

Different Tasks:

Semantic segmentation (**UNet / DeepLab**)

Detection / Instance Segmentation (**RCNN**)

Five Key Components

Task:

Data

Models

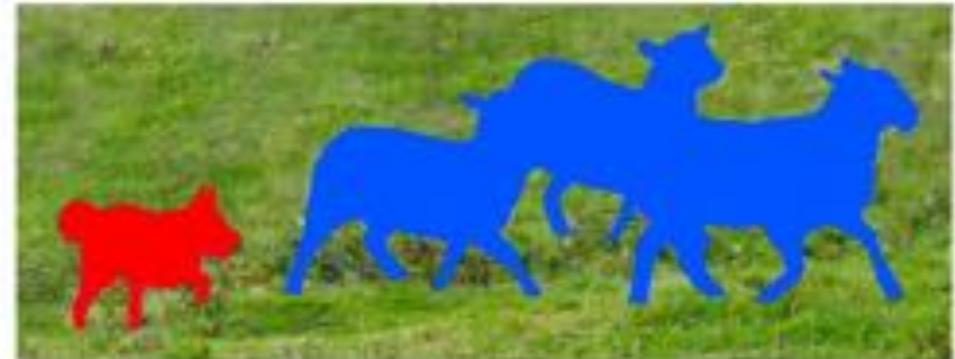
Metrics

Loss

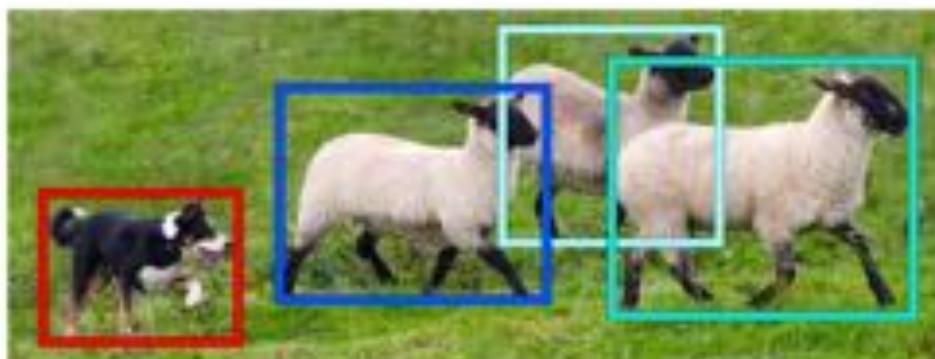
Consider Different Tasks.....



Image Recognition



Semantic Segmentation



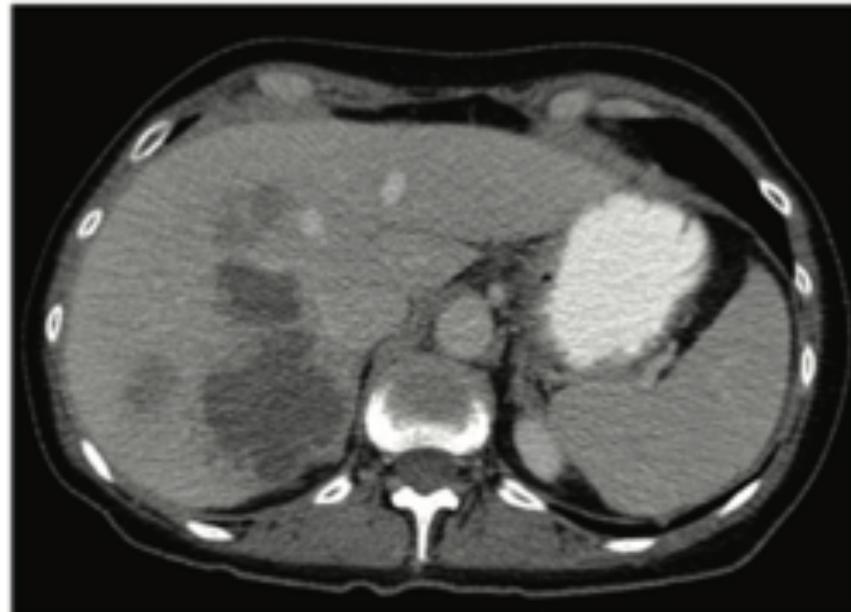
Object Detection



Instance Segmentation

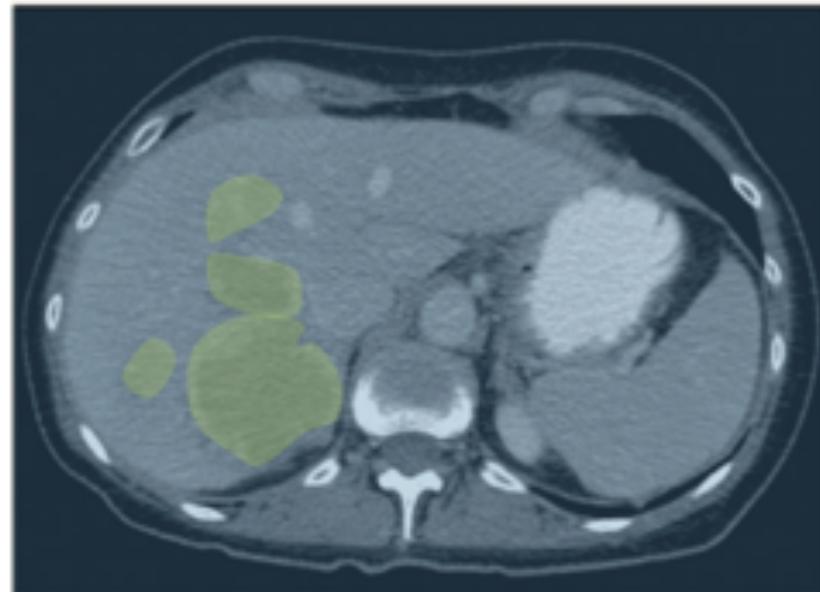
Instance Segmentation is kind of
Semantic segmentation + detection

Classification: liver metastases



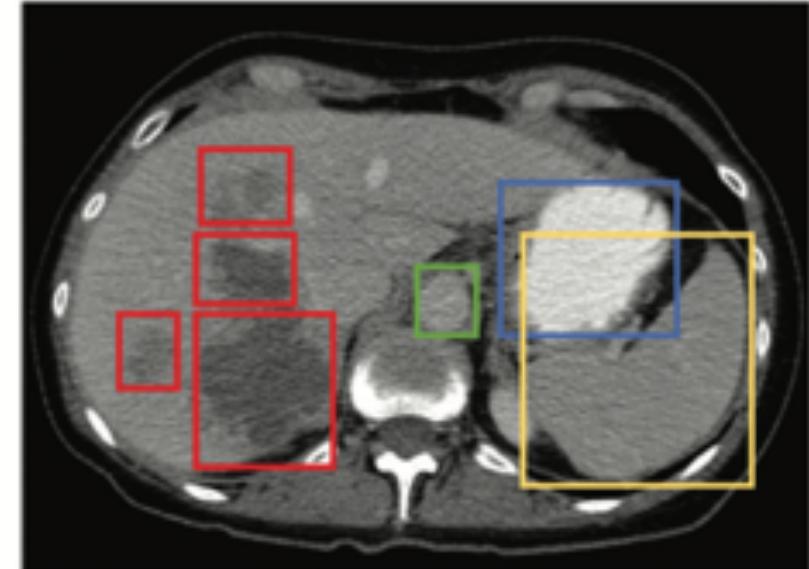
a.

Semantic segmentation



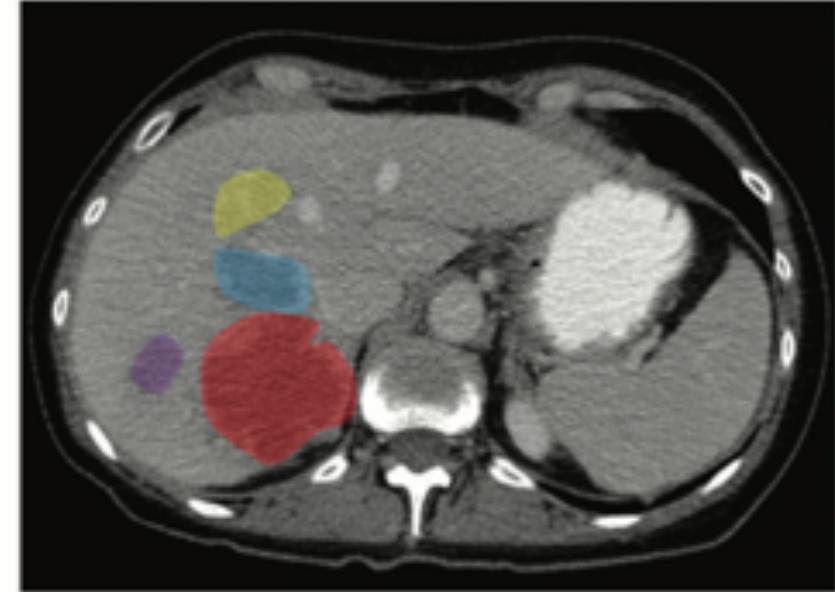
Liver metastases No metastasis

Object detection



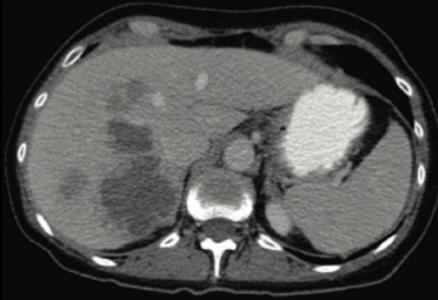
b.

Instance segmentation



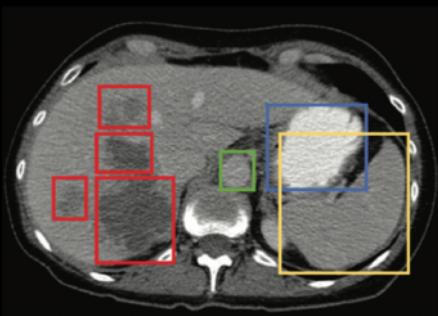
Metastasis 1 Metastasis 2 Metastasis 3 Metastasis 4

Classification



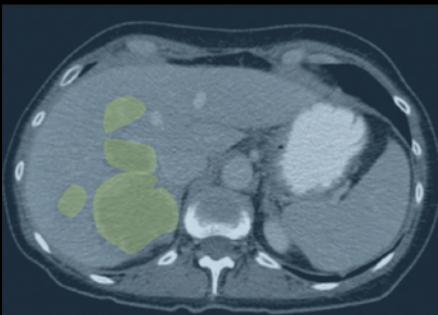
When to do what?

Detection



No annotation
But may need more data

Semantic
Segmentation

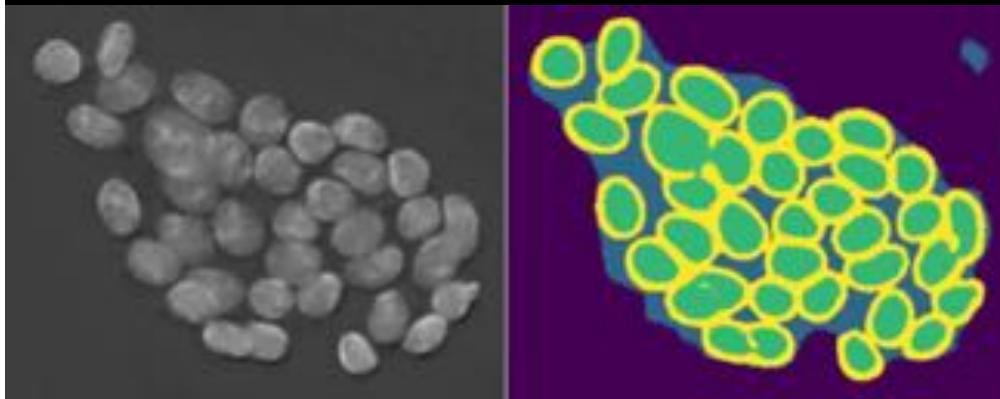


More manual cost
More Information

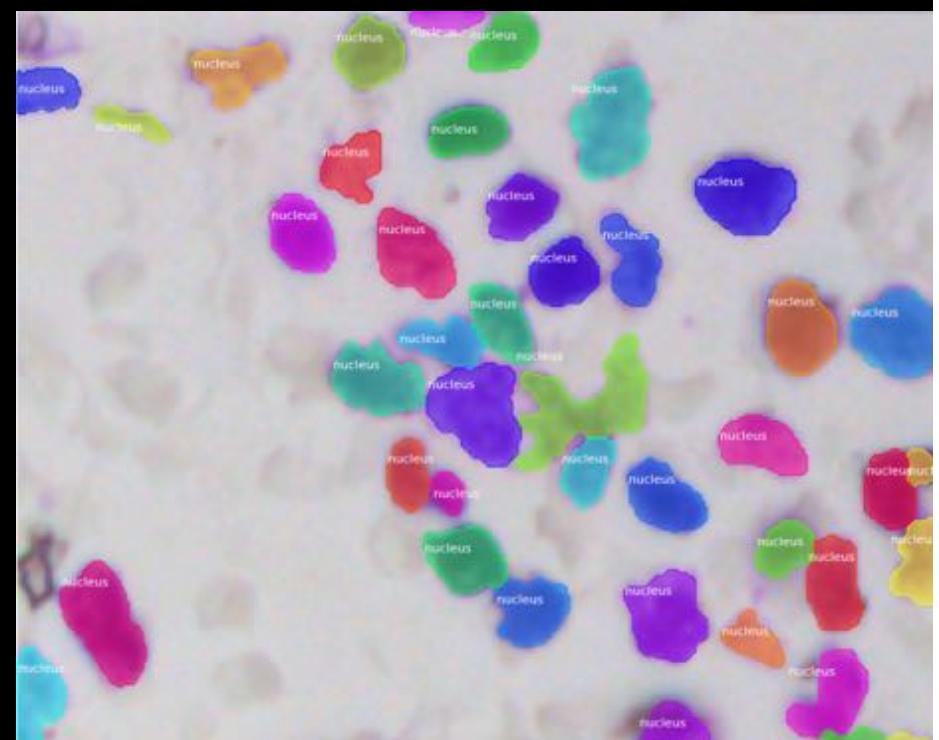
Instance
Segmentation



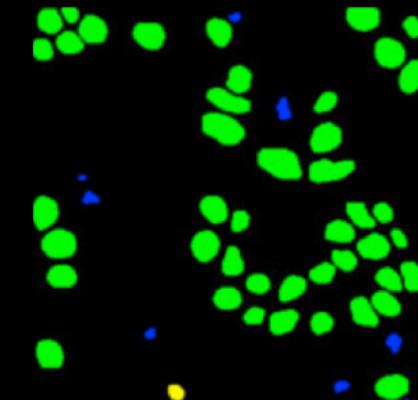
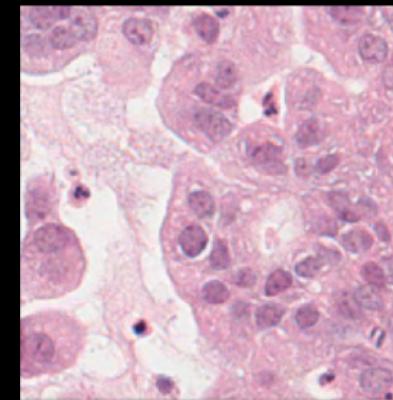
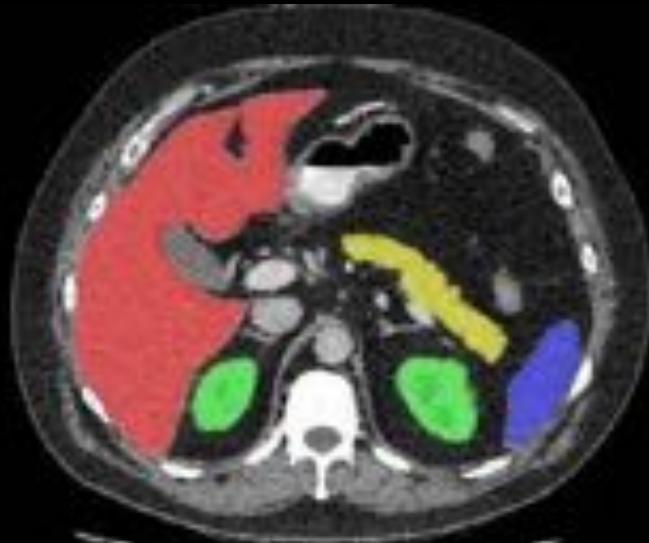
Cell semantic segmentation (cell + boundaries)



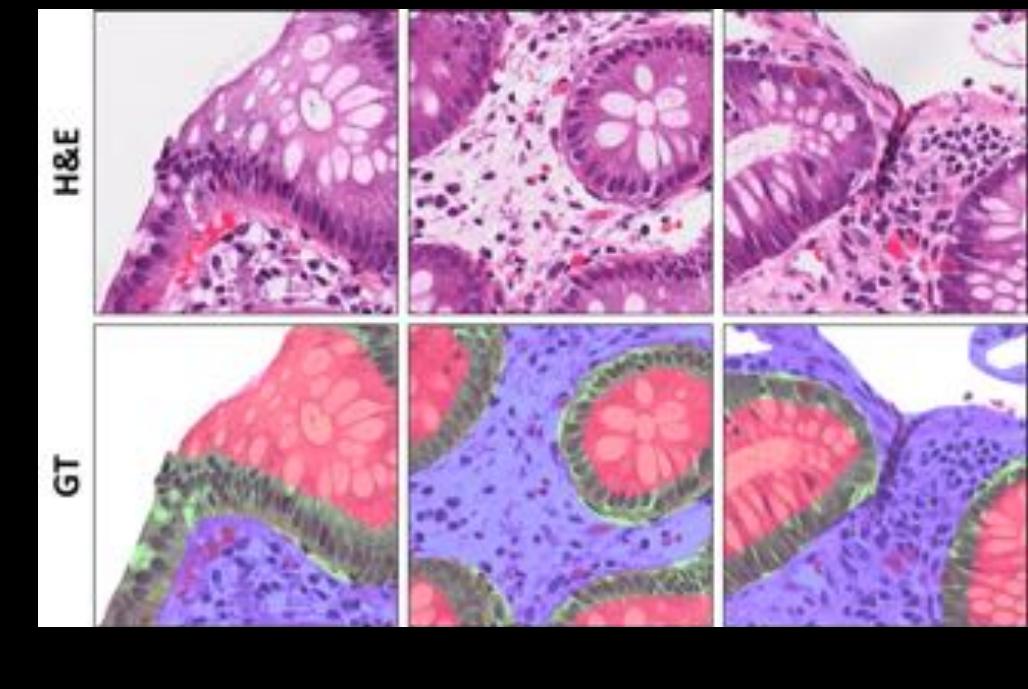
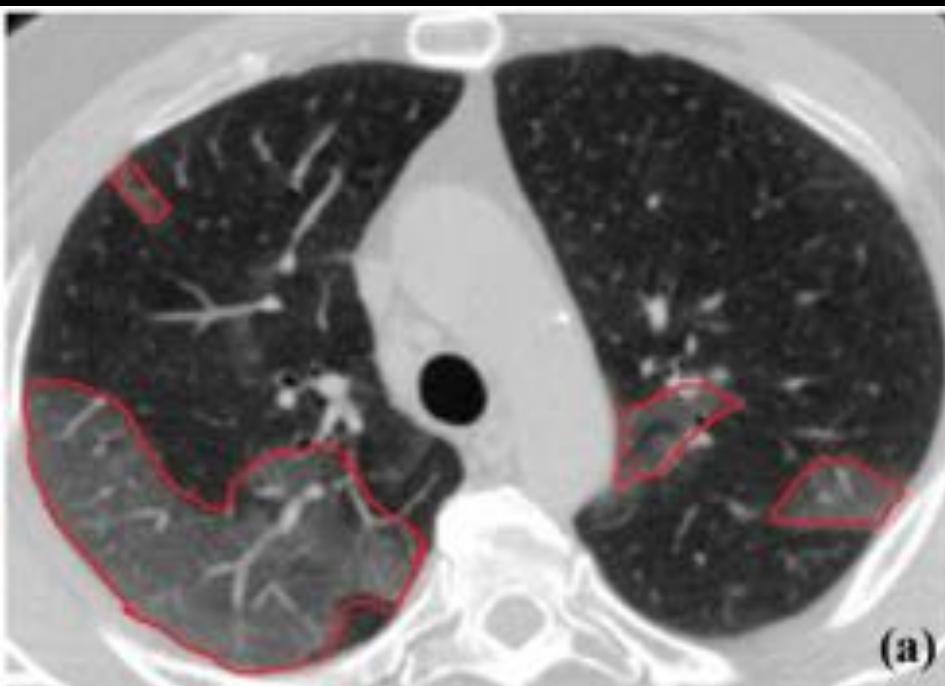
Cell instance segmentation (cell by cell)



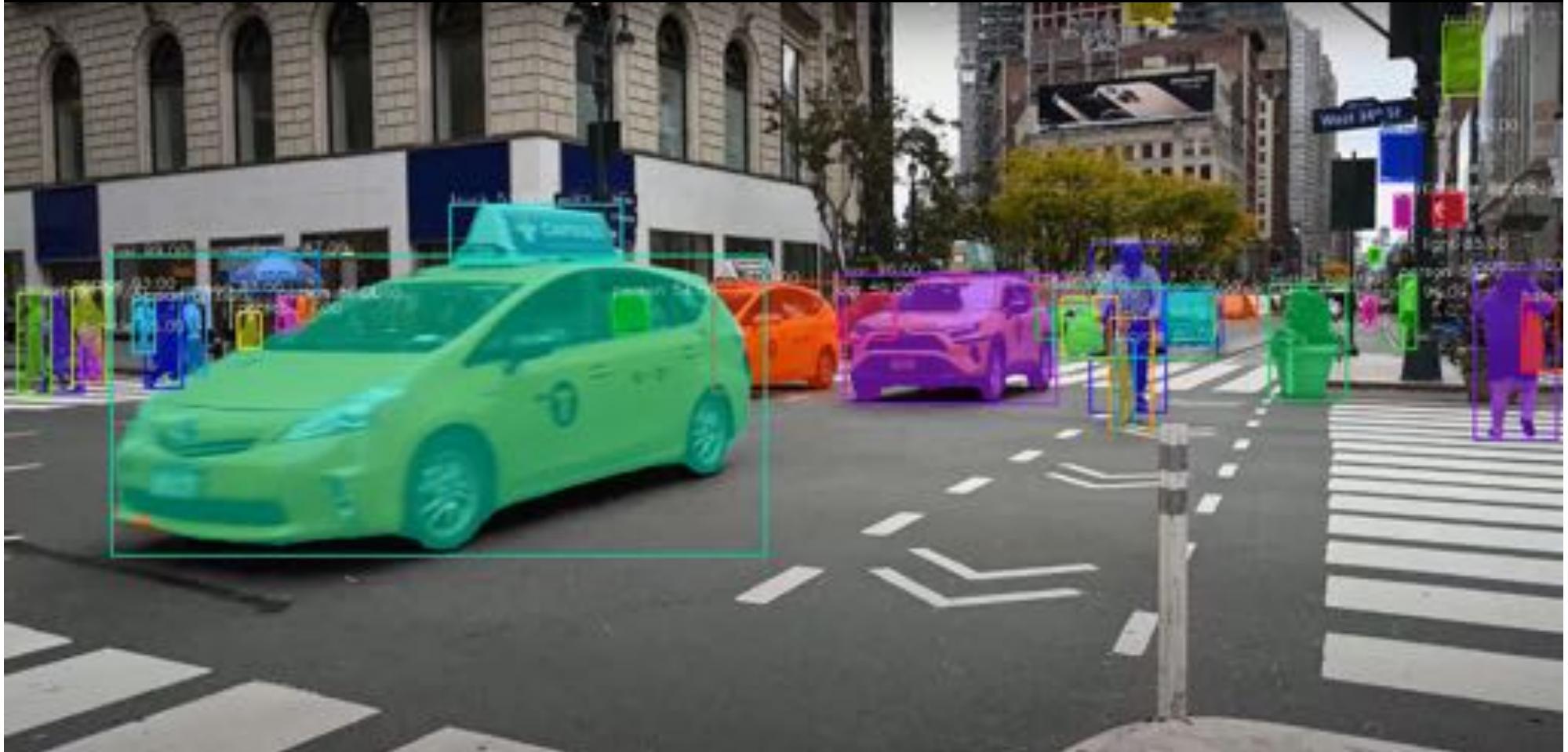
Some instances are clearly defined

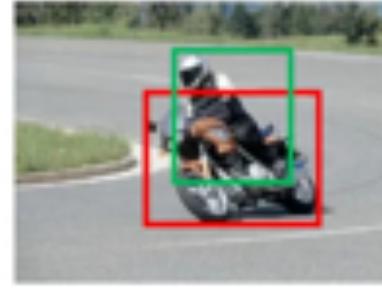


Some are not



Roughly \$1 per object / How much per object for your application?





