

Deep Learning

Lecture 7

Recap

Architectures for classification

- ResNet
- VGG
- MobileNet
- EfficientNet
- ViT

Object detection

- RCNN
- Fast-RCNN
- Faster-RCNN

Semantic segmentation

Motivation



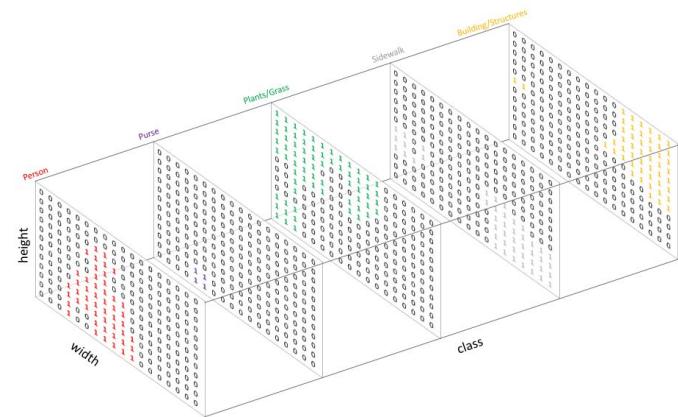
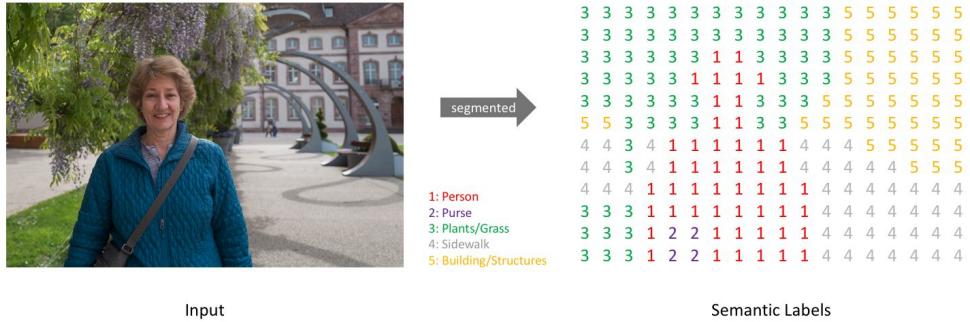
We have car, building on both images, but scenes are different!

Semantic segmentation



Through semantic segmentation we really understand the image

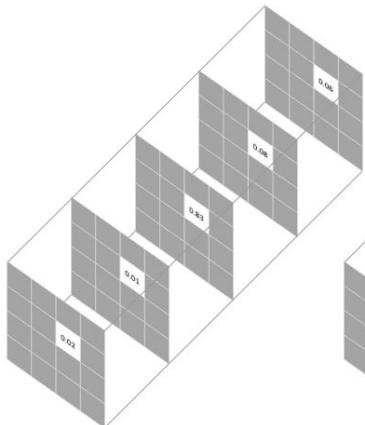
How to make semantic segmentation



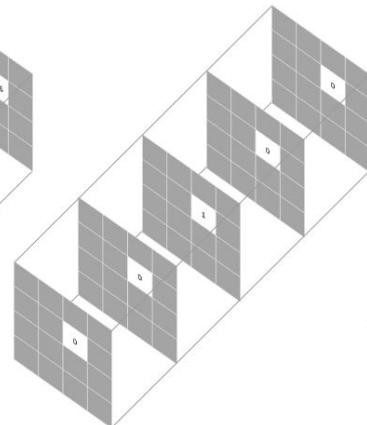
We set every pixel with class value \rightarrow Output $H \times W \times 1$

But neural network work poorly with discrete values? How to parametrize predictions?

Loss functions



Prediction for a selected pixel

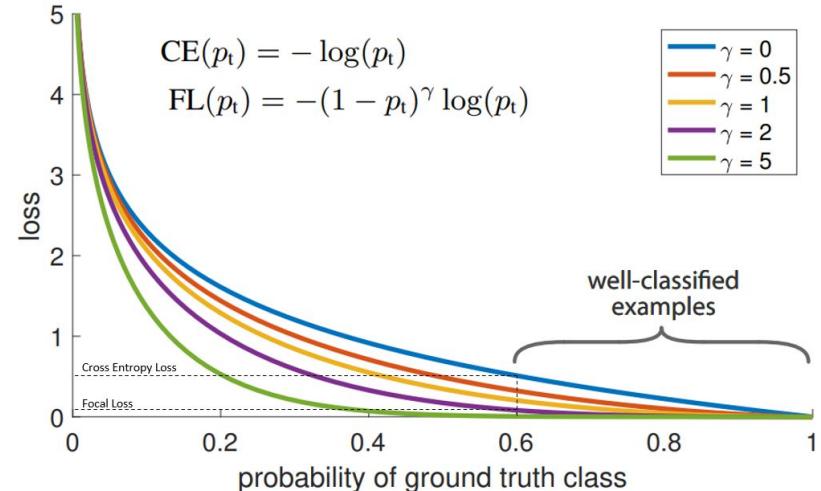


Target for the corresponding pixel

Pixel-wise loss is calculated as the log loss, summed over all possible classes

$$-\sum_{\text{classes}} y_{\text{true}} \log(y_{\text{pred}})$$

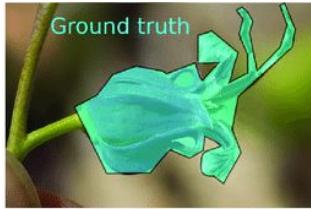
This scoring is repeated over all pixels and averaged



$$\text{FL}(p_t) = -\alpha_t (1 - p_t)^\gamma \log(p_t)$$

Why do we need alphas?

Metrics



Intersection



Union

$$\text{IoU} = \frac{\text{Intersection}}{\text{Union}}$$

Averaging over classes

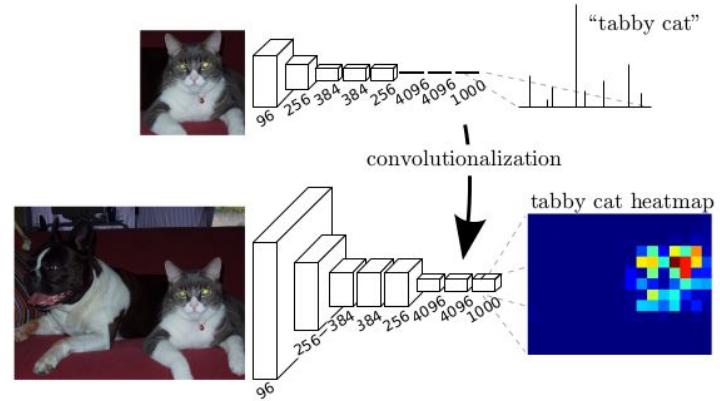
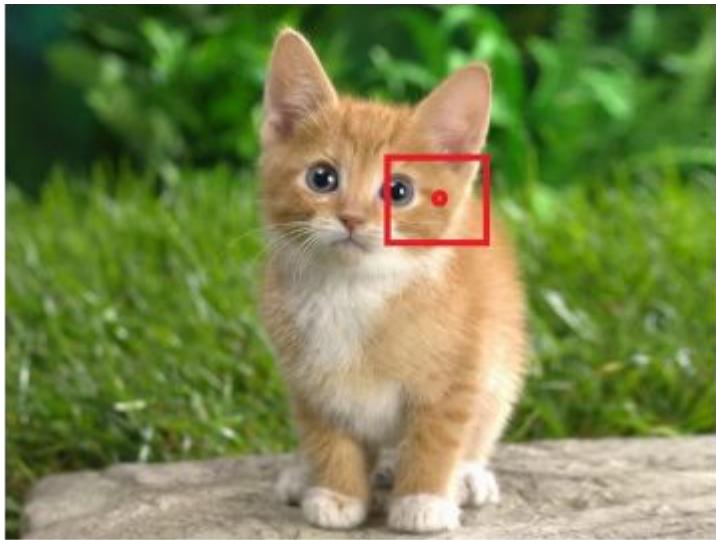
$mIoU$

$$accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$

Averaging over classes

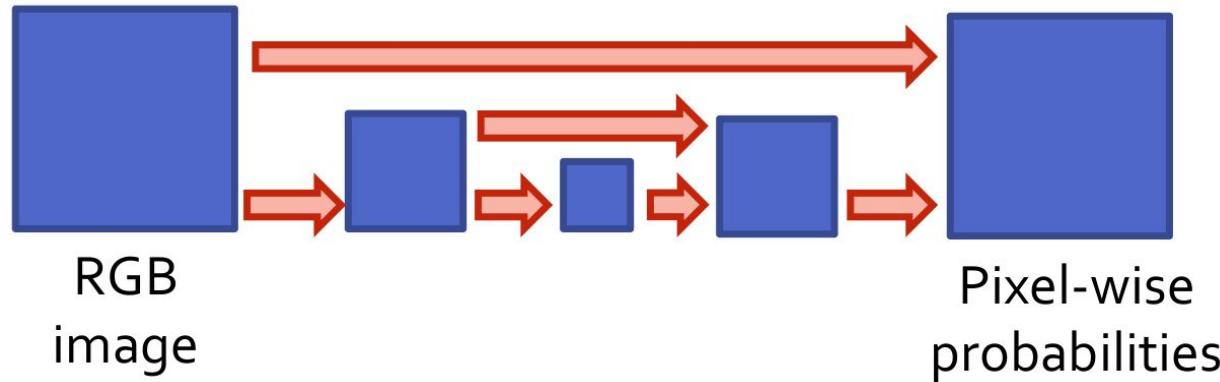
$mAcc$

Simple solution



- Works too long
- Resolution is very bad

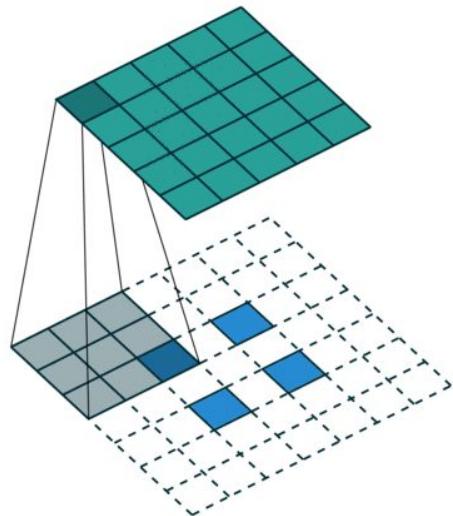
Better solution



This solution solves both problems:

