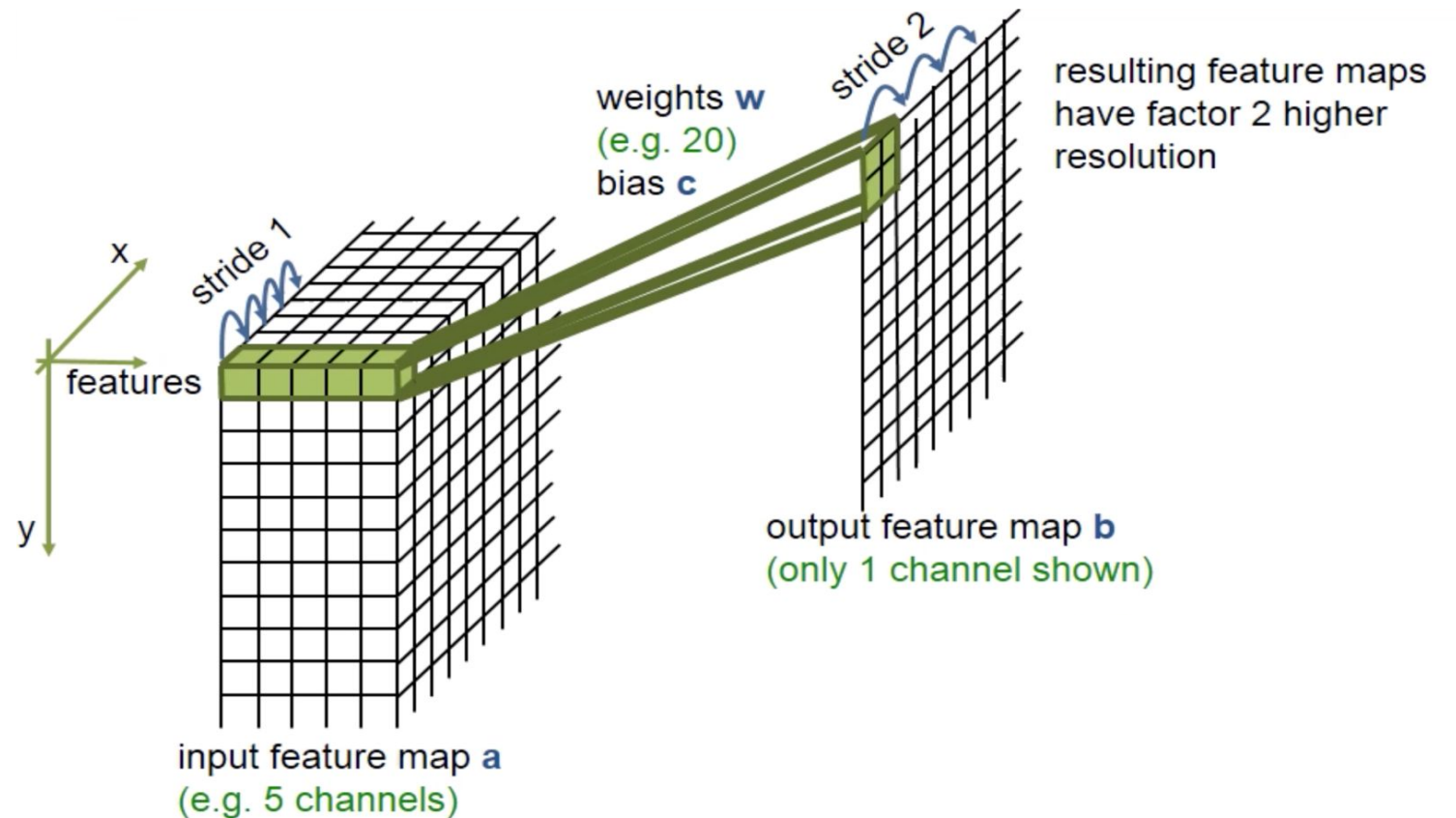
- 2 convolutional layers + ReLU (as before)
- 2x2 up-conv to (spatially upsample the features)

This is new!



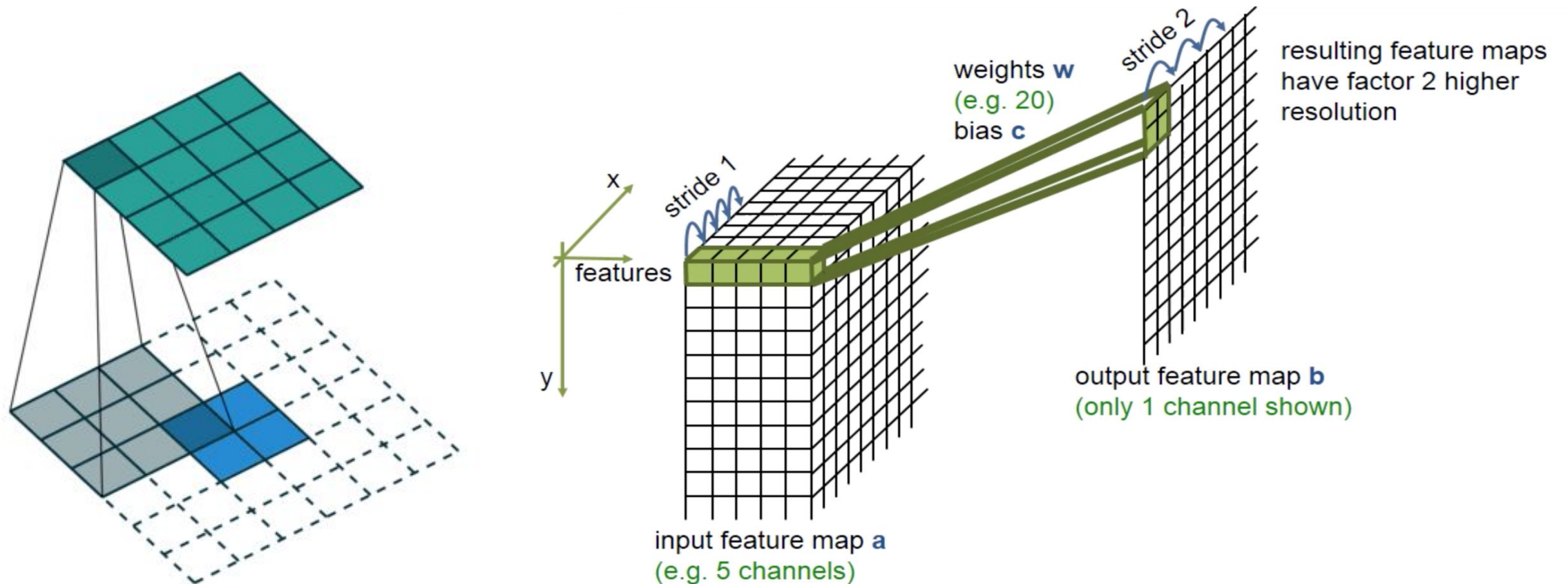
Up-convolutions

Increase the spatial dimension (“resolution”) with a learned convolutional kernel.

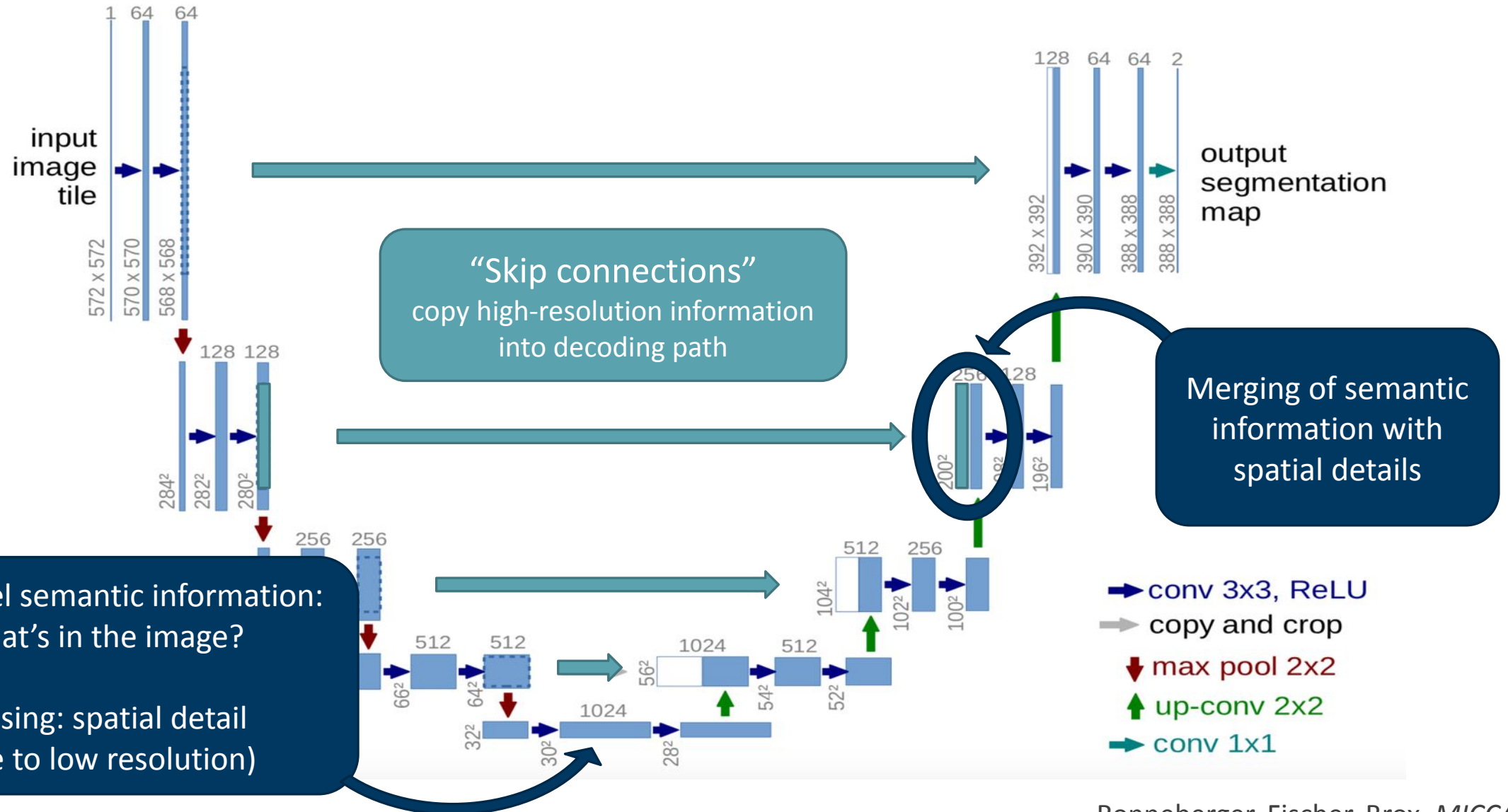


Up-convolutions

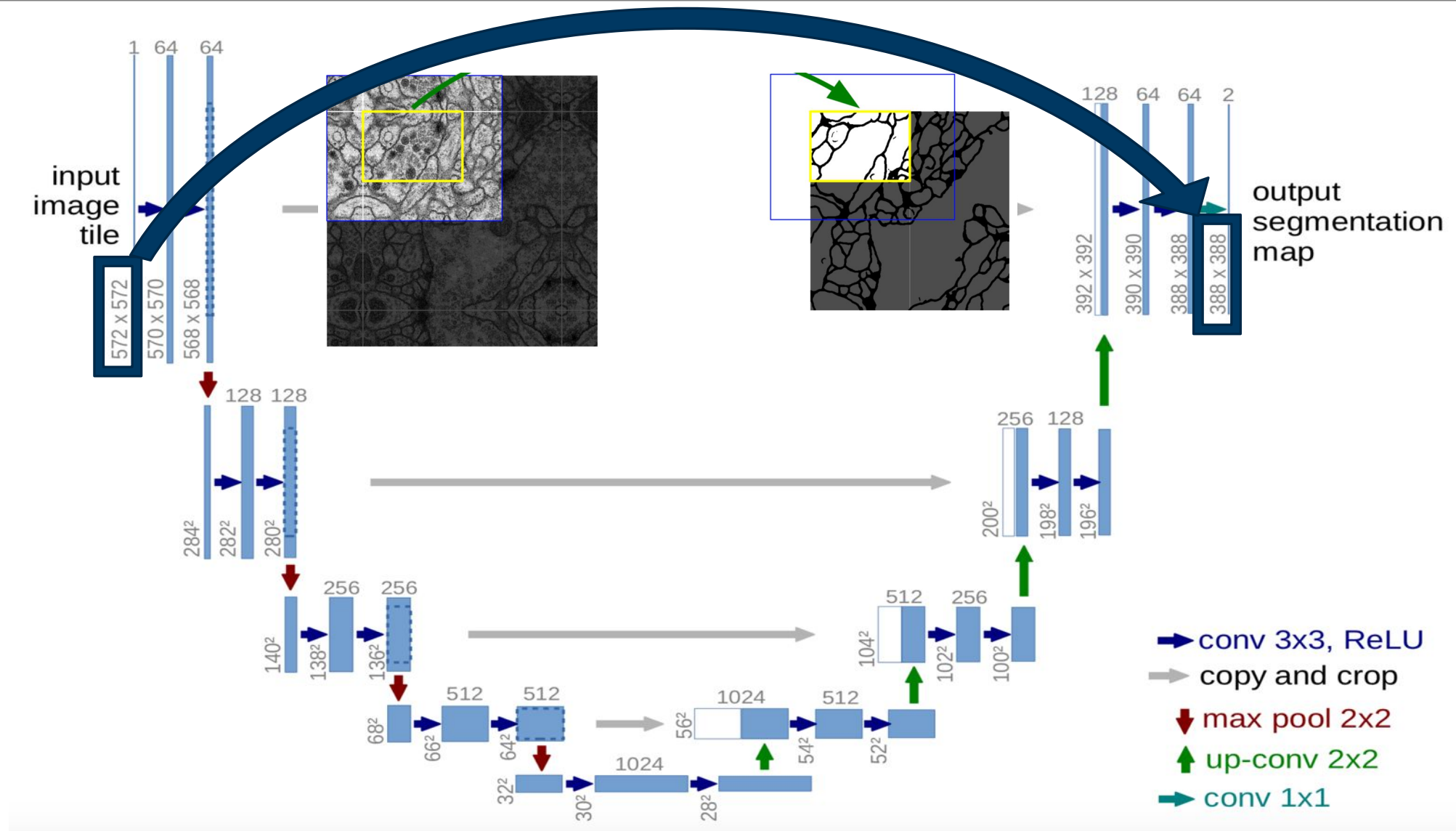
Increase the spatial dimension (“resolution”) with a learned convolutional kernel.