{motorbike, person} {motorbike (point),
person (point)}

{motorbike (b-box),
person (b-box)}

{motorbike (pixel labels),
person (pixel labels)}

1 sec
per class

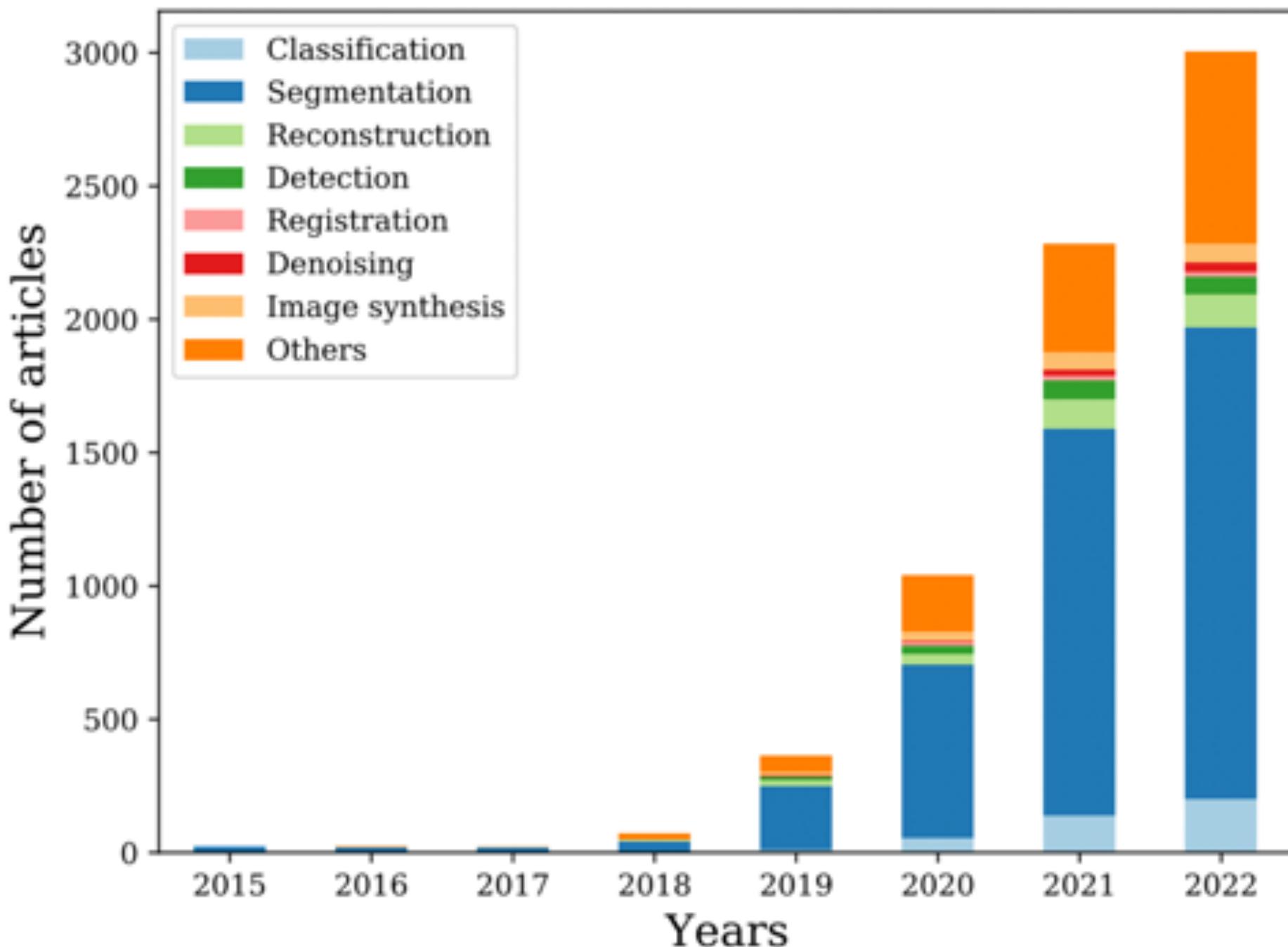
2.4 sec
per instance

10 sec
per instance

78 sec
per instance



Berman et al., What's the Point: Semantic Segmentation with Point Supervision, ECCV 16



Encoder only segmentation (FCN)

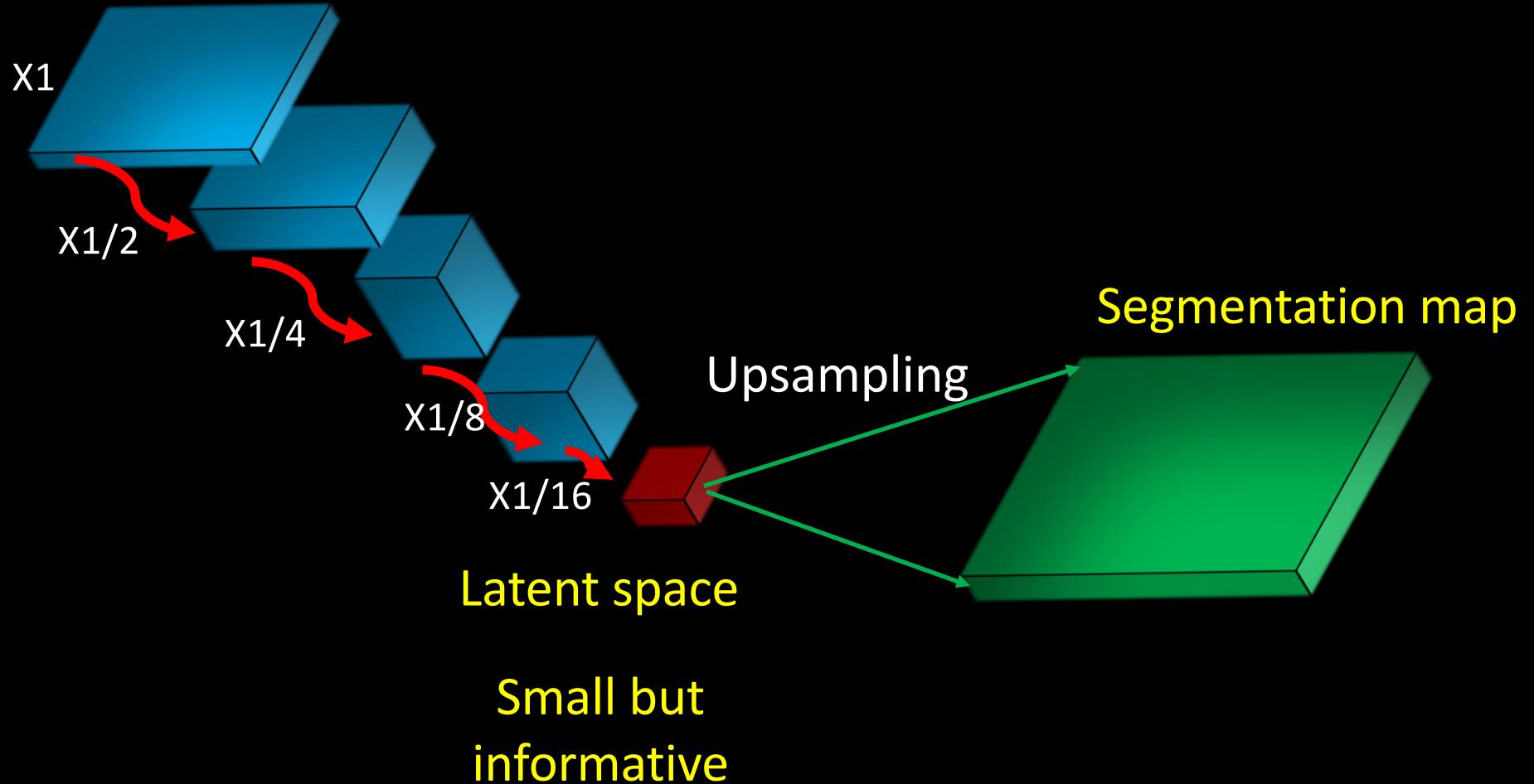
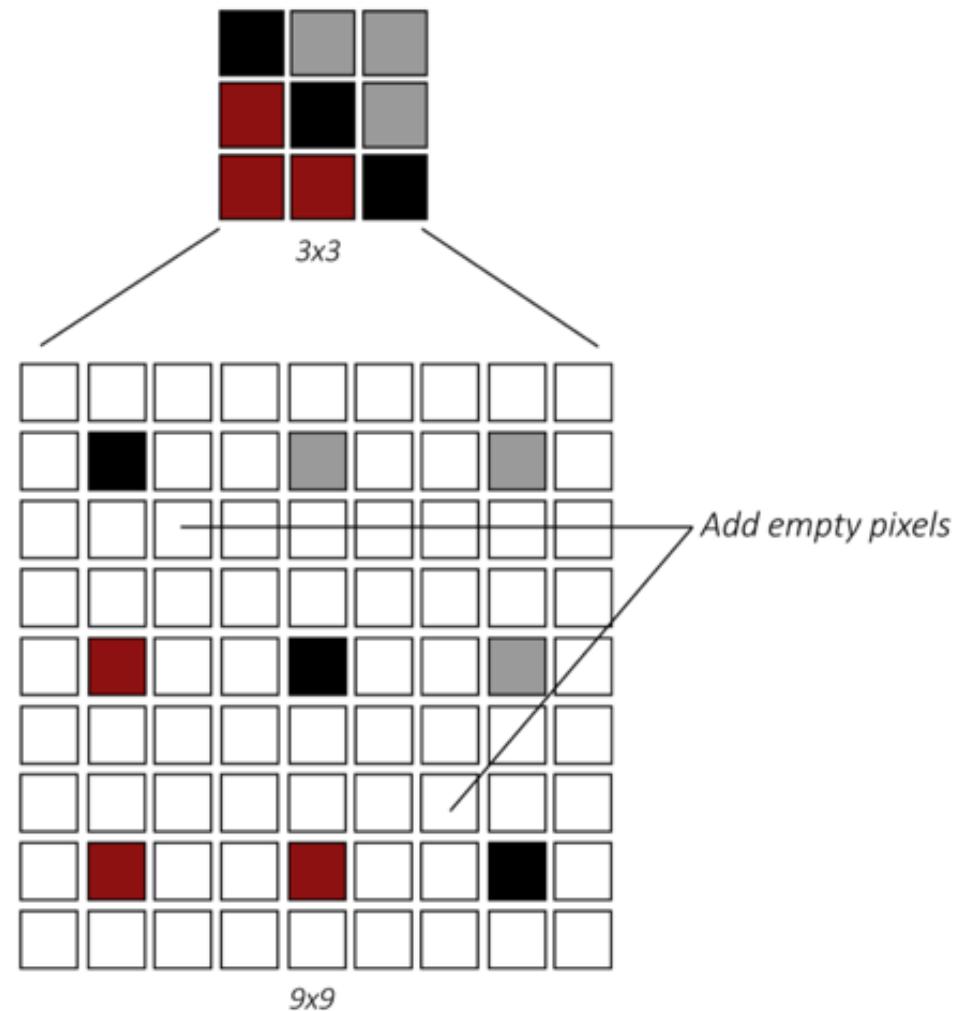
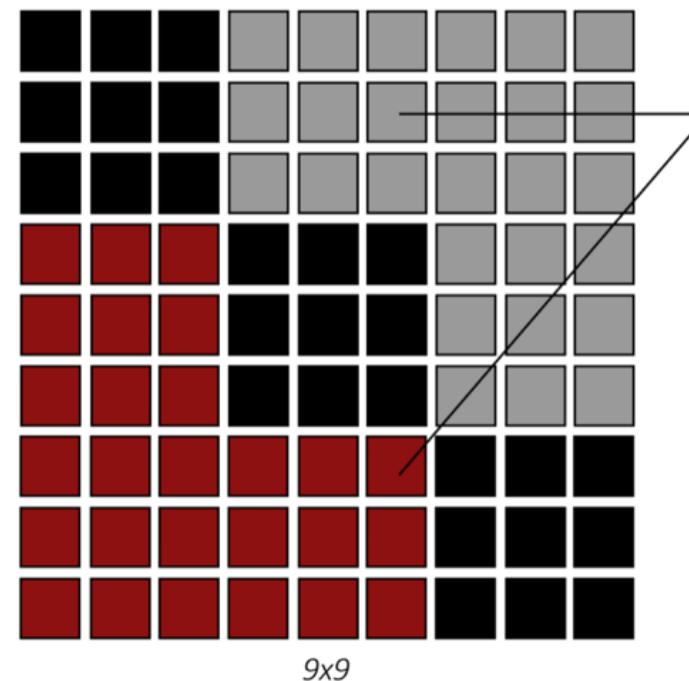
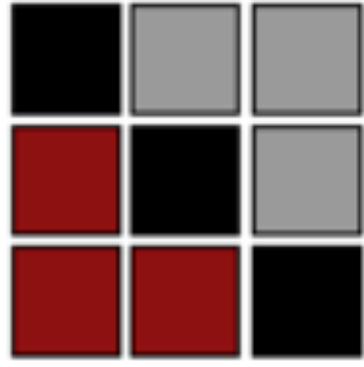


Image Up-sampling Problem

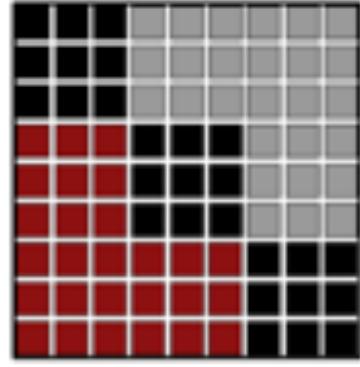


Fill by Nearest
Neighbors

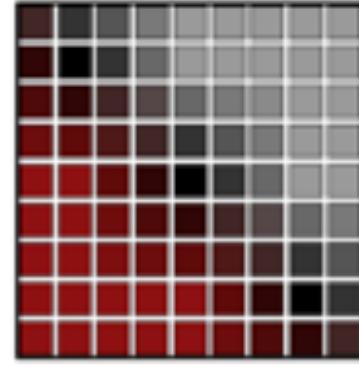




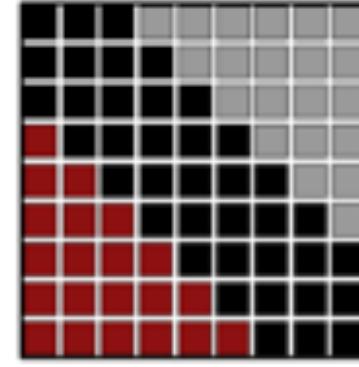
Low Resolution



*Nearest Neighbor
Interpolation*



*Bilinear
Interpolation*



Super-Resolution

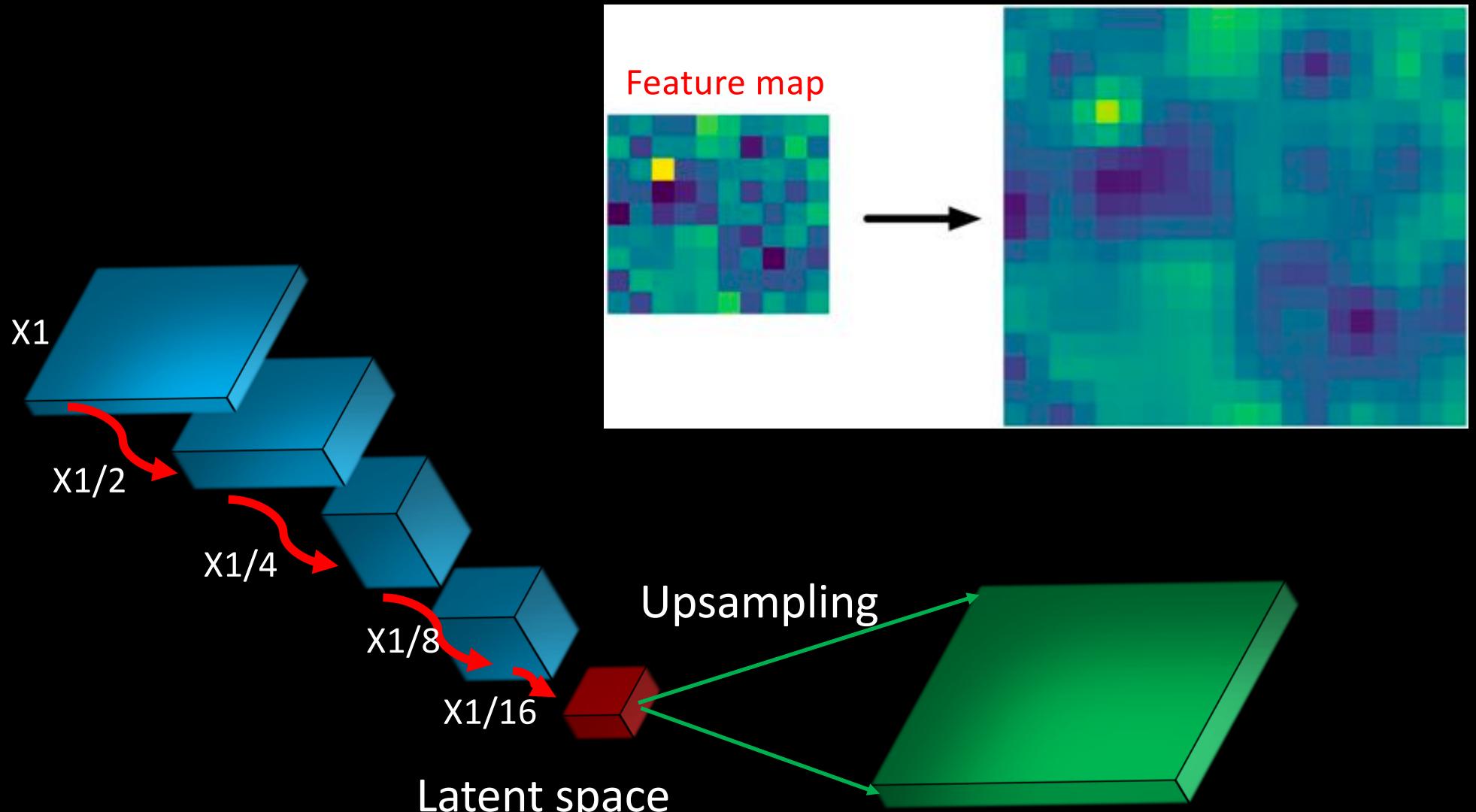
Accurate
but not
smart

Vague
Boundaries

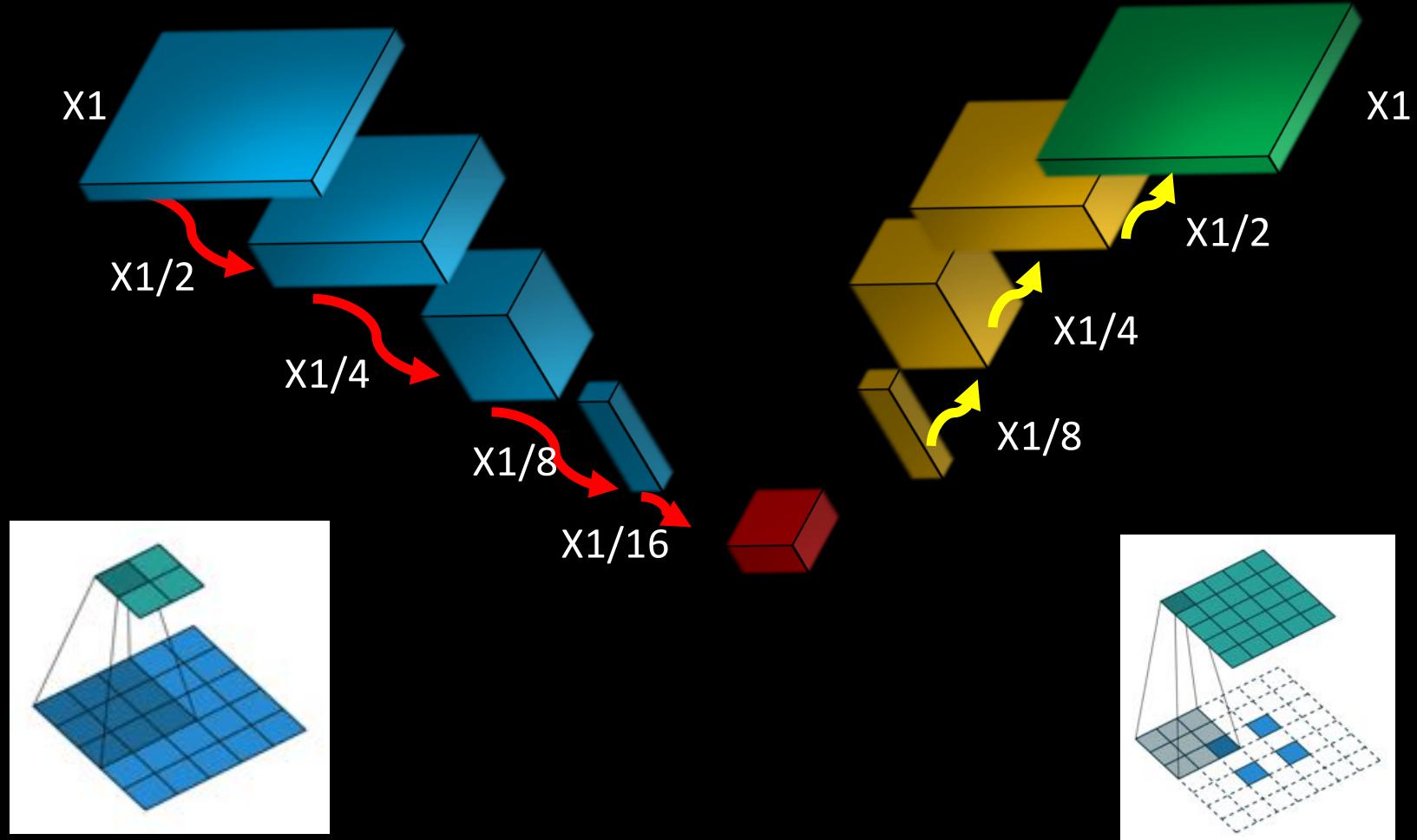
“smart”
interpolation

Encoder only segmentation (FCN)

Upsampled (not very smart!)



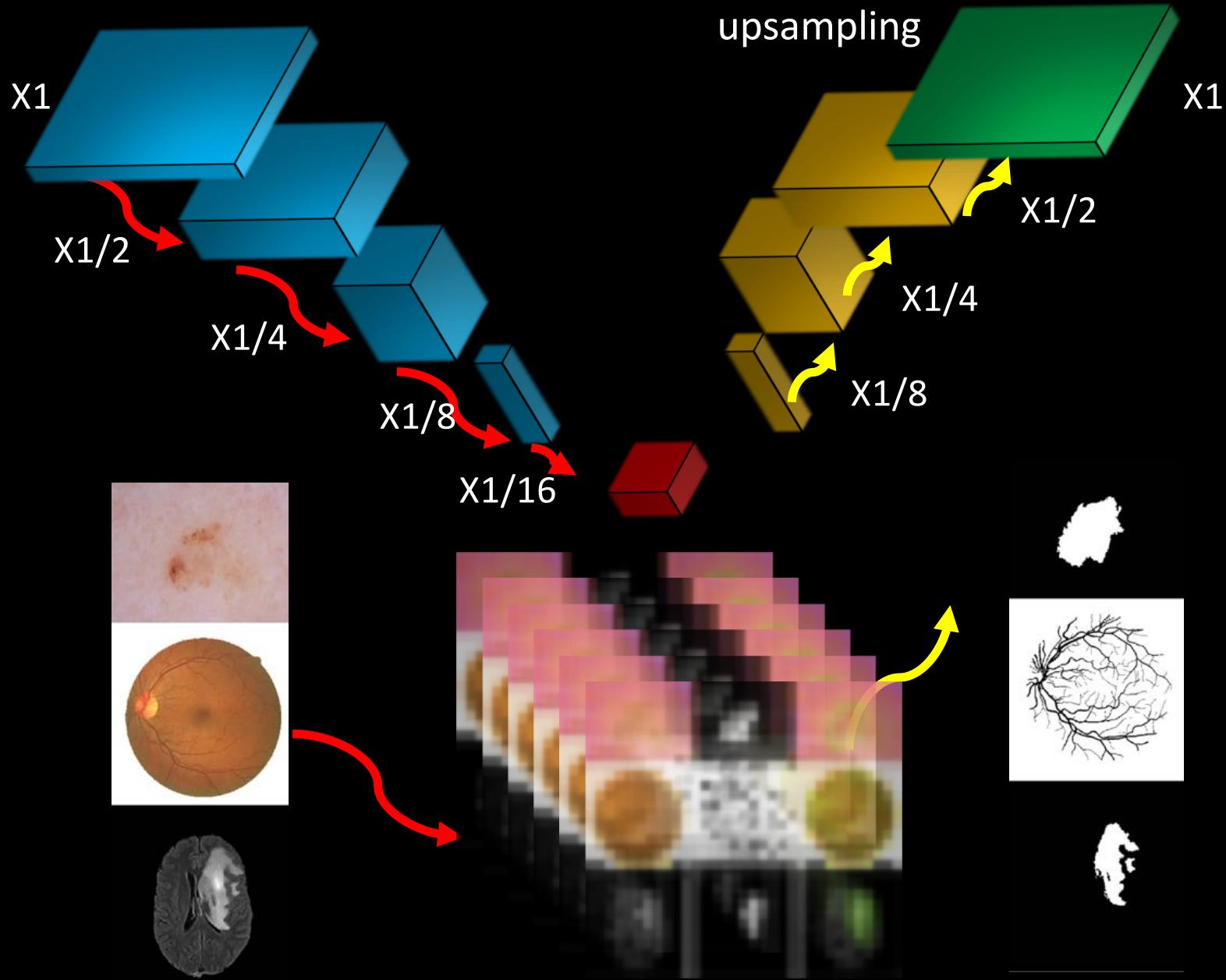
Decoding branch for processing while upsampling