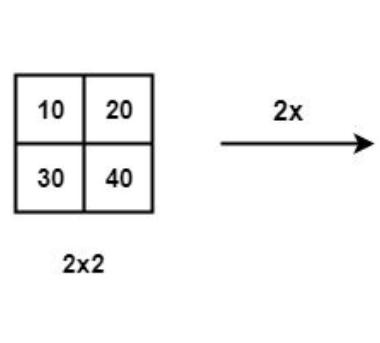
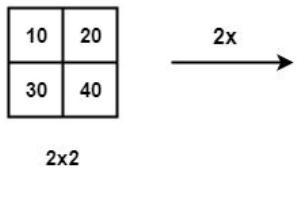
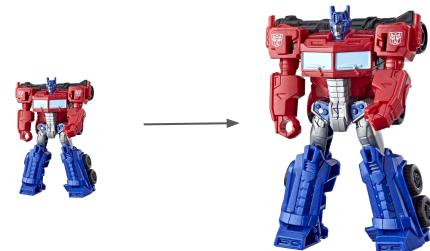
- Big receptive field
- Fine details due to top information

Upsampling



Transposed convolution can be used as an upsampling function. But it can produce checkerboard artifacts

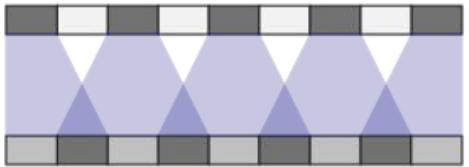
Upsampling



Nearest neighbor upsampling

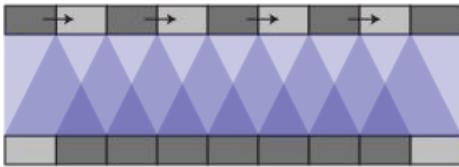
Bilinear interpolation upsampling

Upsampling



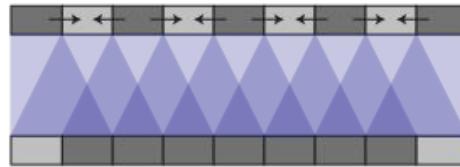
$$\begin{bmatrix} a & c \\ b & \\ a & c \\ & b \end{bmatrix}$$

Deconvolution



$$\begin{bmatrix} a+b & c \\ a & b+c \\ a+b & c \\ a & b+c \end{bmatrix}$$

NN-Resize Convolution



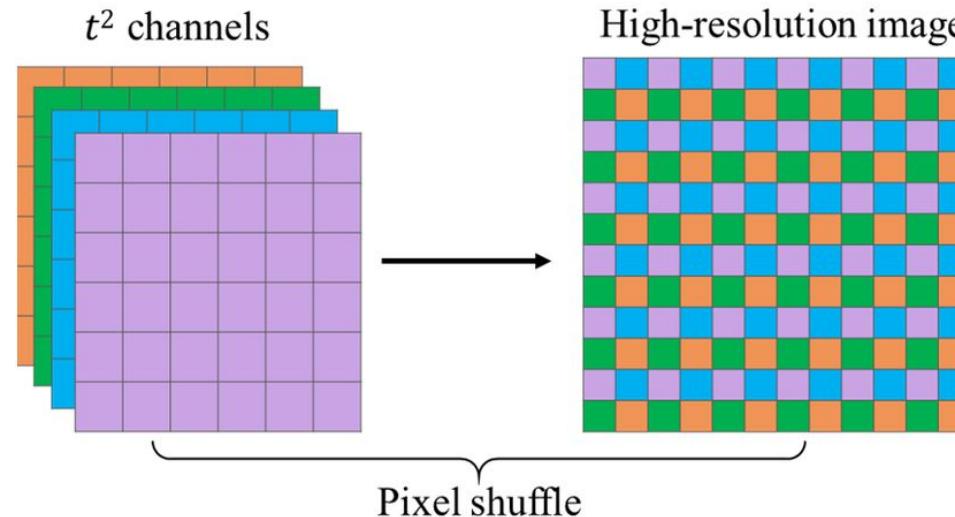
$$\begin{bmatrix} a + \frac{1}{2}b & \frac{1}{2}b+c & & \\ \frac{1}{2}a & \frac{1}{2}a+b+\frac{1}{2}c & \frac{1}{2}c & \\ a + \frac{1}{2}b & & \frac{1}{2}b+c & \\ \frac{1}{2}a & & \frac{1}{2}a+b+\frac{1}{2}c & \frac{1}{2}c \end{bmatrix}$$

Bilinear-Resize Convolution



Voilà!

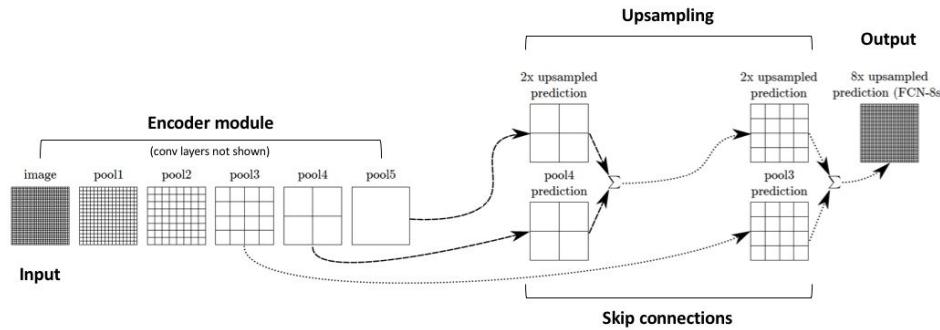
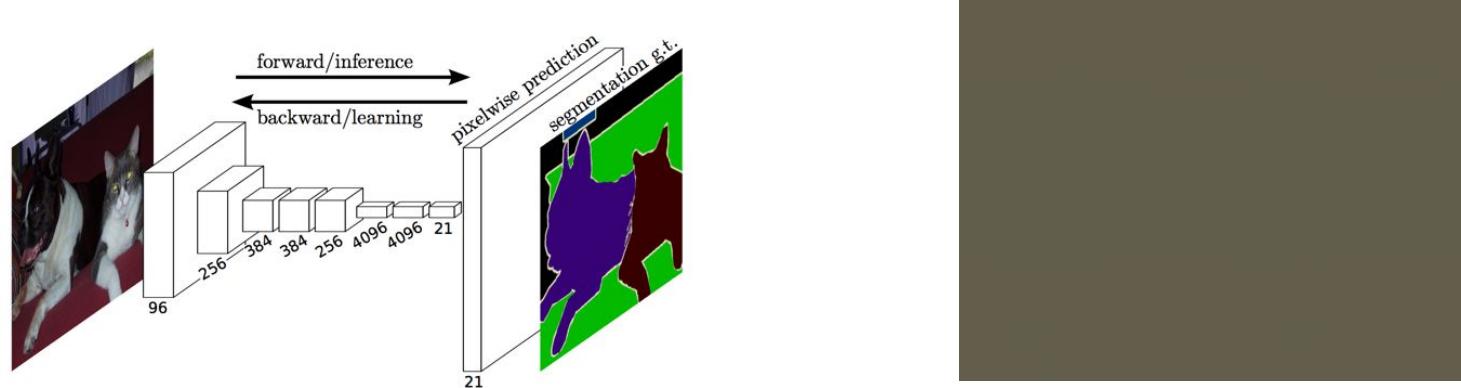
Pixel Shuffle



Transforms additional channels to spatial resolution (usually we do reverse)!