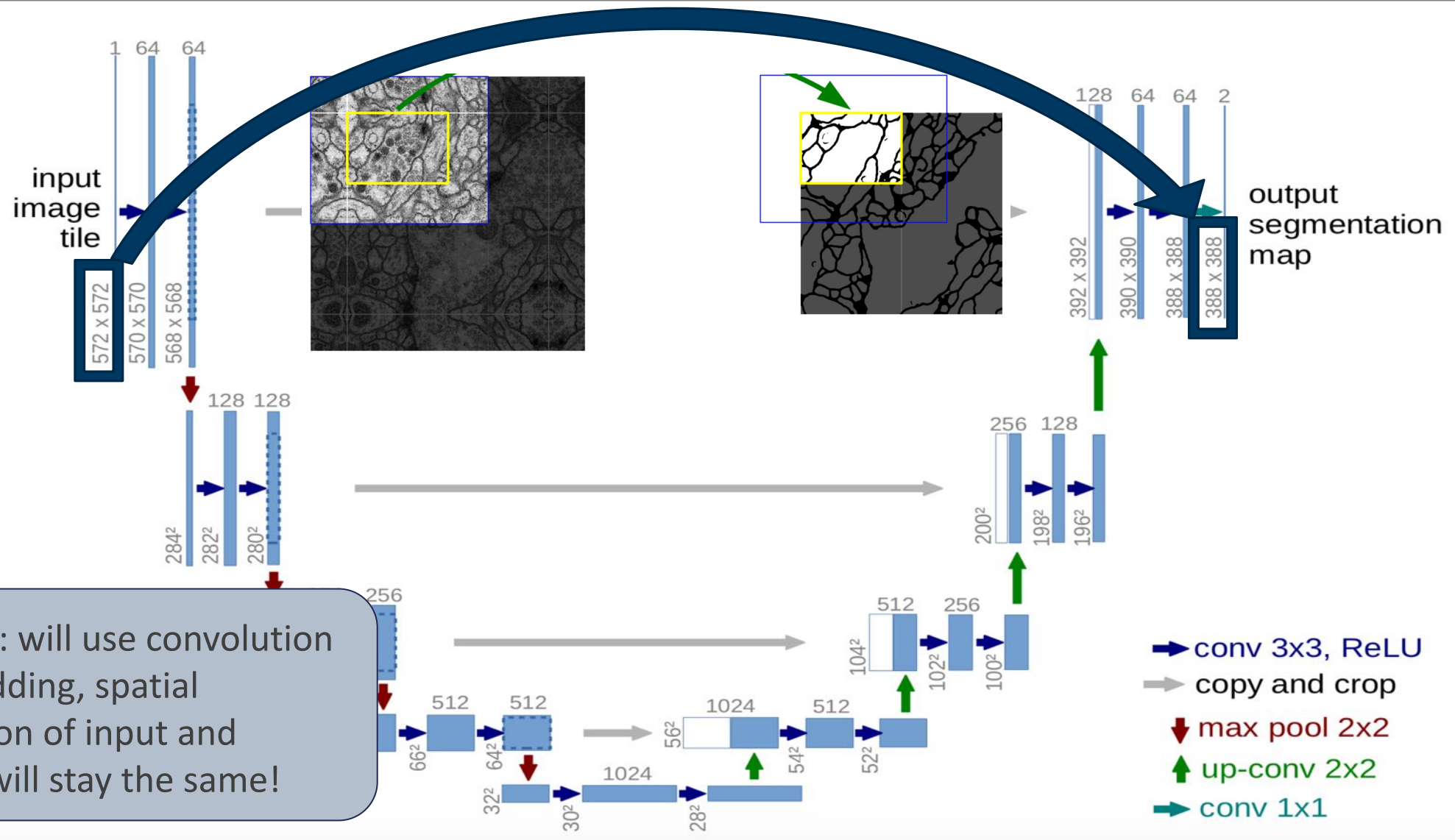
Skip connections



Spatial dimensions

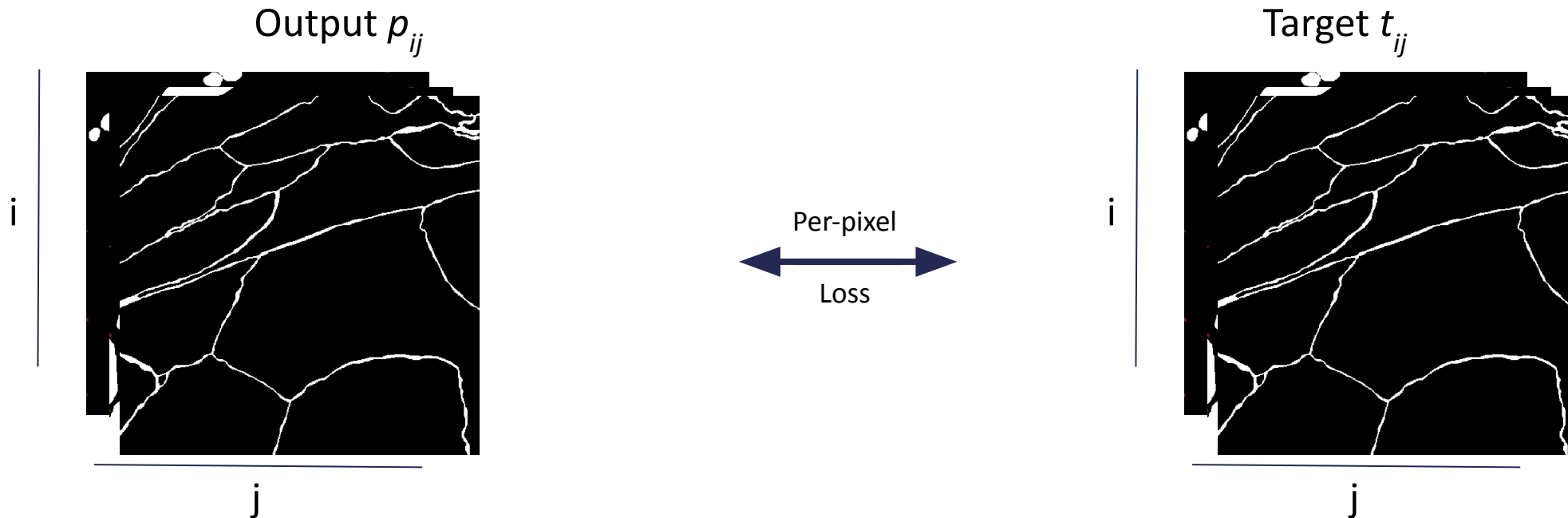


Spatial dimensions



Loss function

Pixel-wise cross entropy: $-\frac{1}{N} \sum_{ij} \sum_k \log p_{ij}^k t_{ij}^k$

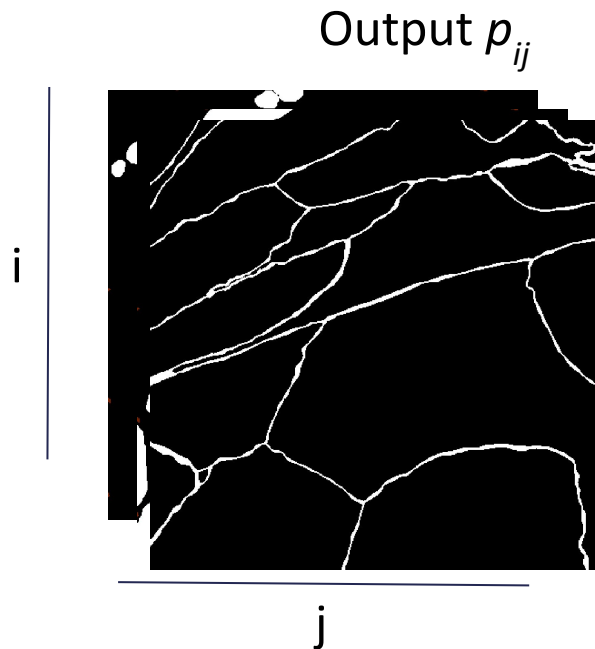


Loss function

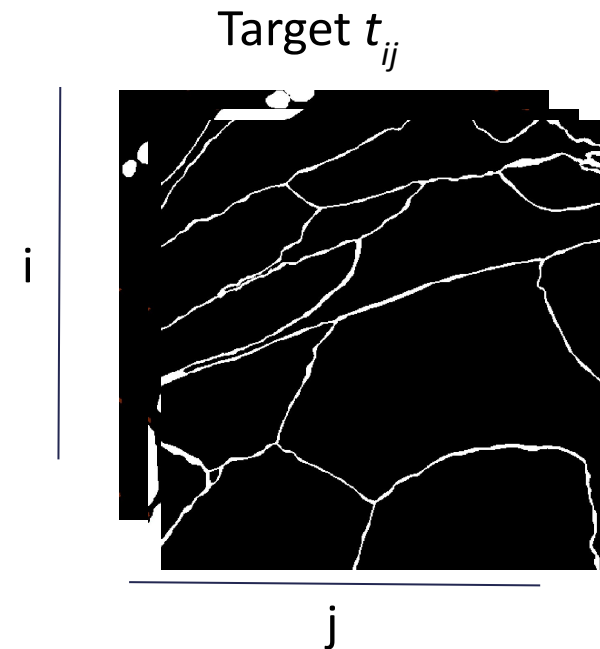
Pixel-wise cross entropy:

$$-\frac{1}{N} \sum_{ij} \sum_k \log p_{ij}^k t_{ij}^k$$

Sum pixels
Sum classes



Per-pixel
Loss



Loss function

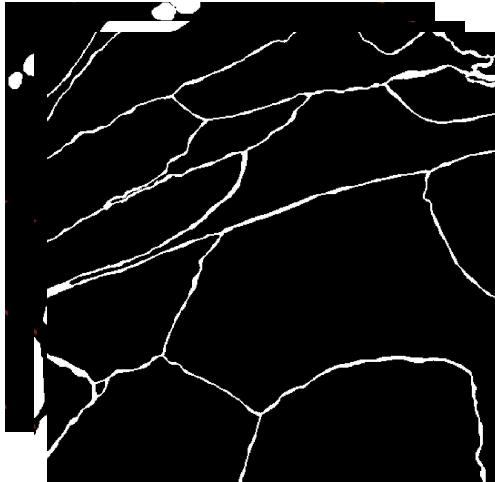
Pixel-wise cross entropy:

$$-\frac{1}{N} \sum_{ij} \sum_k \log p_{ij}^k t_{ij}^k$$

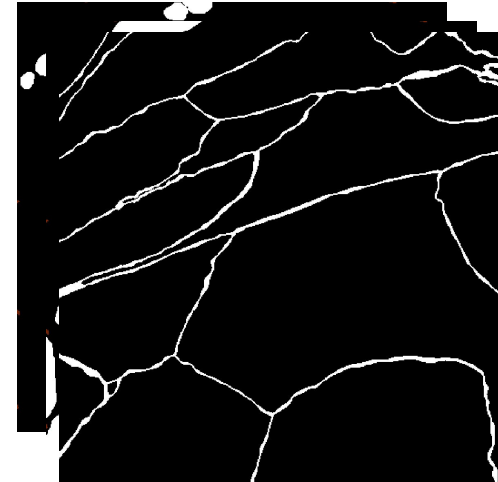
Sum pixels
Sum classes

Output p_{ij}

Target t_{ij}



Problem: target is often **unbalanced**.
E.g. much more background than membrane pixels!
Cross Entropy will not work well



Solutions for imbalanced target

- Weight the (Binary) Cross Entropy with class frequencies.
 - How to weight exactly can be tricky.
- **Use loss function that is not sensitive to imbalance:
Dice Coefficient**

$$\text{DICE} = 1 - \frac{2 \sum_n^N t_n p_n}{\sum_n^N t_n^2 + \sum_n^N p_n^2}$$

Solutions for imbalanced target

- Weight the (Binary) Cross Entropy with class frequencies.
 - How to weight exactly can be tricky.
- **Use loss function that is not sensitive to imbalance:
Dice Coefficient**

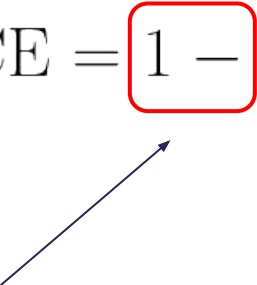
$$\text{DICE} = 1 - \frac{2 \sum_n^N t_n p_n}{\sum_n^N t_n^2 + \sum_n^N p_n^2}$$

Intersection of target and prediction

Union of target and prediction
(the denominator makes it insensitive to class imbalance!)

Solutions for imbalanced target

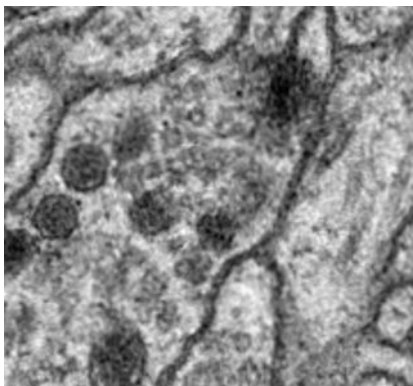
- Weight the (Binary) Cross Entropy with class frequencies.
 - How to weight exactly can be tricky.
- **Use loss function that is not sensitive to imbalance:
Dice Coefficient**

$$\text{DICE} = 1 - \frac{2 \sum_n^N t_n p_n}{\sum_n^N t_n^2 + \sum_n^N p_n^2}$$


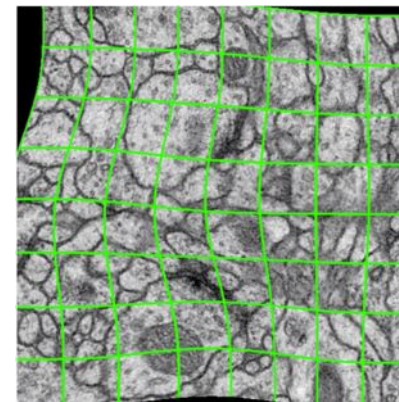
“1 - ” because low values must
correspond to good solution

U-net: training

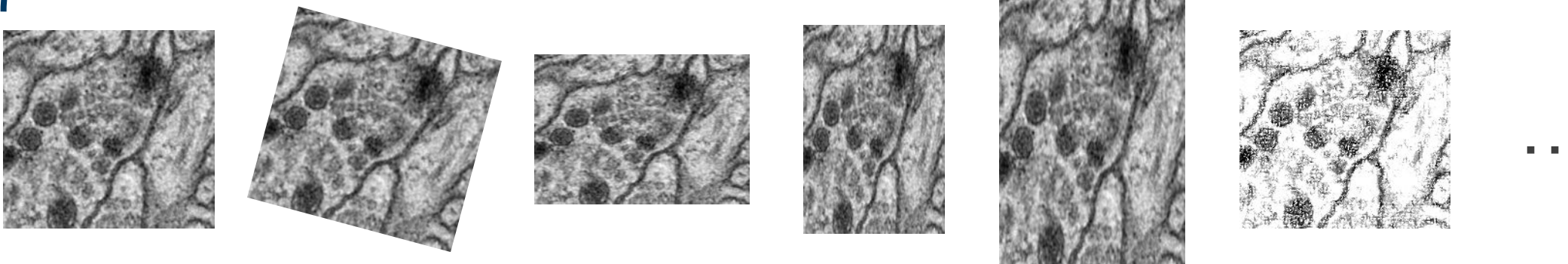
Bioimage datasets are often small!



Elastic image deformations



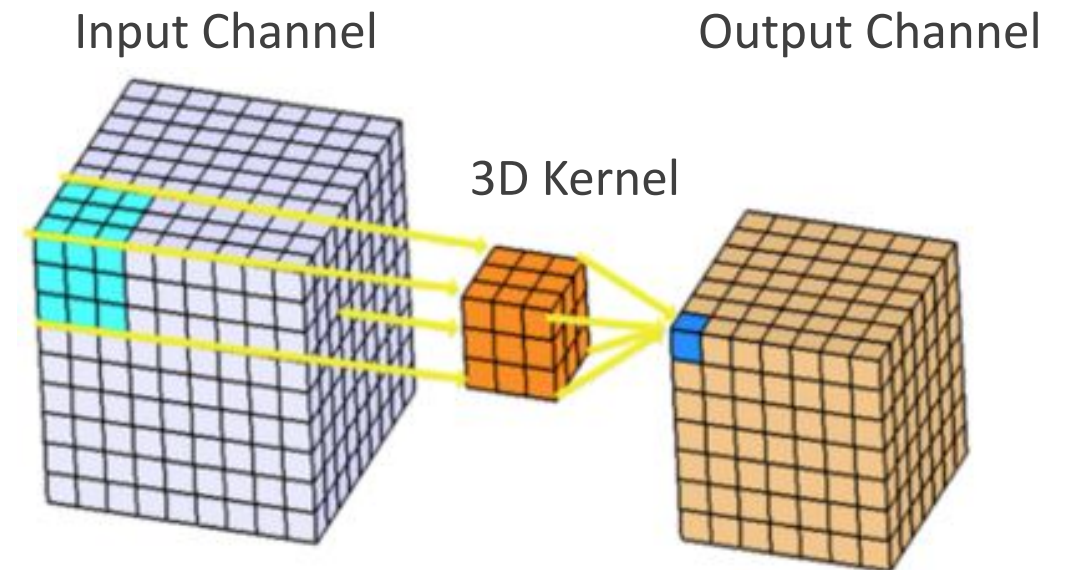
For best results: use data augmentation!



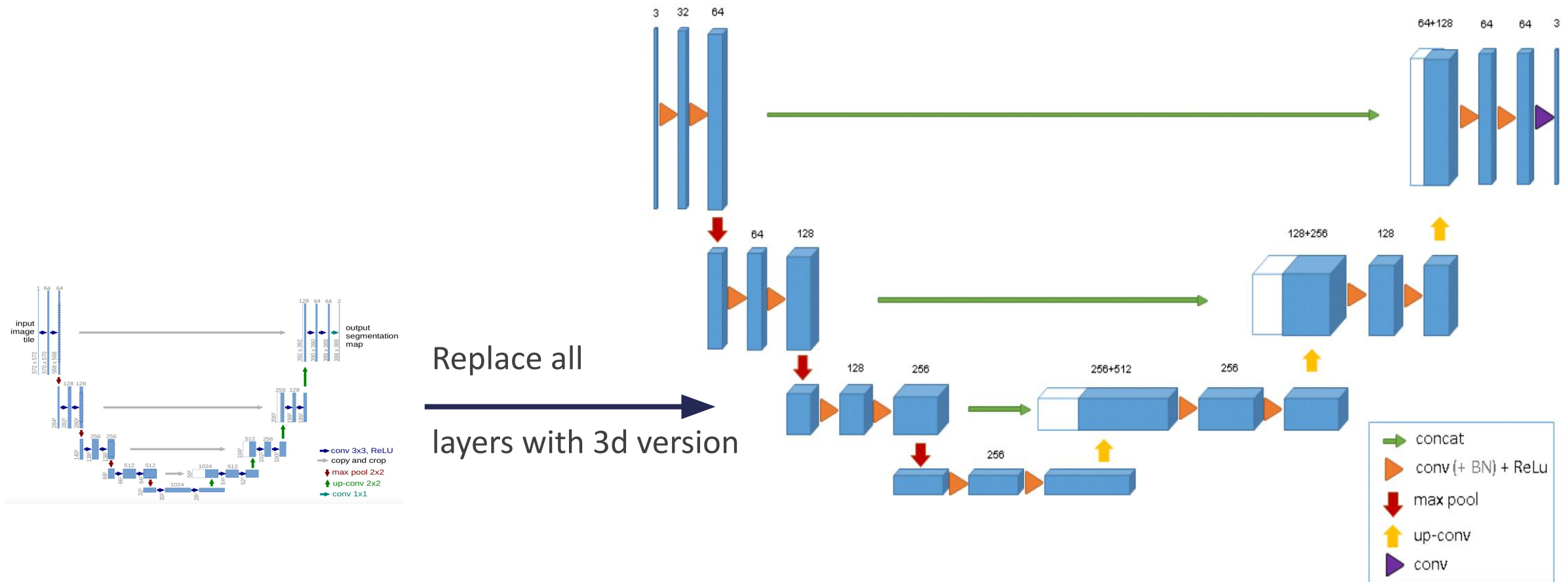
Volumetric data? Use a 3D U-Net!

We have 3d versions for all layers used in U-Net.

- 2d convolution -> 3d convolution
 - 2d kernel $X \times Y$ -> 3d kernel $X \times Y \times Z$
- Needs more computation, but otherwise the same



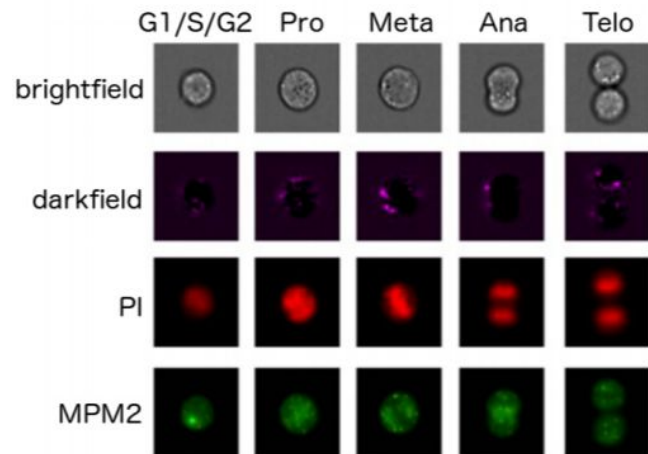
Volumetric data? Use a 3D U-Net!