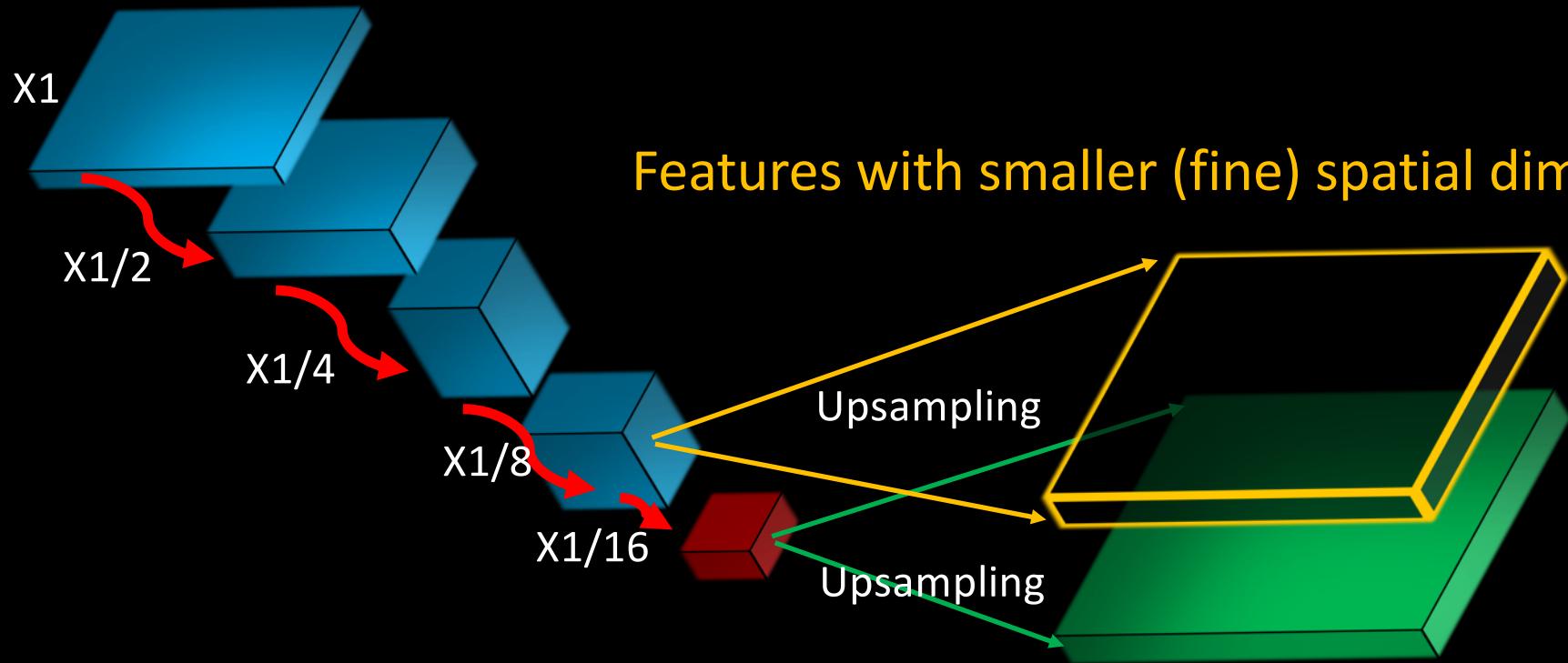


Convolution
(increase features, decrease spatial)

Deconvolution
(decrease features, increase spatial)

Decoding branch for
processing while
upsampling

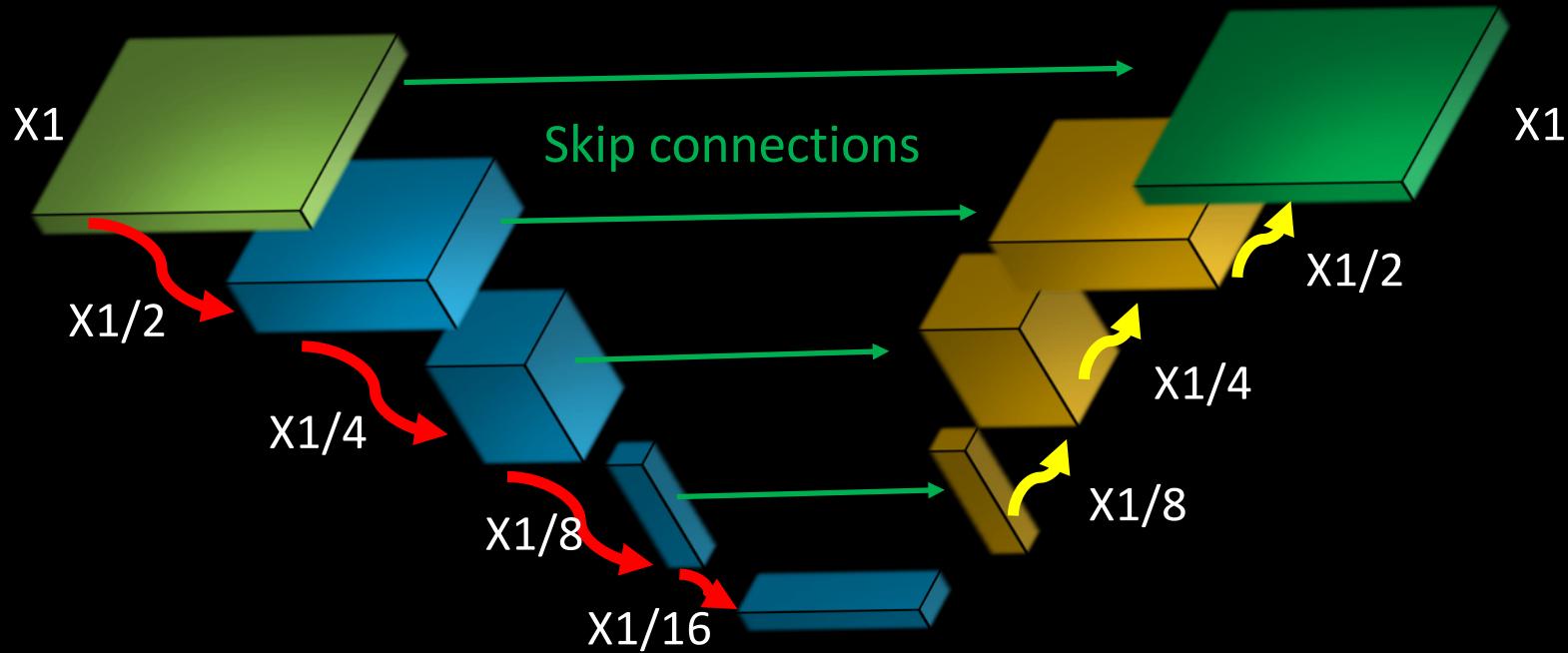




Features with larger (coarse) spatial dimension

Turn out we need features from all the scale level

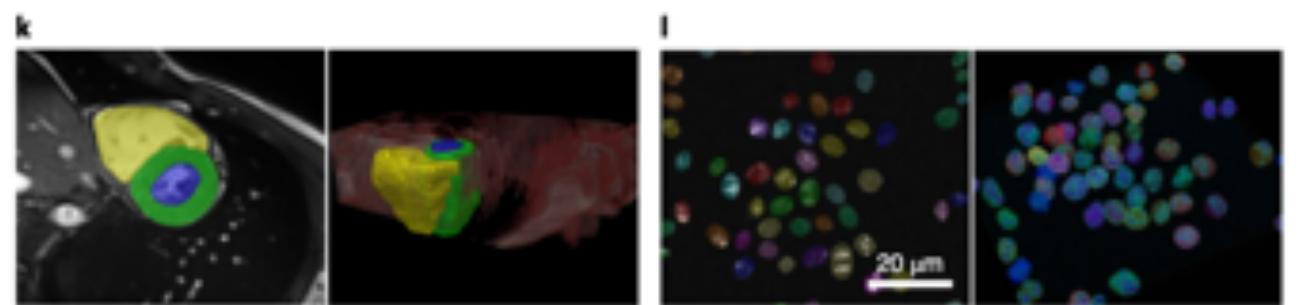
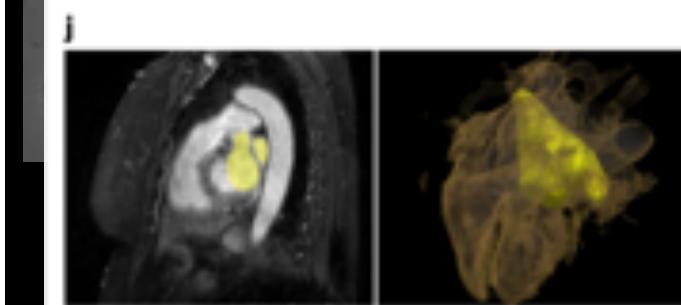
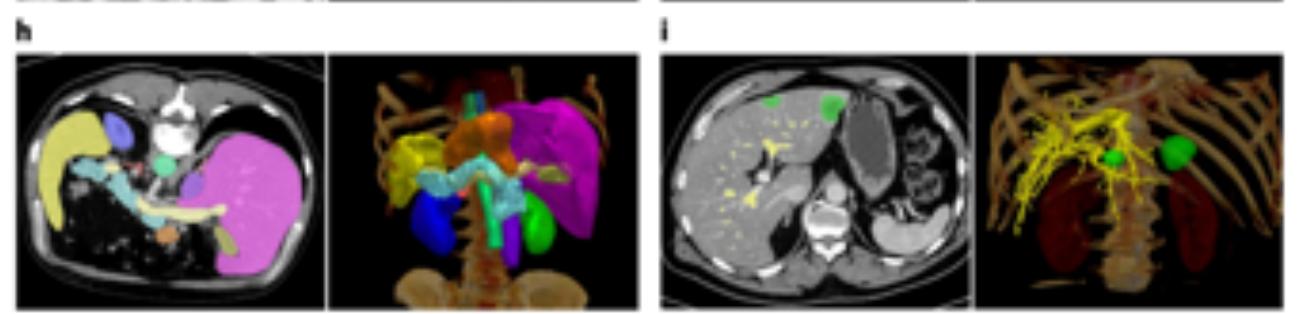
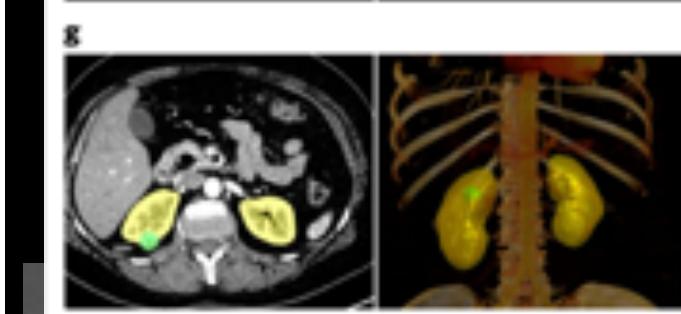
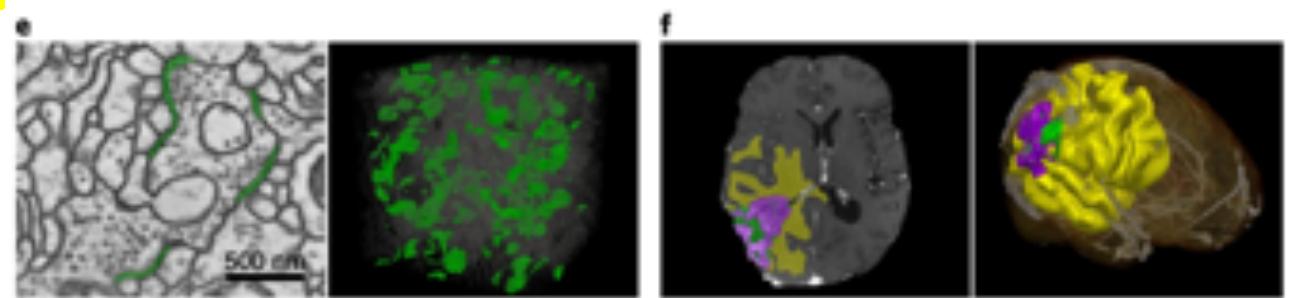
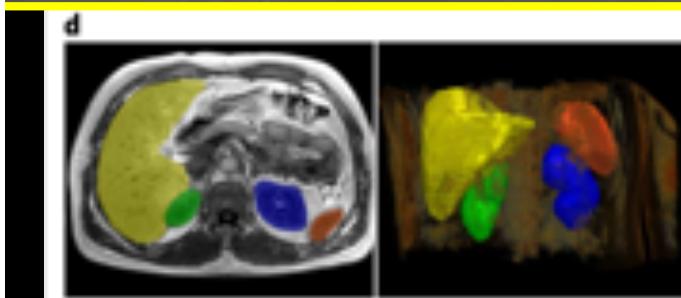
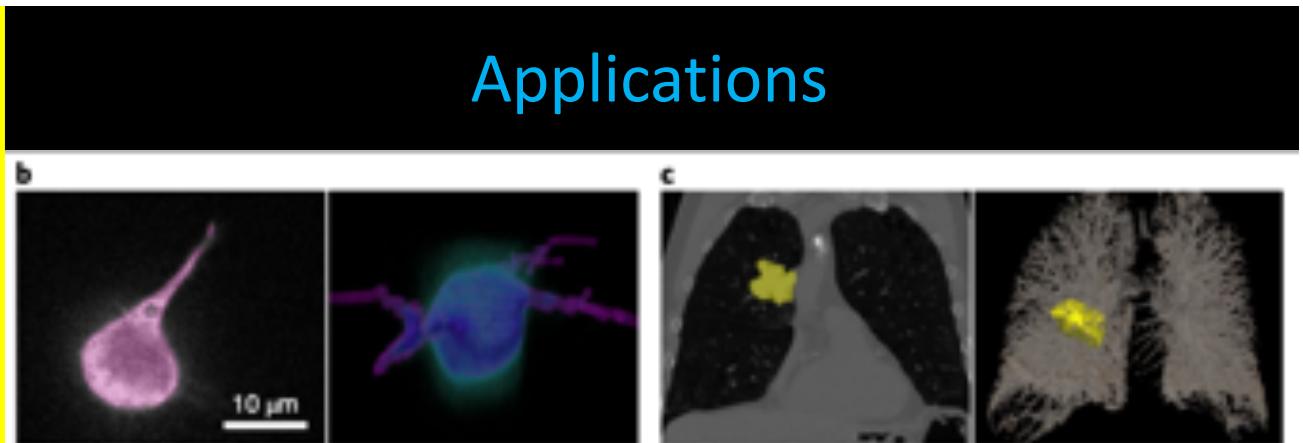
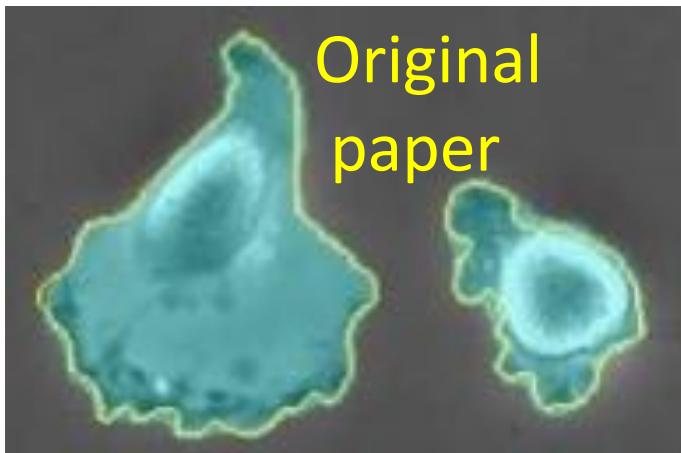
U-Net (encoder-decoder w/ skip connection)



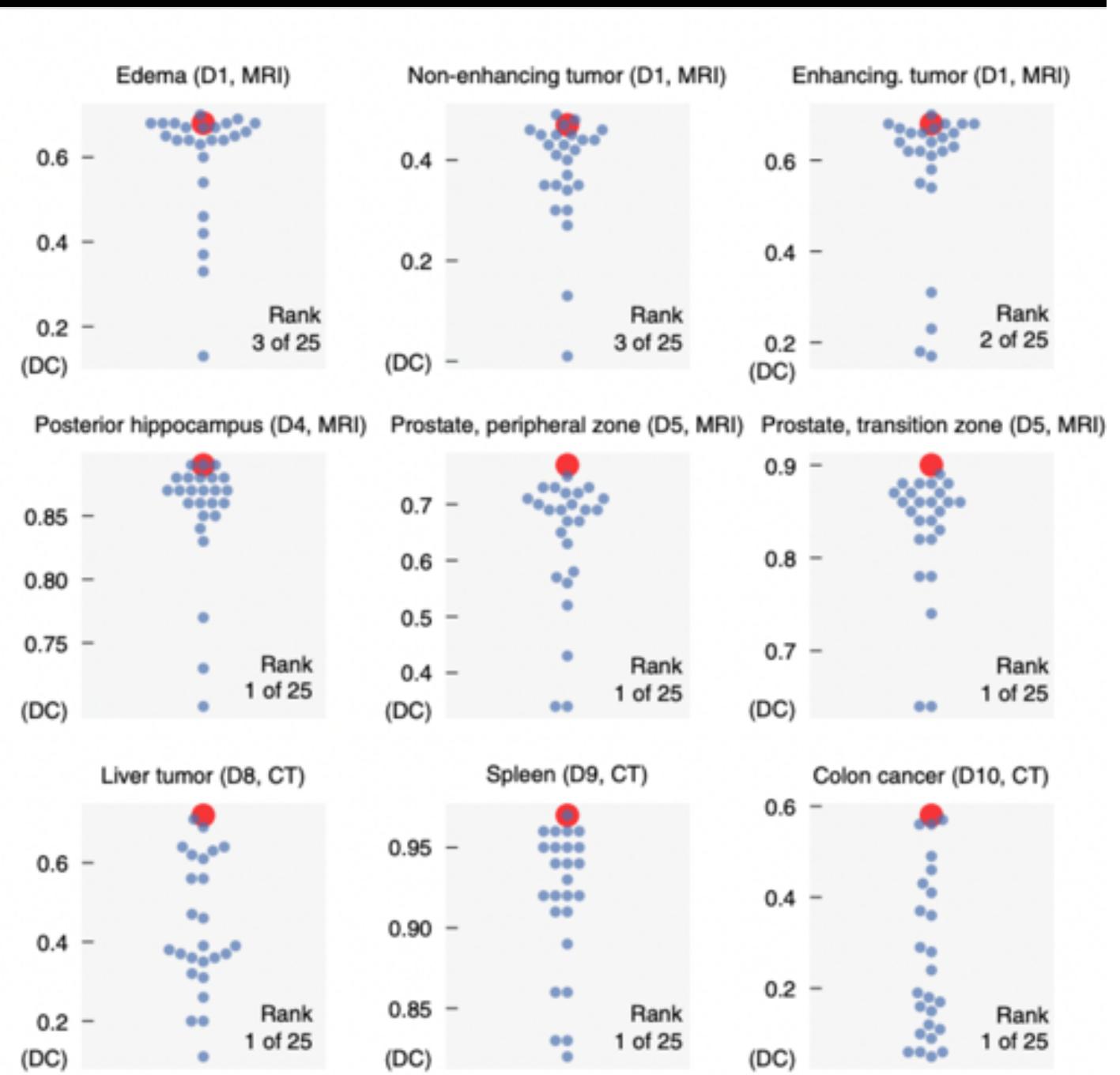
Skip connections: translate full size details
(fine and coarse) from left to right

Original paper

Applications

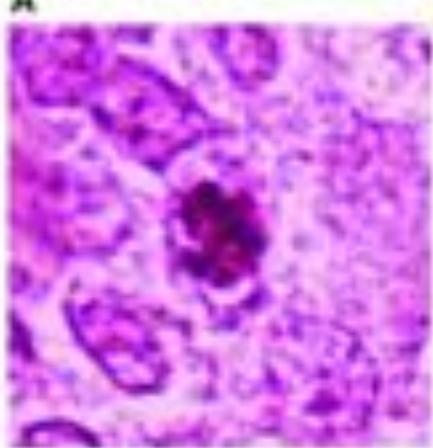


nn (not new) unet

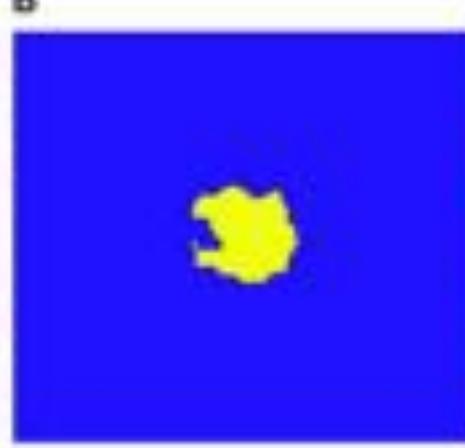


Loss and metrics of segmentation

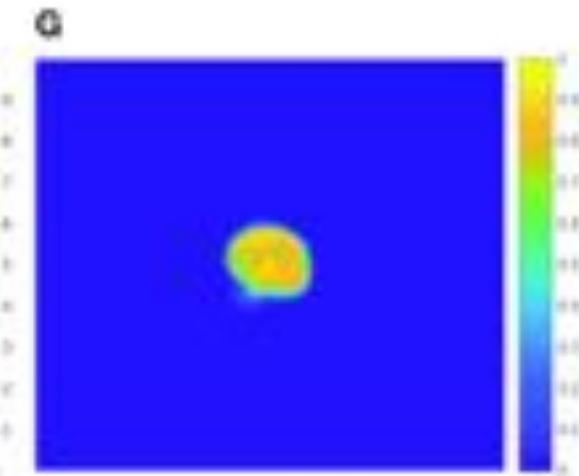
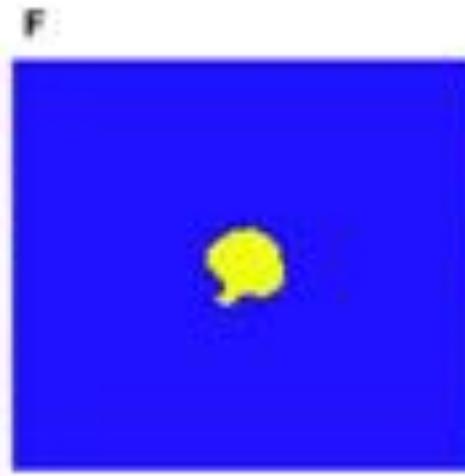
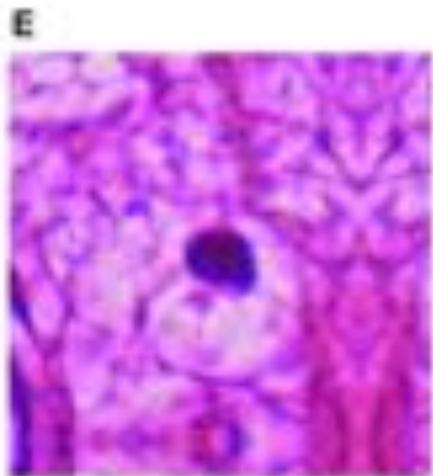
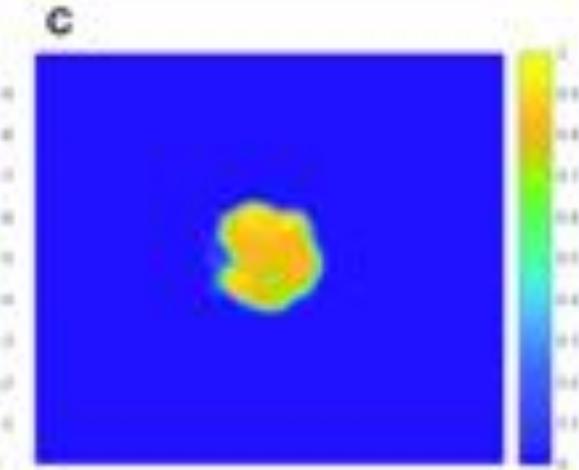
Original



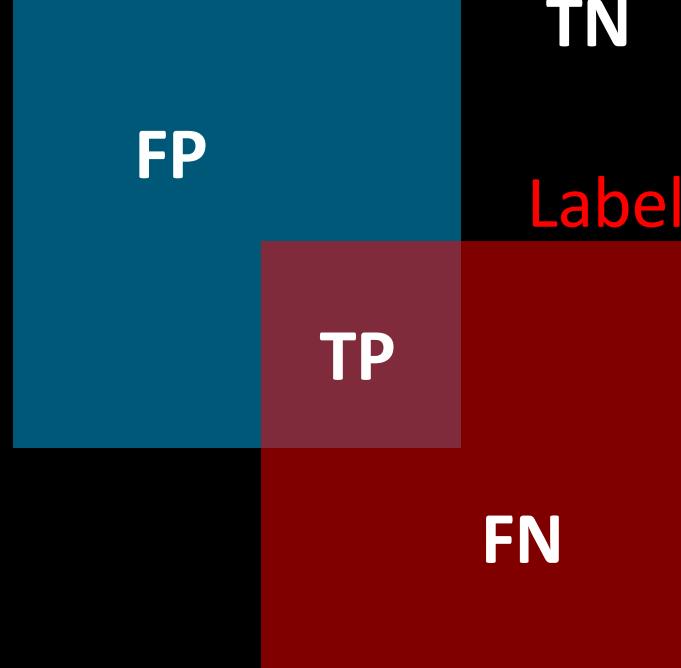
Ground-truth



Prediction
(probability)



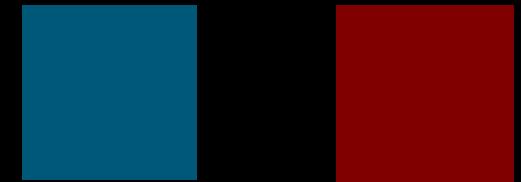
Prediction



Metrics

Dice Coefficient
(F1)