Architectures

FCN



FCN

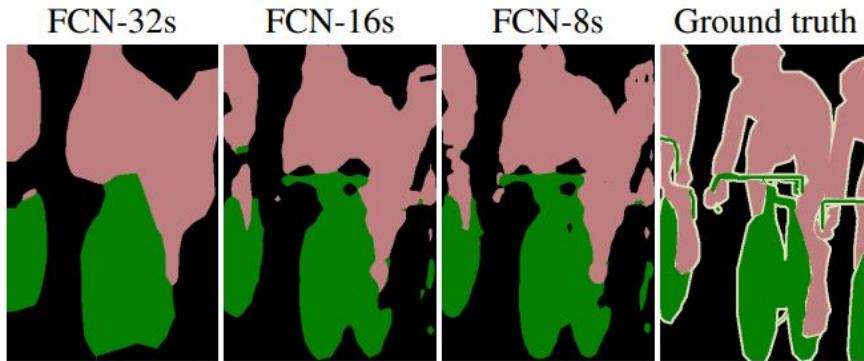
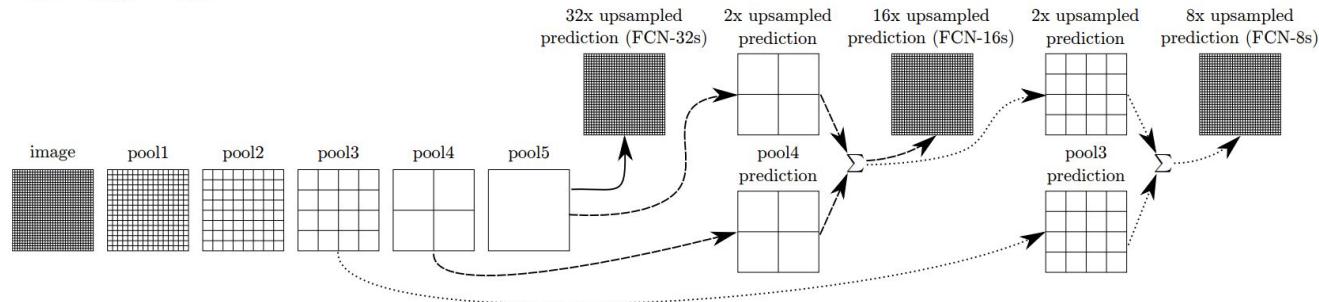
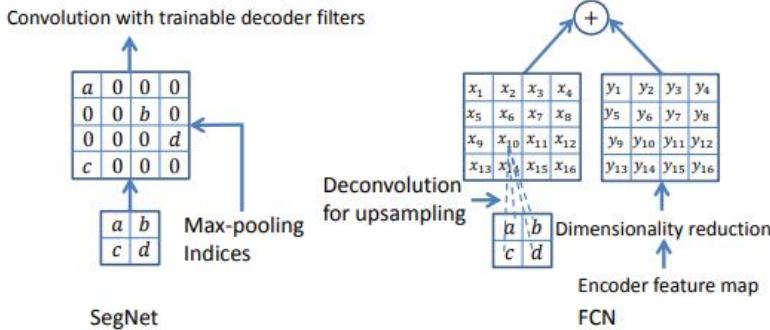
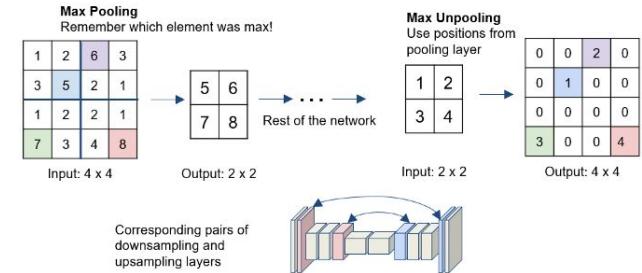
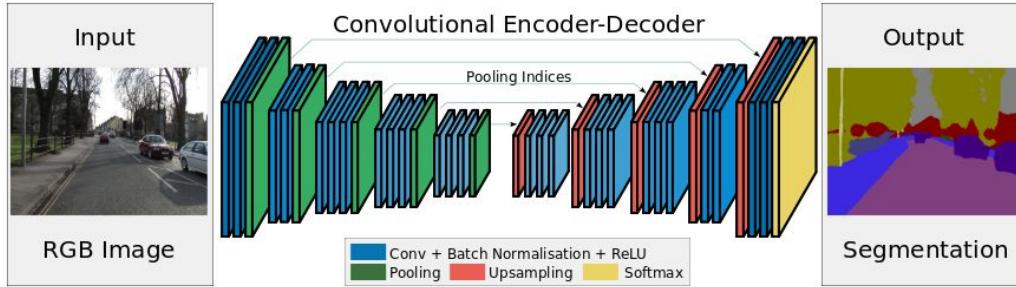


Figure 4. Refining fully convolutional nets by fusing information from layers with different strides improves segmentation detail. The first three images show the output from our 32, 16, and 8 pixel stride nets (see Figure 3).

If we use more top information we get the better quality

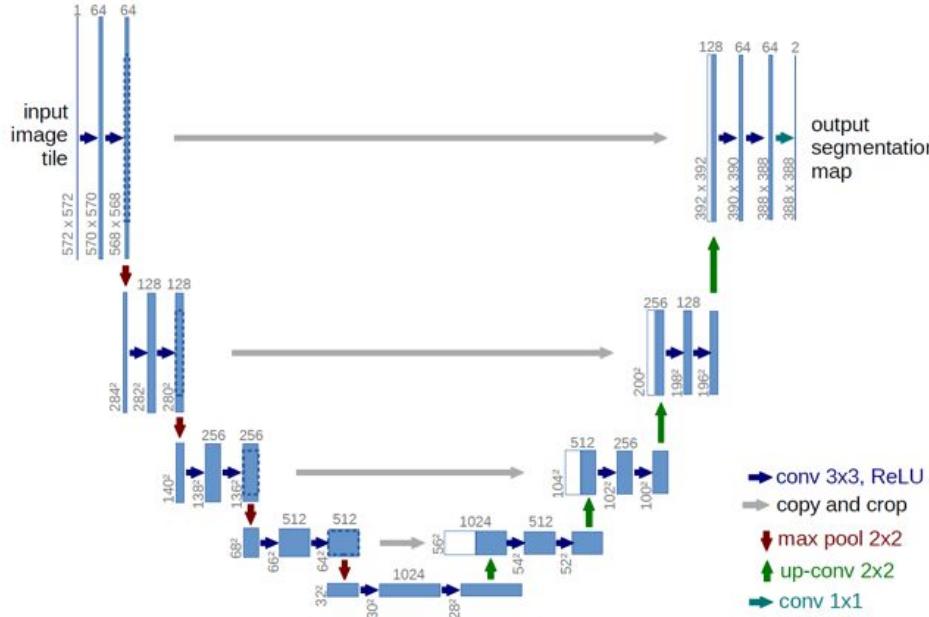


SegNet



- Making encoder and decoder equal in terms of parameters
- Saving more information from encoder

UNet



UNet

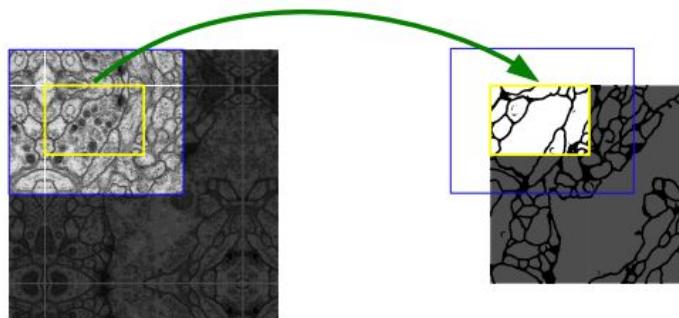
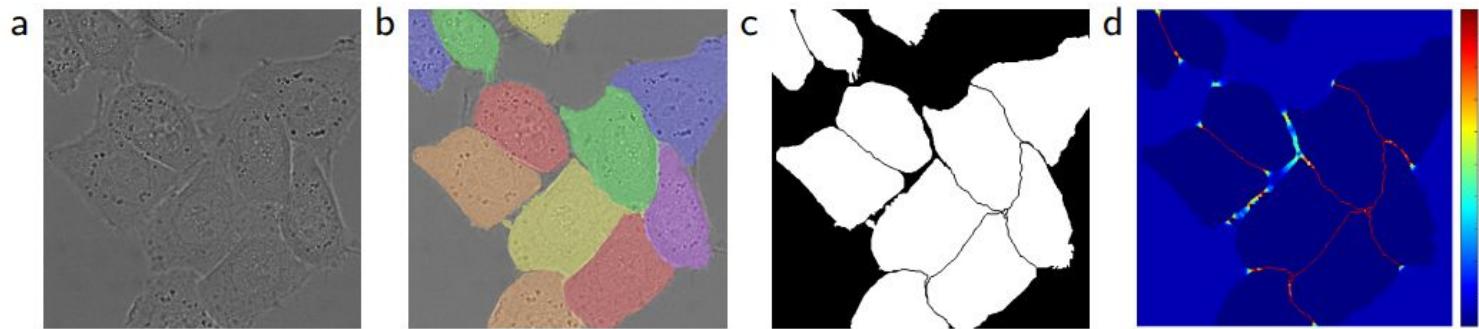


Fig. 2. Overlap-tile strategy for seamless segmentation of arbitrary large images (here segmentation of neuronal structures in EM stacks). Prediction of the segmentation in the yellow area, requires image data within the blue area as input. Missing input data is extrapolated by mirroring



HRNet

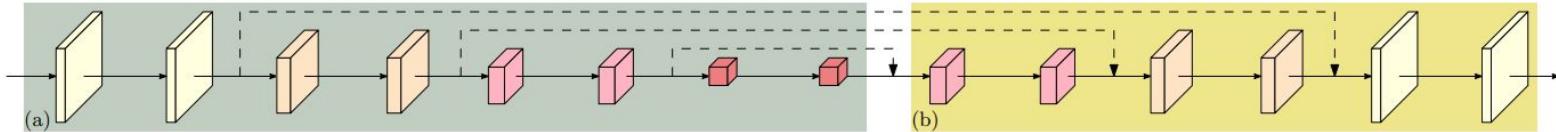
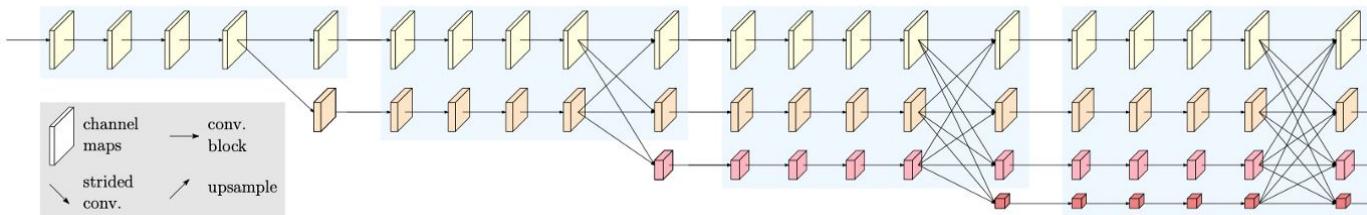


Fig. 1. The structure of recovering high resolution from low resolution. (a) A low-resolution representation learning subnetwork (such as VGGNet [126], ResNet [54]), which is formed by connecting high-to-low convolutions in series. (b) A high-resolution representation recovering subnetwork, which is formed by connecting low-to-high convolutions in series. Representative examples include SegNet [3], DeconvNet [107], U-Net [119] and Hourglass [105], encoder-decoder [112], and SimpleBaseline [152].