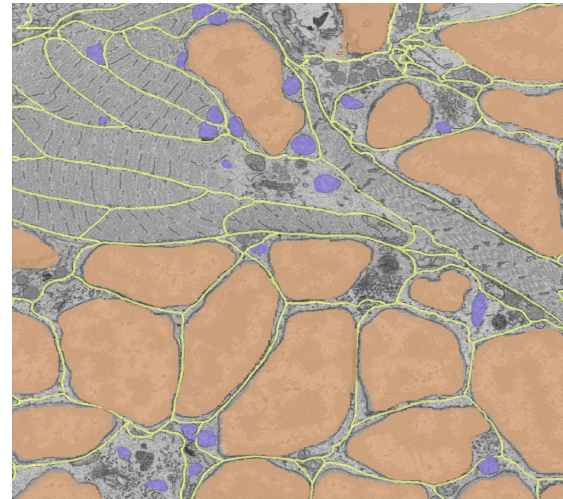
Instance Segmentation

Segmentation for biomedical images

Image
Classification

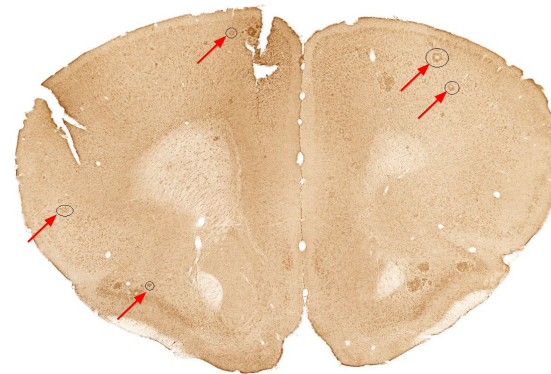


Semantic
Segmentation



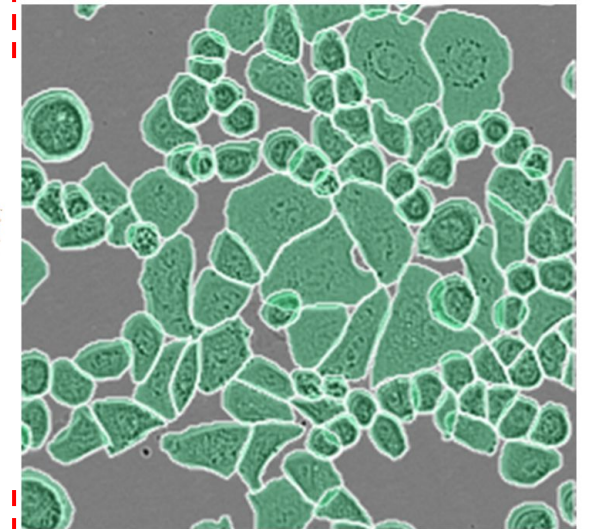
Cell Membrane
Nuclei
Mitochondria

Object
Detection



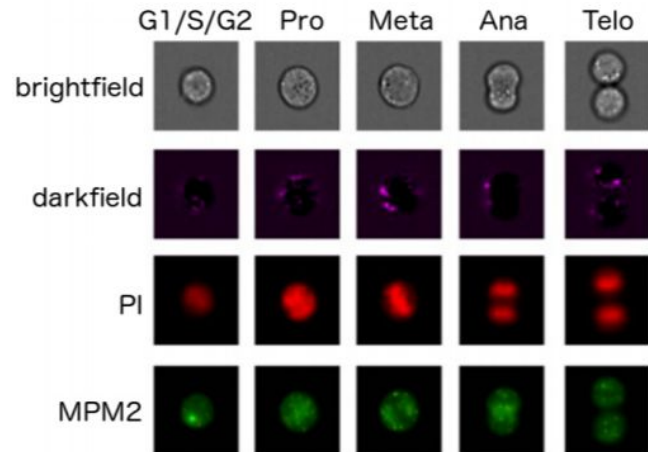
Object detection = find
bounding boxes for each
object -> **we skip this.**

Instance
Segmentation

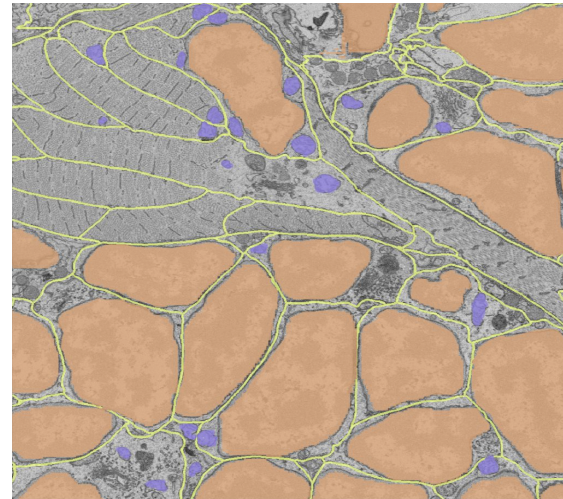


Segmentation for biomedical images

Image
Classification



Semantic
Segmentation

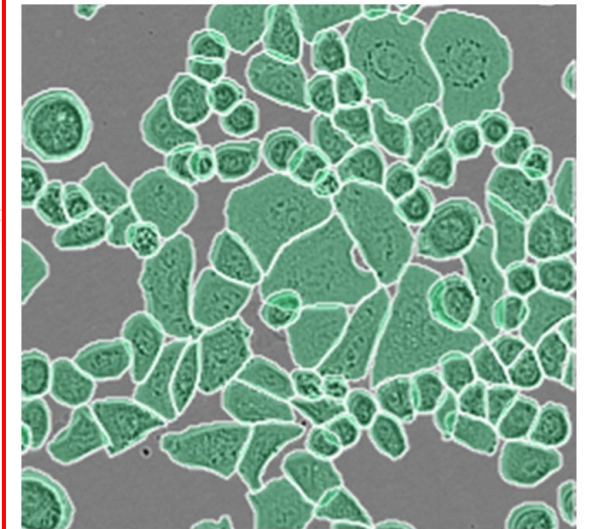


Cell Membrane
Nuclei
Mitochondria

Object
Detection

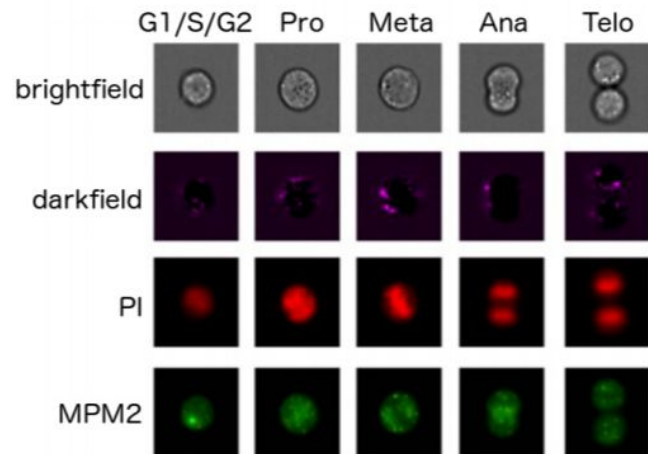


Instance
Segmentation

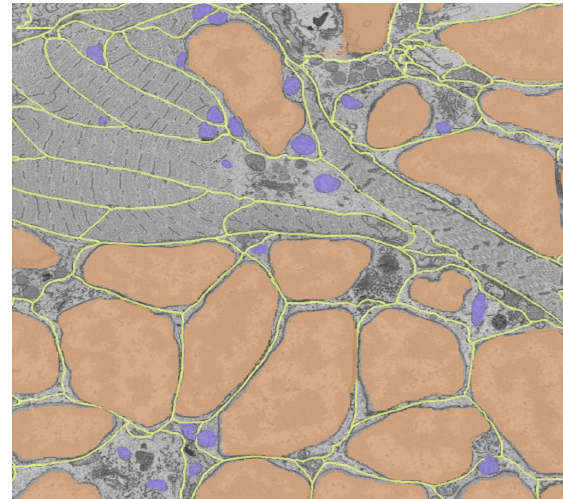


Segmentation for biomedical images

Image
Classification

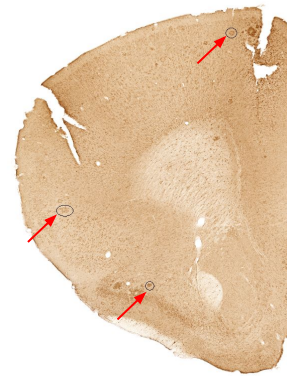


Semantic
Segmentation

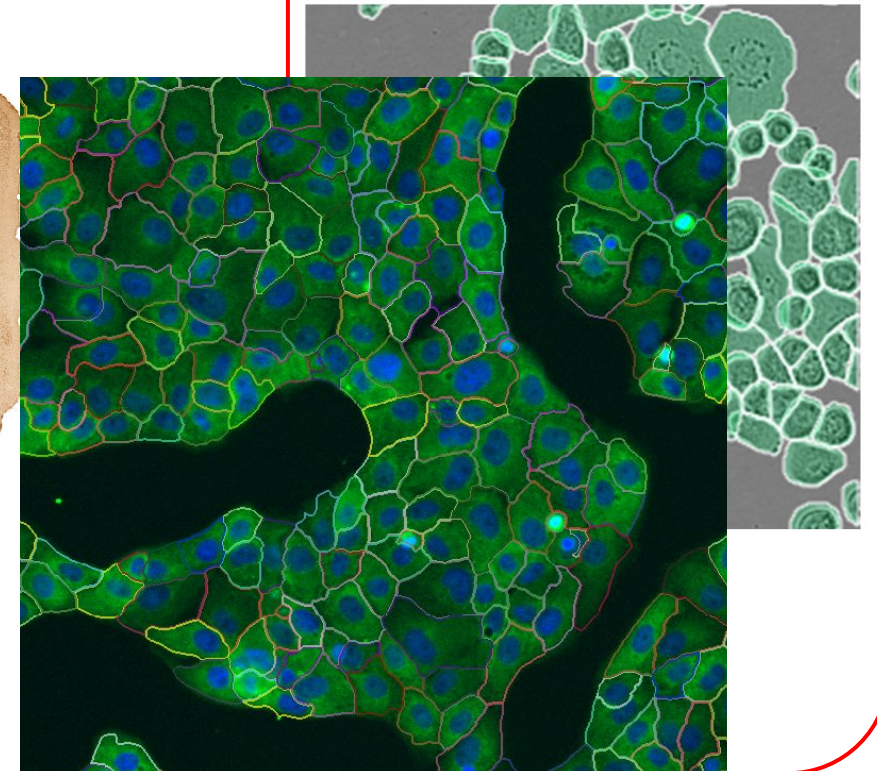


Cell Membrane
Nuclei
Mitochondria

Object
Detection



Instance
Segmentation

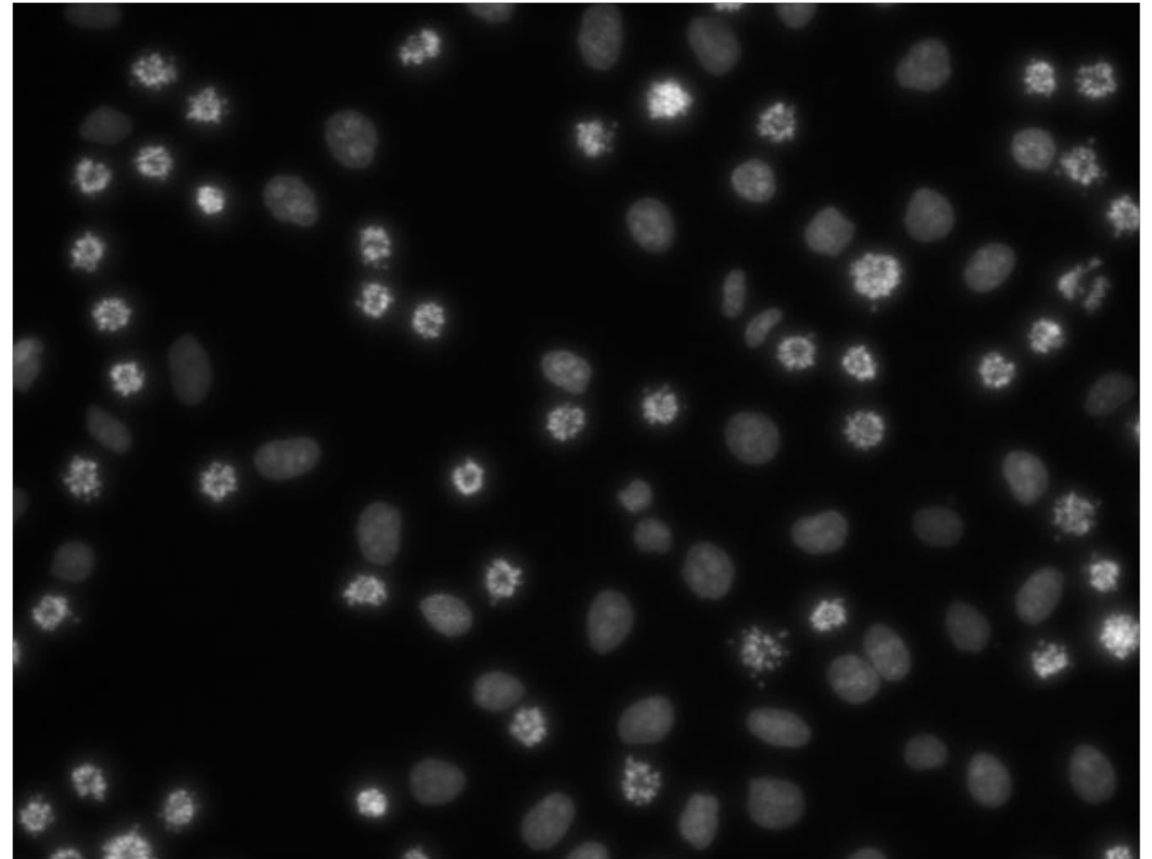


Instance segmentation

Goal: Segment individual nuclei in the image (1 mask per object)

Problem: How many objects are there?
We don't know before solving the task!

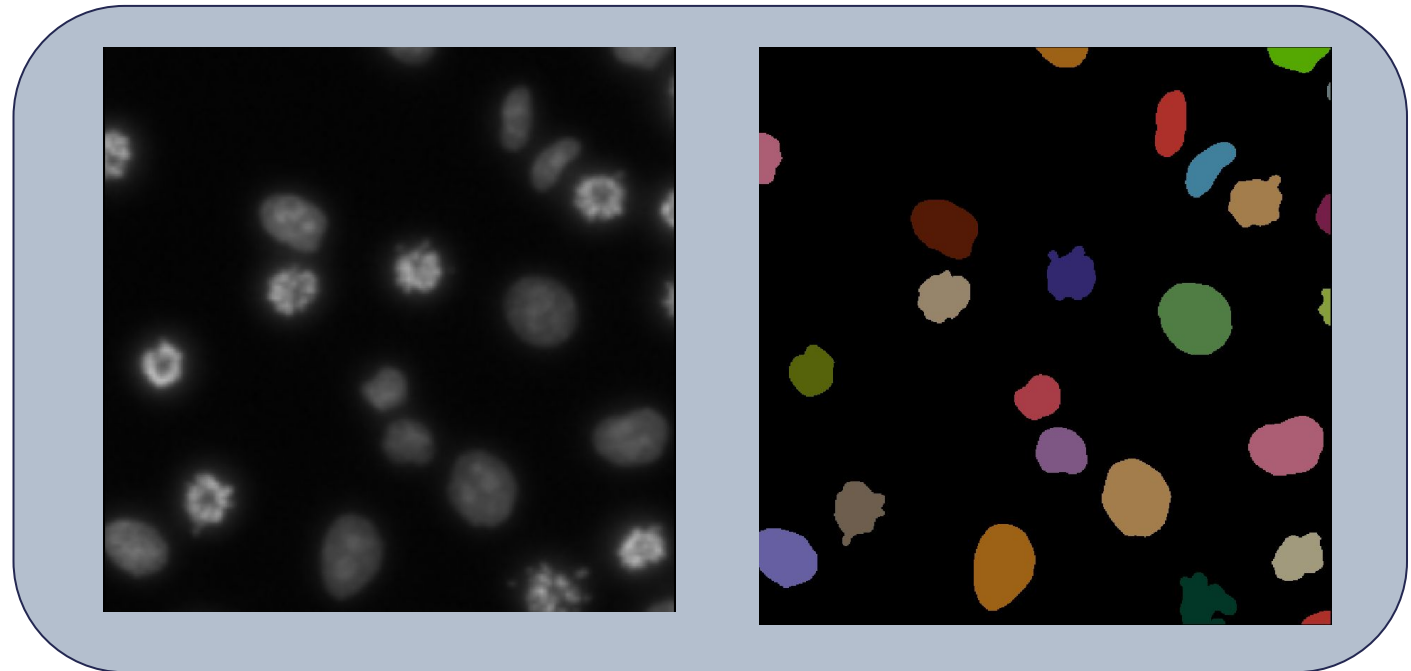
Different in semantic segmentation:
we know how many different categories
to predict in the beginning.



Instance segmentation

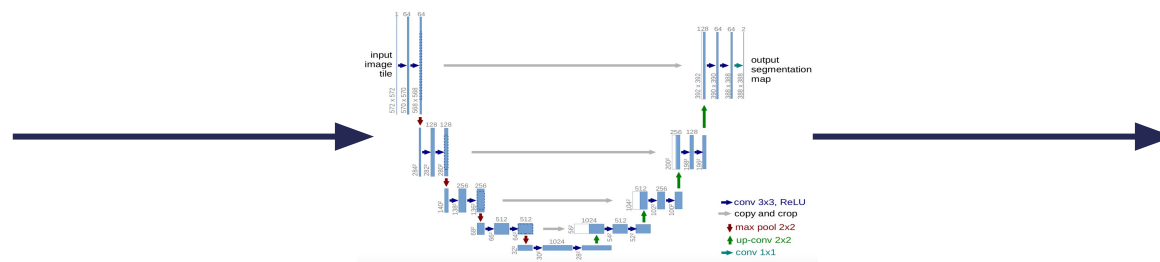
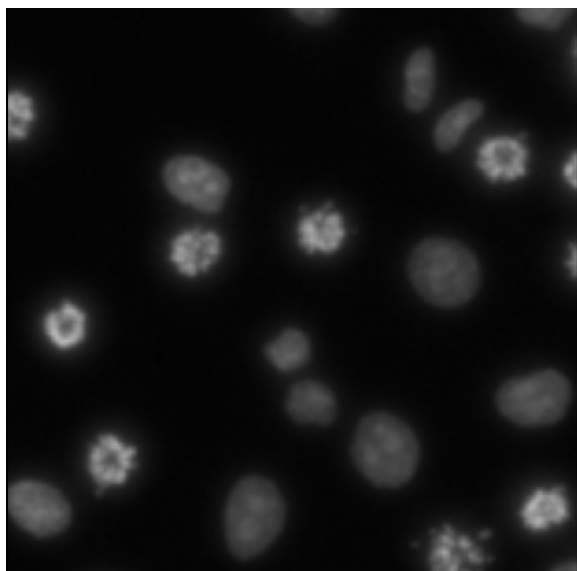
- Each instance has a unique id, all pixels belonging to instance have that id
 - Background is 0
- Visualization: random colors -> each id gets a random color

```
00000000000000000000000000000000
000000001111111111100000
000000011111111111100000
000000011111111222222222
000000000000000000222222
00333333300000000022200
0000333333330000002000
0000003333000000000000
0000000333000000000000
```



Instance segmentation

Predict instances directly?