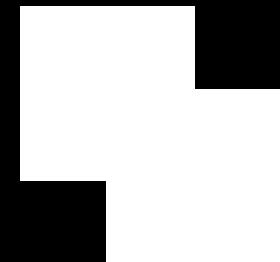
$$\frac{2 * \text{TP}}{t + p}$$



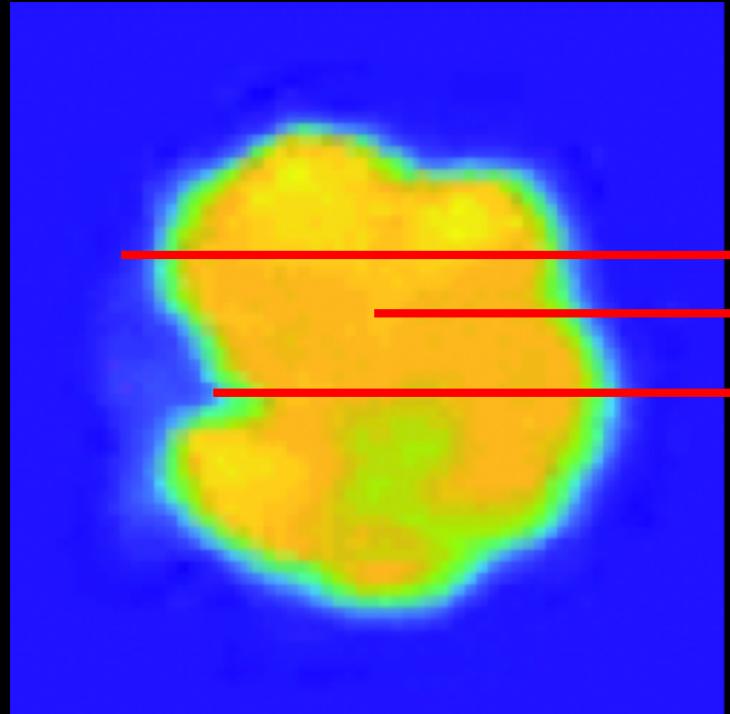
IoU



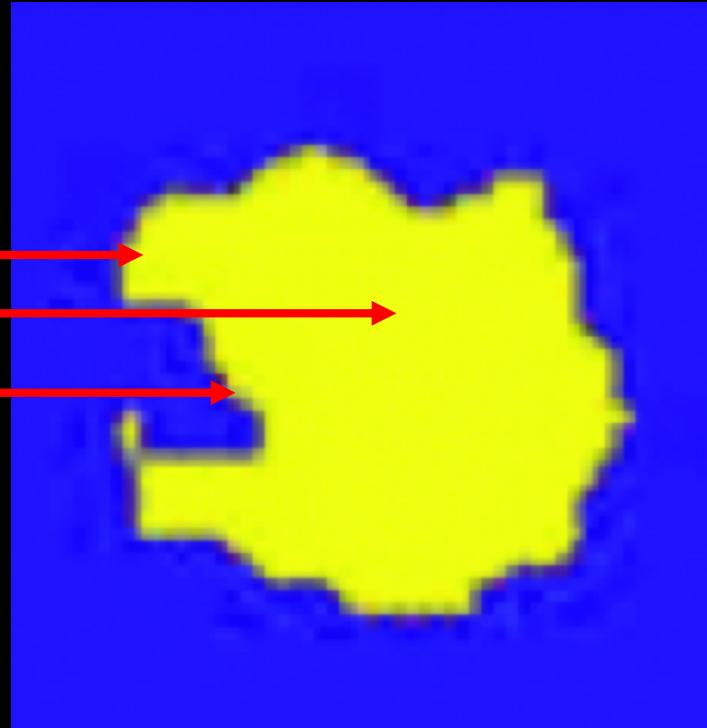
$$\frac{\text{TP}}{t + p - \text{TP}}$$



Prediction (probability)

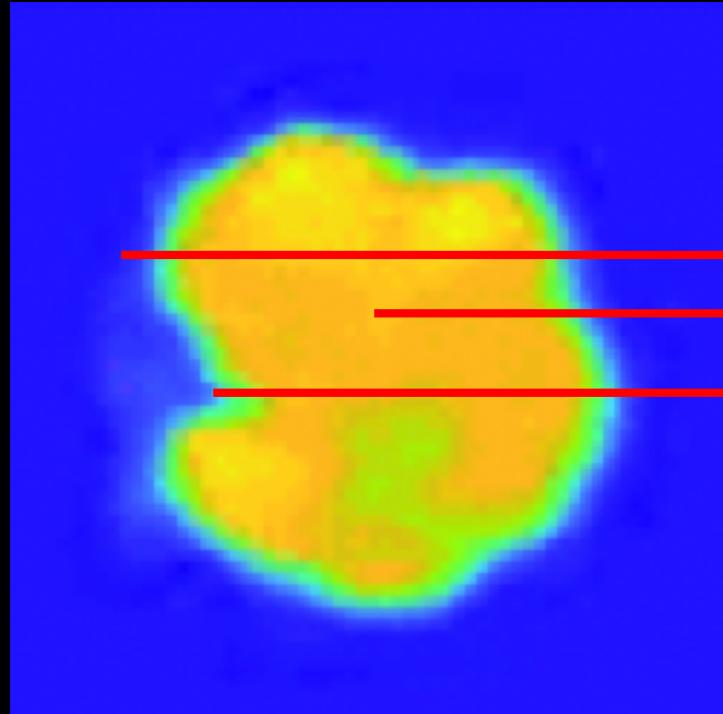


Ground-truth

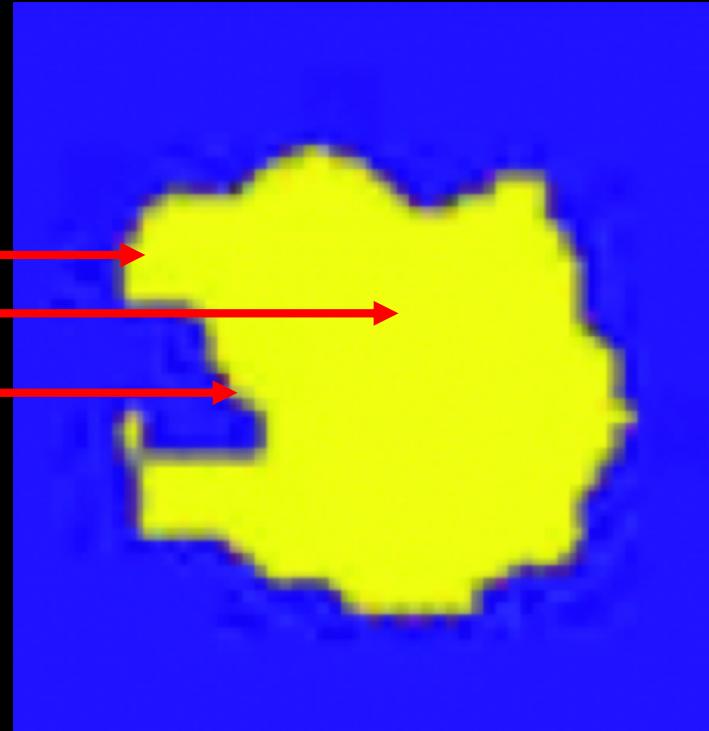


Its basically pixel-by-pixel
classification !!

Prediction (probability)



Ground-truth



Prob.

Label

Cross entropy = $-\log(q) * p$

0.01

1

$-\log(0.01) * 1 = 2$

False Negative. (FN)

0.95

1

$-\log(0.95) * 1 = 0.02$

True Positive (TP)

0.6

0

$-\log(1-0.60) * (1-0) = 0.53$

False Positive. (FP)



Nested UNet

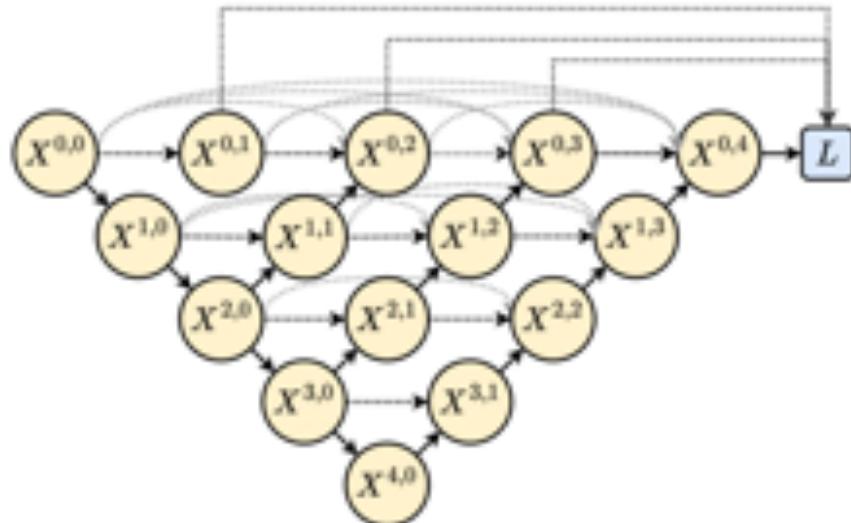


Fig. 6: Architecture of the U-Net++ [67]. The U-Net++ uses a nested structure of convolutional layers $X^{i,j}$ (indicate j -th convolution layer in the i -th scale) connected through the skip connection paths to model multi-scale representation.

Fully-connected Skip Connection

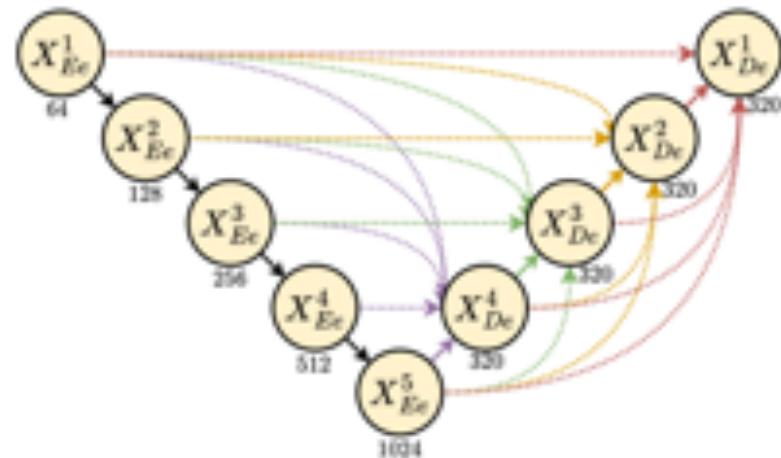
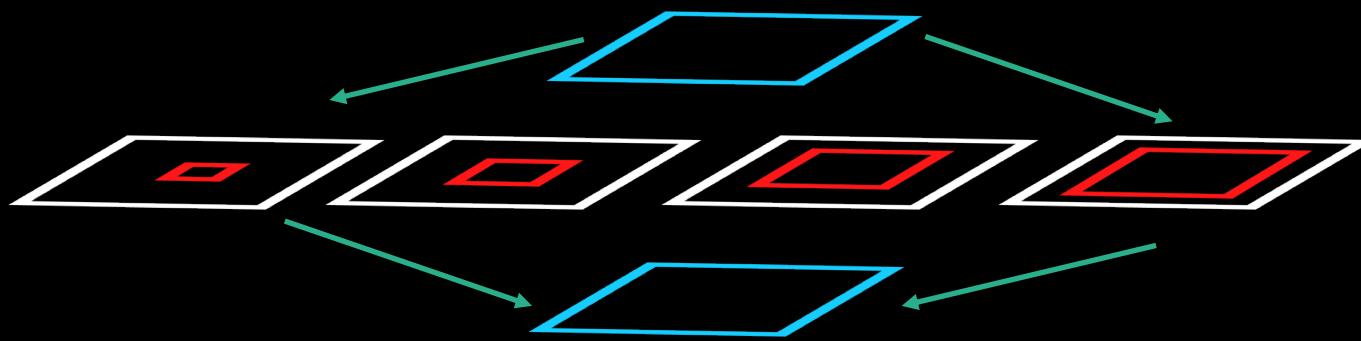
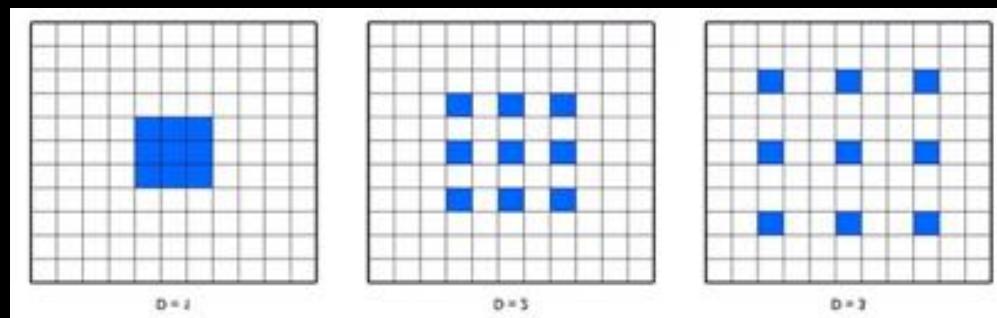


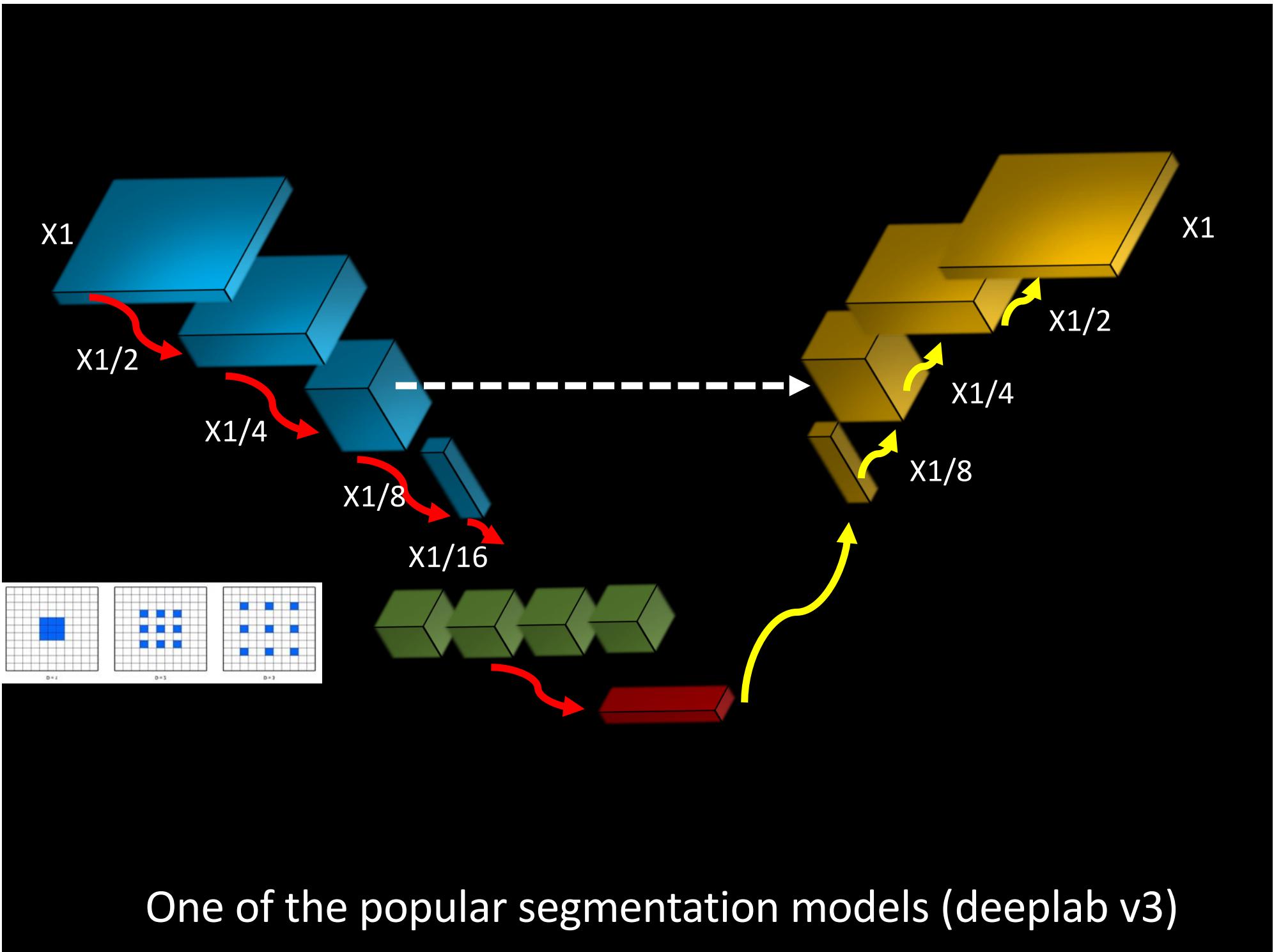
Fig. 7: Architecture showing the full-scale skip connections of the UNet3+ [68]. X_{Ee}^i and X_{De}^i denote the encoder and decoder path feature maps at i -th scale, respectively. In addition, each subscript under each feature map symbol represents the number of feature maps to the corresponding block.

Another popular segmentation model: Deeplab Families

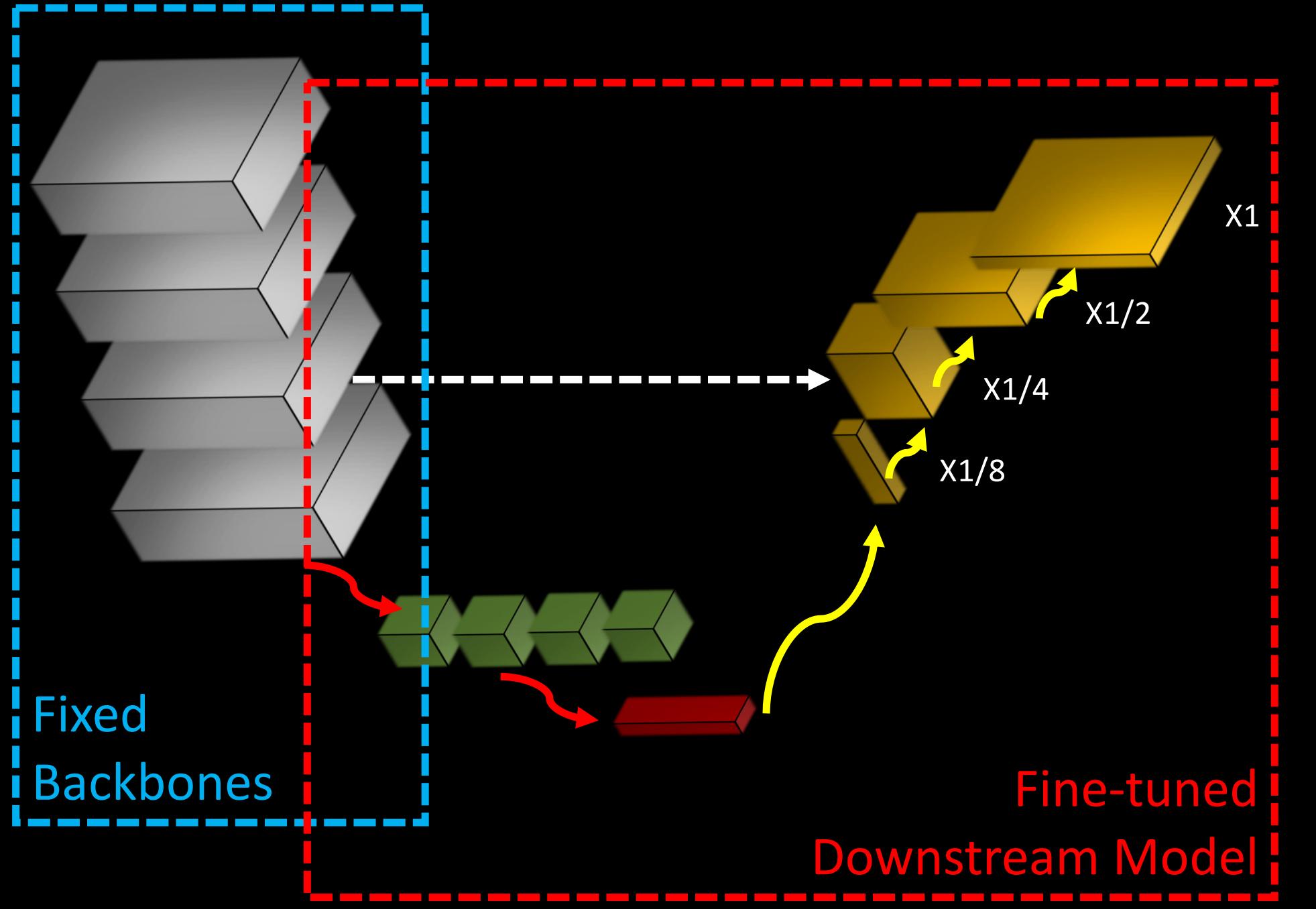


Spatial Pyramid Pooling



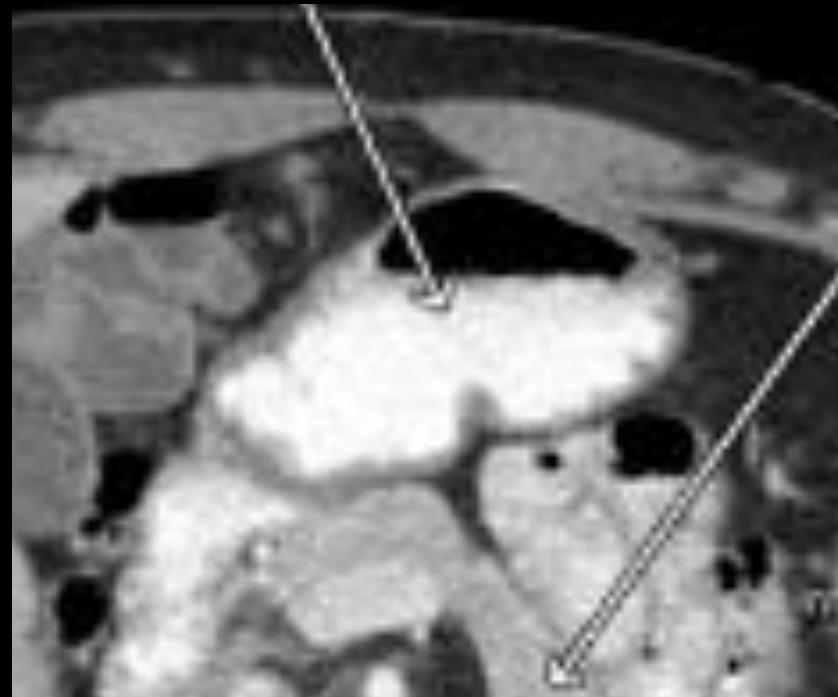


We can also utilize pre-trained CNN and their information





Backbones:
features extractors pretrained with a massive dataset for classification, and useful for “downstream” tasks



Q1: In stomach in this picture?

Q2: Can you point it out? —————> Downstream
tasks

Q3: Can you draw the outline? ↗

Deep Learning researchers (Mildly exaggerated)



How to use an existing model with bare-minimum understanding:

Welcome to the new nnU-Net!

Click [here](#) if you were looking for the old one instead.

Coming from V1? Check out the [TLDR Migration Guide](#). Reading the rest of the documentation is still strongly recommended ;-)

<https://github.com/MIC-DKFZ/nnUNet>

Make sure the code is trustworthy &
can be reproduced



[Notifications](#)



[Fork](#) 1.8k



[Star](#) 6k



Issues 218

Usually, there is a pre-defined data format

```
nnUNet_results/
└─ Dataset002_Heart
    ├─ nnUNetTrainer_nnUNetPlans_2d
    │   ├─ fold_0
    │   ├─ fold_1
    │   ├─ fold_2
    │   ├─ fold_3
    │   ├─ fold_4
    │   └─ dataset.json
    │   └─ dataset_fingerprint.json
    └─ plans.json
    └─ nnUNetTrainer_nnUNetPlans_3d_fullres
        ├─ fold_0
        ├─ fold_1
        ├─ fold_2
        ├─ fold_3
        ├─ fold_4
        └─ dataset.json
        └─ dataset_fingerprint.json
        └─ plans.json
```

First, understand the pre-defined format

uuUNet_raw

Dataset001_heart

imagesTr

heart_0001_0000_nii_gz

imagesTs

heart_0001_0000_nii_gz



Load the XXX.nii.gz file for me and
Print out the datatype and dimension

```
imagesTr  
    heart_0001_0000_nii_gz  
imagesTs  
    heart_0001_0000_nii_gz
```

```
>>> # Load the NIfTI file  
>>> img = nib.load('heart_0001_0000.nii.gz')  
>>> data = img.get_fdata()  
>>> # Print shape and data type  
>>> data.shape  
(256, 256, 120) # Example dimensions - height,  
width, and number of slices  
>>> data.dtype  
dtype('unit16')
```

Our Data

input image.npy (XYZ, float32)
seg mask.npy (XYZ, unit8)

?????

imagesTr
input image (XYZ, uint16)
imagesTs
seg mask. (XYZ, unit8)

Our Data

input image.npy (XYZ, float32)
seg mask.npy (XYZ, unit8)

?????

imagesTr
input image (XYZ, uint16)
imagesTs
seg mask. (XYZ, unit8)



I have a 3D (MRI, mask) pair of (float32 and uint8), help me to convert them to .nii.gz of uint 16 and uint 8

```
>>> import numpy as np, nibabel as nib

>>> mri = np.load('mri.npy').astype(np.uint16) # float32 -> uint16

>>> mask = np.load('mask.npy').astype(np.uint8) # keep uint8

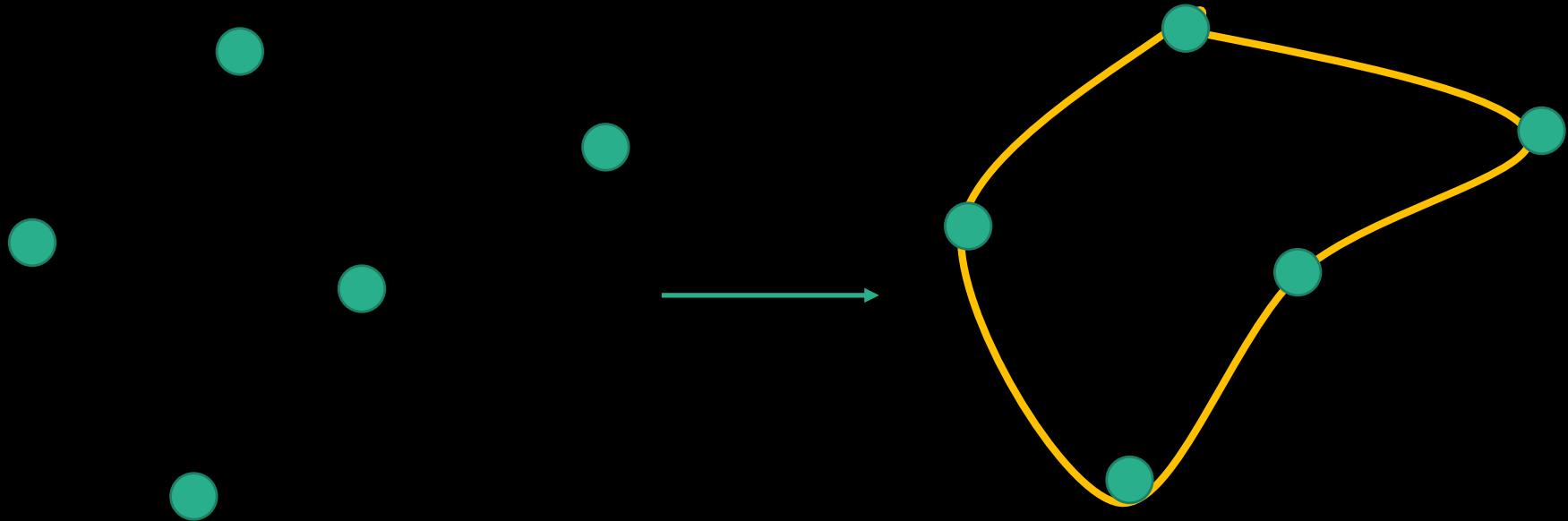
>>> nib.save(nib.Nifti1Image(mri, affine=np.eye(4)), 'mri_uint16.nii.gz')

>>> nib.save(nib.Nifti1Image(mask, affine=np.eye(4)), 'mask_uint8.nii.gz')

>>> print(f"Saved shapes: {mri.shape}, {mask.shape}")
```

Saved shapes: (256, 256, 120), (256, 256, 120)

Change annotation format:



I have a key point files of (x,y,z),
Help me to convert them to closed area...