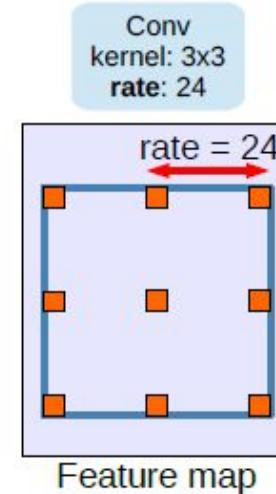
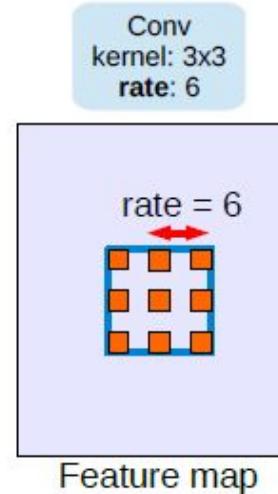
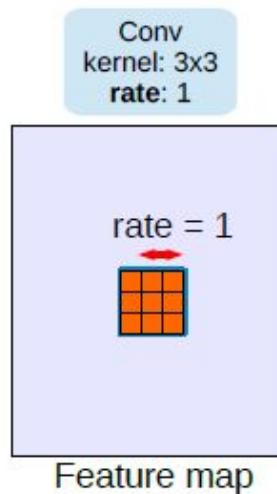
Base template for segmentation architecture



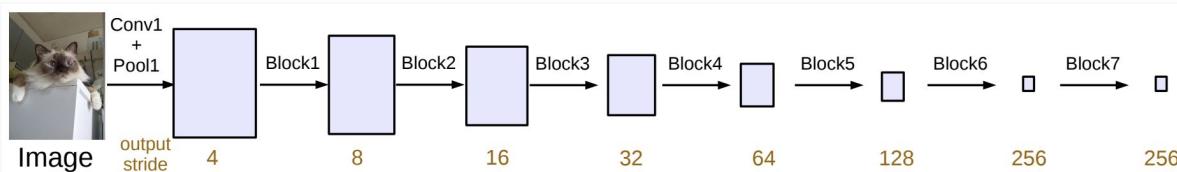
HRNet

Atrous convolution

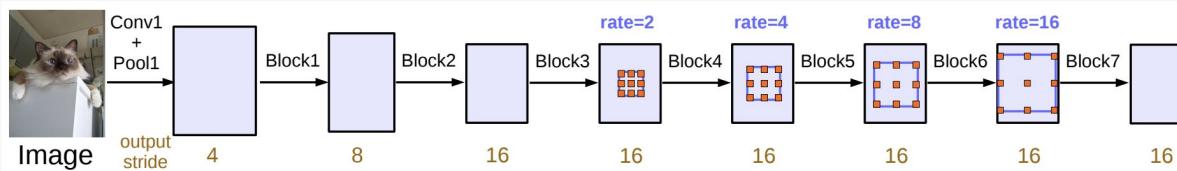
How to increase receptive field w/o lowering image resolution



Atrous convolution



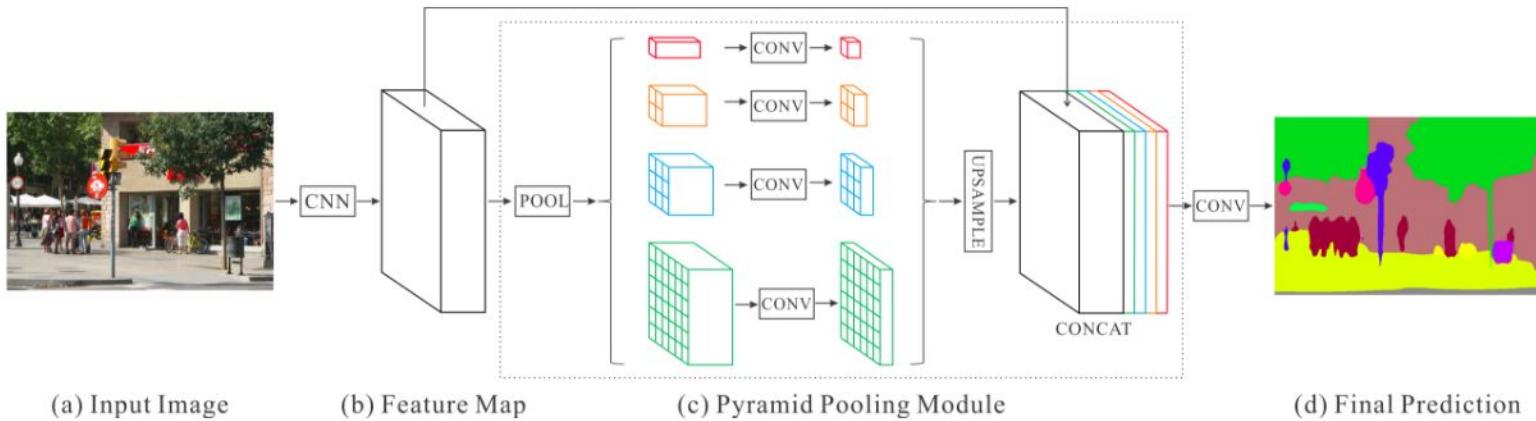
(a) Going deeper without atrous convolution.



(b) Going deeper with atrous convolution. Atrous convolution with $rate > 1$ is applied after block3 when $output_stride = 16$.

Figure 3. Cascaded modules without and with atrous convolution.

PSPNet

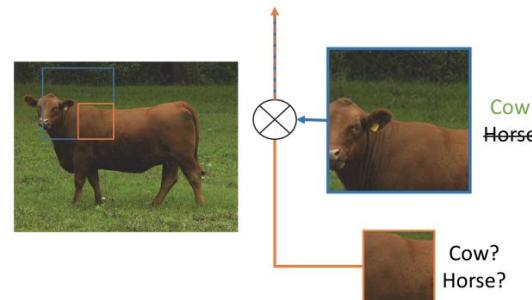


(a) Input Image

(b) Feature Map

(c) Pyramid Pooling Module

(d) Final Prediction



DeepLab

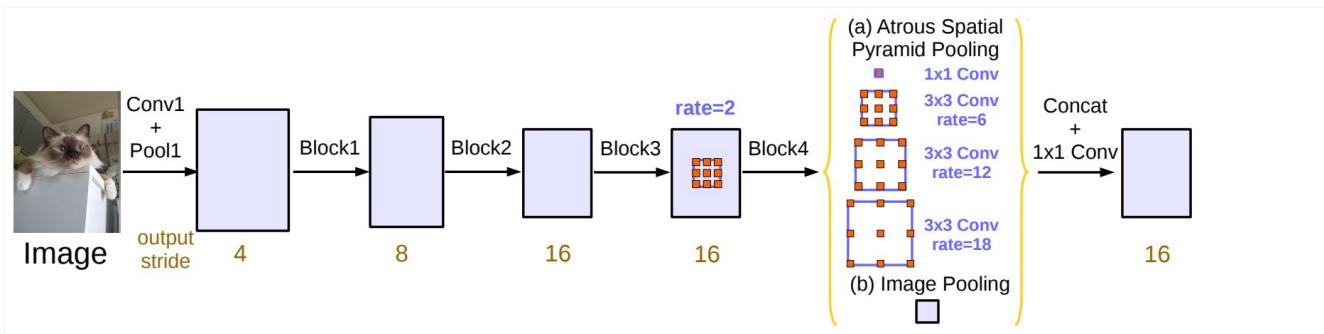
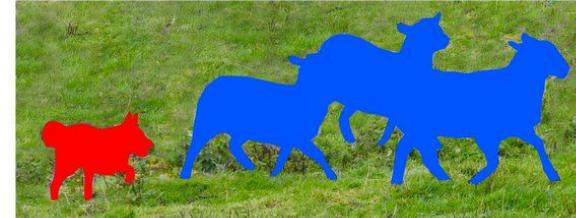


Figure 5. Parallel modules with atrous convolution (ASPP), augmented with image-level features.