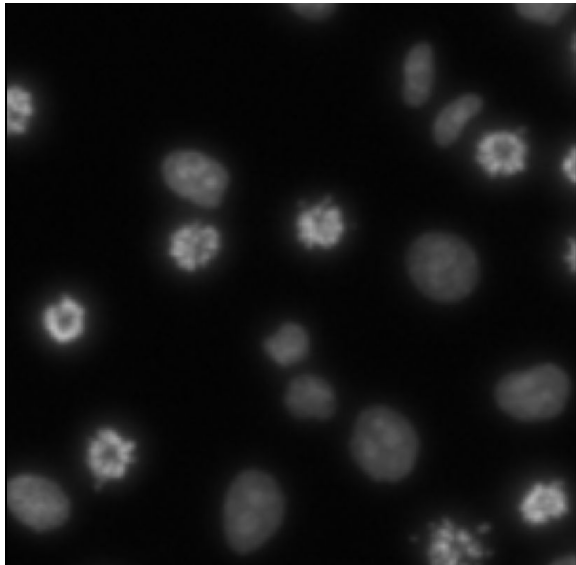
Instance segmentation

Problem: instance segmentation is invariant under relabeling:
Swapping instance labels *1* and *2* does not change the meaning!

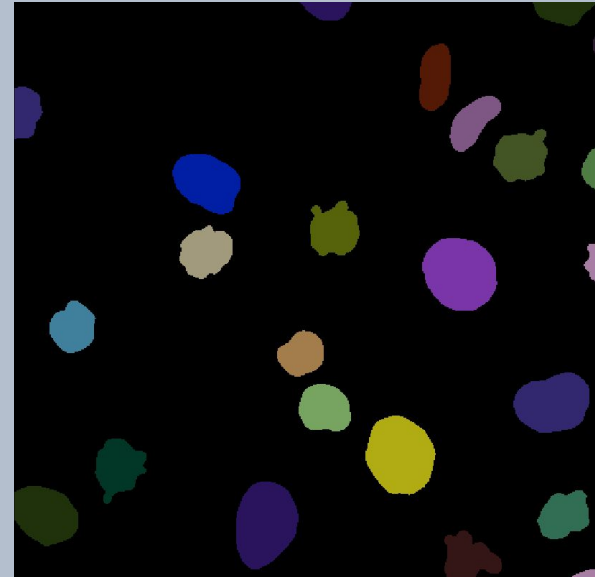
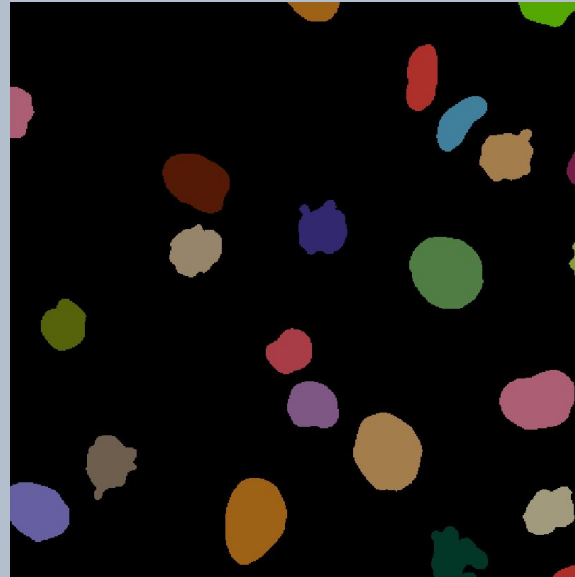
000000000000000000000000		000000000000000000000000
0000000111111111100000		0000000222222222200000
0000001111111111100000		0000002222222222200000
0000001111111112222222	Swap IDs 1, 2	0000002222222221111111
0000000000000000222222	←→	0000000000000000111111
003333330000000022200		003333330000000011100
0000333333330000002000		0000333333330000001000
0000003333000000000000		0000003333000000000000
0000000333000000000000		0000000333000000000000

Instance segmentation

Problem: instance segmentation is invariant under relabeling:
Swapping instance labels 1 and 2 does not change the meaning!



Equivalent instance segmentations!

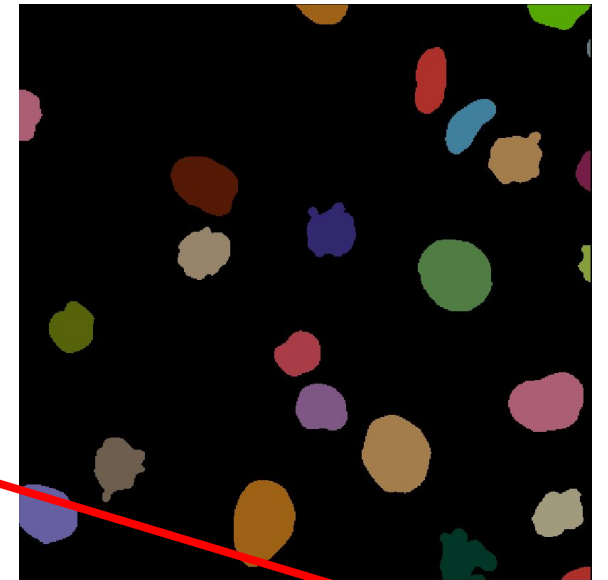
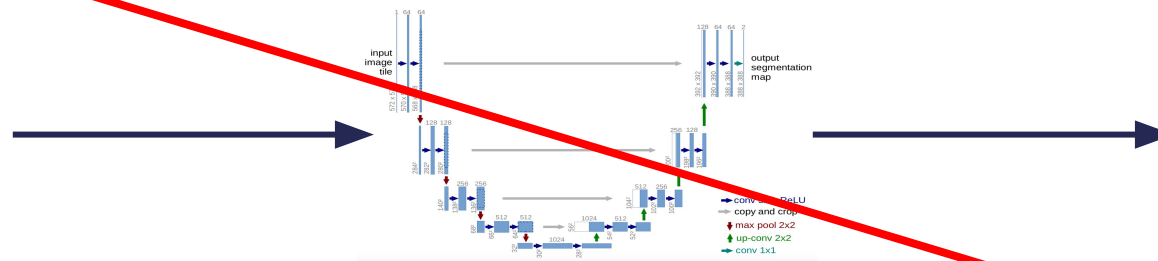
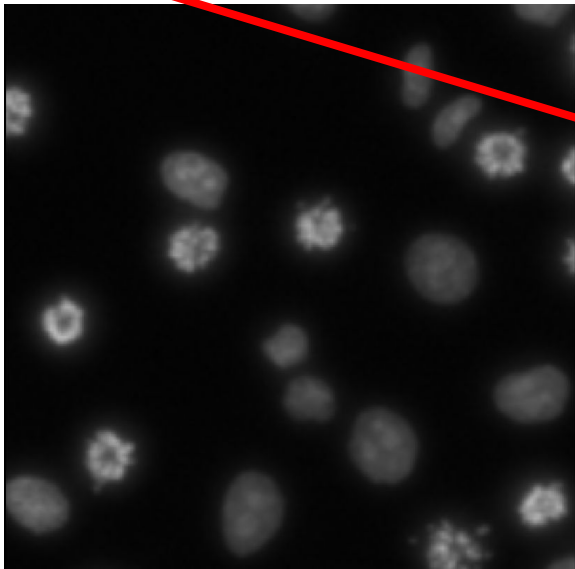


*On-going research:
differentiable instance
segmentation

Instance segmentation

Predict instances directly?

Not possible due to relabeling invariance, can't formulate a differentiable loss*



Instance segmentation

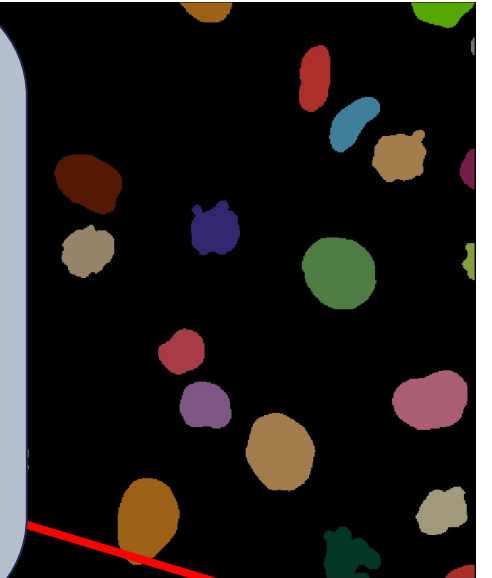
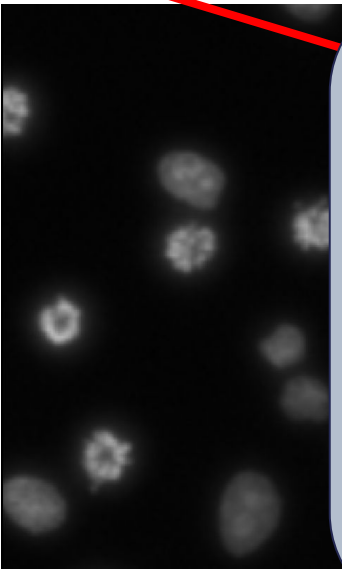
Predict instances directly?

Not possible due to relabeling invariance, can't formulate a differentiable loss

Instead:

- Semantic segmentation with deep learning (or regression)
- Post-processing to get the objects

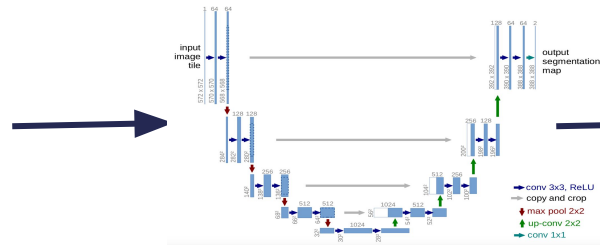
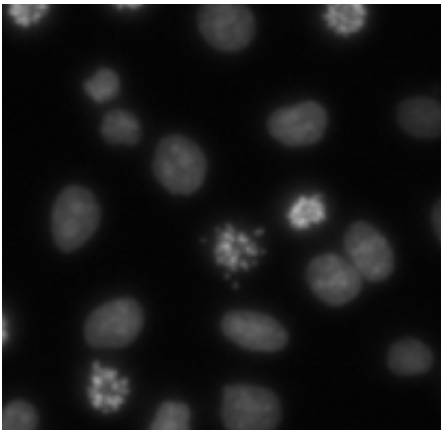
Many different approaches, will discuss three now.



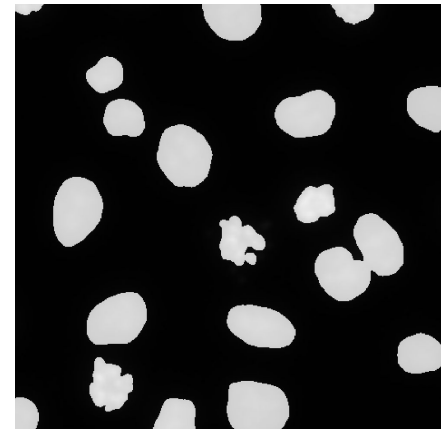
Foreground prediction

Predict foreground (binary segmentation)

Input

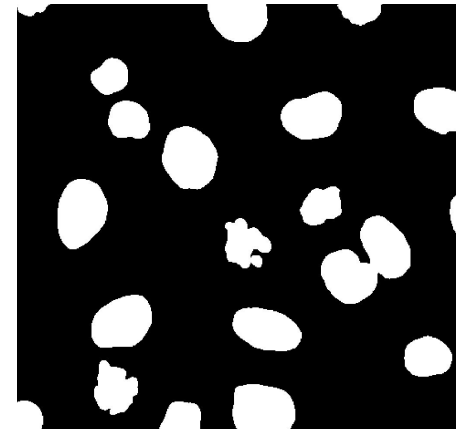


Prediction



Loss

Target

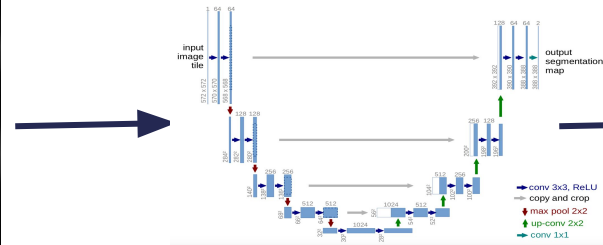
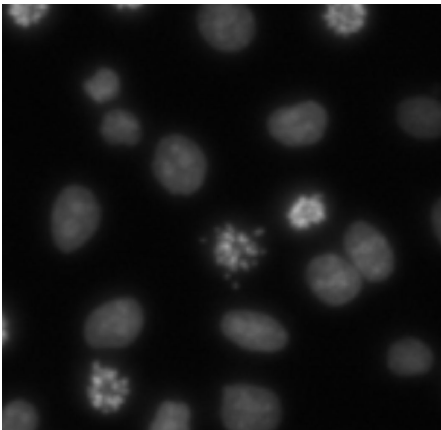


Training

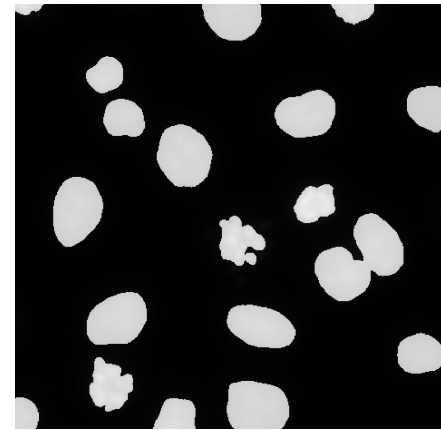
Foreground prediction

Predict foreground (binary segmentation), apply connected components

Input



Prediction



Threshold
Connected
Comp.

Instance Seg.

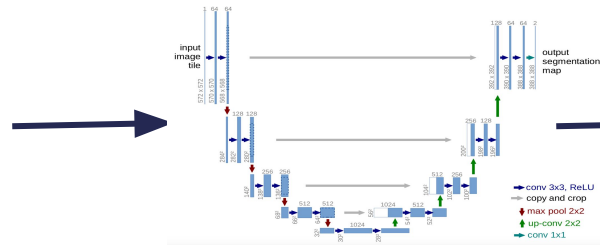
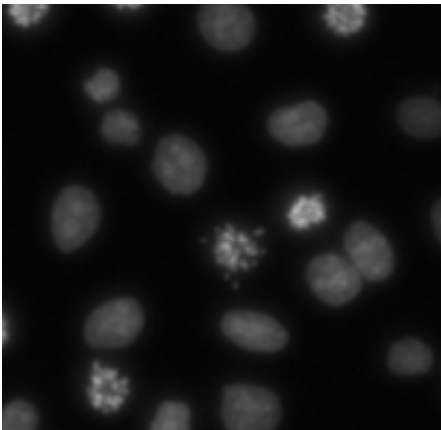


Prediction

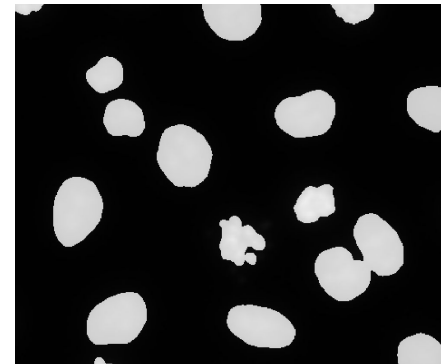
Foreground prediction

Predict foreground (binary segmentation), apply connected components

Input

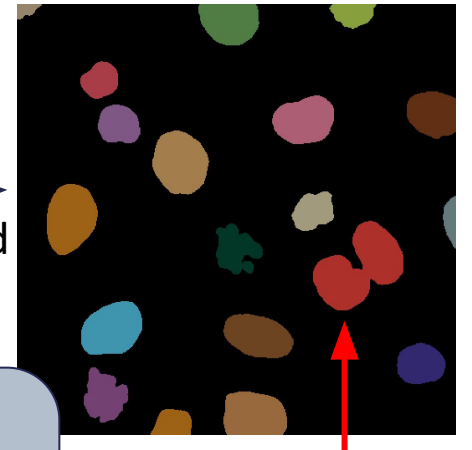


Prediction



Threshold
Connected
Comp.

Instance Seg.

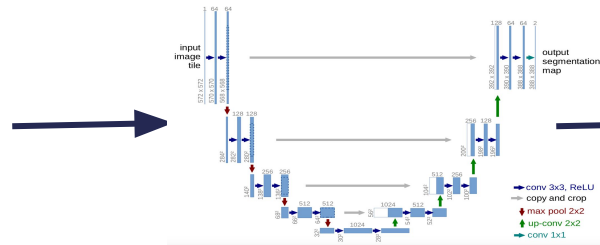
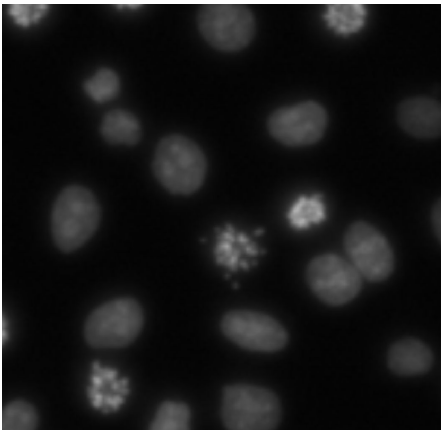


Problem: touching objects!