Your Code &
Data



Data Formation

Environment

Output formation

Model Definition

Fixing the interface

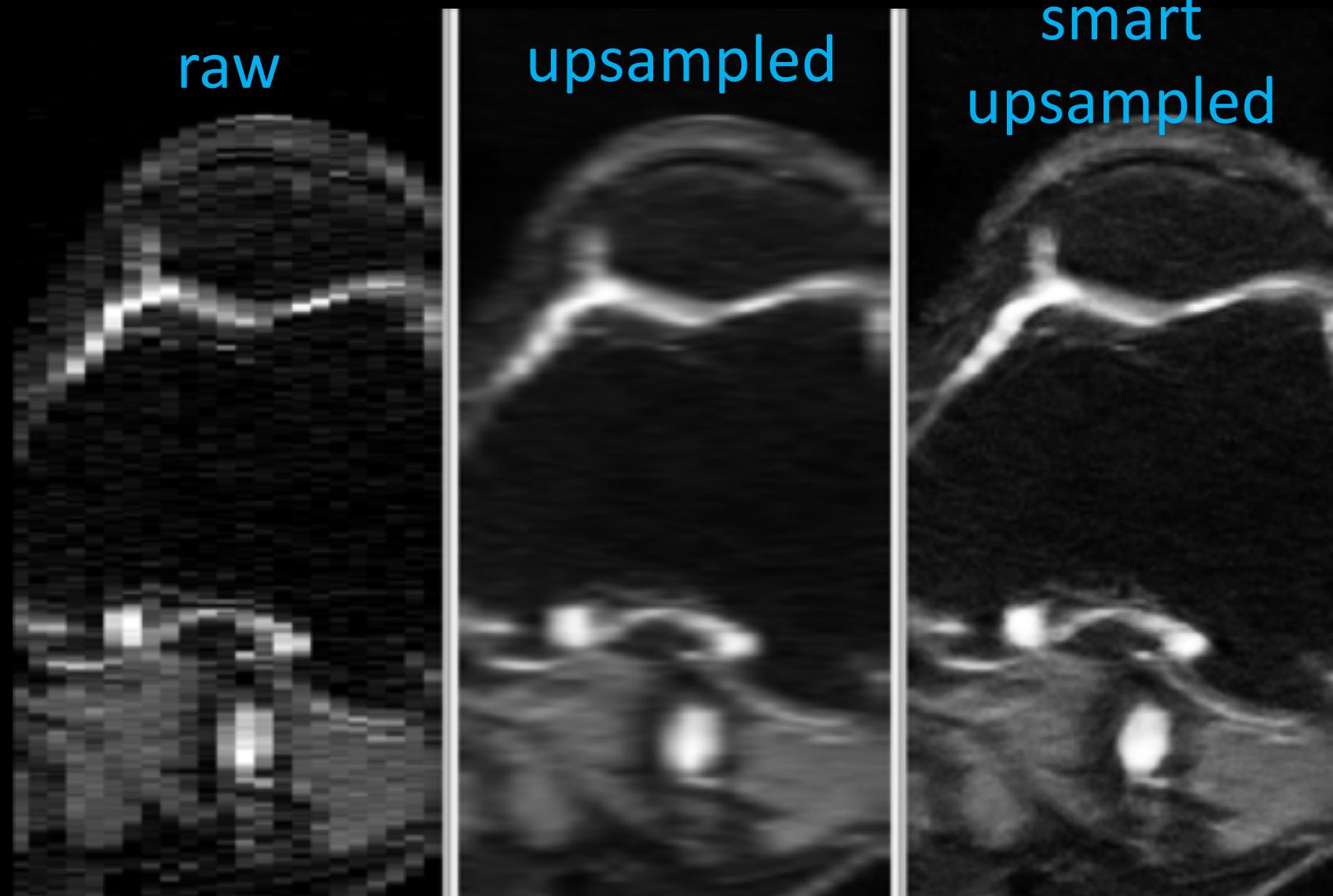
Existing Code

Think about what is
routine



What required
intelligence!

We don't usually operate in a 3D high-resolution space (i.e. anisotropic resolution):



Front
Annotation

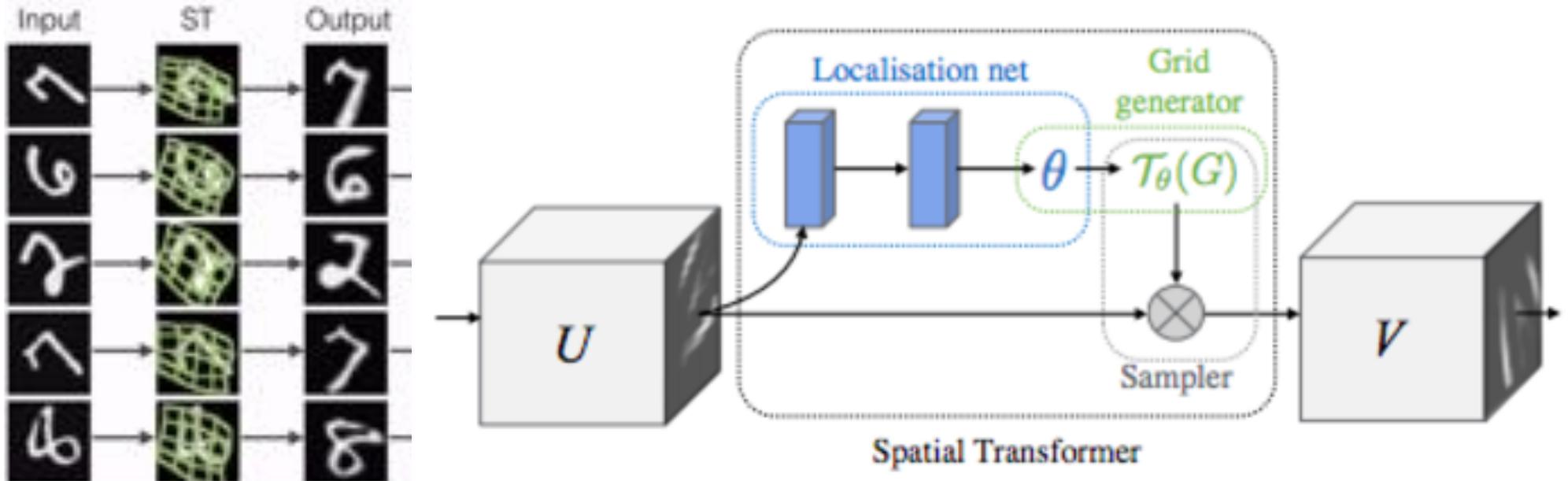


Side
Annotation

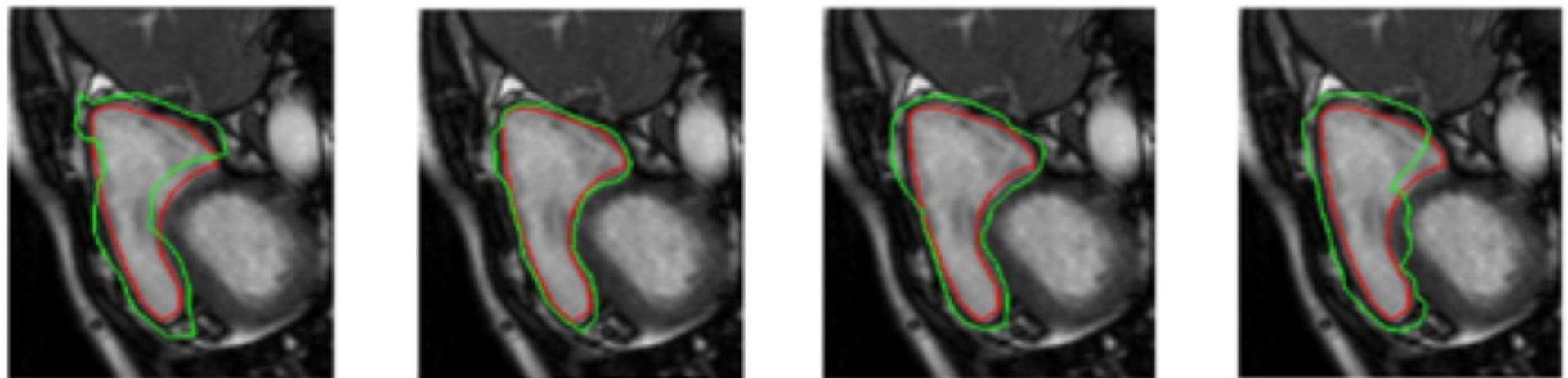
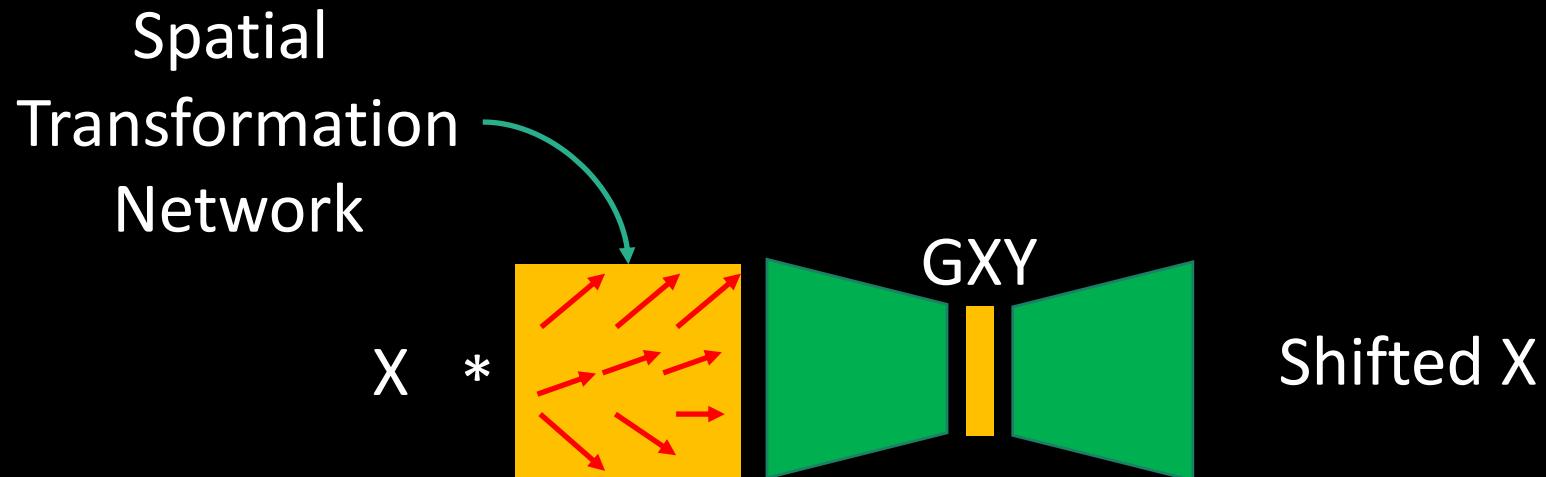


Resample?
Smoothing?
Multi-view
annotation?

Imaging Registration by Spatial Transformation

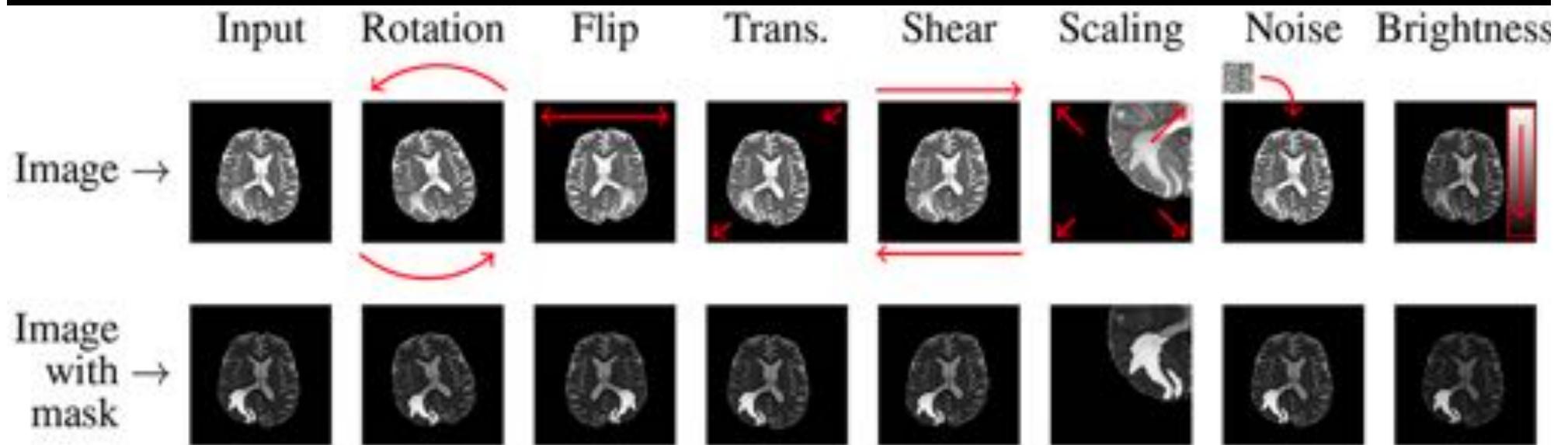


Imaging Registration by Spatial Transformation



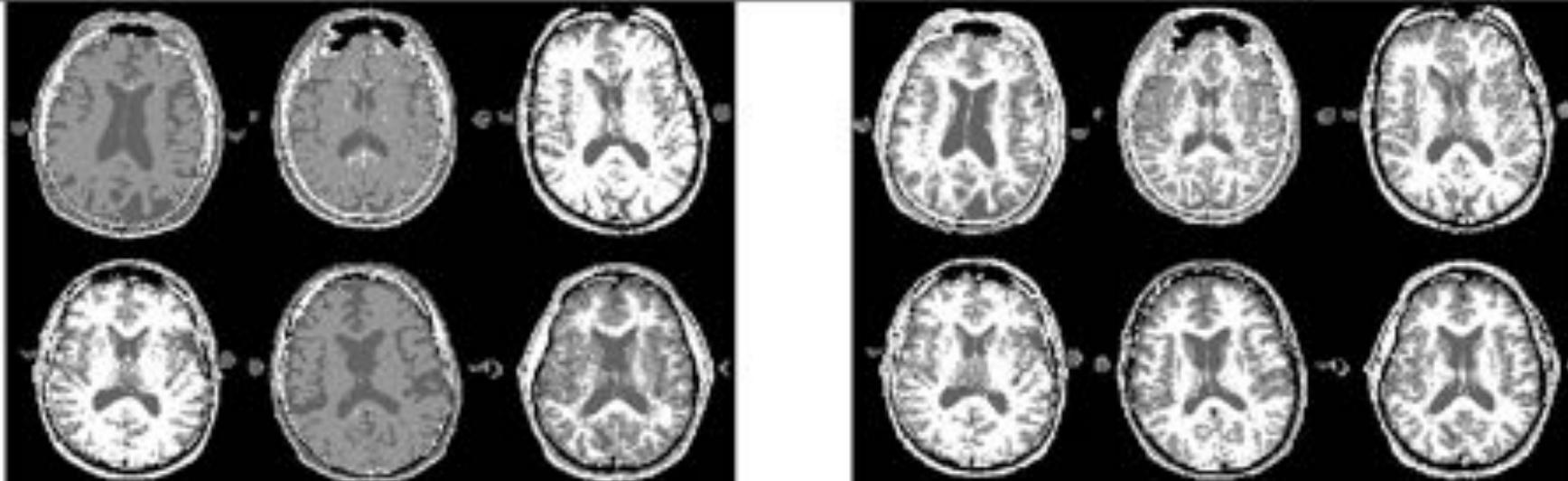
ELASTIC REGISTRATION OF MEDICAL IMAGES WITH GANS

Data Augmentation:

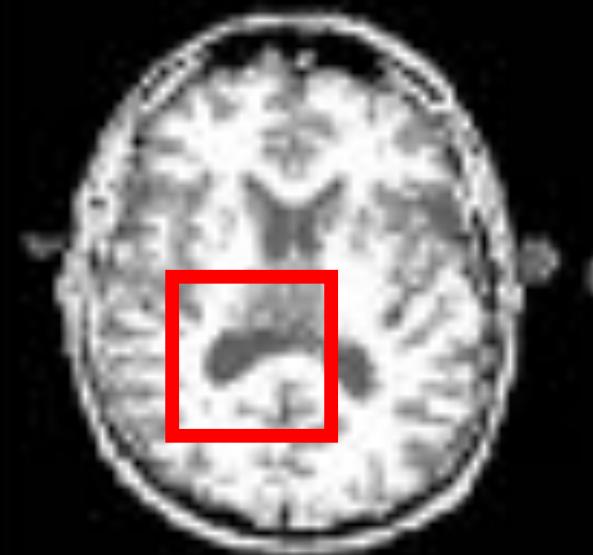


Intensity Normalization / augmentation

paper read at IEEE ISBI 2004

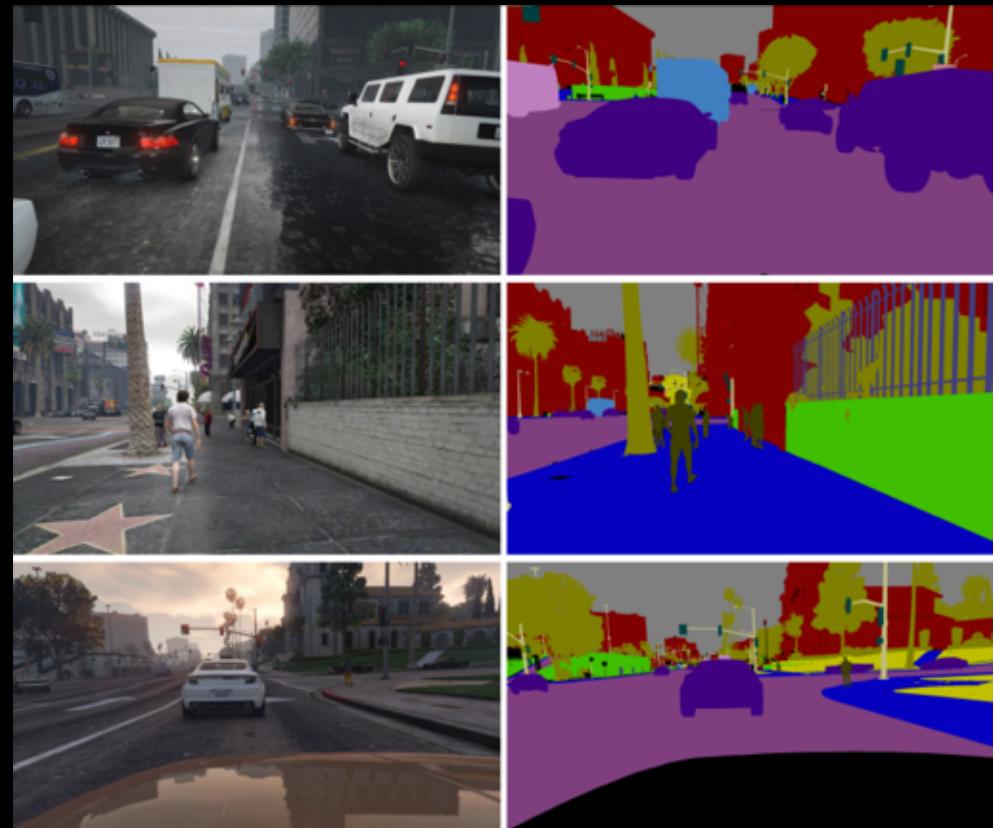
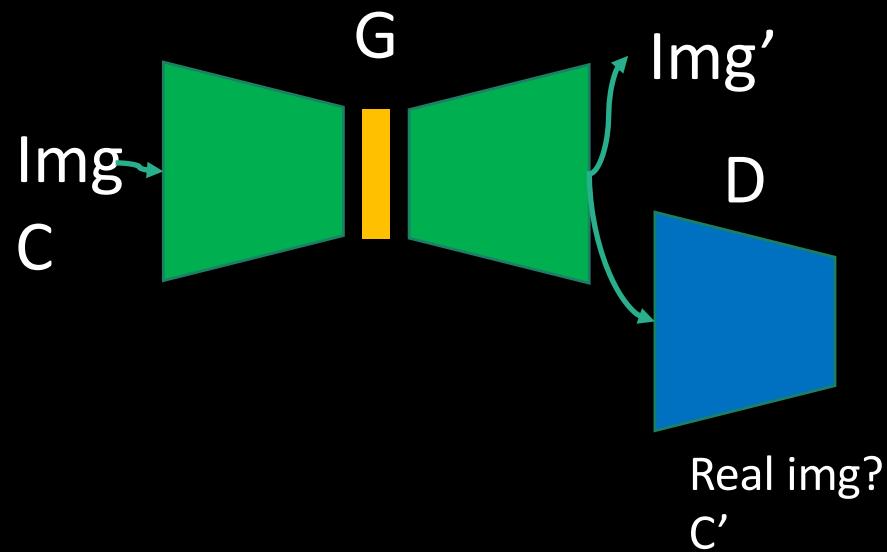


Patch size

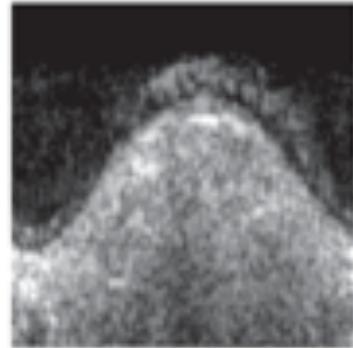


“Smart” Data Augmentation:

