# Semantic Segmentation

## ISTD 50.035

## Computer Vision

# Semantic Segmentation



Person
Bicycle
Background

Label each pixel of an image with a class value -> dense prediction
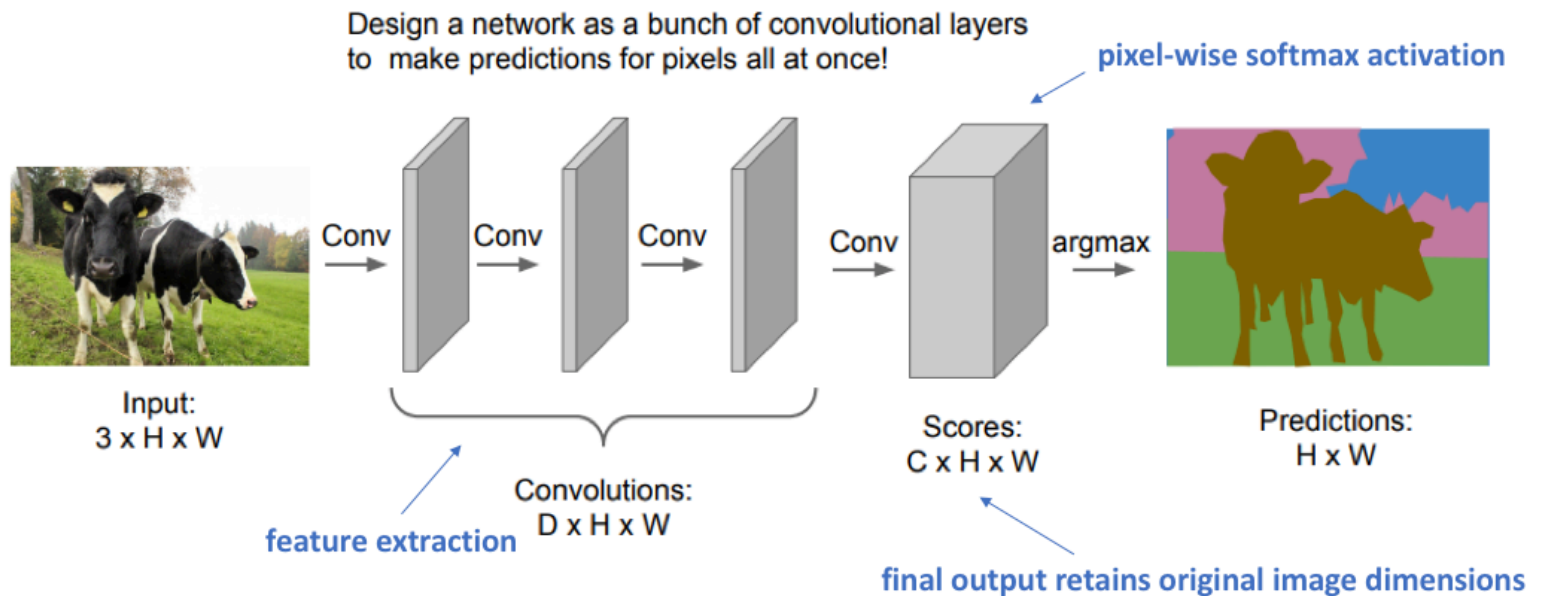
# Semantic Segmentation



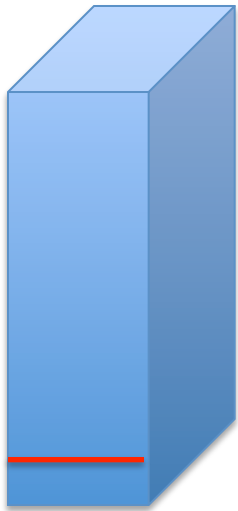Input Image          Semantic Segmentation          Instance Segmentation

# Semantic Segmentation with ConvNet



Design a network as a bunch of convolutional layers to make predictions for pixels all at once!