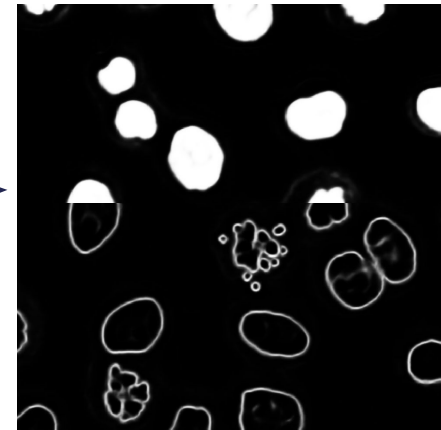
Foreground and boundary prediction

Predict foreground and boundaries

Input



Prediction



Loss

Target

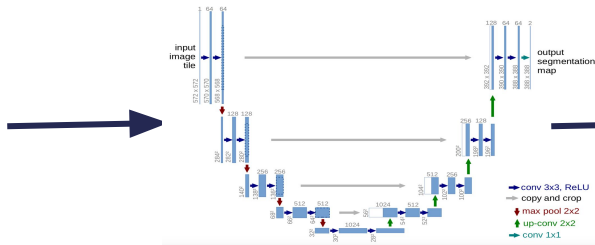
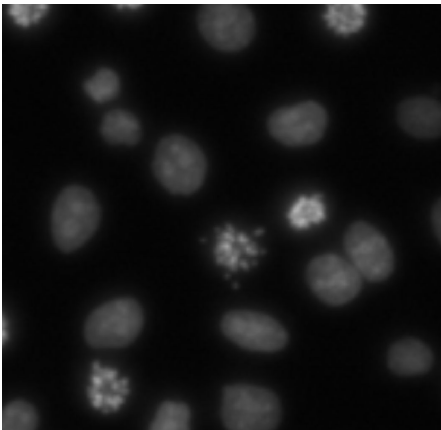


Training

Foreground and boundary prediction

Subtract boundaries from foreground before threshold, apply watershed to get back full size

Input



Prediction



Instance Seg.

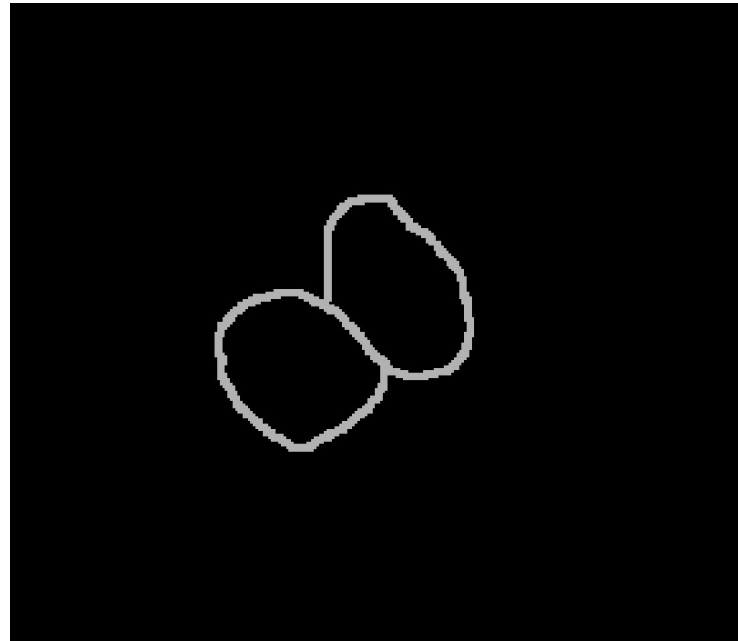
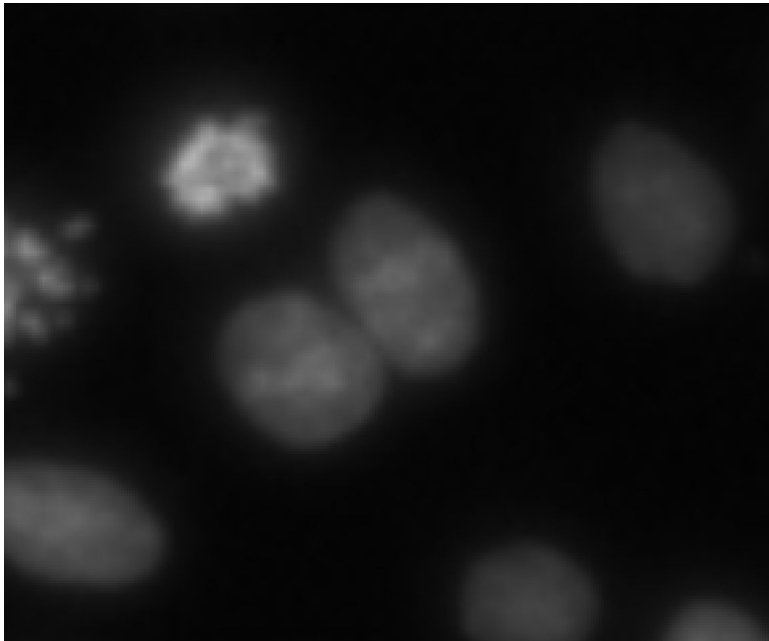


Prediction

Foreground and boundary prediction

So we can just predict foreground and boundaries and solve any segmentation problem?

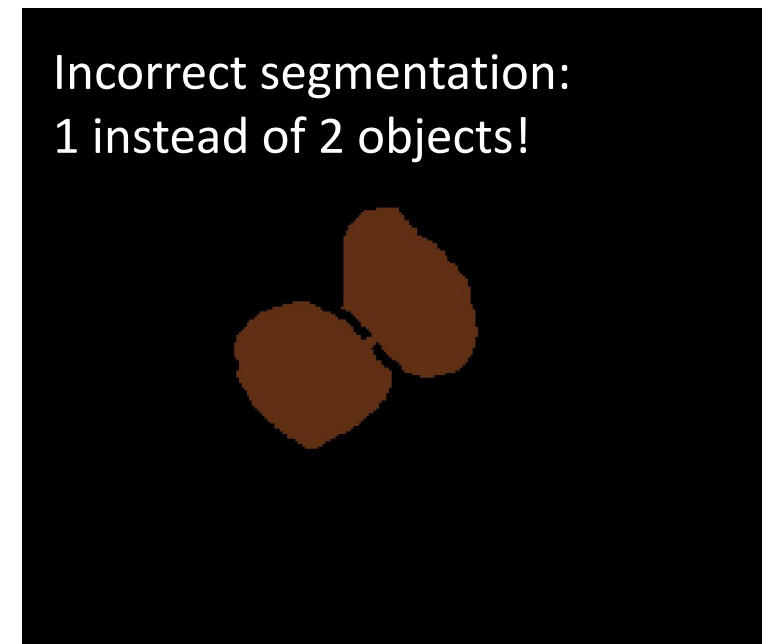
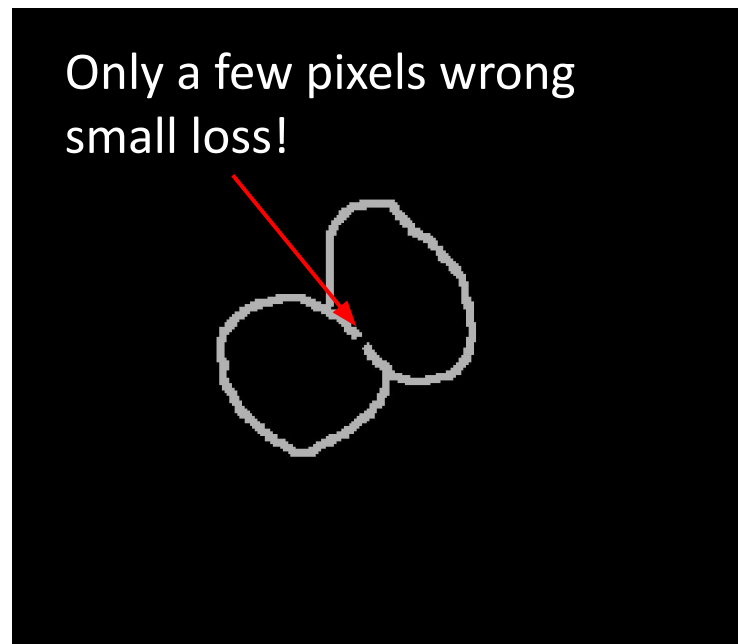
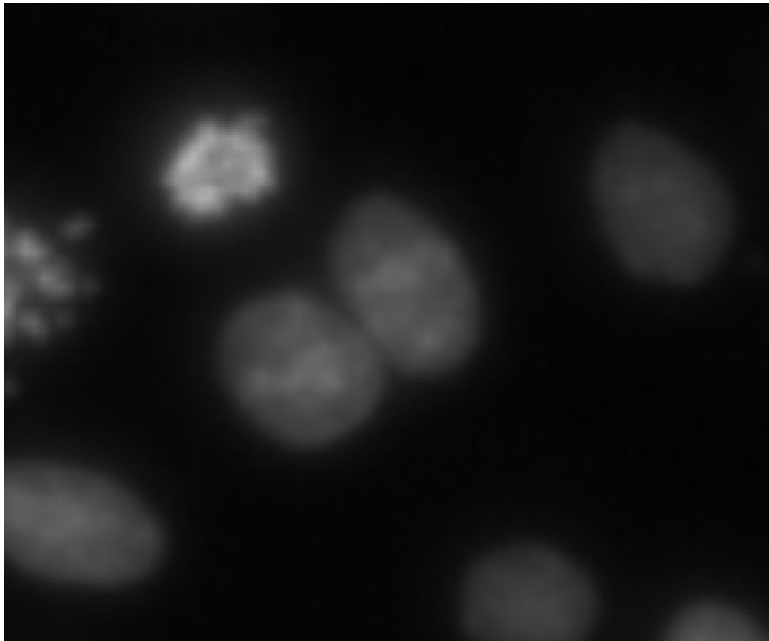
Correct boundary prediction, correct segmentation.



Foreground and boundary prediction

So we can just predict foreground and boundaries and solve any segmentation problem?

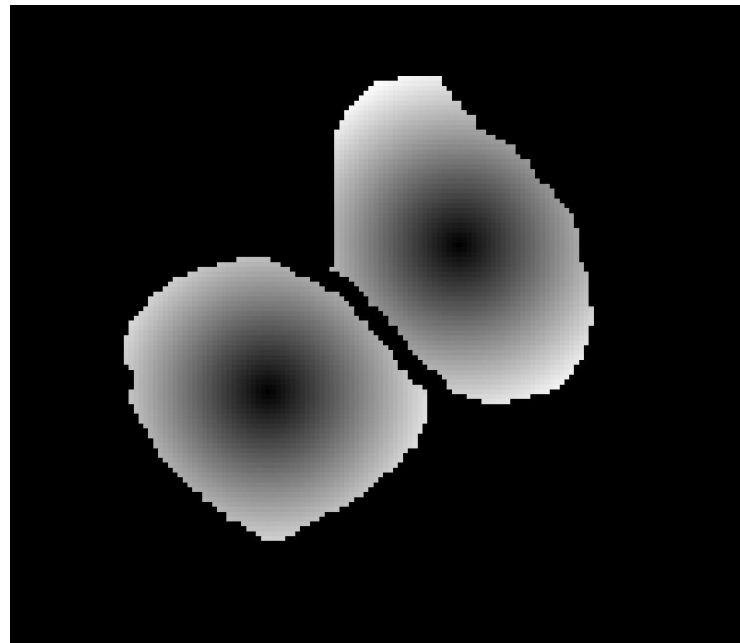
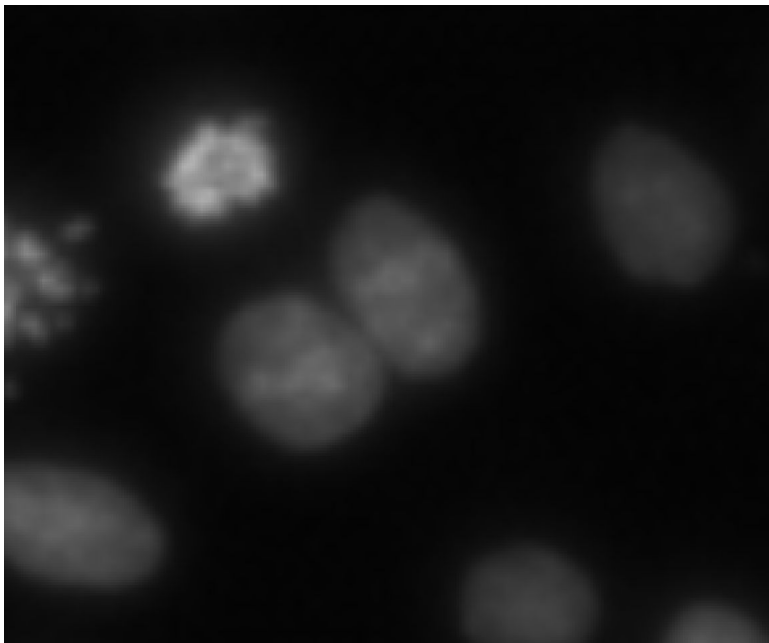
Small error in boundary prediction, very incorrect segmentation.



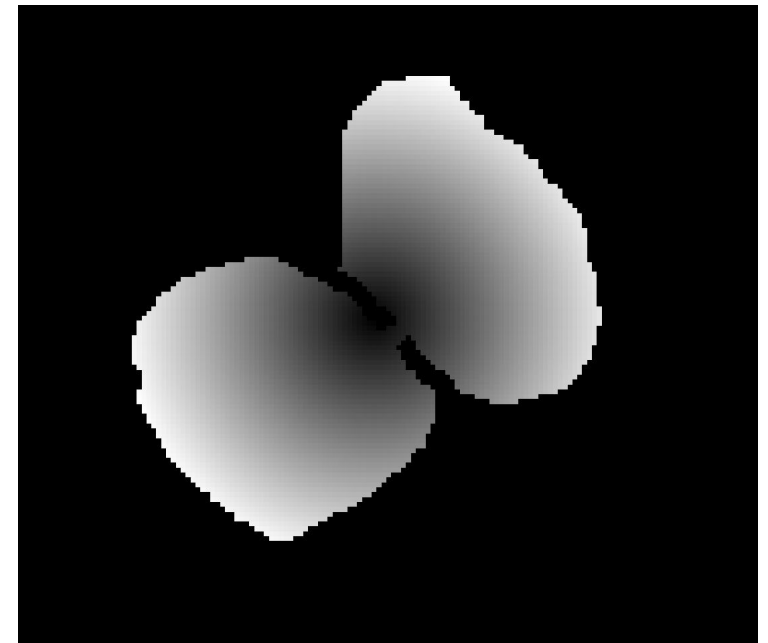
Foreground and boundary prediction

Instead: predict representation with network that changes a lot when segmentation changes.

Example: predict distance to center point.



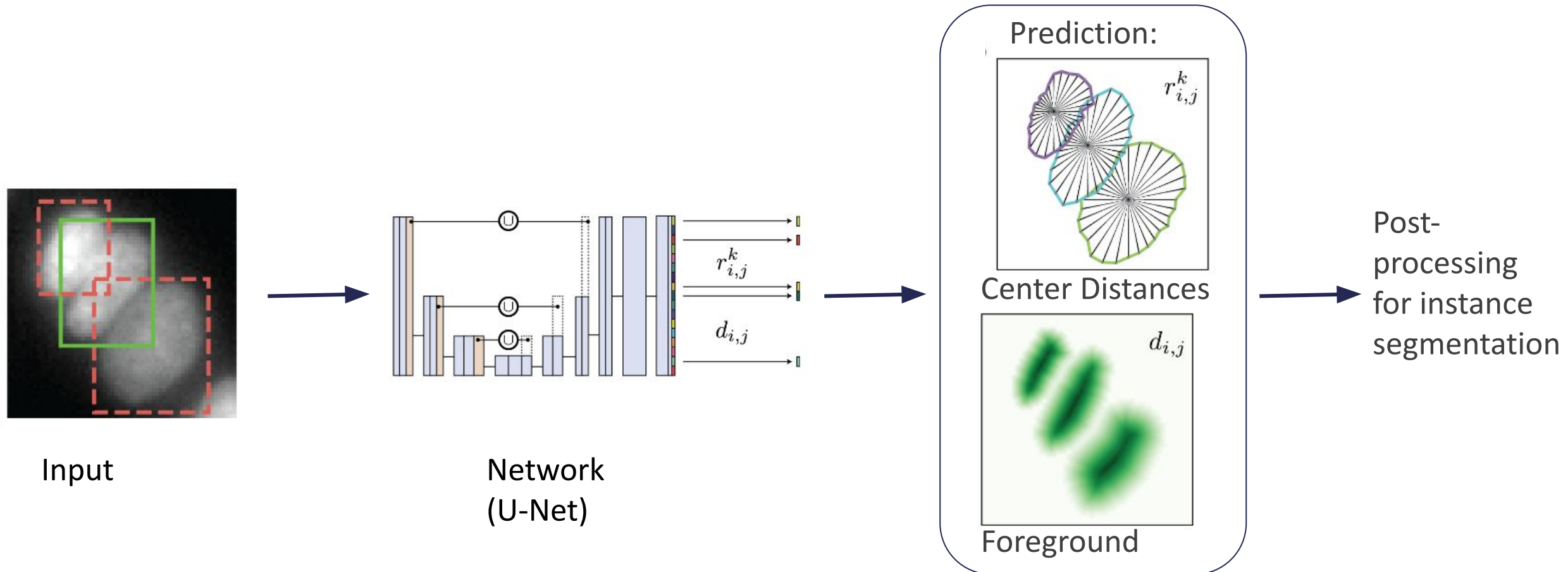
Distance to center (correct)



Distance to center (incorrect)