Image + condition GAN



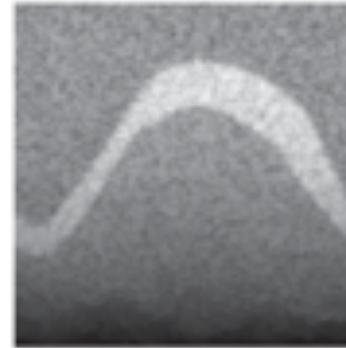
Mask Conditioned Image Generation



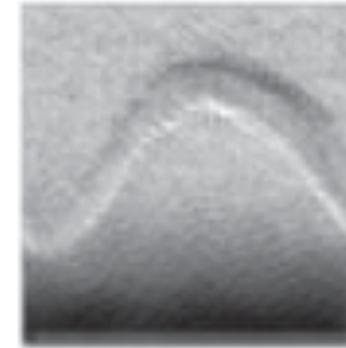
(a) Real image



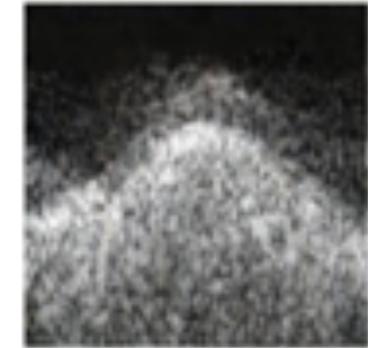
(b) Tissue map



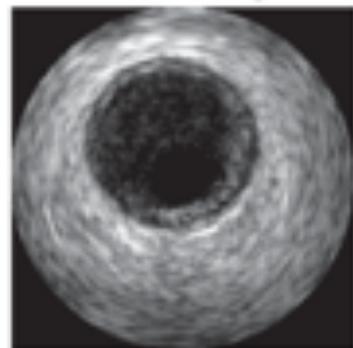
(c) Pseudo B-Mode



(d) Stage I GAN



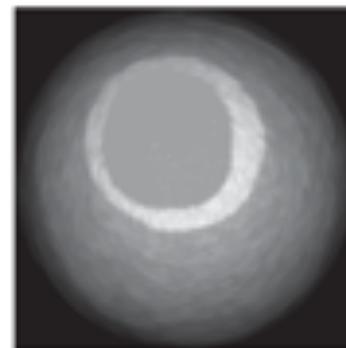
(e) Stage II GAN



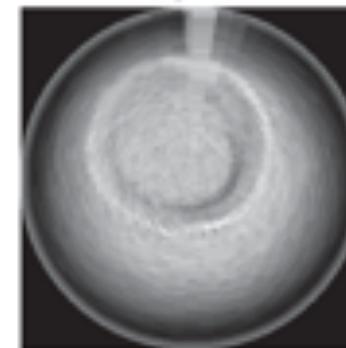
(f) Real image



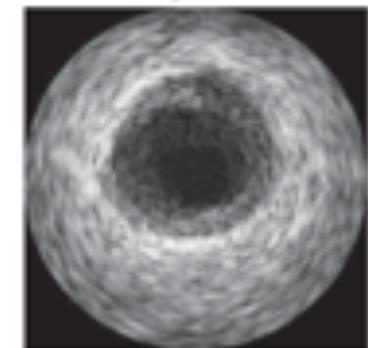
(g) Tissue map



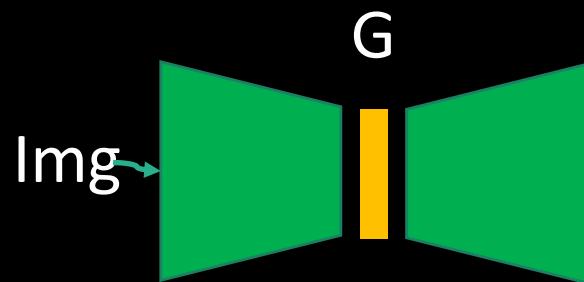
(h) Pseudo B-Mode



(i) Stage I GAN

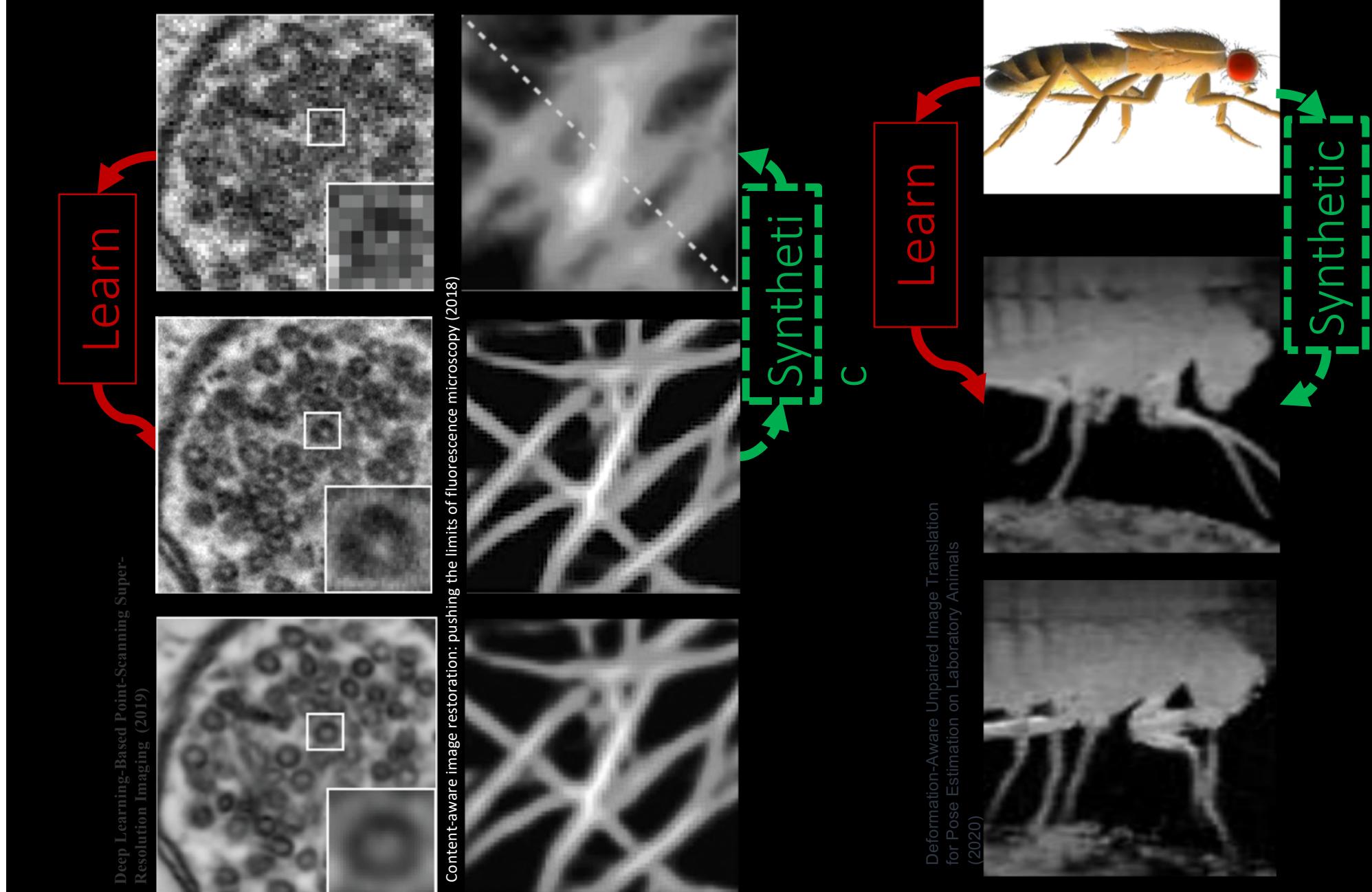


(j) Stage II GAN

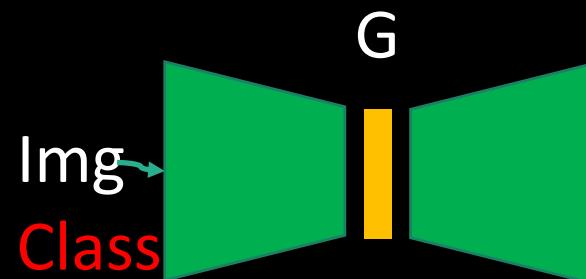
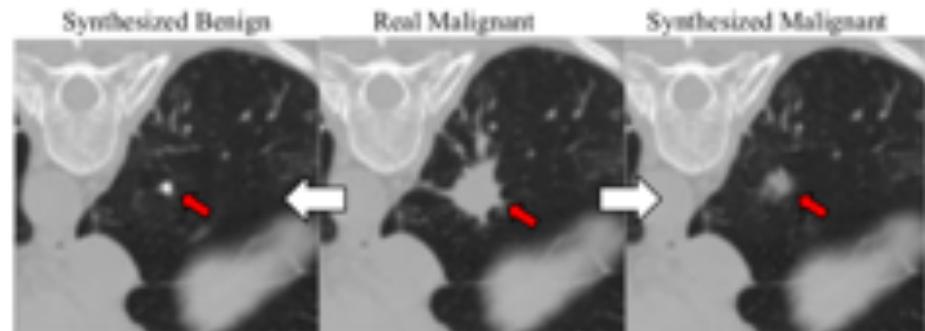
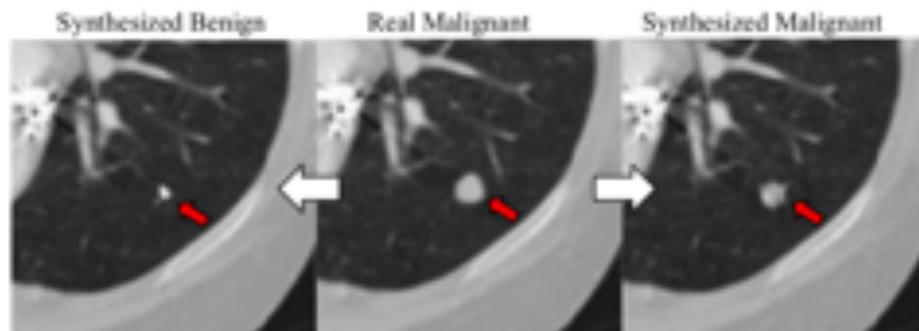
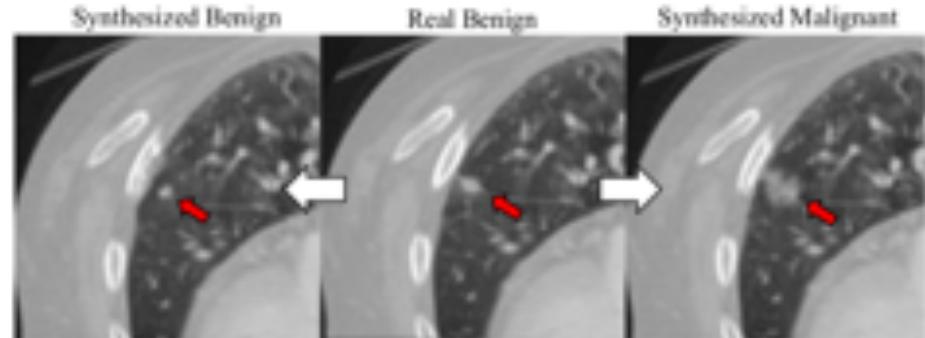
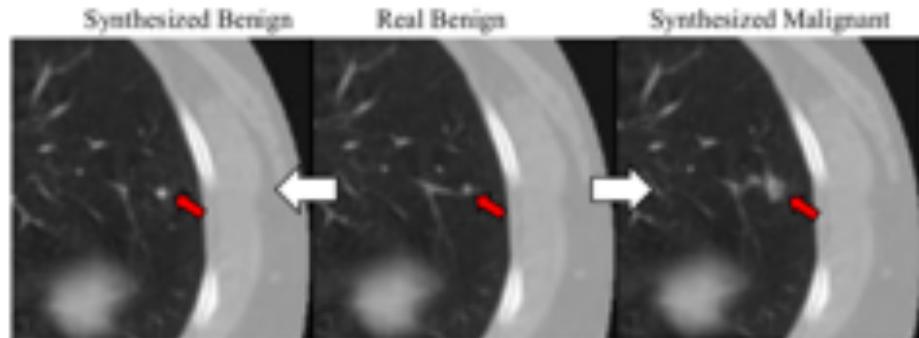


You got free images with free mask!

Mask Conditioned Image Generation (2)



Class + Mask Conditioned Image Generation



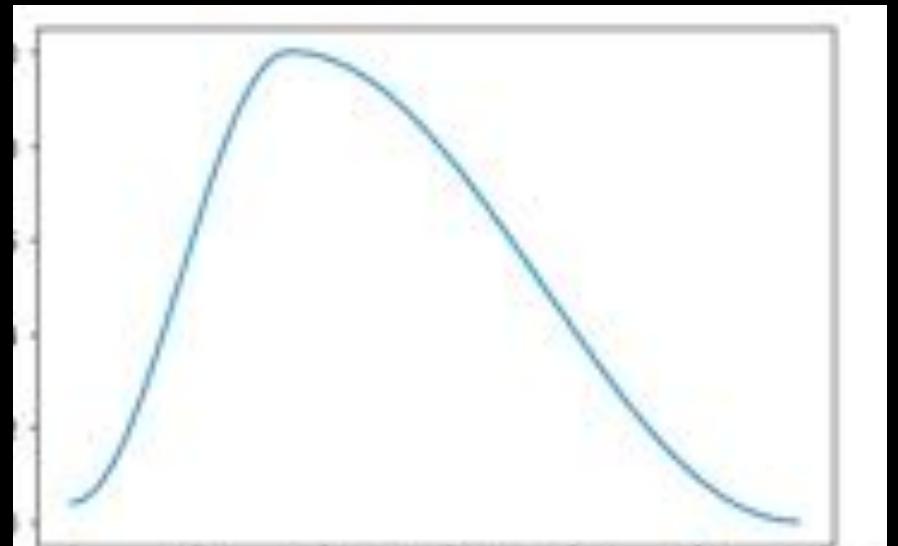
CLASS-AWARE ADVERSARIAL LUNG NODULE SYNTHESIS IN CT IMAGES (MICCAI 2020)

Finally, Hyper-parameters Tuning

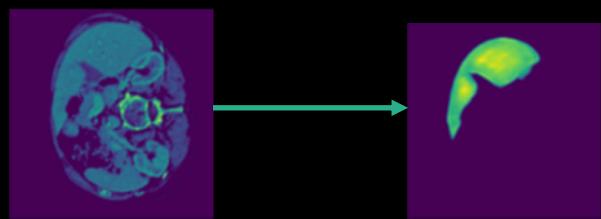
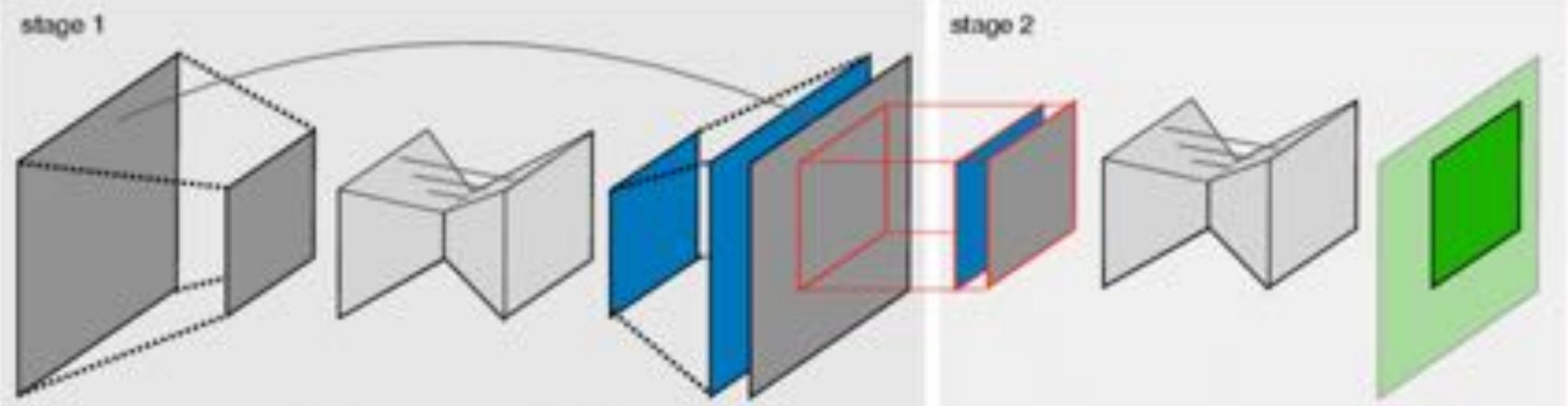
Optimizer: SGD with momentum

Loss: Dice & Cross Entropy

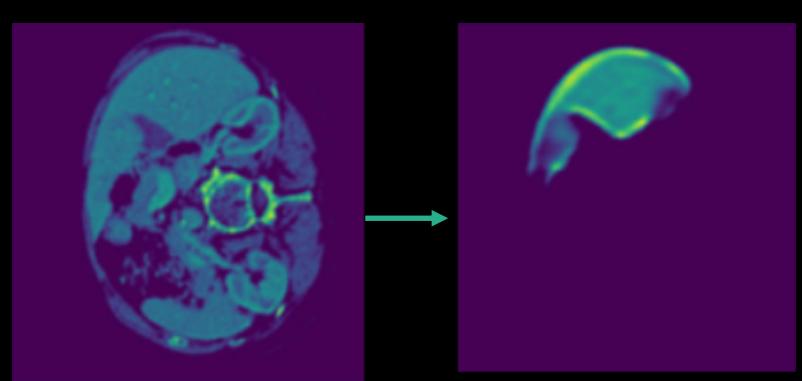
Learning Rate: Polynomial Learning Rate



“Cascade” Segmentation Model

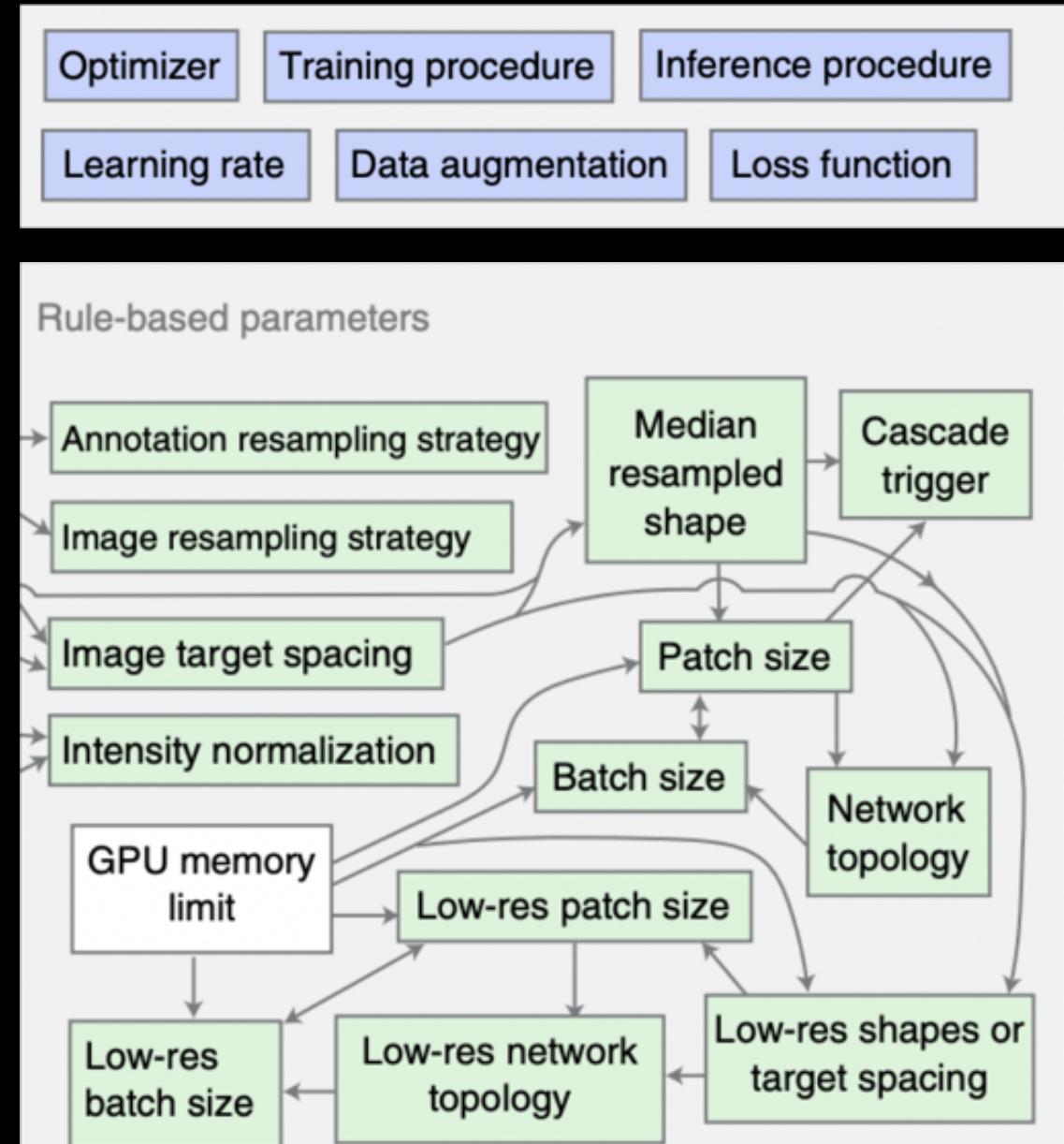


One stage for overall shape



One stage for details

nn (not new) unet



Detection / Instance Segmentation



Class: cell

X = 60 pts

Y = 40 pts

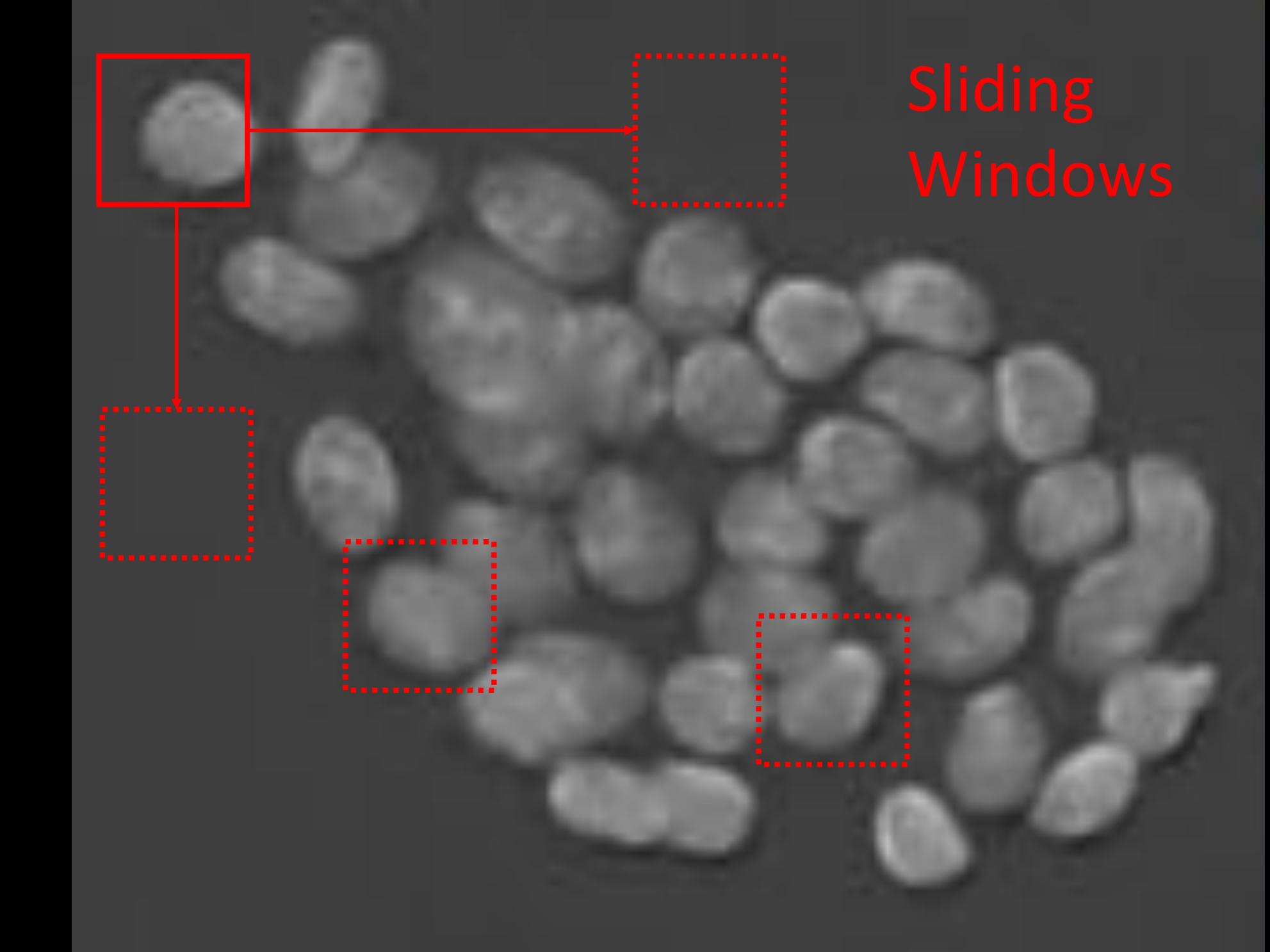
W = 5.2 pts

H = 4.3 Pts

Multi Class
Grade A 30%
Grade B 30%
Grade C 40%

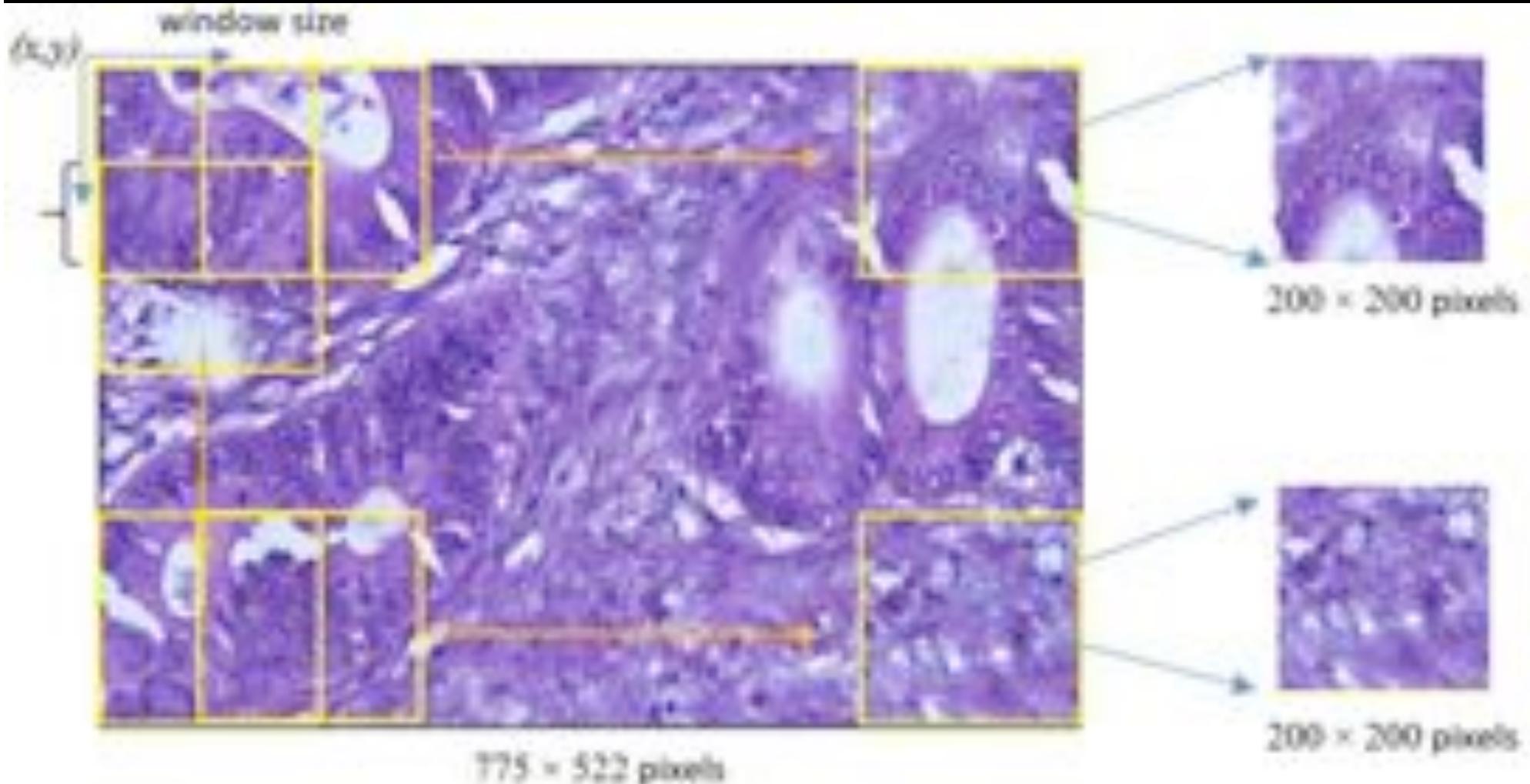
X = 60 pts
Y = 40 pts
W = 5.2 pts
H = 4.3 Pts