pixel-wise softmax activation

Input:
3 x H x W

Conv → Conv → Conv → Conv → argmax

feature extraction

Convolutions:
D x H x W

Scores:
C x H x W

final output retains original image dimensions

Predictions:
H x W

**Downside:** Preserving image dimensions throughout entire network will be computationally expensive.

Probability vector of C classes at each pixel location
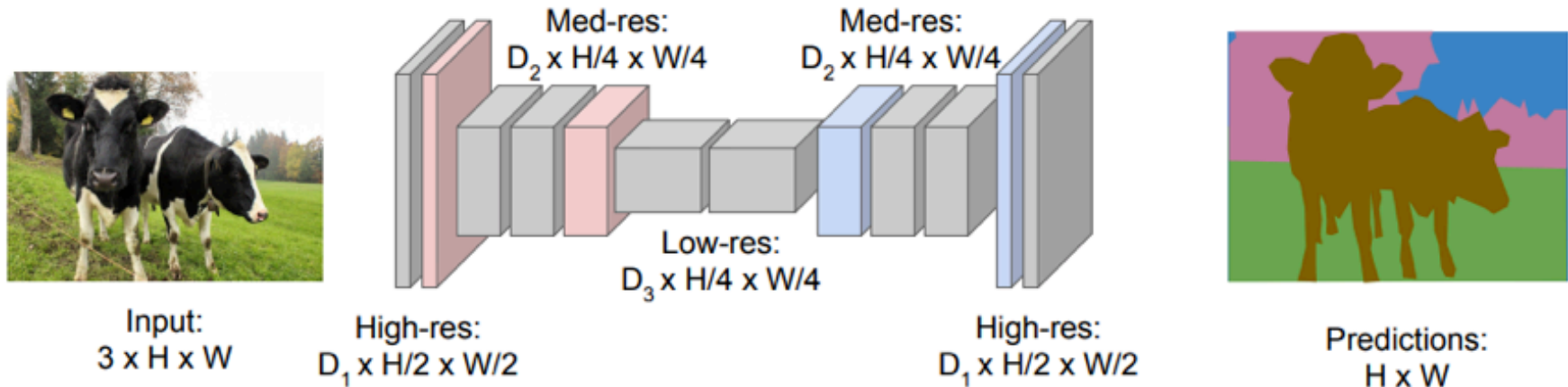
# Pixel wise softmax loss

Probability vector of C classes at each pixel location

Loss at each pixel location = -log $p_y$

Loss = (spatial) sum of loss at each pixel

C

# Semantic Segmentation with ConvNet



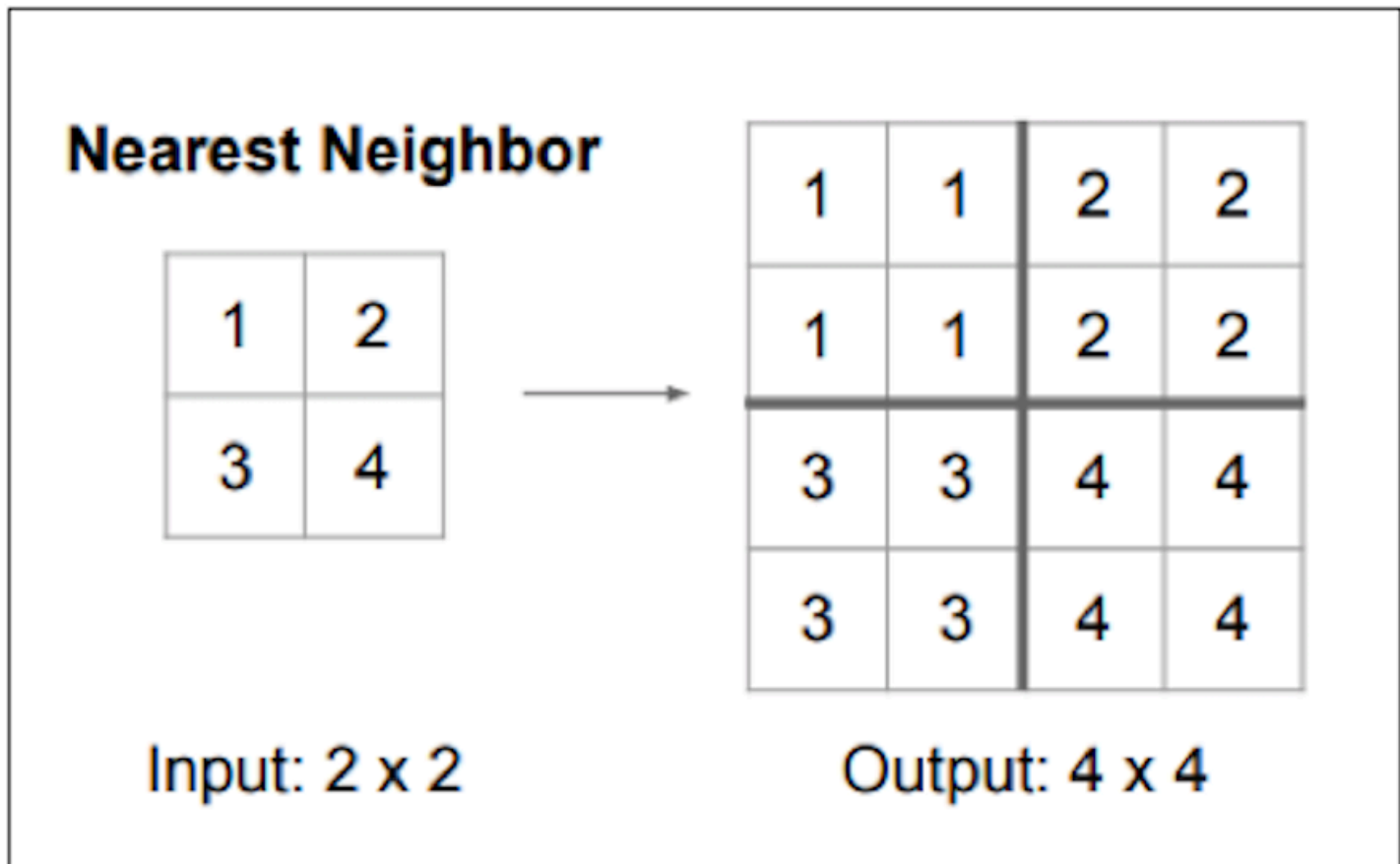Design network as a bunch of convolutional layers, with **downsampling** and **upsampling** inside the network!