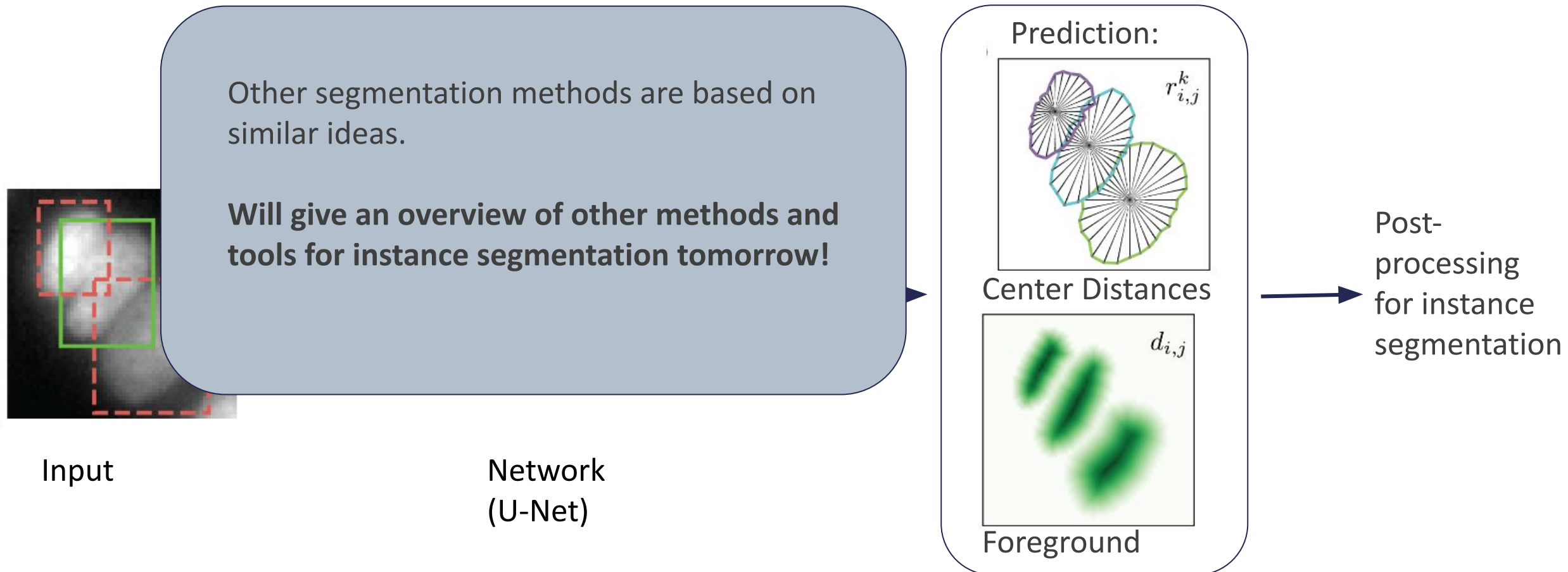
StarDist: Popular for Nucleus Segmentation

Similar idea: predict distance to center along different direction for each pixel (and foreground)



StarDist: Popular for Nucleus Segmentation

Similar idea: predict distance to center along different direction for each pixel (and foreground)



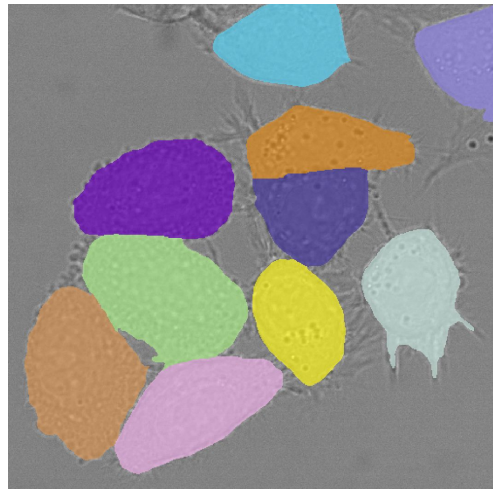
Cell Tracking

Following cells over time

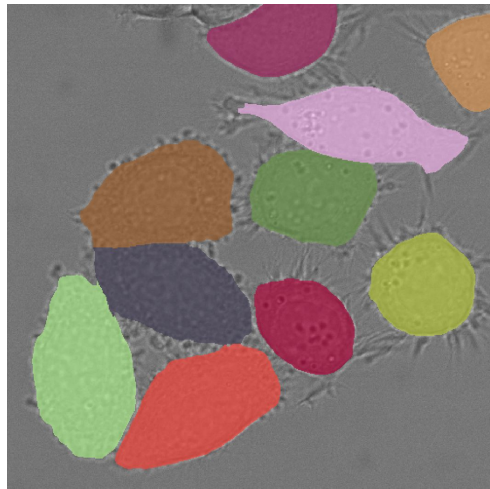
Tracking by detection

Cell tracking: follow cells *and* cell divisions over time to build cell lineage graph.

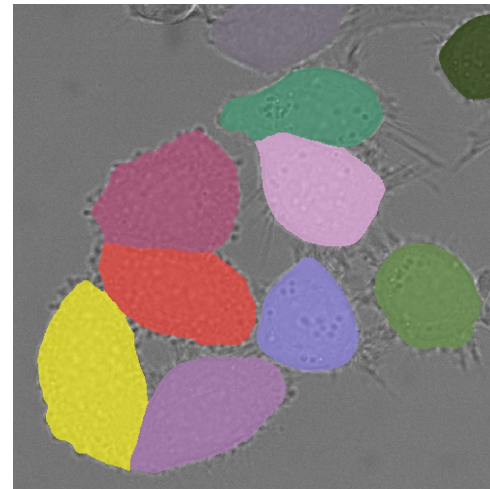
Tracking by detection: first segment individual cells per frame



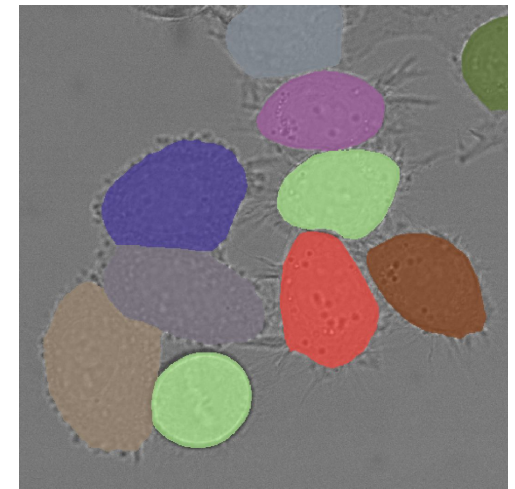
$t = 0$



$t = 1$



$t = 2$

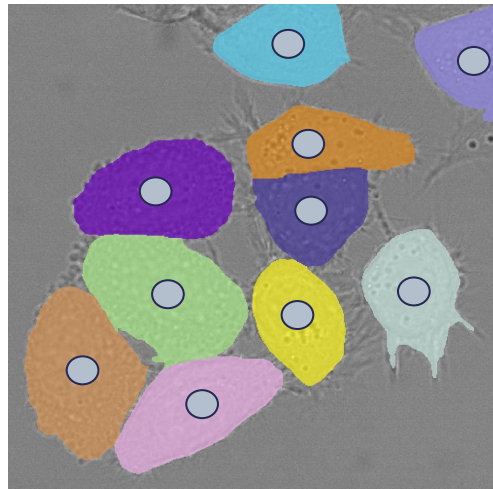


$t = 3$

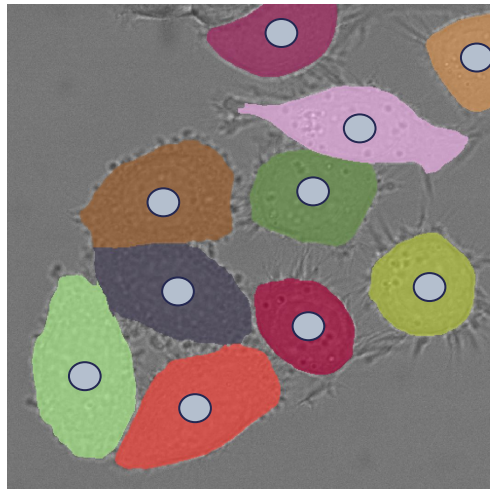
Tracking by detection

Cell tracking: follow cells *and* cell divisions over time to build cell lineage graph.

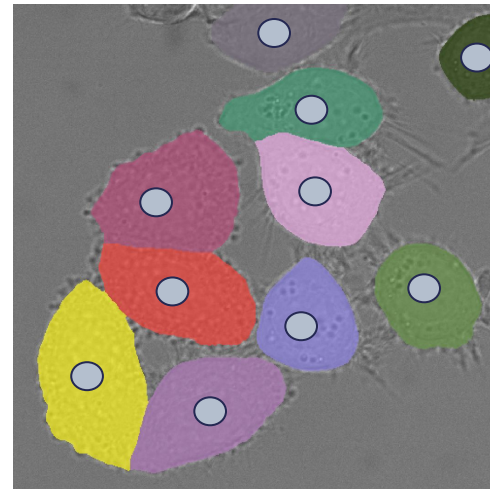
Tracking by detection: build a graph with cells as nodes possible connections (= same cell) across frames as edges



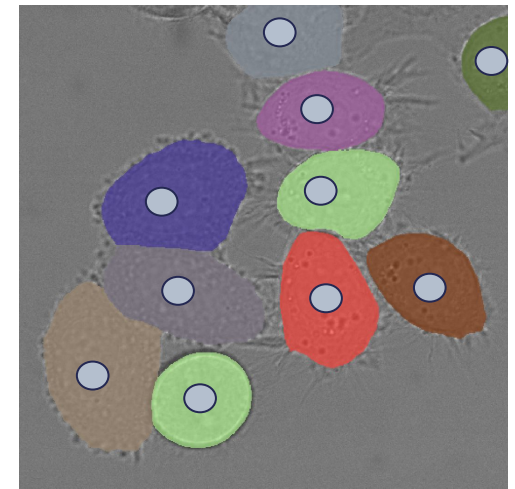
$t = 0$



$t = 1$



$t = 2$

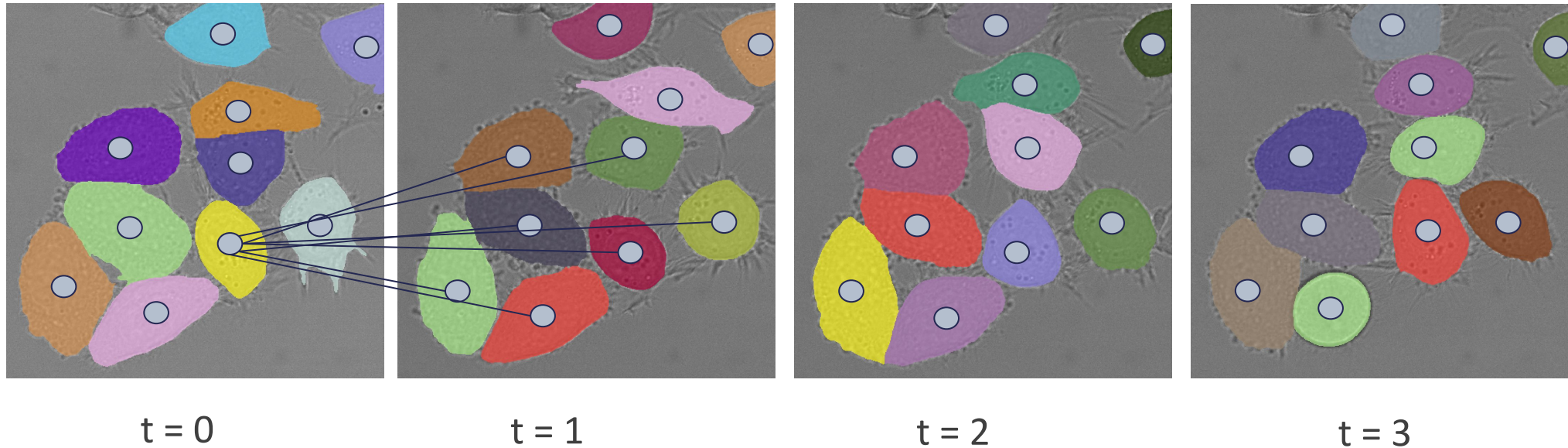


$t = 3$

Tracking by detection

Cell tracking: follow cells *and* cell divisions over time to build cell lineage graph.

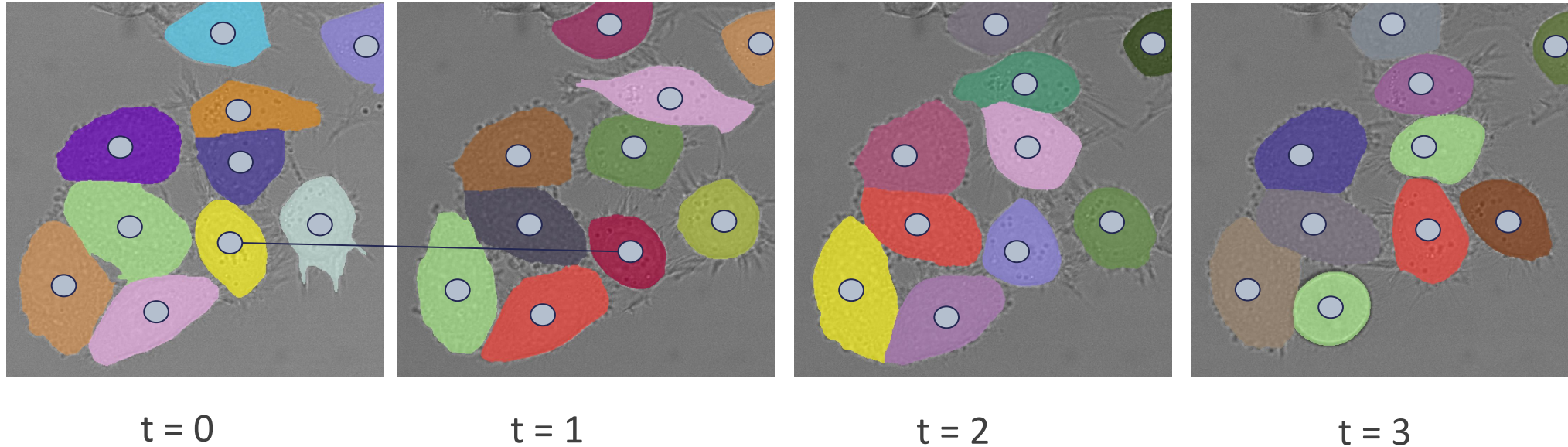
Tracking by detection: build a graph with cells as nodes possible connections (= same cell) across frames as edges



Tracking by detection

Cell tracking: follow cells *and* cell divisions over time to build cell lineage graph.

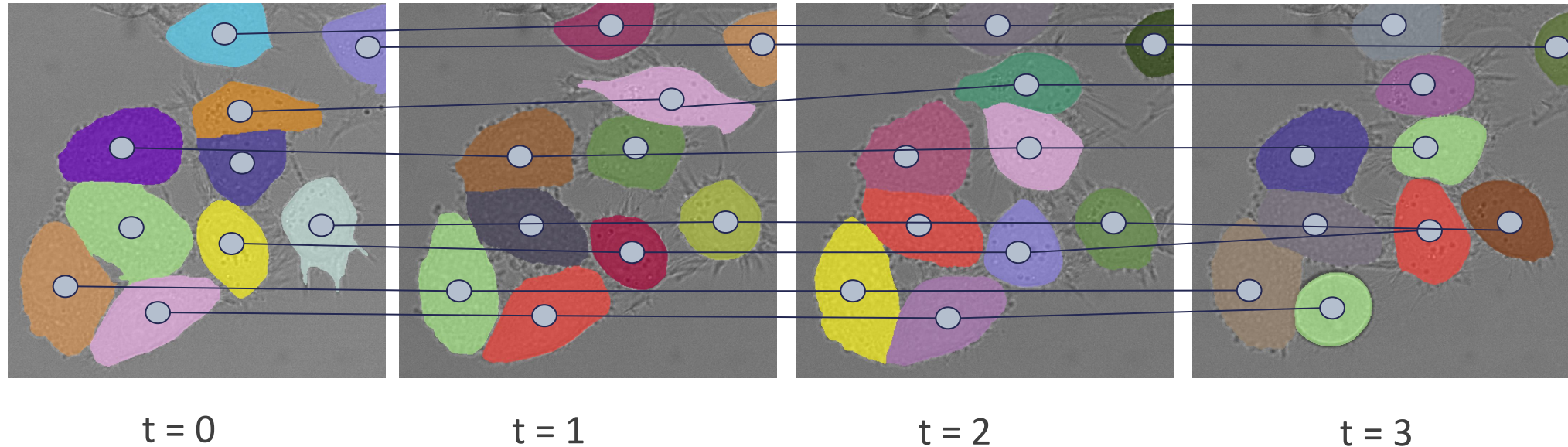
Tracking by detection: select the most likely edges / connections given the constraints that a cell may only link to one or two cells (cell division) in the next frame



Tracking by detection

Cell tracking: follow cells *and* cell divisions over time to build cell lineage graph.

Tracking by detection: select the most likely edges / connections given the constraints that a cell may only link to one or two cells (cell division) in the next frame **for all edges** (shown for a subset here)



layer controls

color by: track_id

colormap: turbo

blending: additive

opacity: 1.00

tail width: [slider]

tail length: [slider]

head length: [slider]

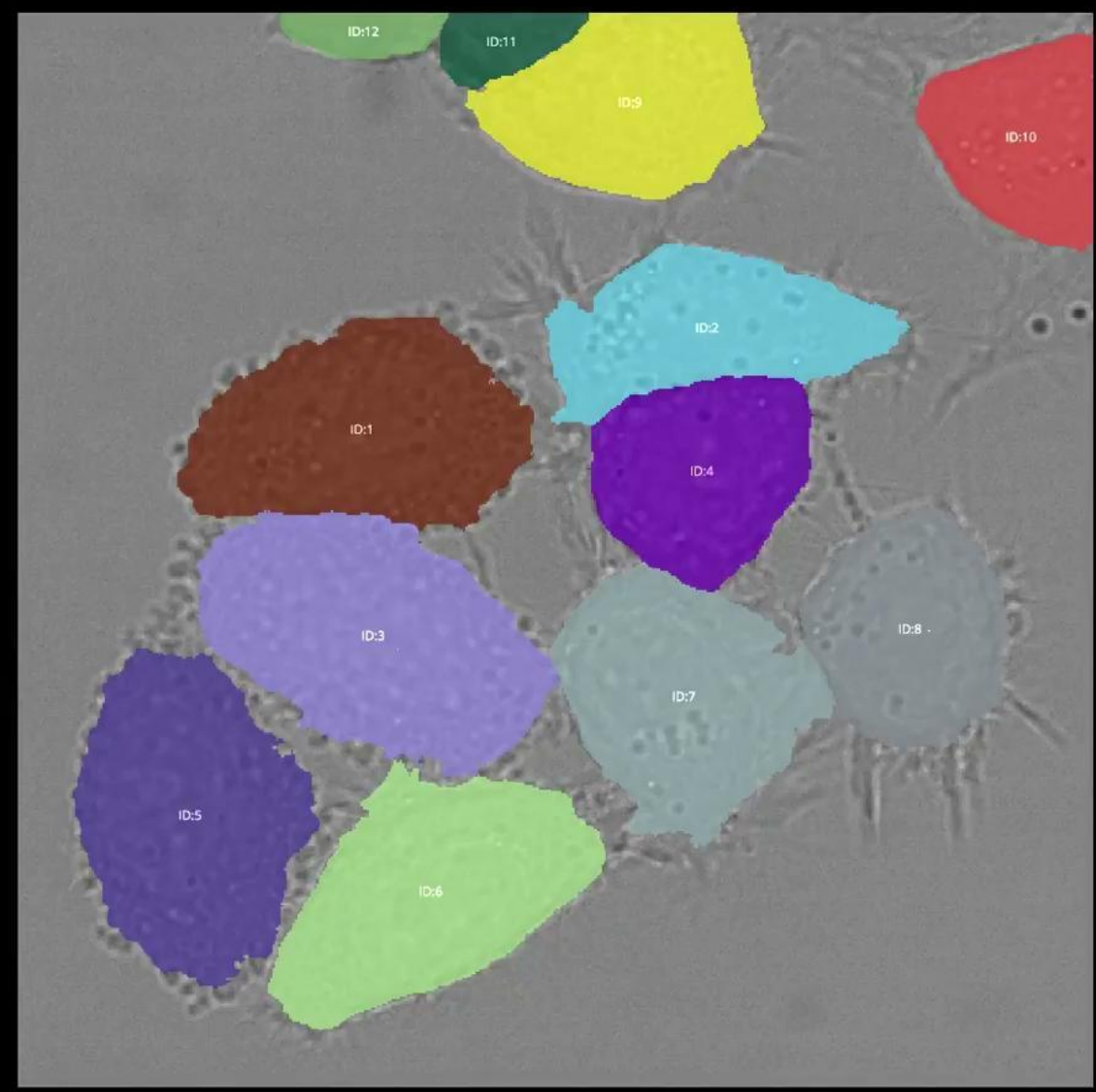
tail: ☒

show ID: ☒

graph: ☐

layer list

- Tracks
- lineage_segmentation
- segmentation
- images



Tracking by detection – How?