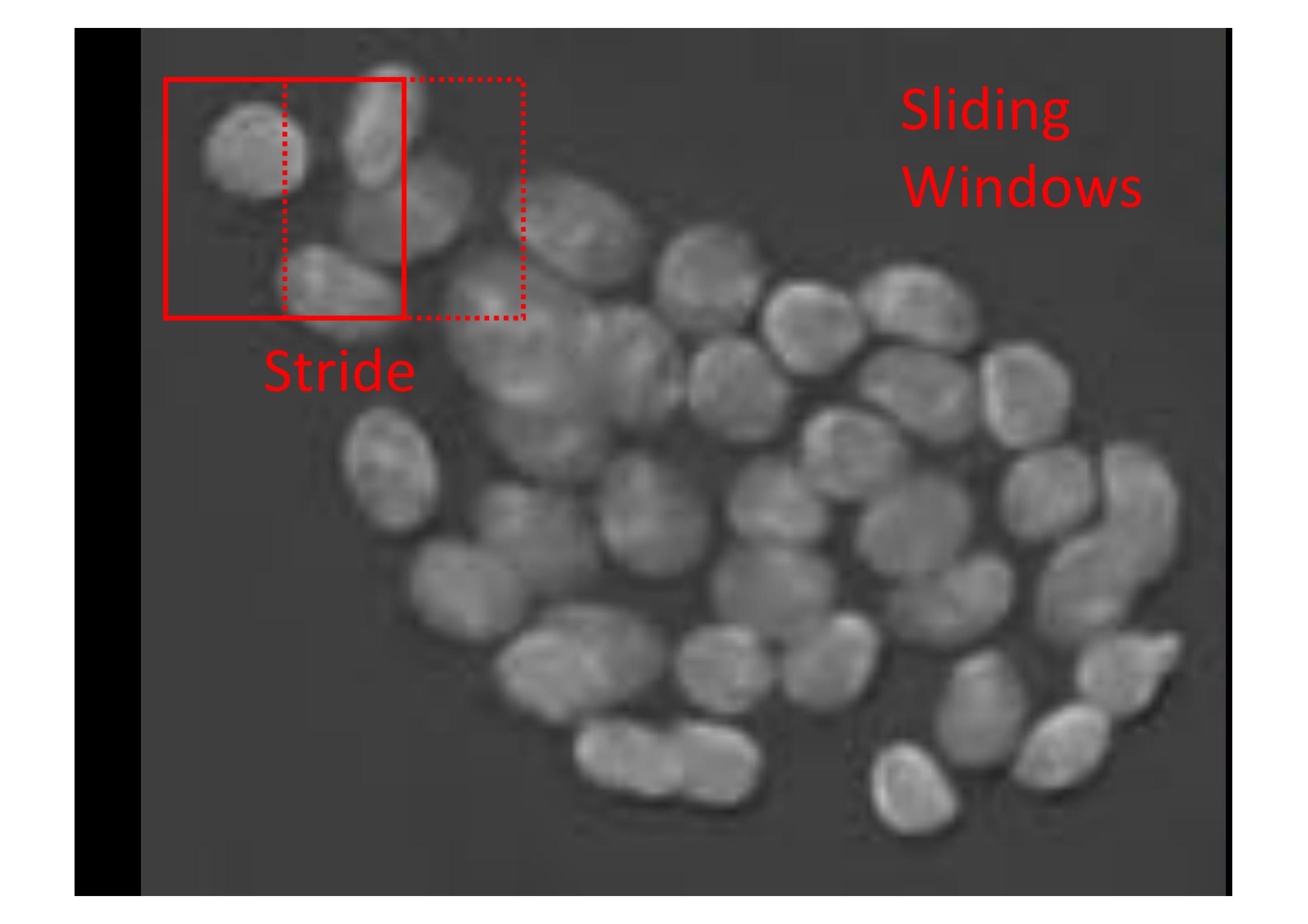


Sliding
Windows

Sliding Windows



Conditional sliding windows: An approach for handling data limitation in colorectal histopathology image classification

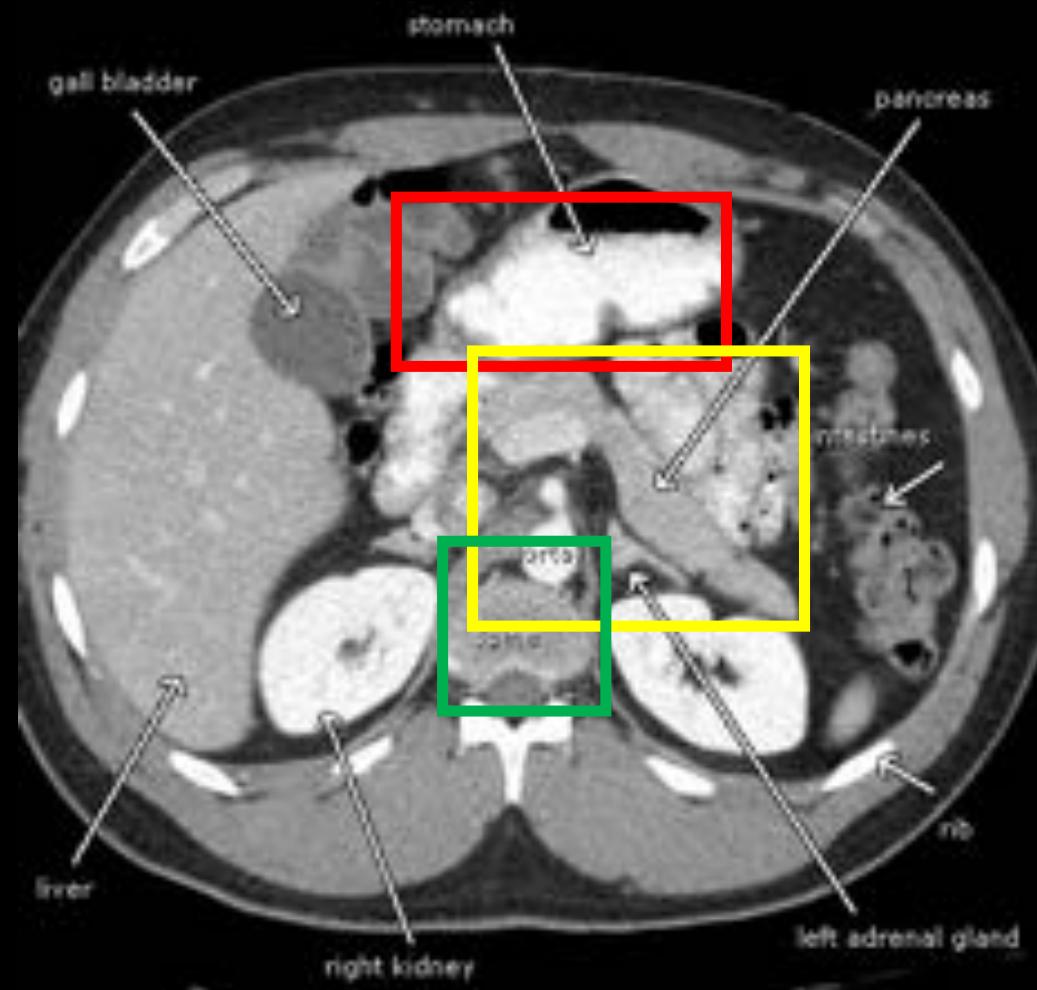


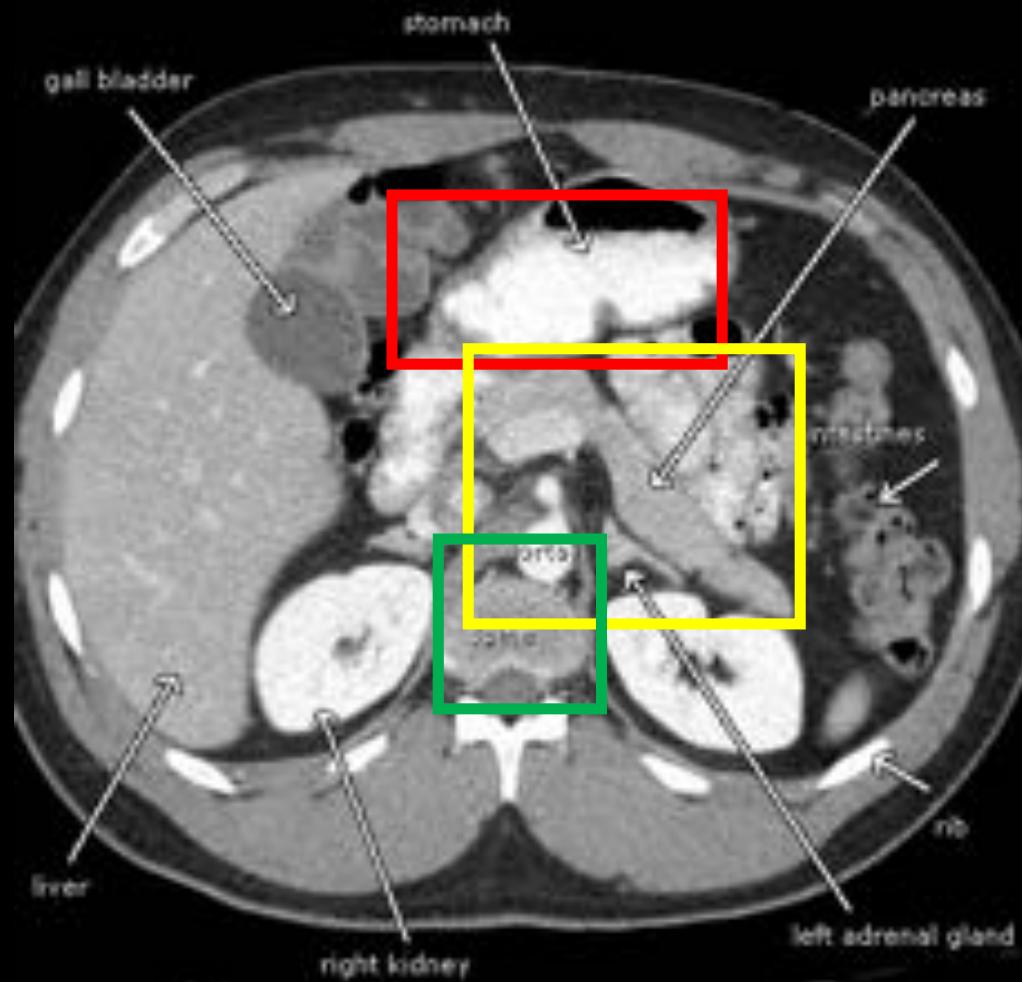
Sliding
Windows

Stride

Sliding Windows Shortcomings:

- slow
- no global context



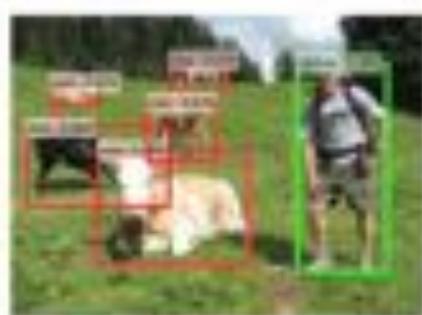
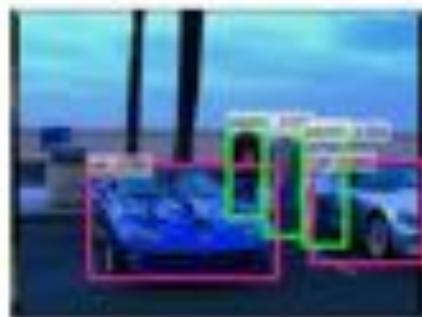
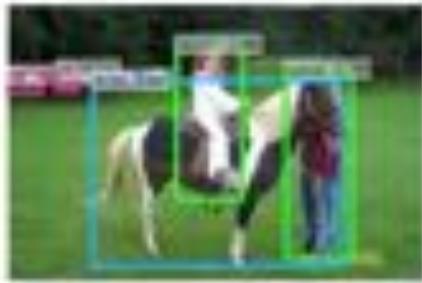


Different size

Different aspect ratio

May be overlapping

(Sometimes) need
global context

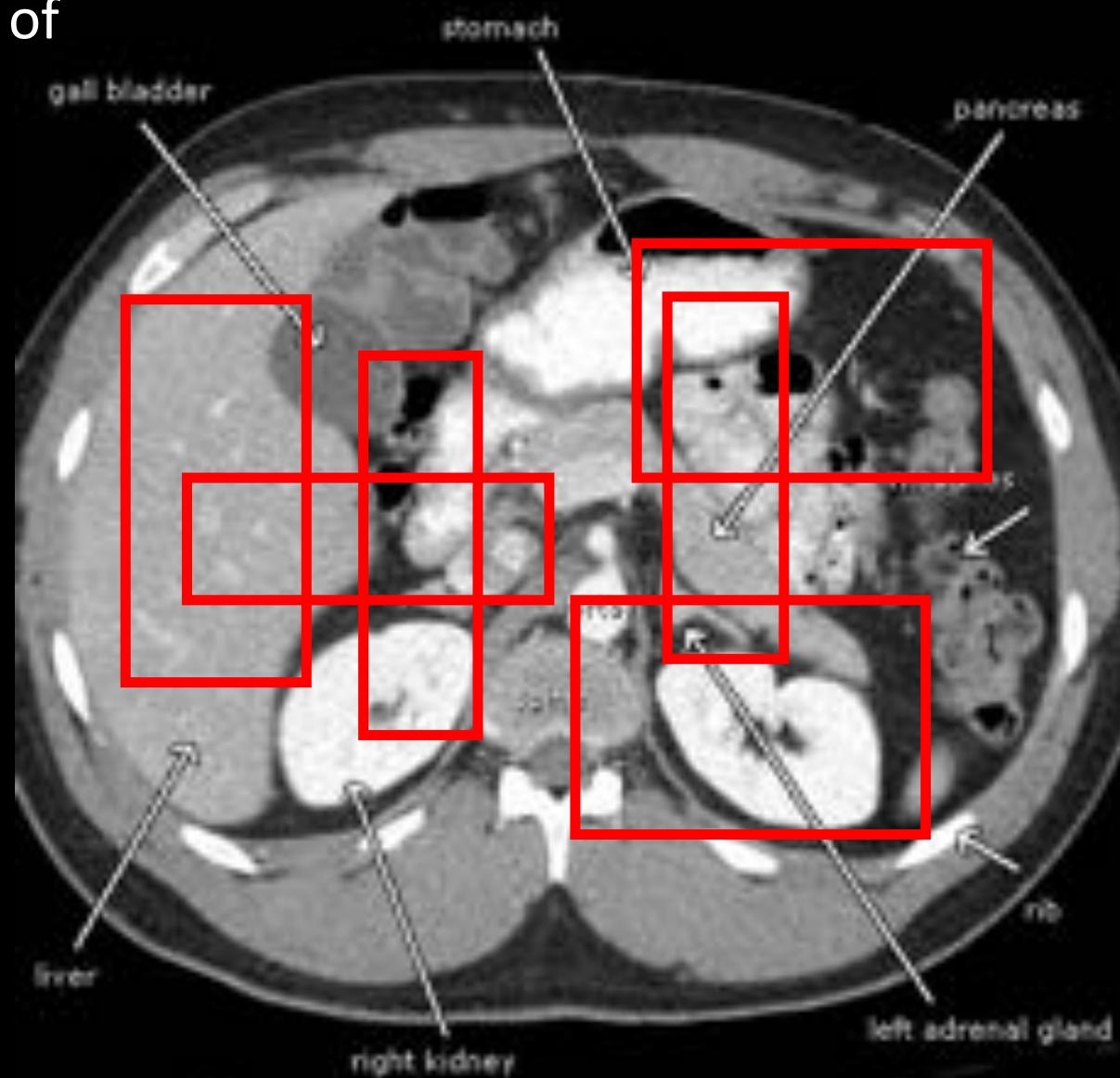


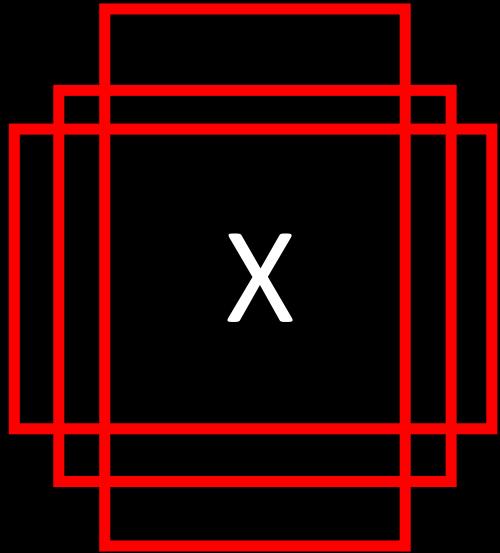
We can propose a bunch of Regions with different

Location,

Size,

And aspect ratio

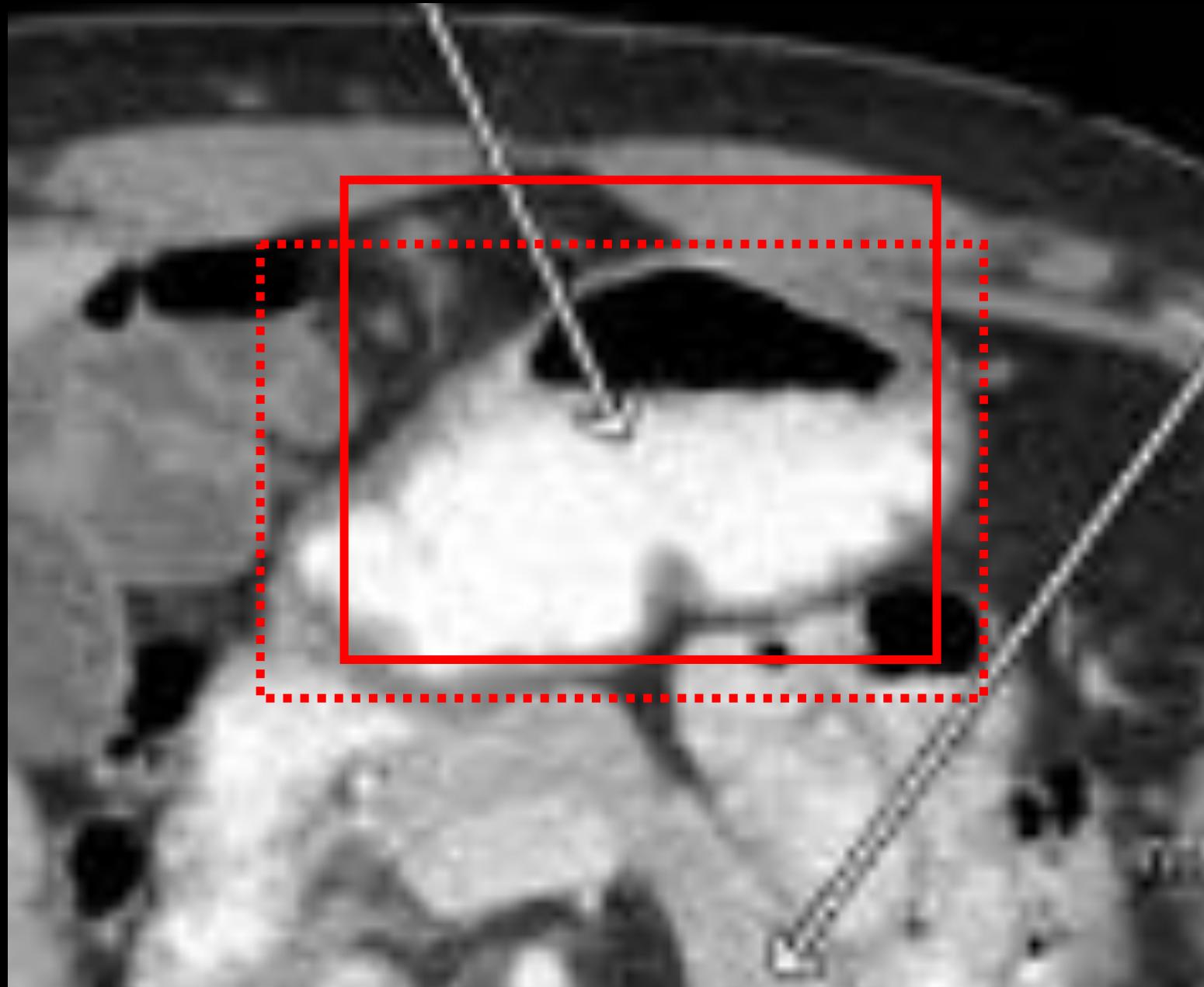




All the bounding boxes around the same regions should share some similar information

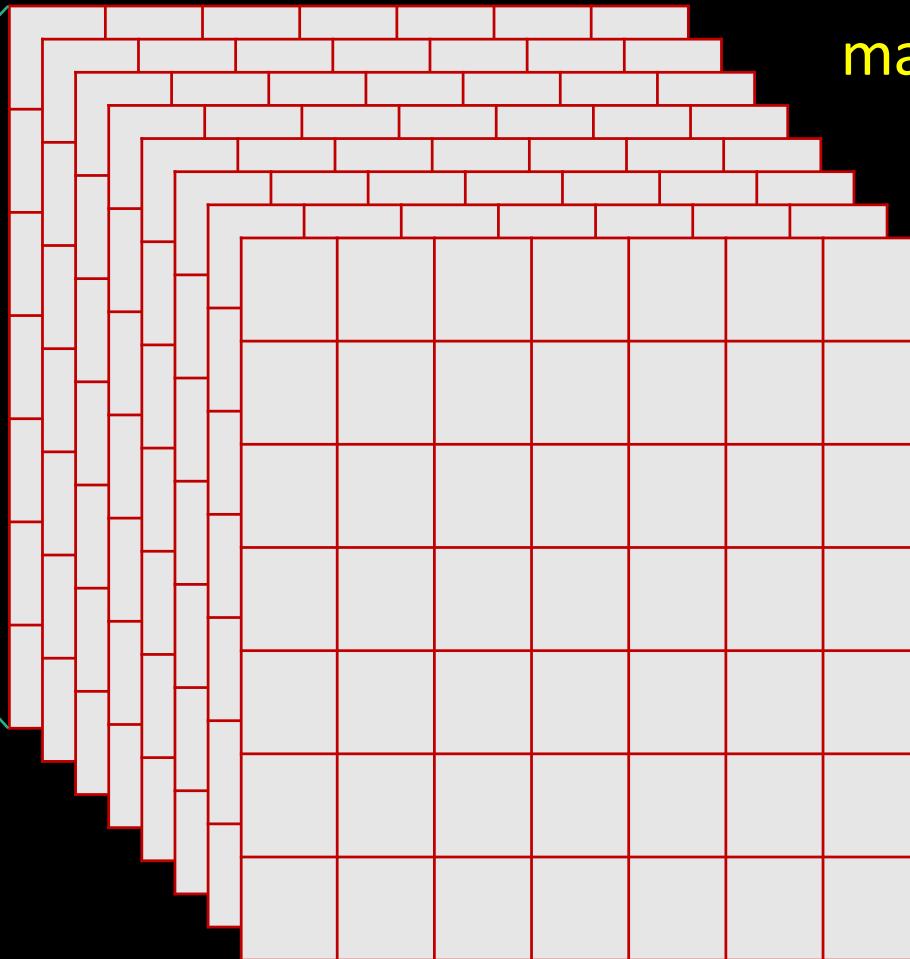
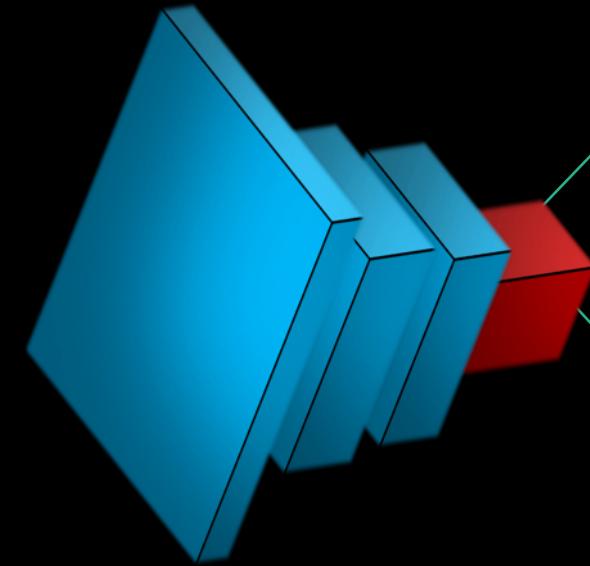
We call the center an “Anchor”

We can propose “rough” bounding boxes then refine it



Latent Space often contained a rough spatial information (for example, in a 7x7 grid)