Med-res:
$D_2 \times H/4 \times W/4$

Med-res:
$D_2 \times H/4 \times W/4$

Low-res:
$D_3 \times H/4 \times W/4$

Input:
$3 \times H \times W$

High-res:
$D_1 \times H/2 \times W/2$

High-res:
$D_1 \times H/2 \times W/2$

Predictions:
$H \times W$

**Solution:** Make network deep and *work at a lower spatial resolution* for many of the layers.
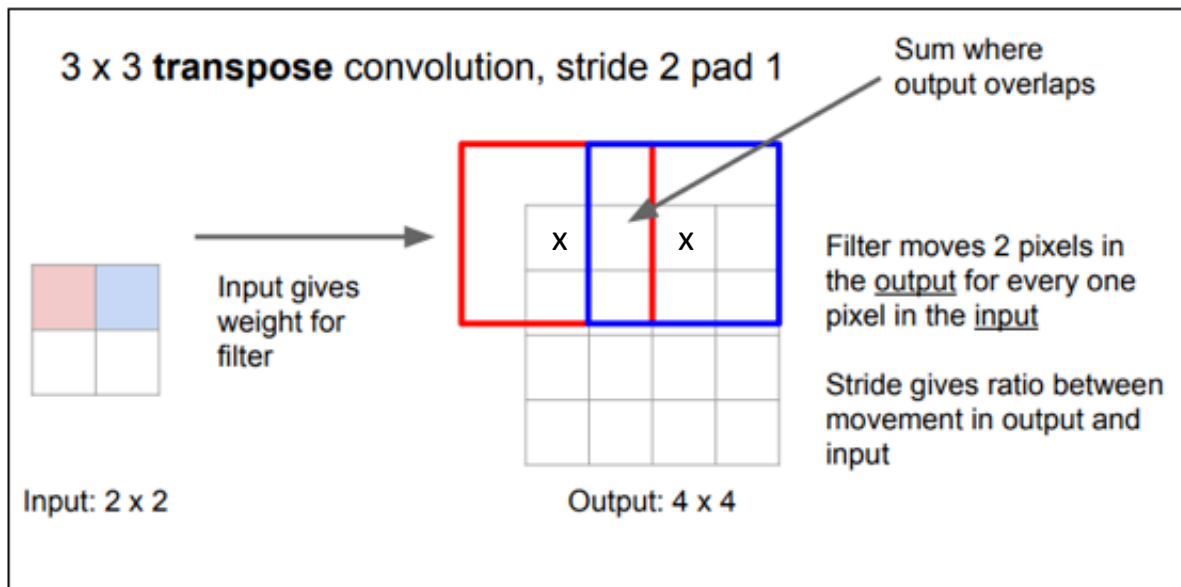
Probability vector of C classes at each pixel location

# Upsampling

Unpooling: inverse of max/average pooling

**Nearest Neighbor**

| 1 | 2 |
|---|---|
| 3 | 4 |

→

| 1 | 1 | 2 | 2 |
|---|---|---|---|
| 1 | 1 | 2 | 2 |
| 3 | 3 | 4 | 4 |
| 3 | 3 | 4 | 4 |