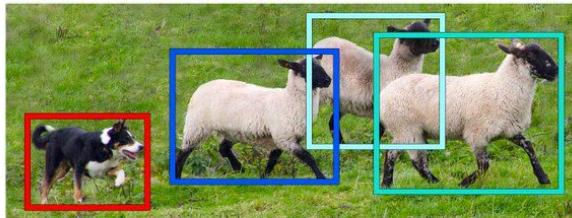
Instance Segmentation



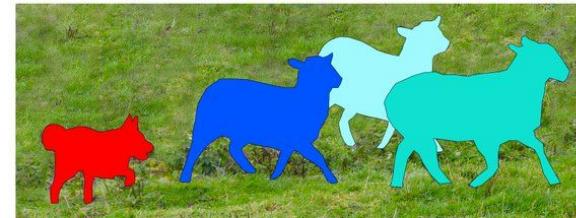
Image Recognition



Semantic Segmentation



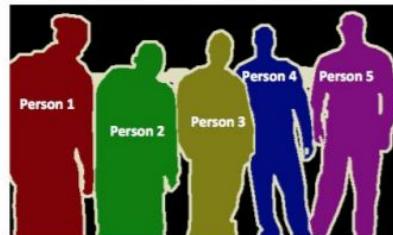
Object Detection



Instance Segmentation

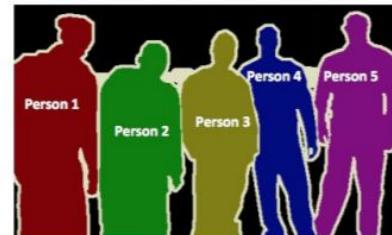
Instance Segmentation

R-CNN driven



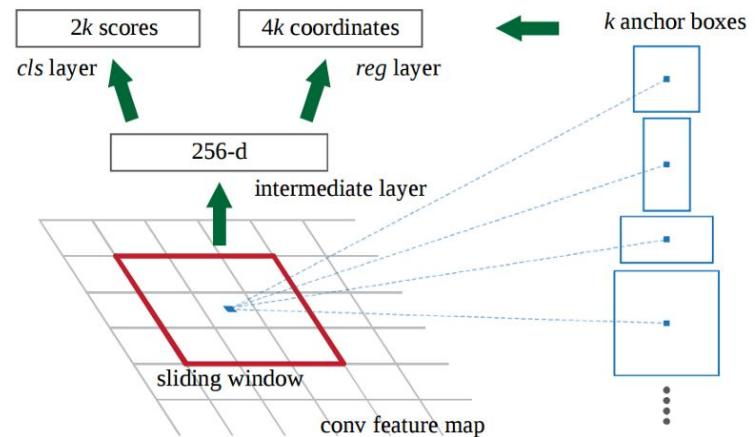
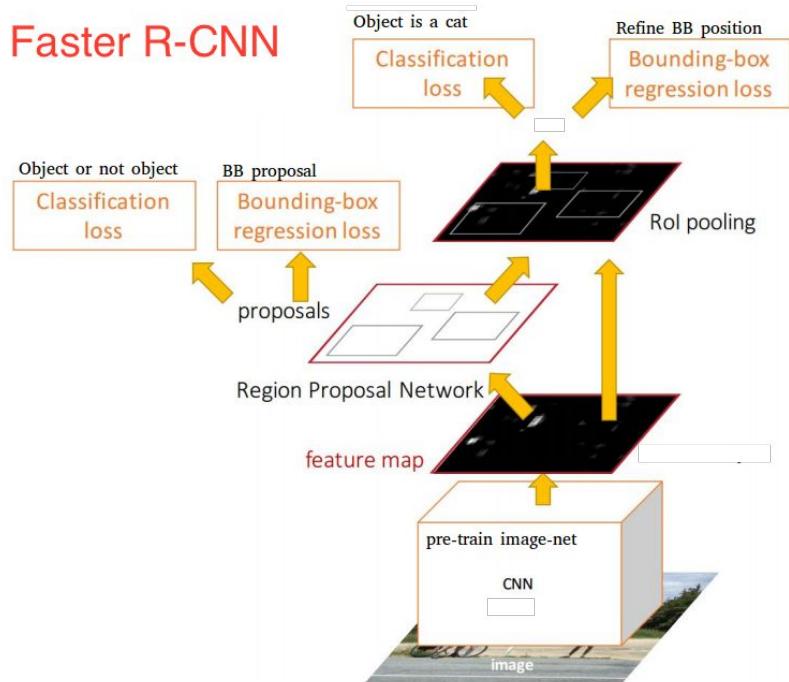
dit: Kaiming He

FCN driven

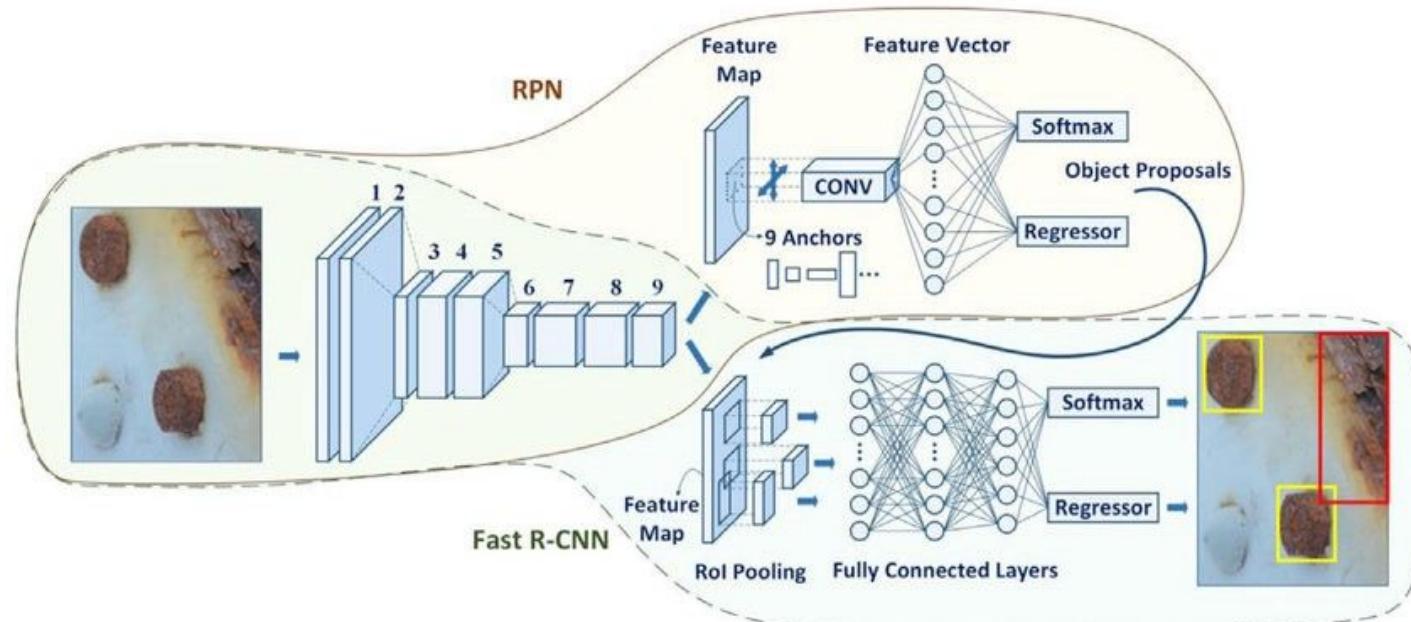


Reminder:

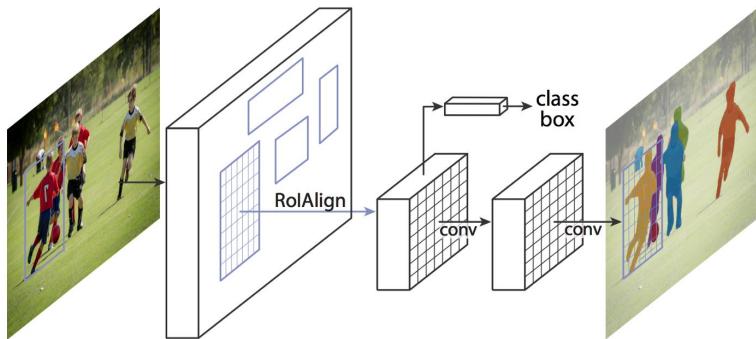
Faster R-CNN



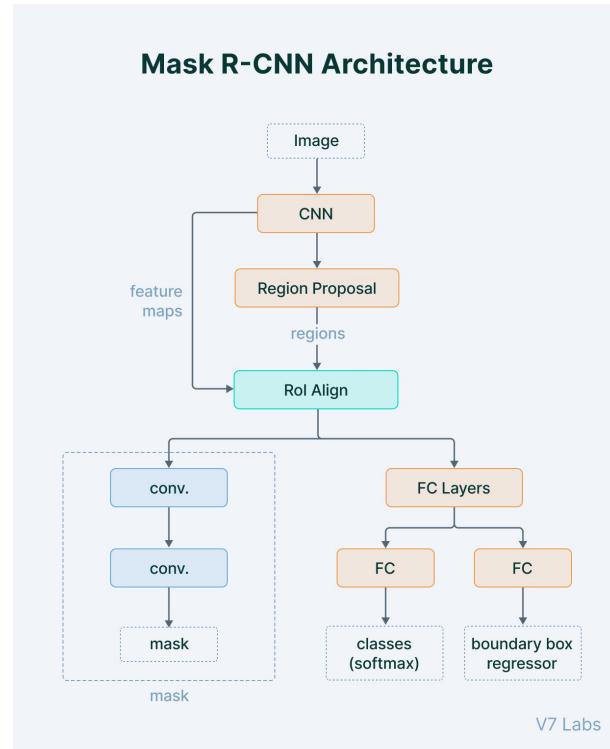
Faster RCNN



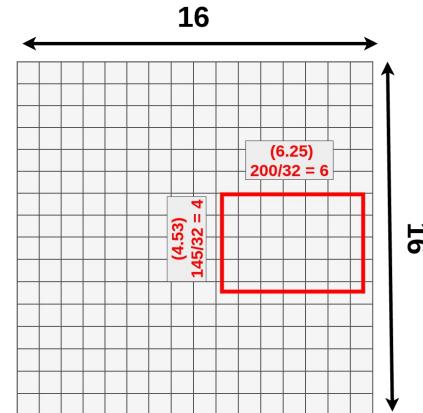
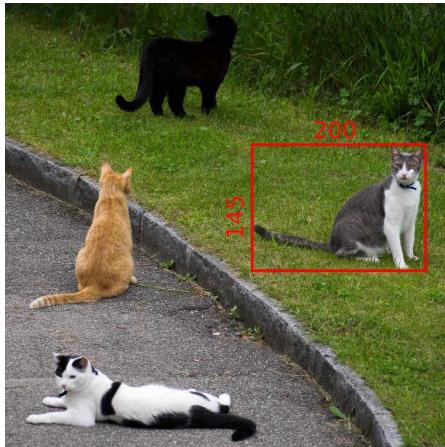
Mask R-CNN