Latent
feature
map

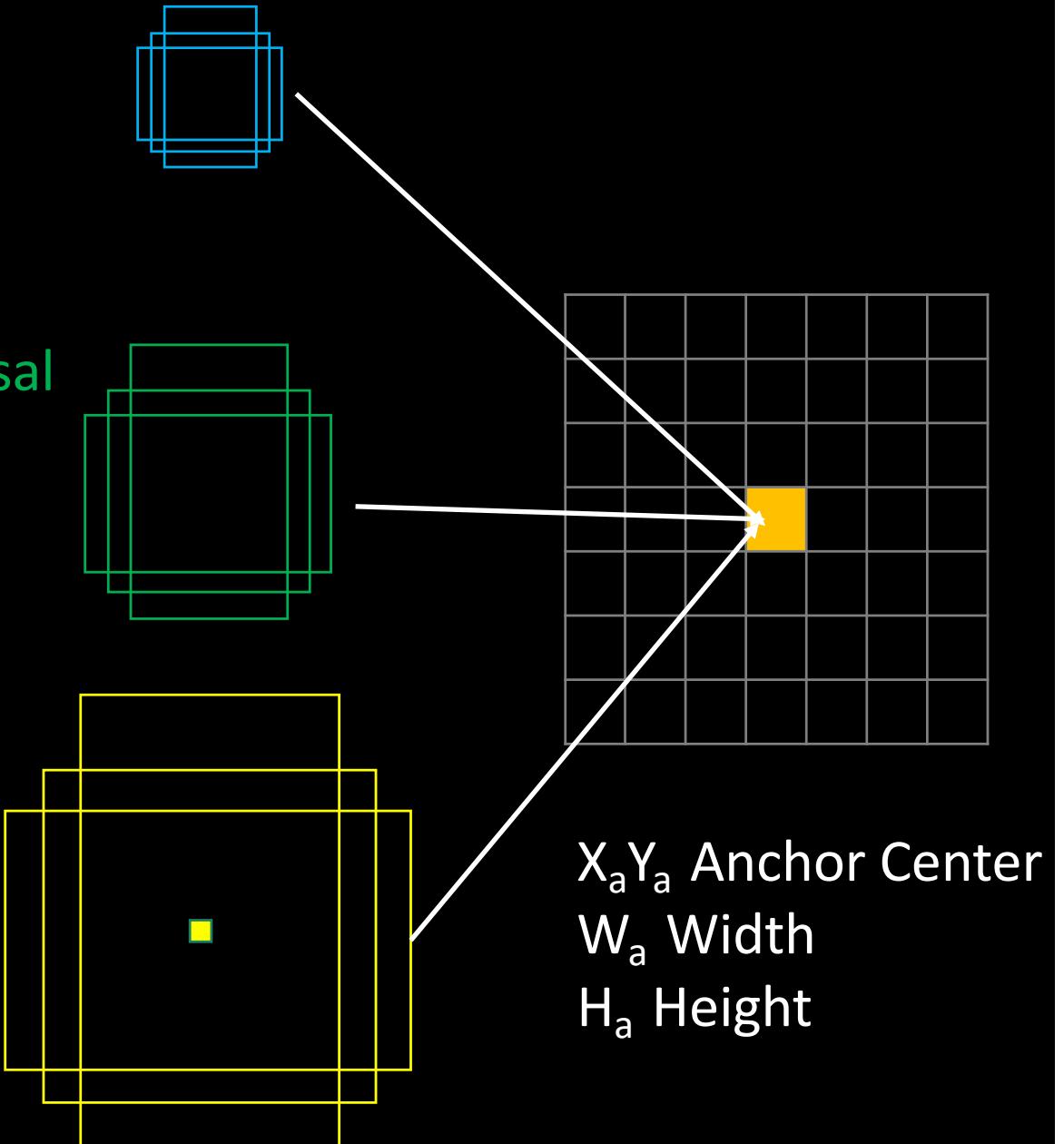
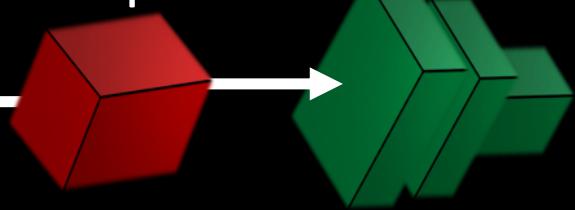


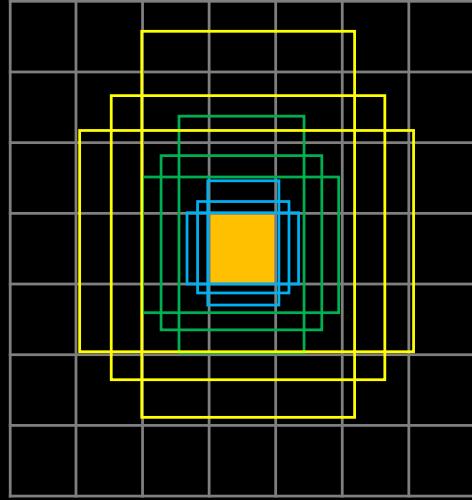
From a
backbone
model

Encoder
(Backbone)

Latent
feature
map

Regional Proposal
Network(RPN)





For every anchor point:



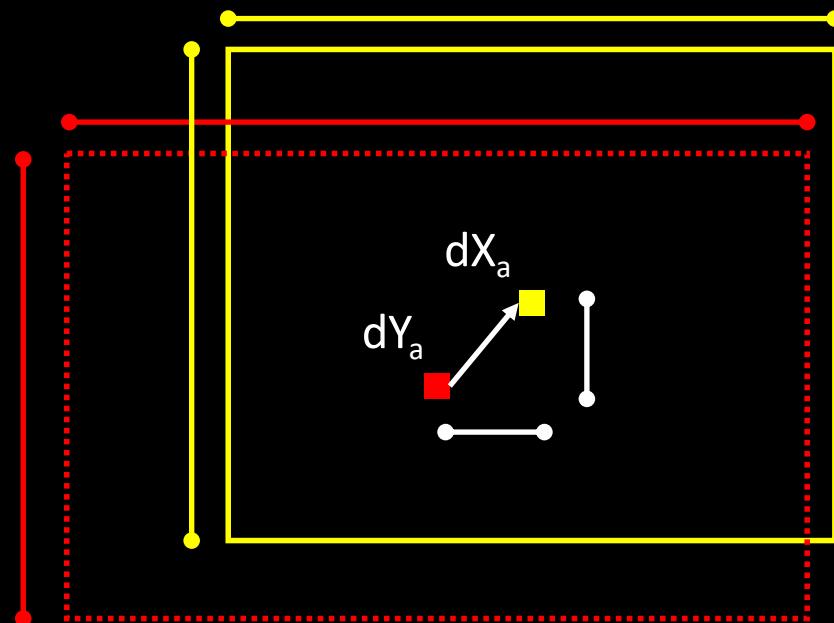
Q1: Is it an object?:

9 (boxes) * c (classes) output probability

Bounding box regression

Q2: How far should the box be adjusted?

9 * 4 (x, y, w, h) regression values



dX_a dY_a Anchor Center

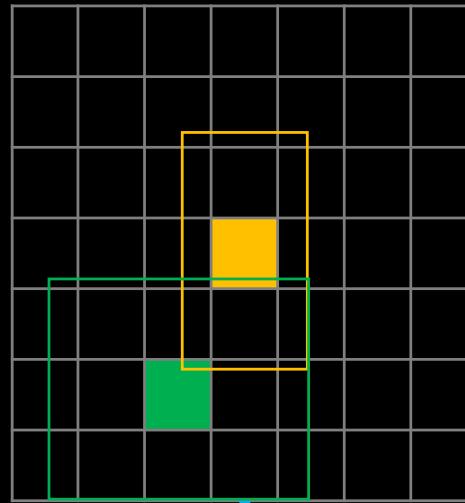
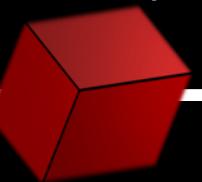
dW_a Width

dH_a Height

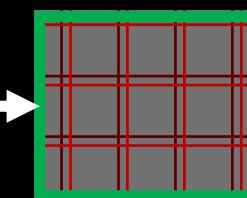
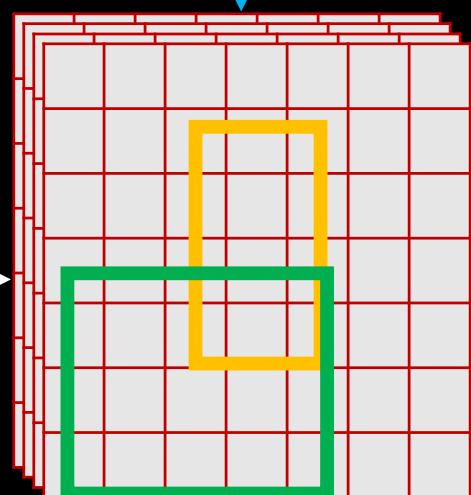
Regional Proposal Network(RPN)

Encoder
(Backbone)

Latent
feature
map

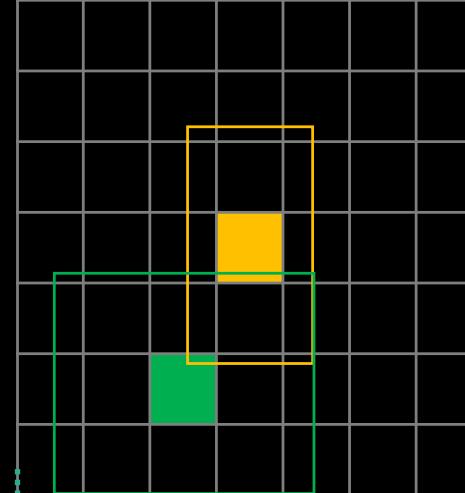
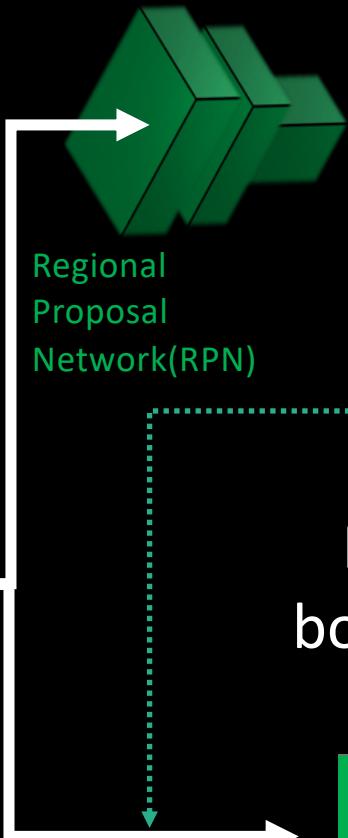


Aligned ROI Features

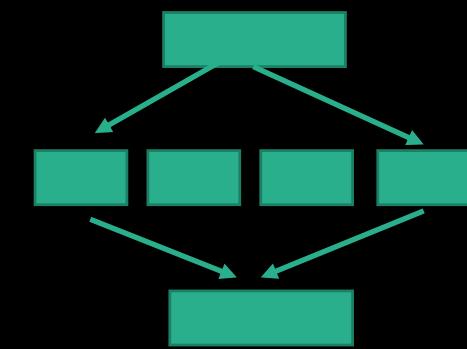
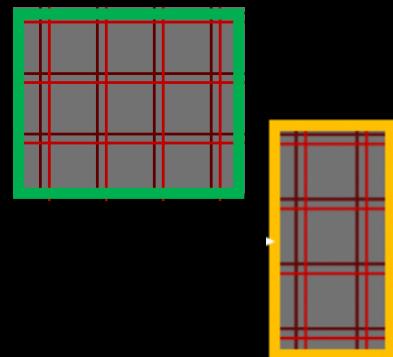


Encoder
(Backbone)

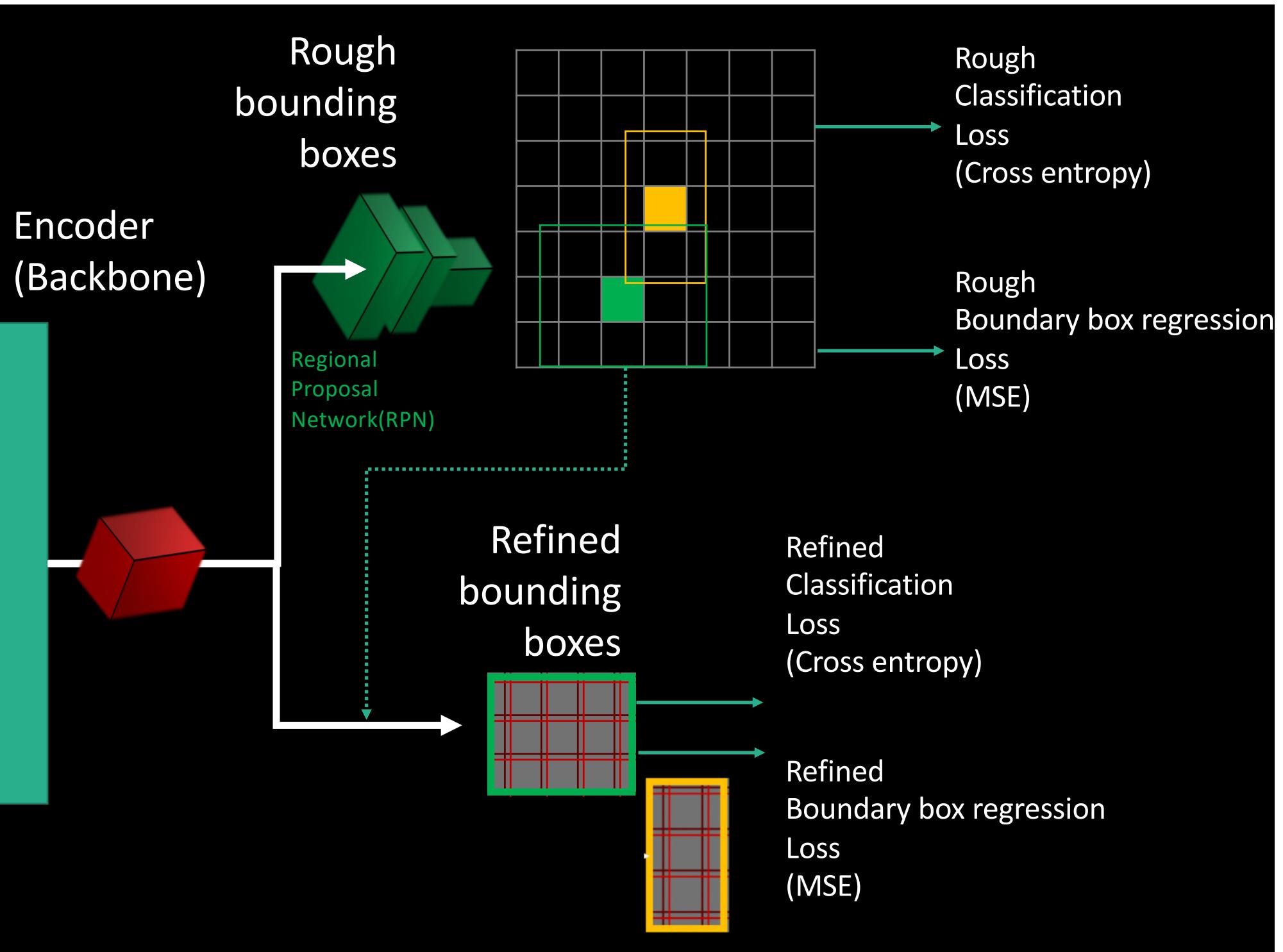
Rough
bounding
boxes

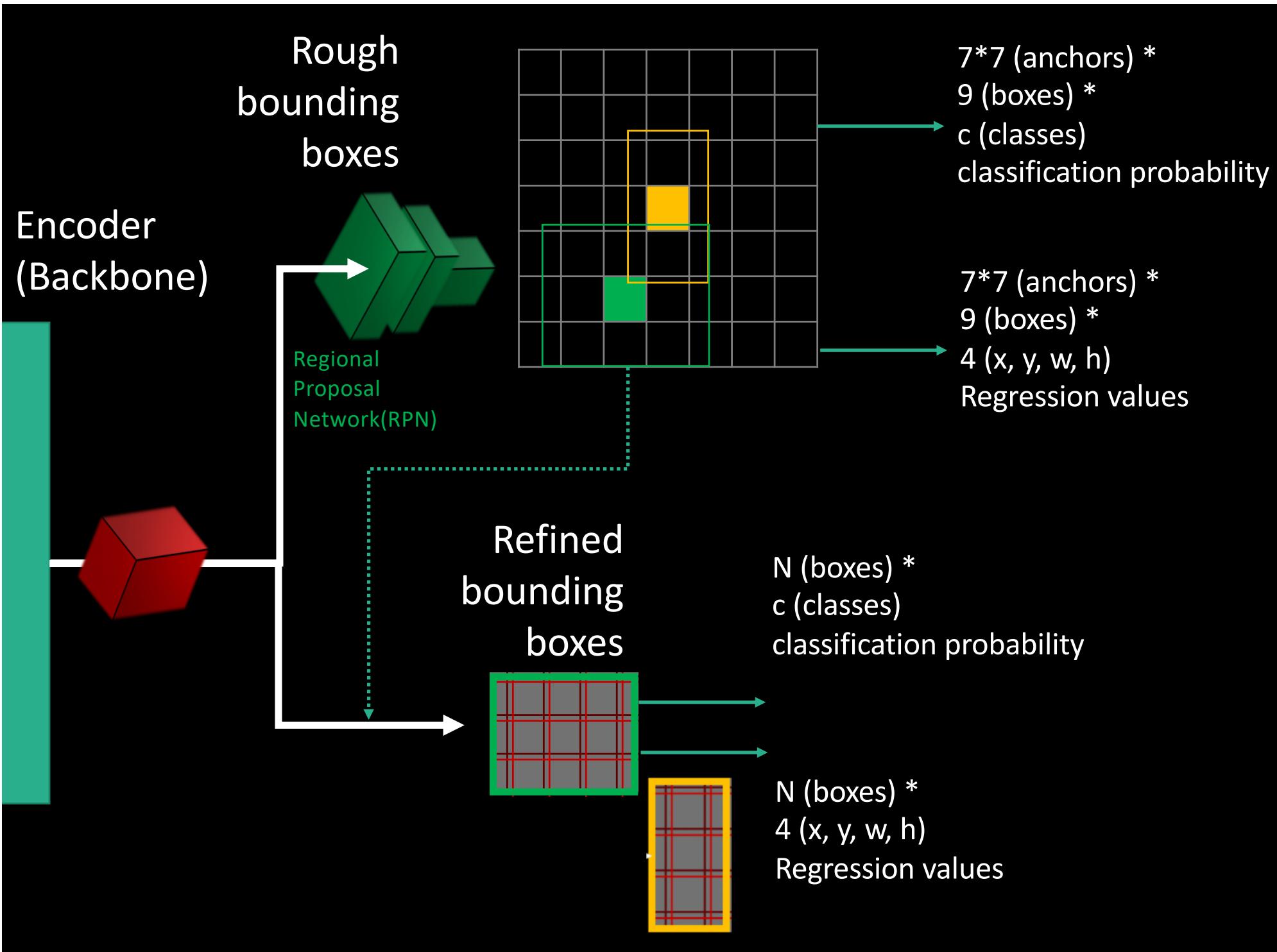


Refined
bounding
boxes



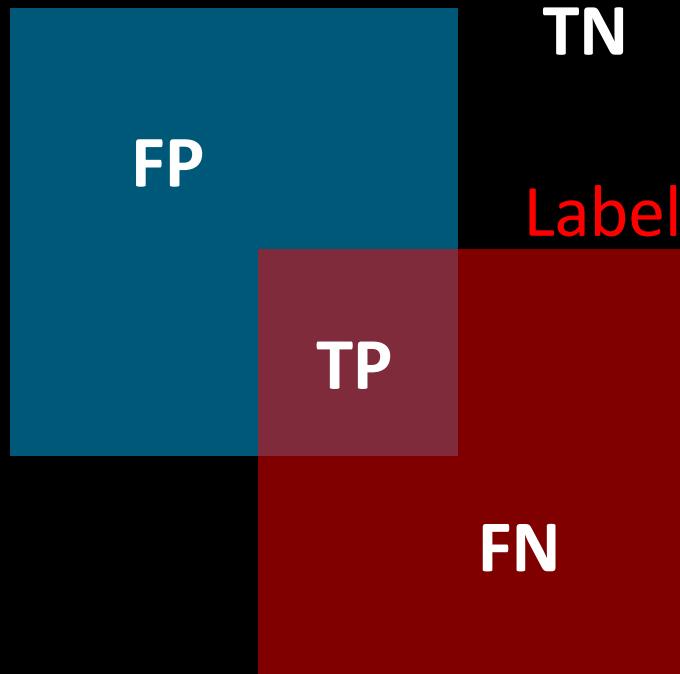
It's kind of multi
branch





Performance Metrics of Segmentation

Prediction

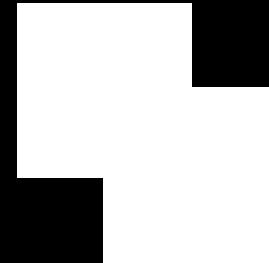


Label

N

IoU

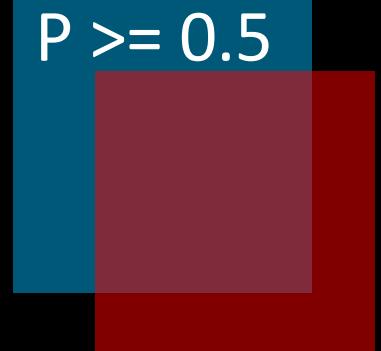
TP



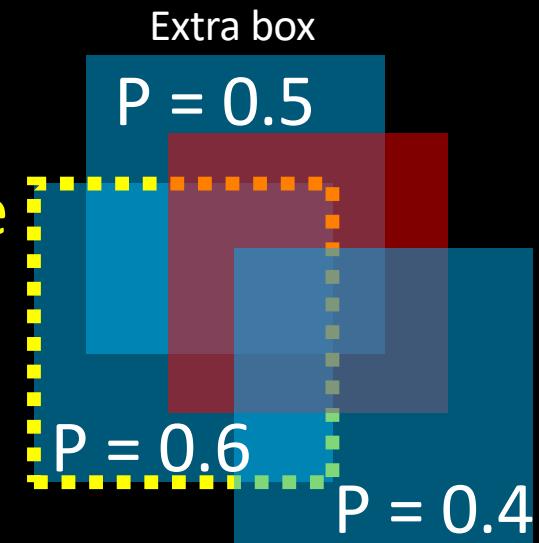
Performance Metrics of Bounding Boxes



True positive



False positive



False negative



Typical Components of a Instance segmentation Project

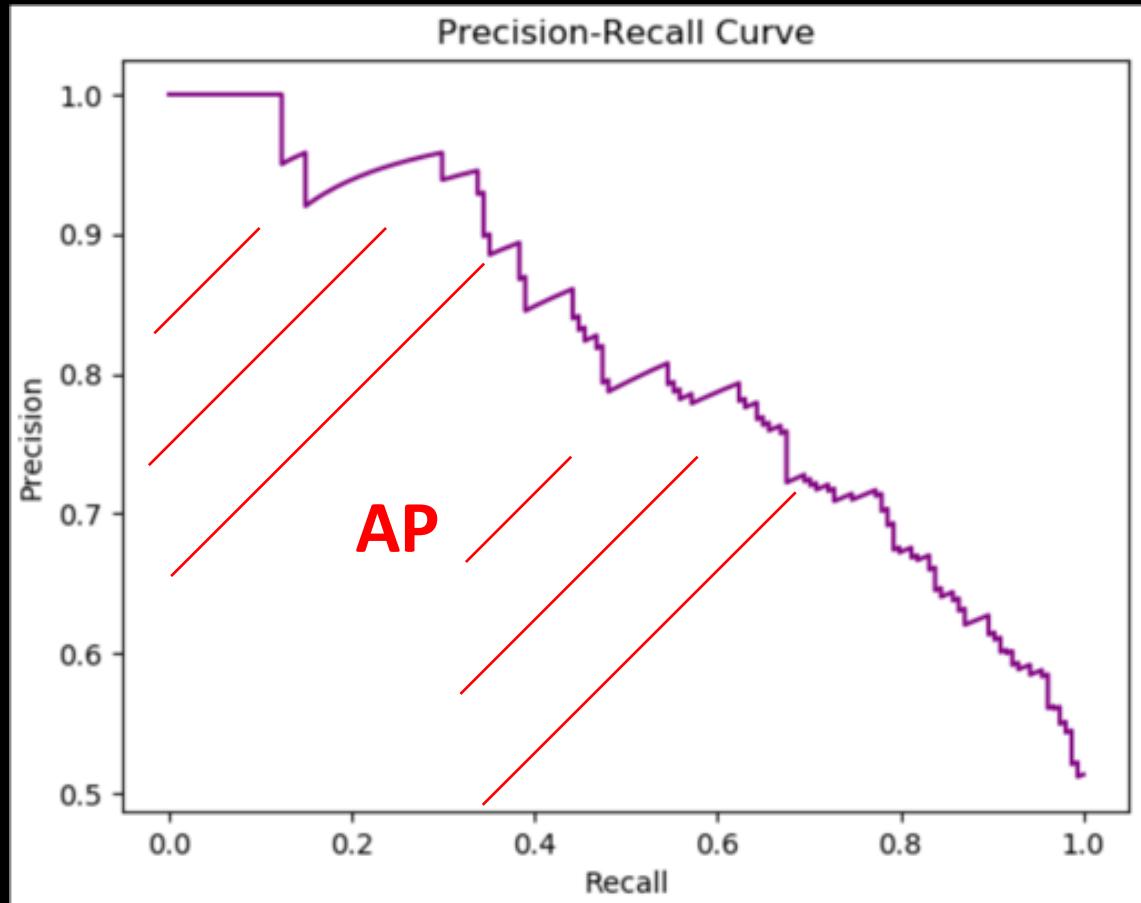
Task:	Data	Models	Metrics	Loss
Instance Segmentation	Image + Bounding boxes/ Segmentation mask + PREPROCESSING	EX: Resnet101 (backbone) + MaskRCNN (architecture)	IoU F1 Dice coeff.	Segmentation Loss Cross-entropy Bounding-box Regression loss

Precision

$$\frac{\text{TP}}{\text{TP} + \text{FP}}$$

Recall

$$\frac{\text{TP}}{\text{TP} + \text{FN}}$$



We often calculate mAP
(mean over all the classes)