Input: 2 x 2          Output: 4 x 4

# Upsampling

Transpose convolution



3 x 3 **transpose** convolution, stride 2 pad 1

Sum where output overlaps

Input gives weight for filter

Filter moves 2 pixels in the output for every one pixel in the input

Stride gives ratio between movement in output and input
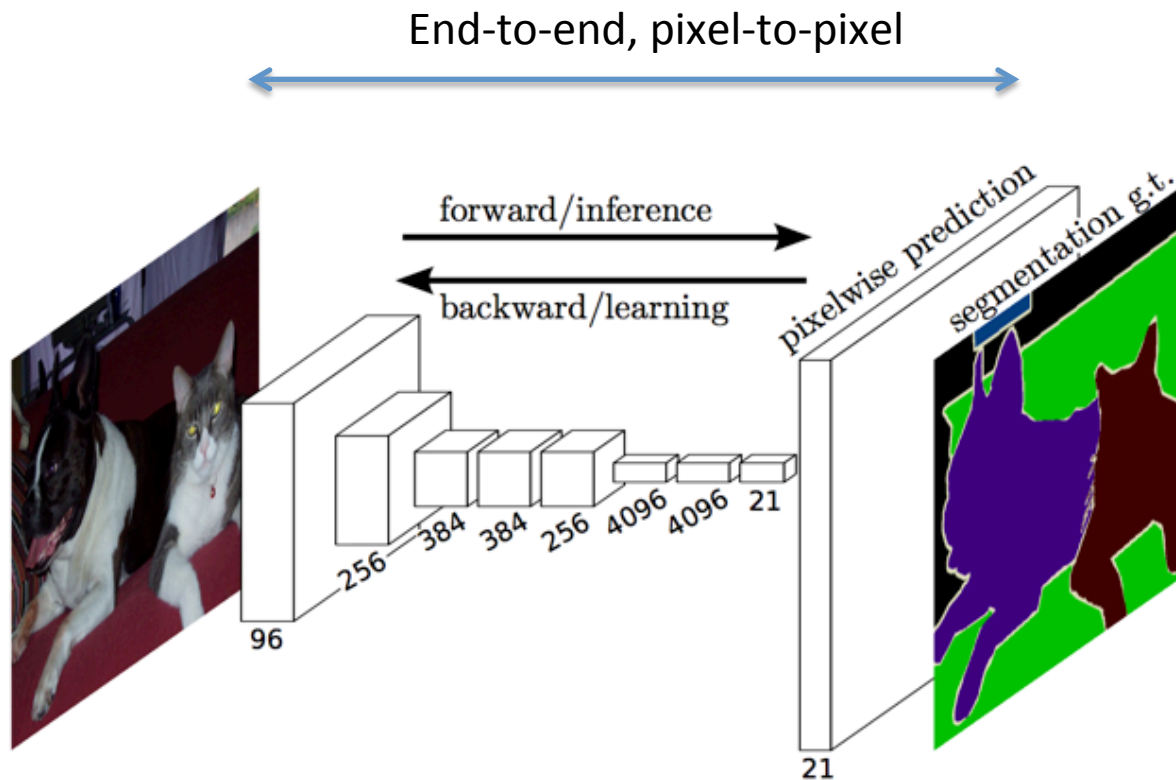
Input: 2 x 2

Output: 4 x 4

Standard convolution: input (matrix) * filter (matrix) -> output (scalar)
Transpose convolution: input (scalar) * filter (matrix) -> output (matrix)