Adding mask prediction head



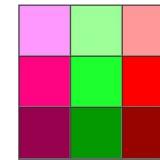
RoI Pooling



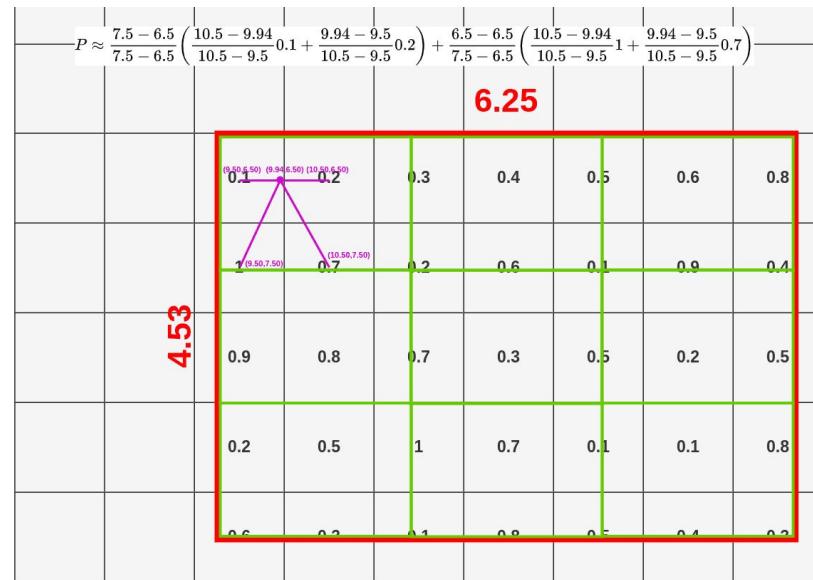
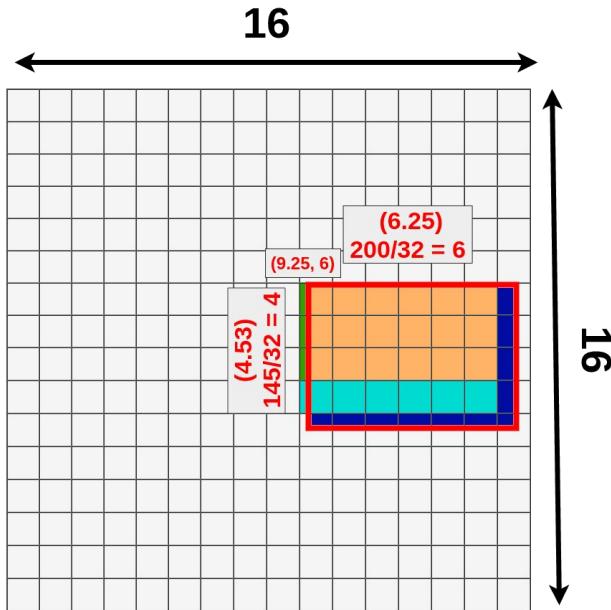
4x6 RoI

0.1	0.2	0.3	0.4	0.5	0.6
1	0.7	0.2	0.6	0.1	0.9
0.9	0.8	0.7	0.3	0.5	0.2

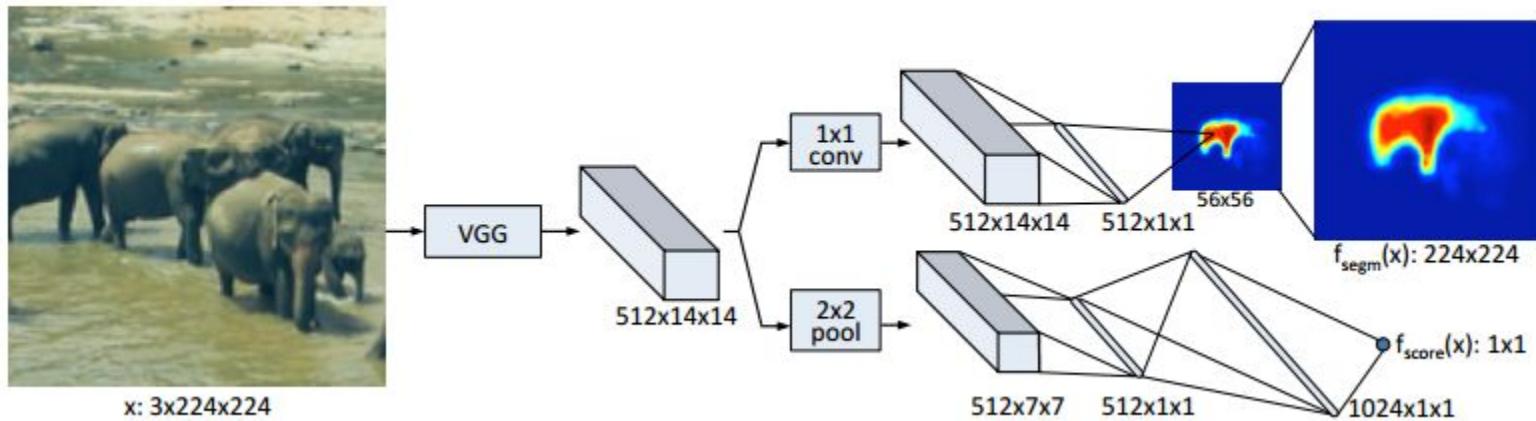
3x3 RoI Pooling



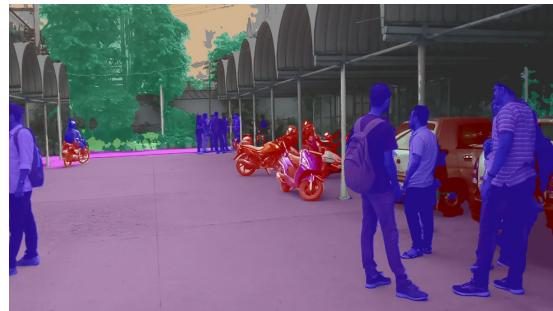
RoI Align



DeepMask



Panoptic Segmentation



Semantic



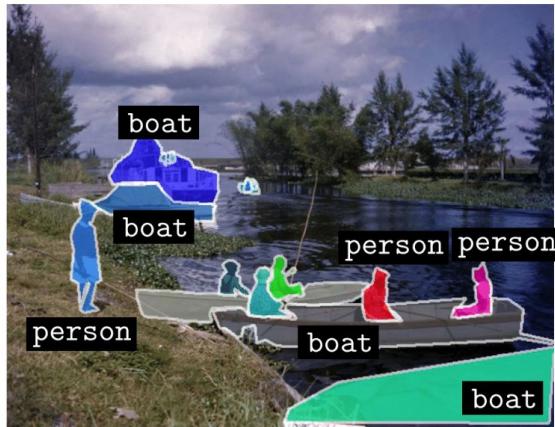
Instance



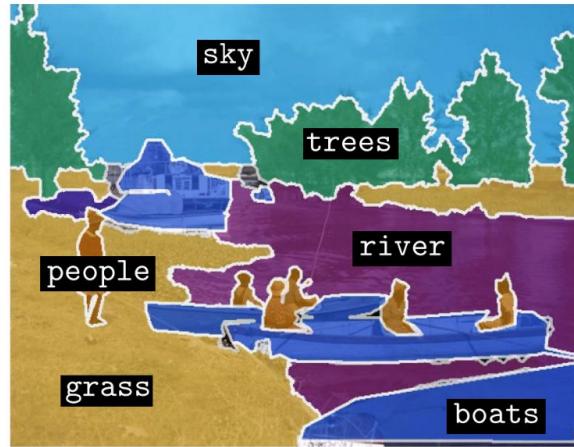
Panoptic

Panoptic Segmentation

What do every segmentation do?



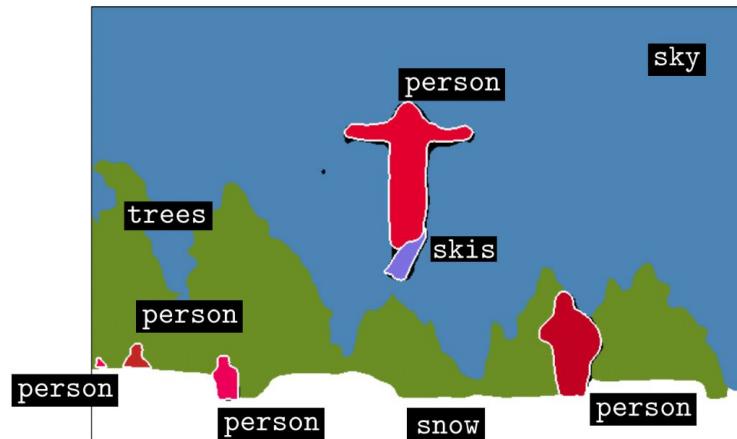
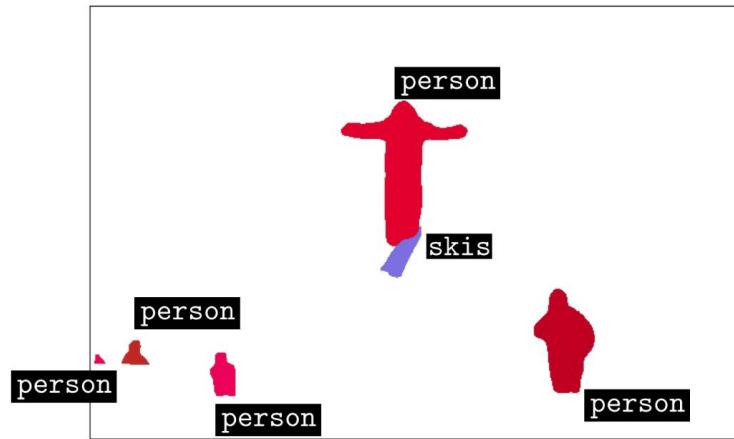
instance segmentation



semantic segmentation

Panoptic Segmentation

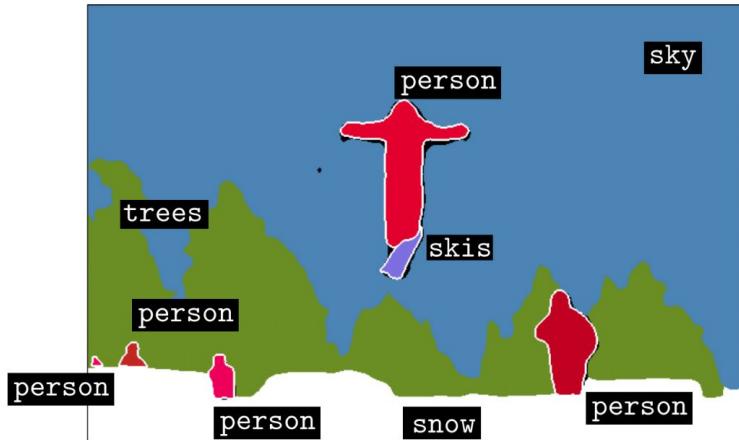
What do instance segmentation see?



Let's add semantic segmentation information

Panoptic Segmentation

How does the real image looks like?