“Classical” approach: Select edges according to:

- Largest mask overlap and/or smallest distance
- + motion model
- and tracking constraints

Implemented in:

- **TrackMate** (Fiji Plugin, <https://www.nature.com/articles/s41592-022-01507-1>)
- btrack (Python / napari plugin, <https://github.com/quantumjot/btrack>)
- motile (Python / napari plugin, <https://github.com/funkelab/motile>)

Tracking by detection – How?

Deep learning-based approach: Learn the likelihood of edge selection.

State-of-the-art: **Trackastra**

- Transformer-based architecture for cell tracking
- Pre-trained models can be applied directly to tracking problems
- Python / napari-plugin
- <https://arxiv.org/abs/2405.15700>

Pretrained models & Transfer learning

Re-using models

Training networks for a new task:

- requires a significant amount of annotated data (large manual effort)
- computational skills to train the network

Pre-trained networks: use network already trained on similar data!

Re-using models

Training networks for a new task:

- requires a significant amount of annotated data (large manual effort)
- computational skills to train the network

Pre-trained networks: use network already trained on similar data!

Tools like StarDist / CellPose / μ SAM (more on this tomorrow):

- Provide networks for specific tasks (Nucleus / Cell Segmentation)
- What if my problem does not fit?

Re-using models

Training networks for a new task:

- requires a significant amount of annotated data (large manual effort)
- computational skills to train the network

Pre-trained networks: use network already trained on similar data!

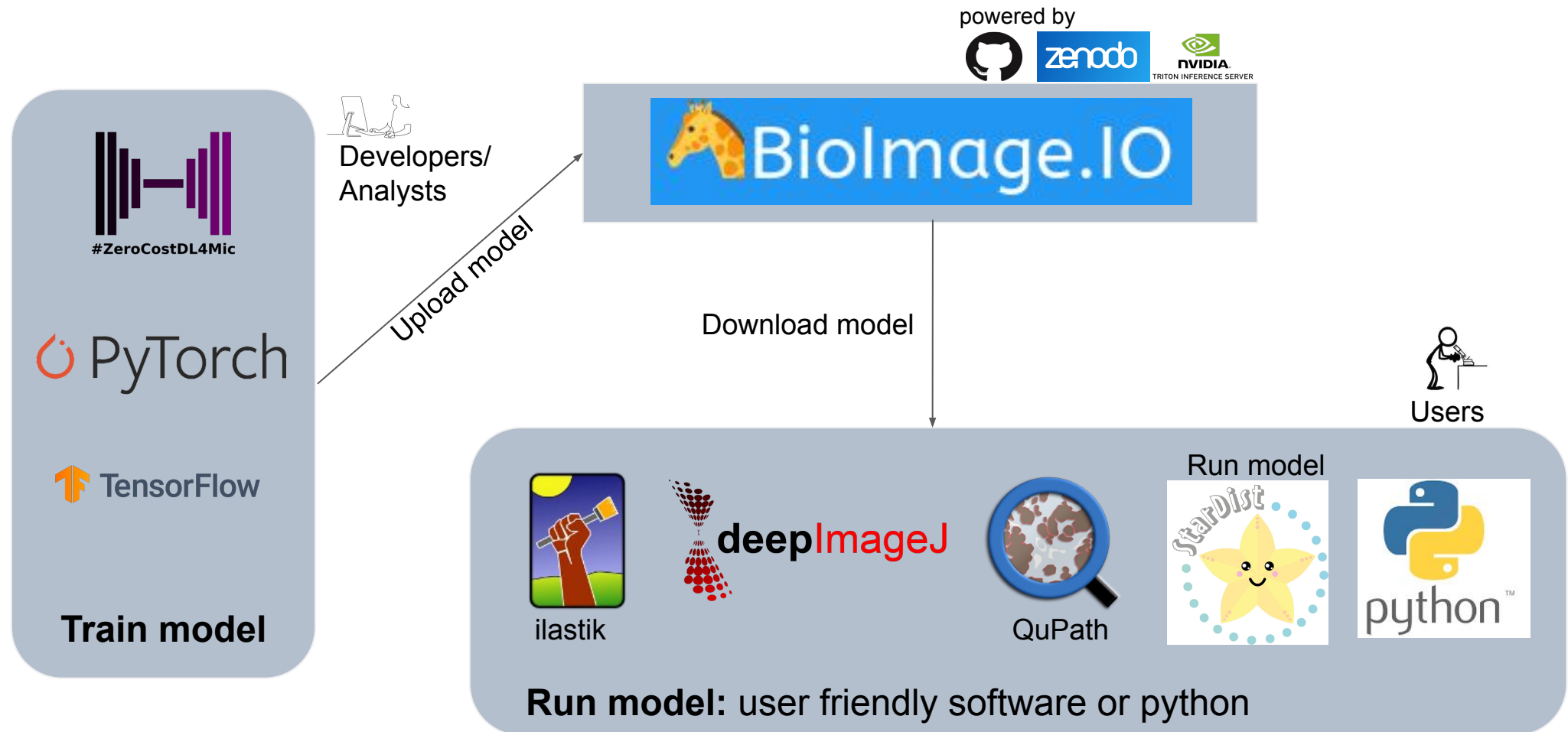
Tools like StarDist / CellPose / μ SAM (more on this tomorrow):

- Provide networks for specific tasks (Nucleus / Cell Segmentation)
- What if my problem does not fit?

BioImage.IO: Ressource for sharing pre-trained models for bio-image analysis.

Biolmage.IO

<https://bioimage.io/#/models>



Limits of pre-trained models

Pre-trained models often not good enough because of differences in training data!

Solution: Fine-tuning

- Annotate some of your data
- Update the model by training on that data

Transfer learning!

Transfer learning: the idea

1

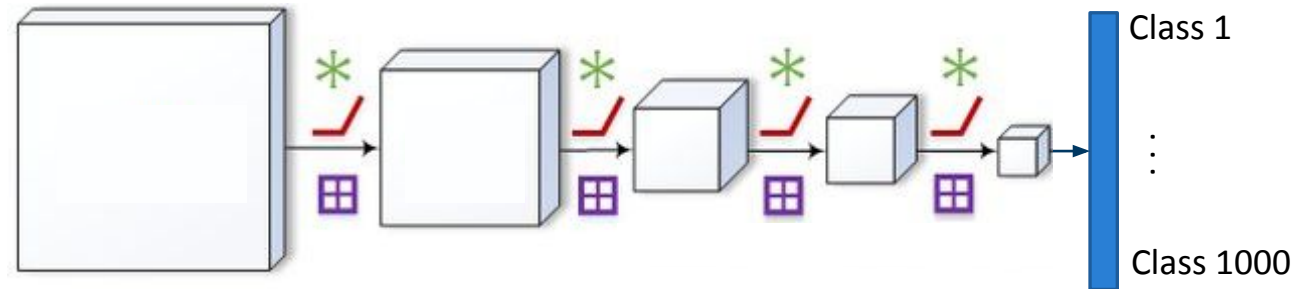
Pre-training

(massive data)

IMAGENET

1 million images

1000 classes



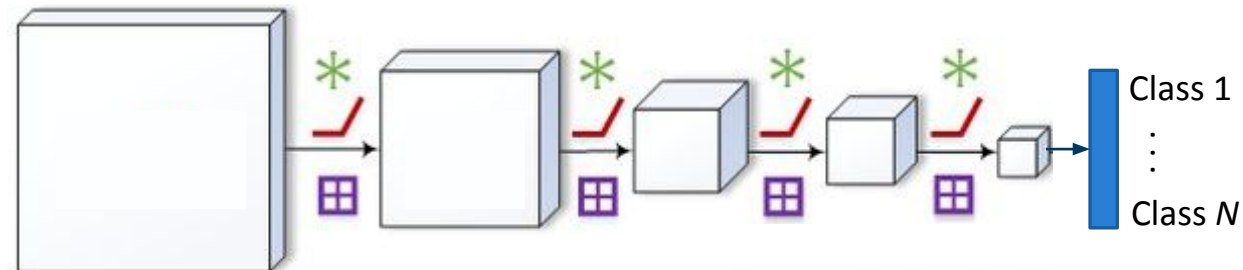
Transfer weights

Replace classifier

2

Fine-tuning

(much less data)



Recognize flowers

Transfer learning: the idea

1

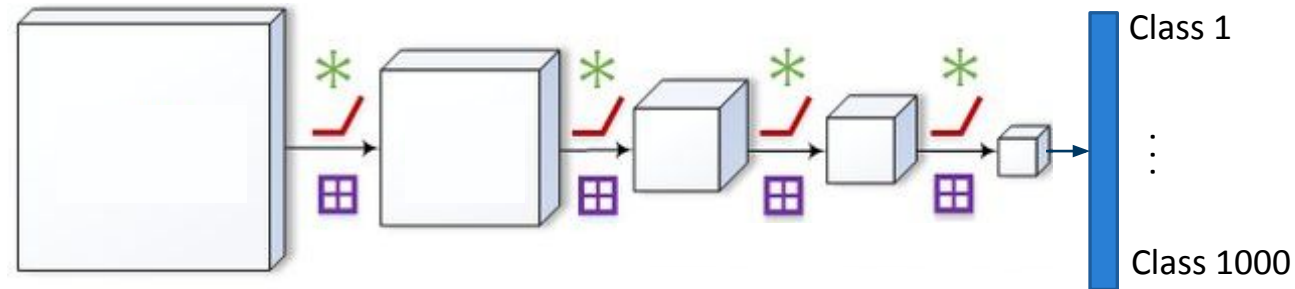
Pre-training

(massive data)

IMAGENET

1 million images

1000 classes



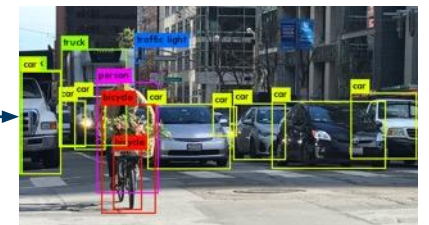
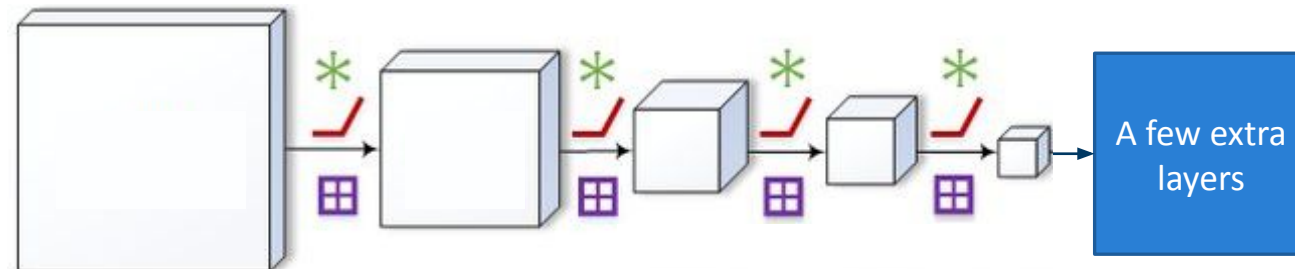
IMAGENET

Transfer weights

2

Fine-tuning

(much less data)

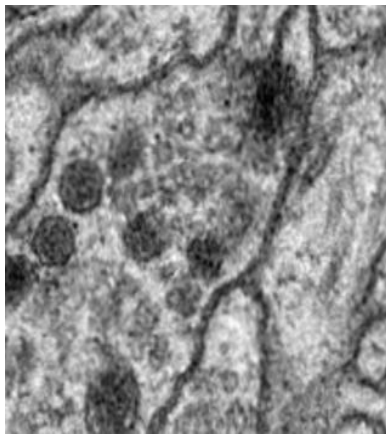


Object detection

Transfer learning for microscopy

But: natural images are quite different from microscopy data:

Benefits of ImageNet pretraining is limited*

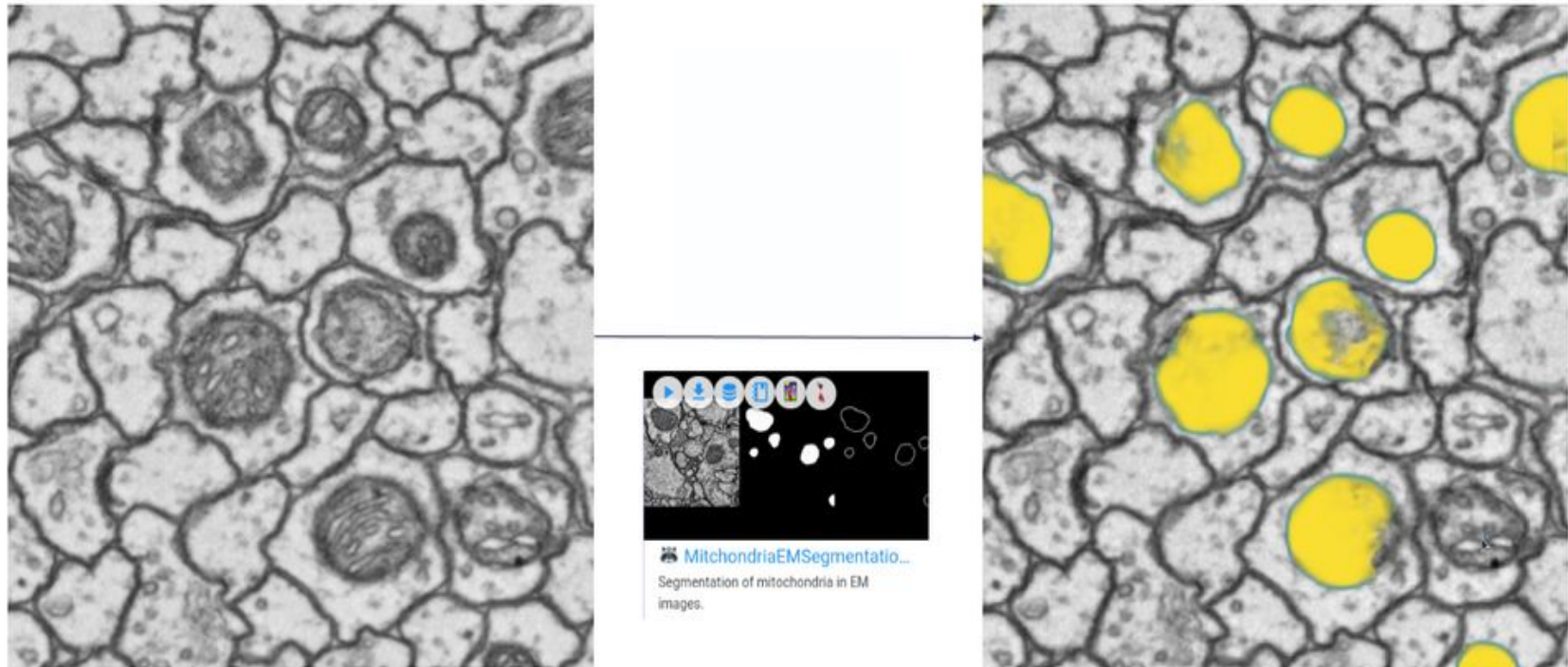


* not true anymore for newer, even bigger and more diverse pre-training datasets



Transfer learning for microscopy

- Find the best model for your data from BiolImage.IO
- Apply to your data, correct result
- Fine-tune the model on the corrected results