Learn filter mask for optimal upsampling

# Fully convolutional network (FCN)



End-to-end, pixel-to-pixel

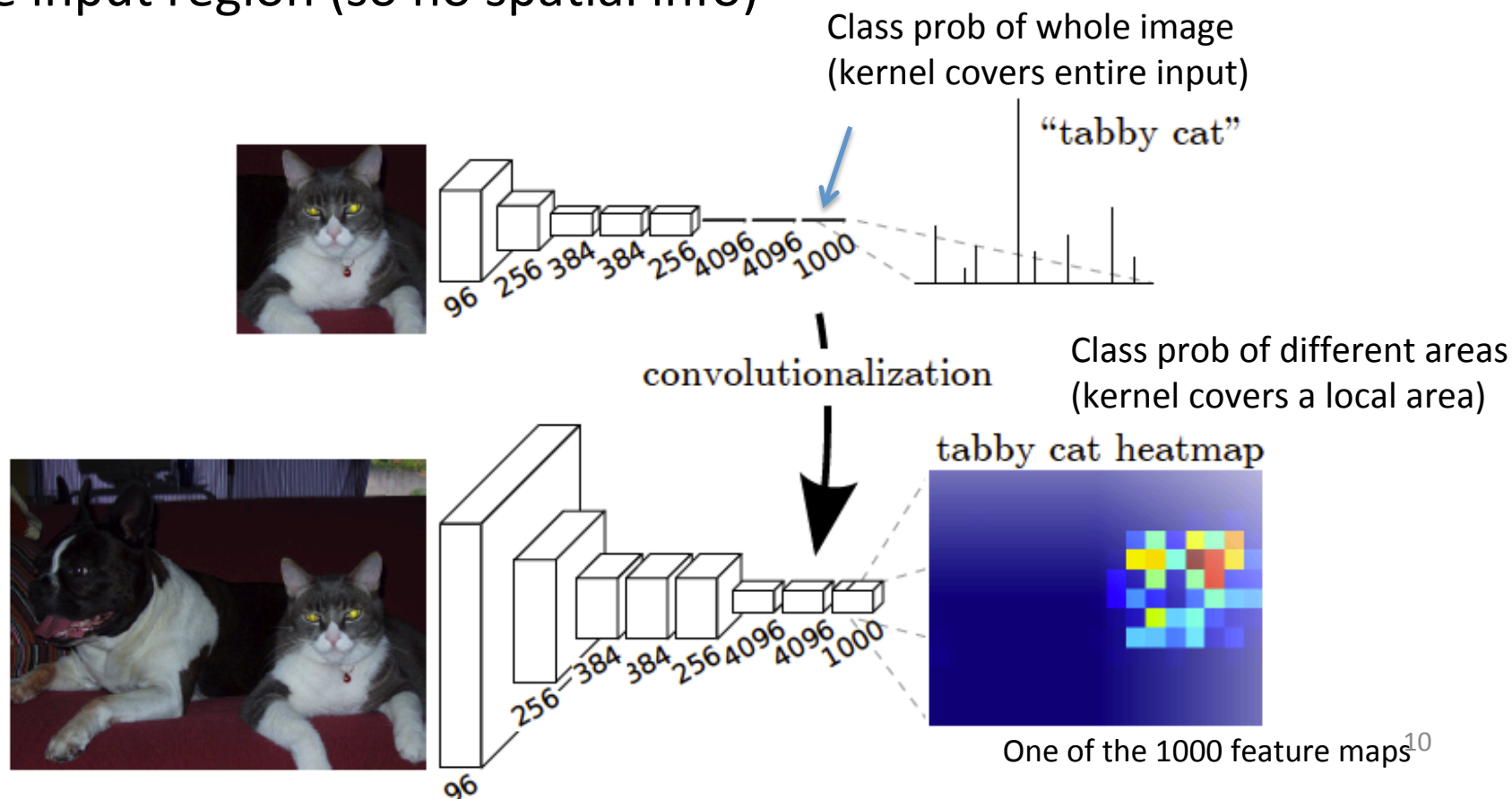[Long et al. 2014]
-Well pre-trained network as 'encoder'
-Transpose convolution layers to upsample the coarse feature map to full-resolution segmentation map
-Trained end-to-end, pixel-to-pixel