Metric

Semantic segmentation

- IoU, per-pixel metric

Instance segmentation

- mAP, object-size agnostic

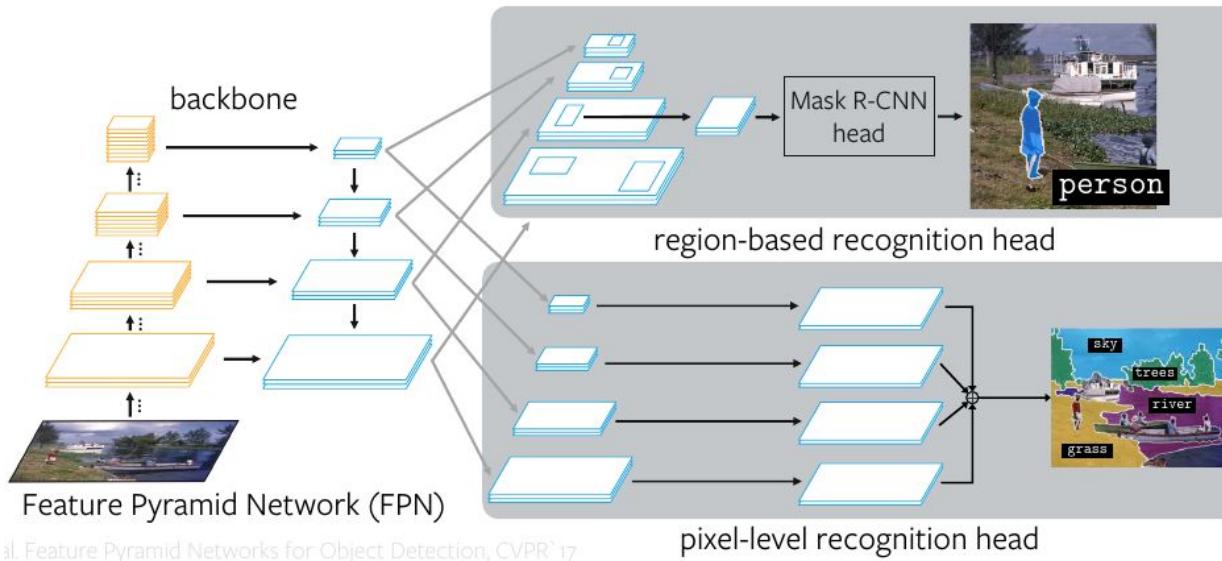
Panoptic segmentation should include both metrics simultaneously

Panoptic Segmentation

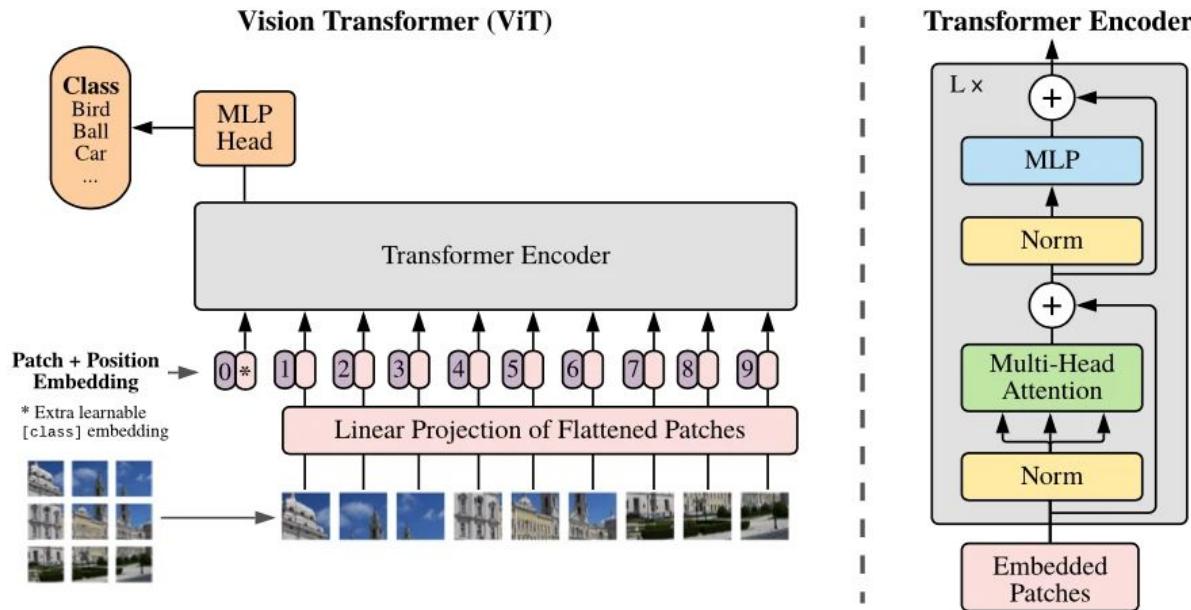
How to evaluate performance of panoptic segmentation model?

$$\begin{aligned} \text{PQ} &= \frac{\sum_{(p,g) \in TP} \text{IoU}(p, g)}{|TP| + \frac{1}{2}|FP| + \frac{1}{2}|FN|} \\ &= \underbrace{\frac{\sum_{(p,g) \in TP} \text{IoU}(p, g)}{|TP|}}_{\text{segmentation quality (SQ)}} \times \underbrace{\frac{|TP|}{|TP| + \frac{1}{2}|FP| + \frac{1}{2}|FN|}}_{\text{recognition quality (RQ)}} \end{aligned}$$

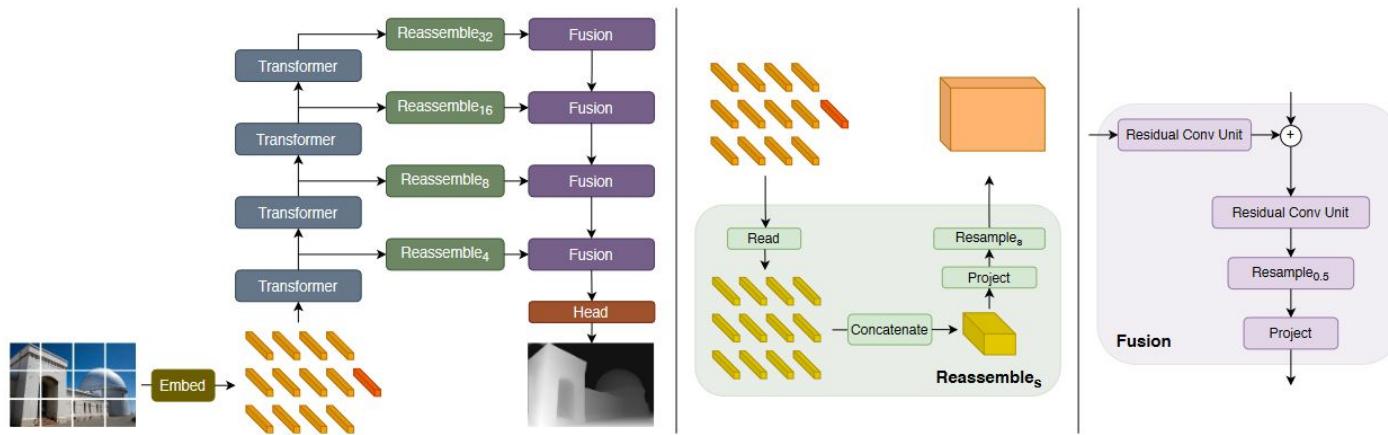
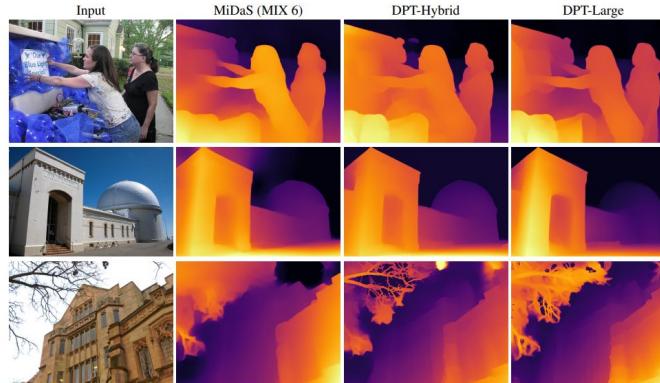
Panoptic architecture



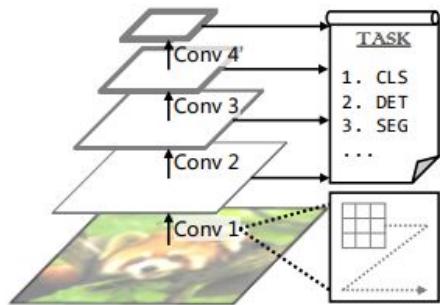
Reminder: ViT



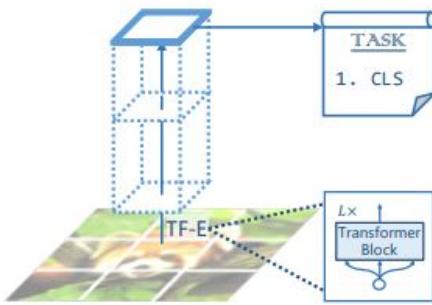
DPT (Vision Transformers for Dense Prediction)



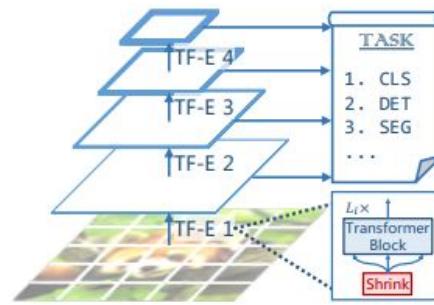
Pyramid Vision Transformer



(a) CNNs: VGG [54], ResNet [22], etc.