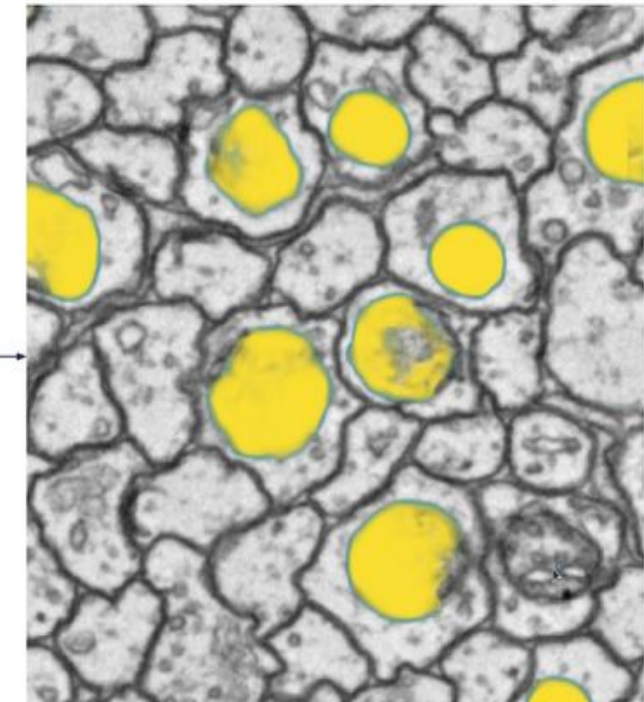
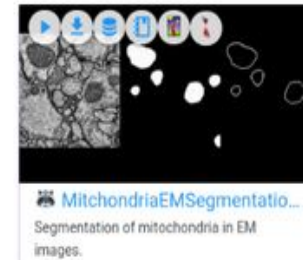
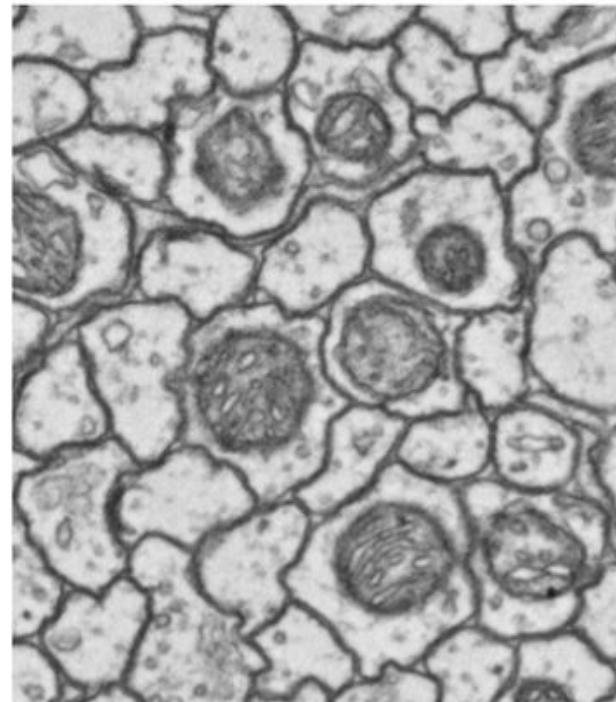


Transfer learning for microscopy

- Find the best model for your data from BiolImage.IO
- Apply to your data, correct result
- Fine-tune the model on the corrected results

Limitations:

- No unified tool for fine-tuning
- Different architectures and frameworks (PyTorch & Tensorflow)
- **Custom python code depending on the model needed**



Deep Learning for Image Restoration

Denoising & Other Tasks

Image-to-image tasks

Many analysis tasks: transform image to a different version:

- May change spatial dimension: $\mathbf{H} \times \mathbf{W} \rightarrow \mathbf{H}' \times \mathbf{W}'$
- May change channels: $\mathbf{C} \rightarrow \mathbf{C}'$

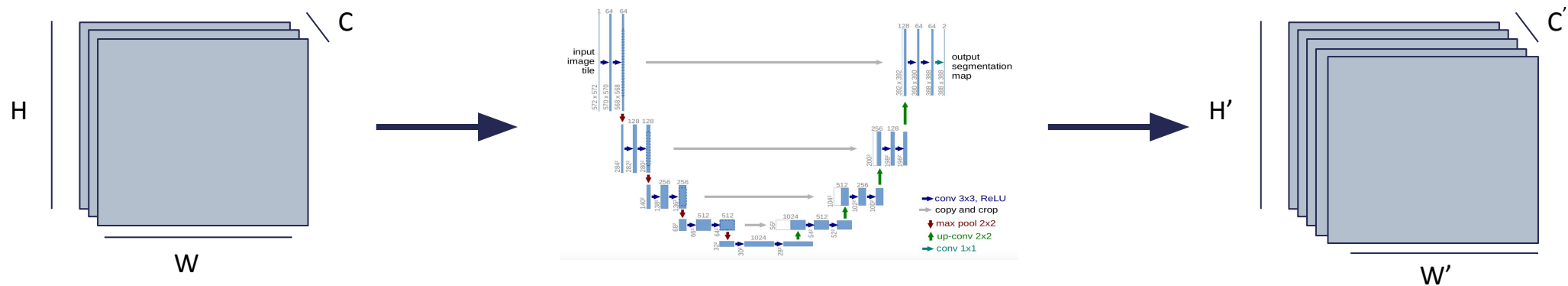


Image-to-image tasks

Many analysis tasks: transform image to a different version:

- May change spatial dimension: $H \times W \rightarrow H' \times W'$
- May change channels: $C \rightarrow C'$

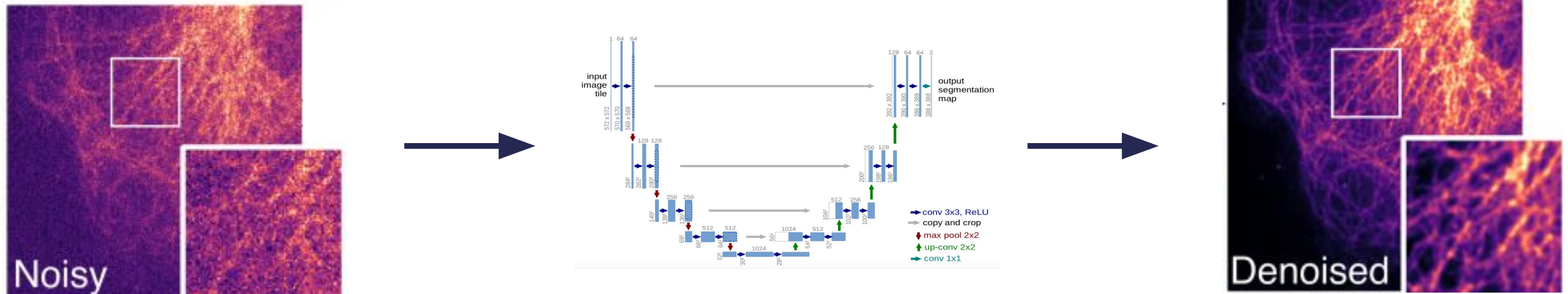
We can use a U-Net for this!



Denoising

Most common image restoration task: Denoising

- Spatial dimensions and channels stay the same
- Improve signal in the image



Denoising: CARE

<https://www.biorxiv.org/content/10.1101/236463v5>

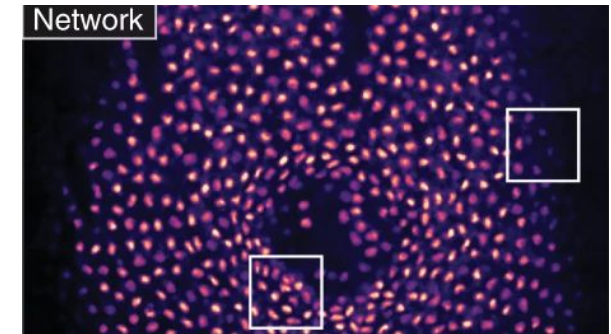
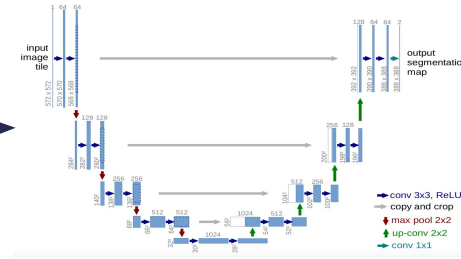
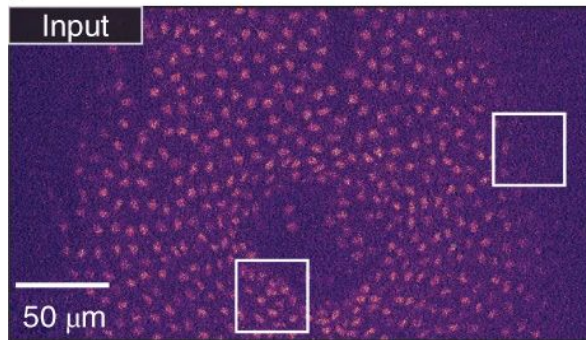
How do we obtain ground-truth for denoising? **Image in noisy and denoised conditions.**

- Fluorescence microscopy: image with different laser intensities.
 - High intensity = better signal
 - Not feasible for long experiment due to bleaching / phototoxicity
- EM: image with longer electron beam dwell times.
 - Longer dwell time = better signal
 - Takes too long for large scale acquisition, may damage sample.

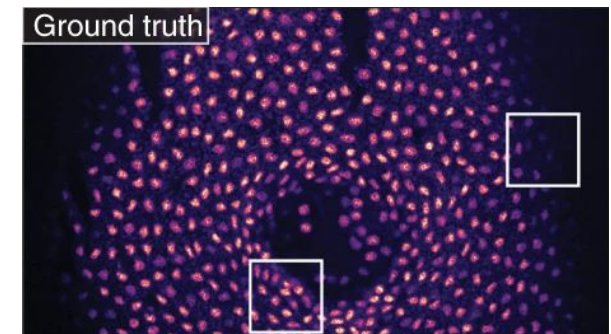
Denoising: CARE

<https://www.biorxiv.org/content/10.1101/236463v5>

How do we obtain ground-truth for denoising? **Image in noisy and denoised conditions.**



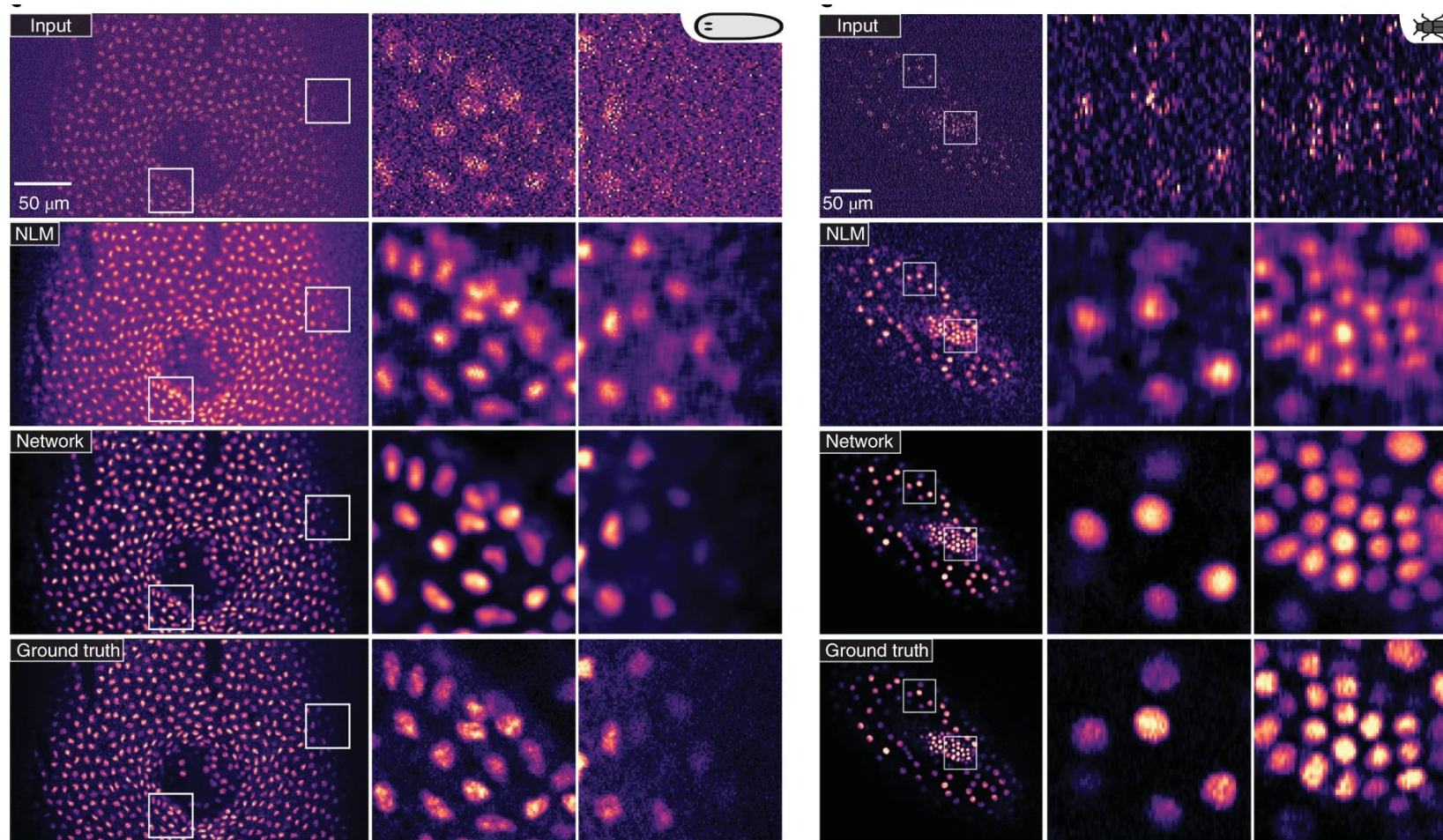
MSE
Loss



MSE = Mean Squared Error (L2 Norm):
Squared difference of pixel intensities,
averaged over all pixels

Denoising: CARE

<https://www.biorxiv.org/content/10.1101/236463v5>

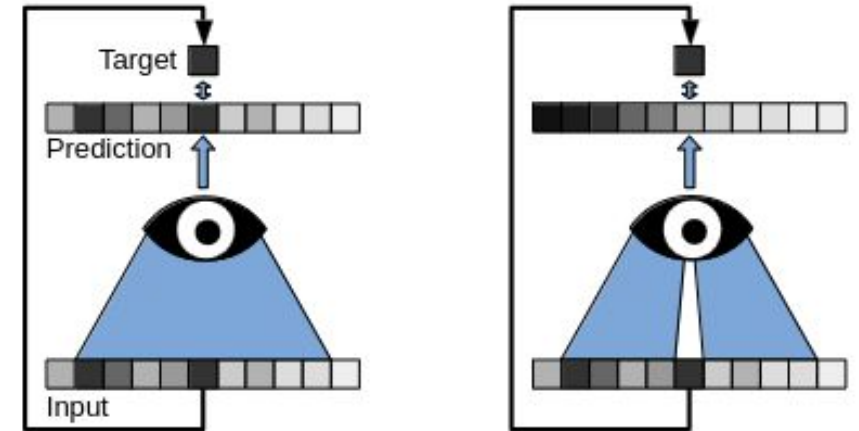


What if we can't take better images?

Can we do anything if clean targets are not available?

Yes! There are tasks we can solve for noisy data that allow denoising!

- Noise2Noise: predict a different noisy version of the image.
- **Noise2Void**: predict pixels from surrounding in noisy data. (Popular for microscopy data!)

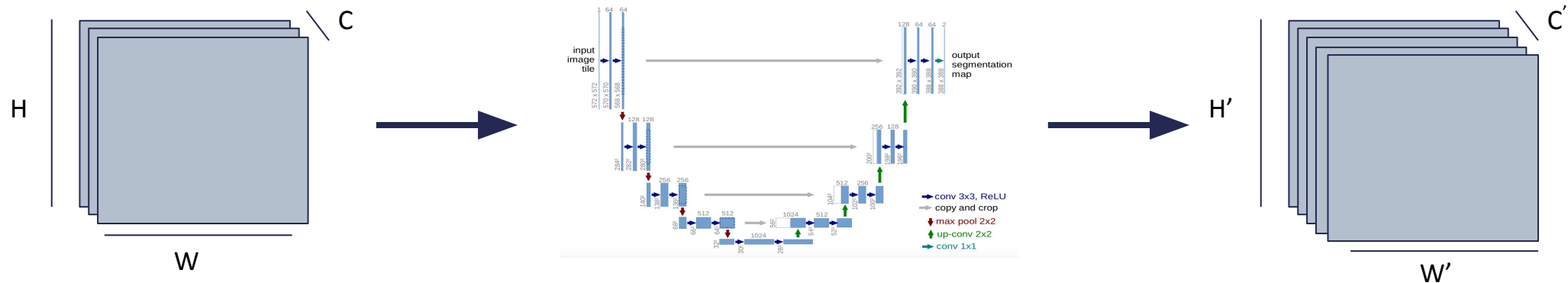


<https://arxiv.org/abs/1811.10980>

Other image restoration tasks

Several other image restoration tasks that can be solved with U-Net (or similar)

- “Super-resolution”: $H \times W \rightarrow (s * H) \times (s * W)$
 - for scale factor $s > 1$
- Slice prediction for anisotropic data: $H \times W \times D \rightarrow H \times W \times (s * D)$
 - e.g. $s = 2$: double the number of slices = double axial resolution



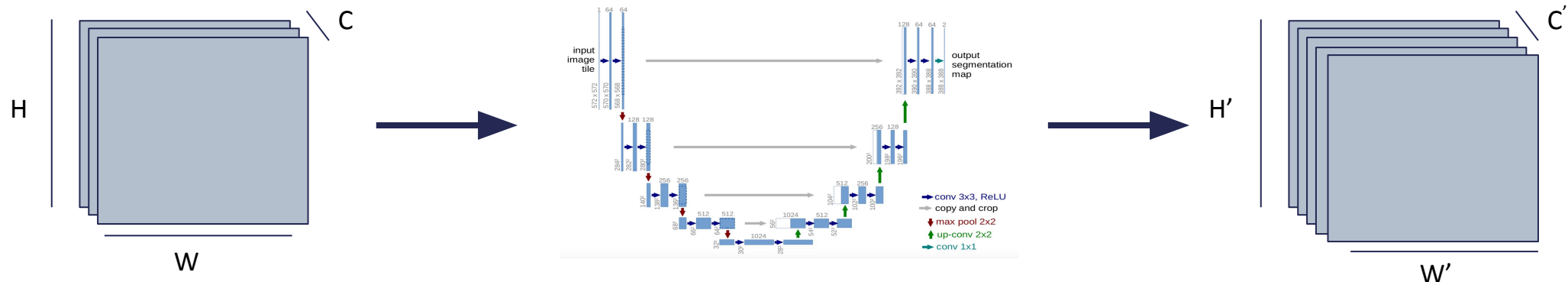
Other image restoration tasks

Several other image restoration tasks

- “Super-resolution”
 - for scale
- Slice prediction
 - e.g. $s = 2$

Framework for image restoration tasks in microscopy:
CAREamics: <https://careamics.github.io>

Provides a python library and napari plugin



Second Exercise

Explain the second exercise

Explain second and third exercise