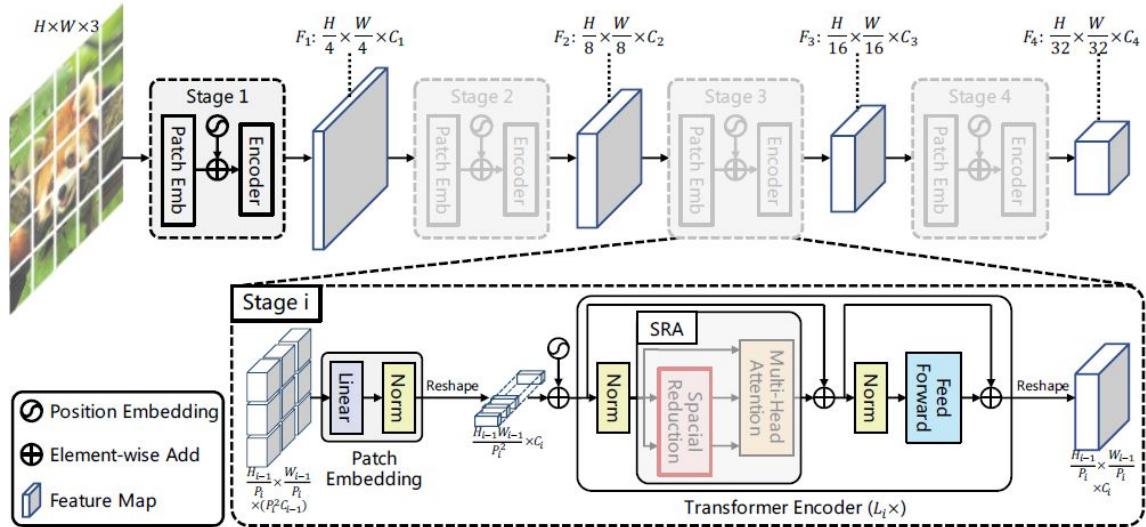


(b) Vision Transformer [13]

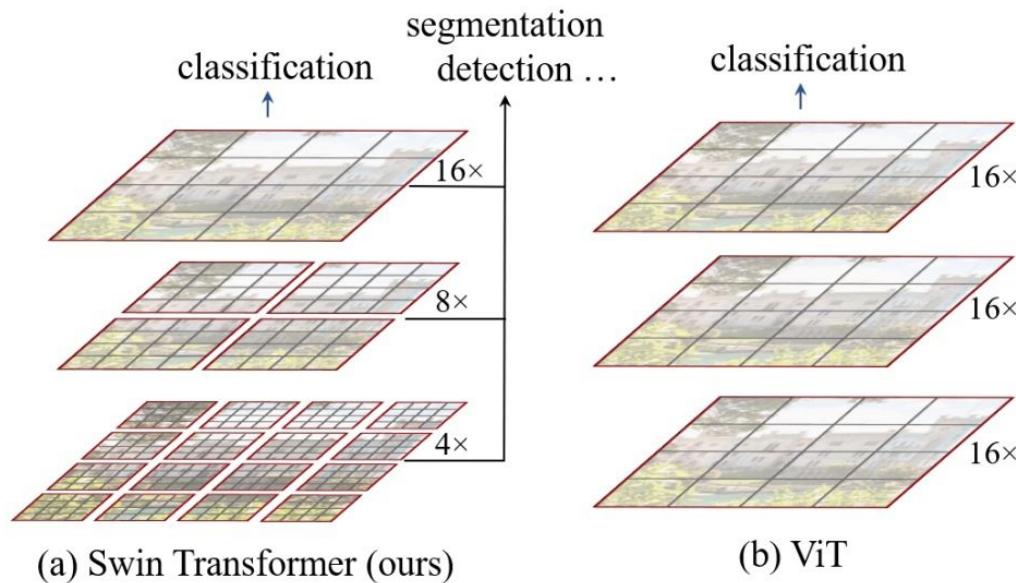


(c) Pyramid Vision Transformer (ours)

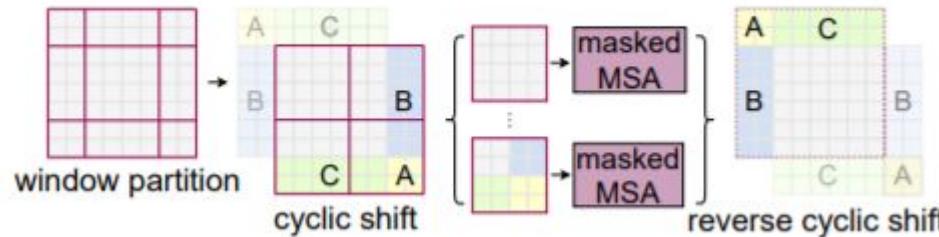
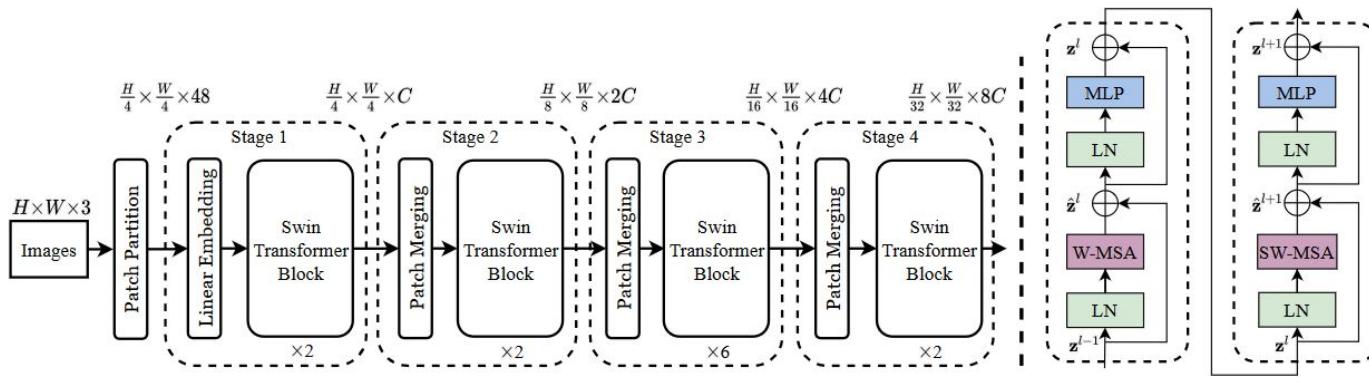
Pyramid Vision Transformer



Swin



Swin



HRFormer

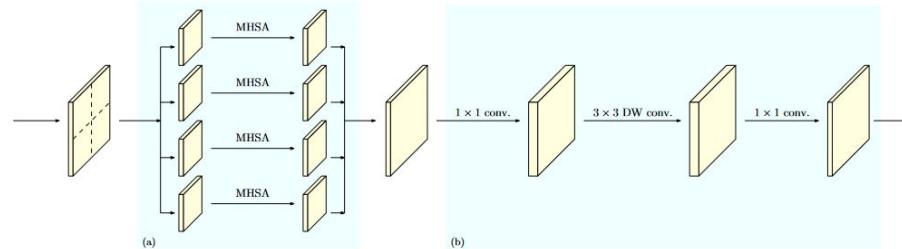
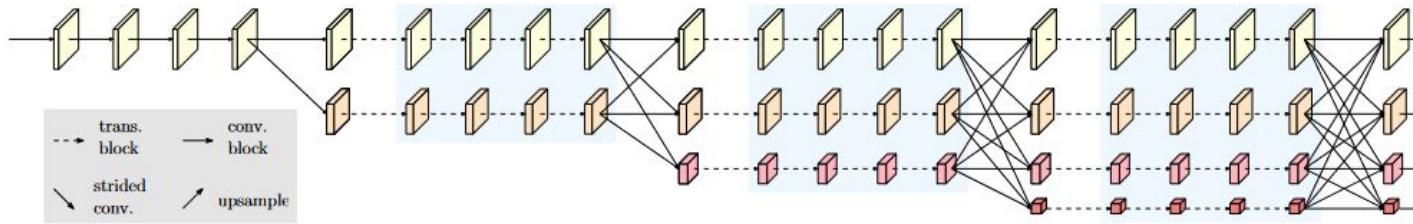


Figure 1: **Illustrating the HRFormer block.** The HRFormer block is composed of (a) local-window self-attention and (b) feed-forward network (FFN) with depth-wise convolution. The local-window self-attention scheme is inspired by the interlaced sparse self-attention [56, 21].

HRFormer

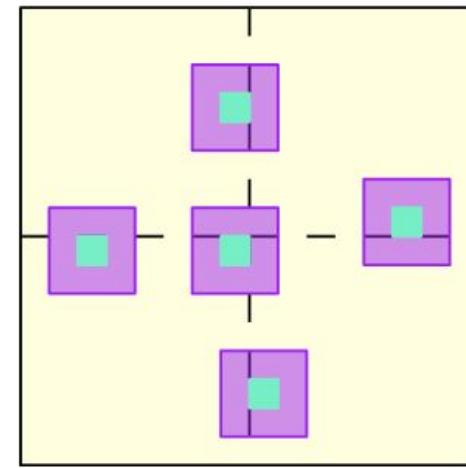
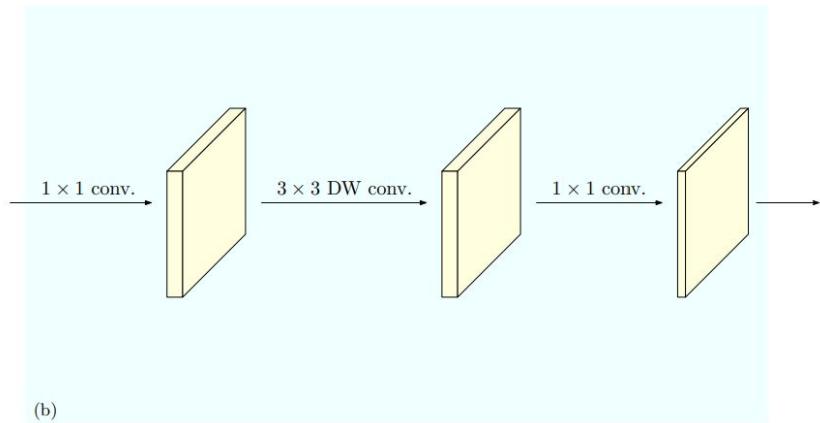


Figure 3: Illustrating that FFN with 3×3 depth-wise convolution connects the non-overlapping windows.

Recap

- Semantic segmentation problem
- Upsampling
- Architectures
- Panoptic / Instance segmentation