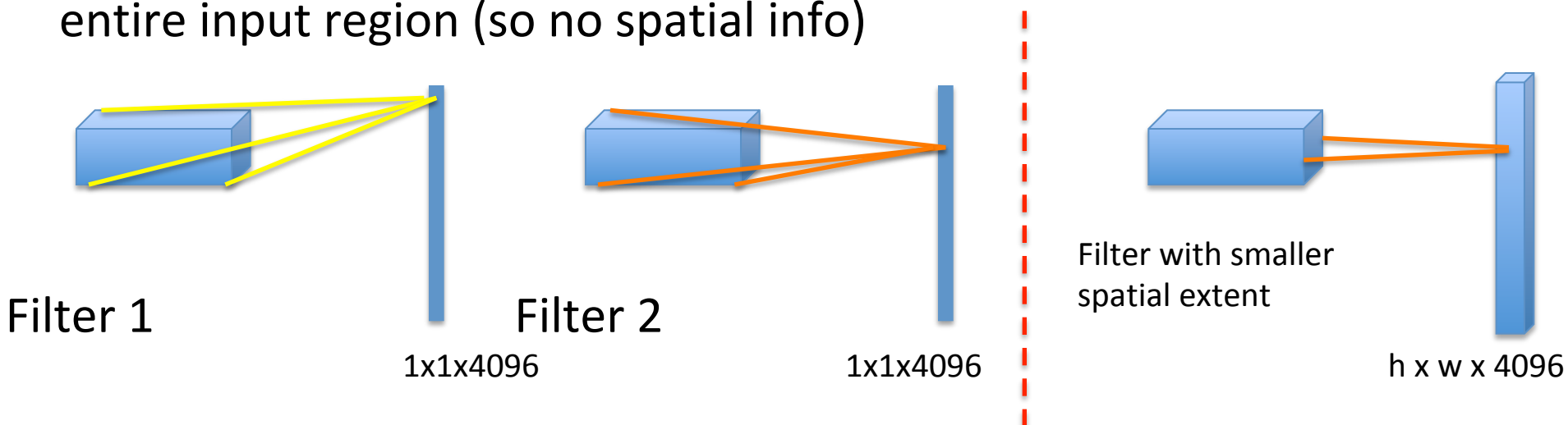# Fully convolutional network (FCN)

Fully connected layer

-Discard spatial information -> not suitable for SS

-Can be viewed as convolutions with kernels that cover the entire input region (so no spatial info)



Class prob of whole image
(kernel covers entire input)

"tabby cat"

convolutionalization

Class prob of different areas
(kernel covers a local area)

tabby cat heatmap

One of the 1000 feature maps

# Fully convolutional network (FCN)

Fully connected layer
-Discard spatial information -> not suitable for SS
-Can be viewed as convolutions with kernels that cover the entire input region (so no spatial info)

Filter 1

1x1x4096

Filter 2

1x1x4096

Filter with smaller spatial extent

h x w x 4096

Fully connected layer: 4096 filter masks; filter size = input
Output of fully connected layer: 1x1 feature map (total 4096 channels)
Fully convolutional: use small filter

# Issue in baseline FCN

Ground truth target

Predicted segmentation

"Semantic segmentation faces an inherent tension between semantics and location: global information resolves what while local information resolves where"

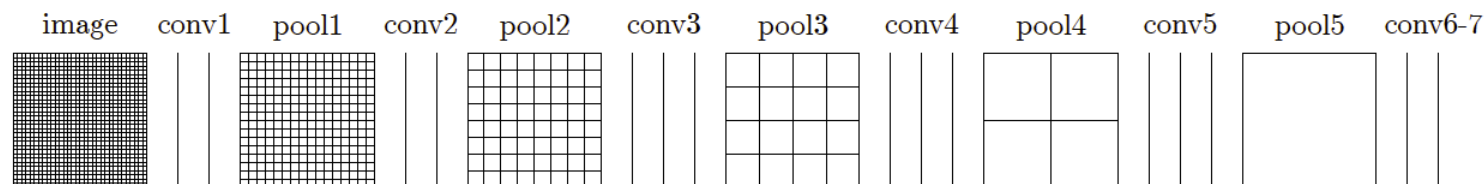"Combining fine layers and coarse layers lets the model make local predictions that respect global structure"

# Combining what and where

- Deep, coarse semantic information
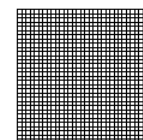- Shallow, fine appearance information


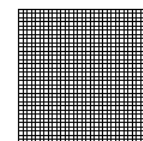
Pooling with stride 2

Pooling with stride 4 equivalently

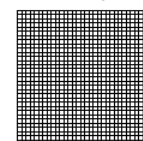Pooling with stride 16 equivalently

image  conv1  pool1  conv2  pool2  conv3  pool3  conv4  pool4  conv5  pool5  conv6-7

32x upsampled prediction (FCN-32s)

2x conv7
pool4

16x upsampled prediction (FCN-16s)

4x conv7
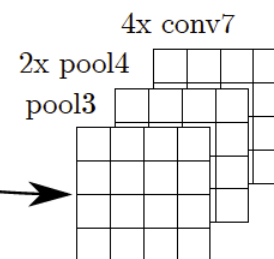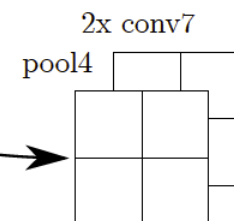2x pool4
pool3

8x upsampled prediction (FCN-8s)

Pool4: C 1x1 filters on pool4 -> class prediction

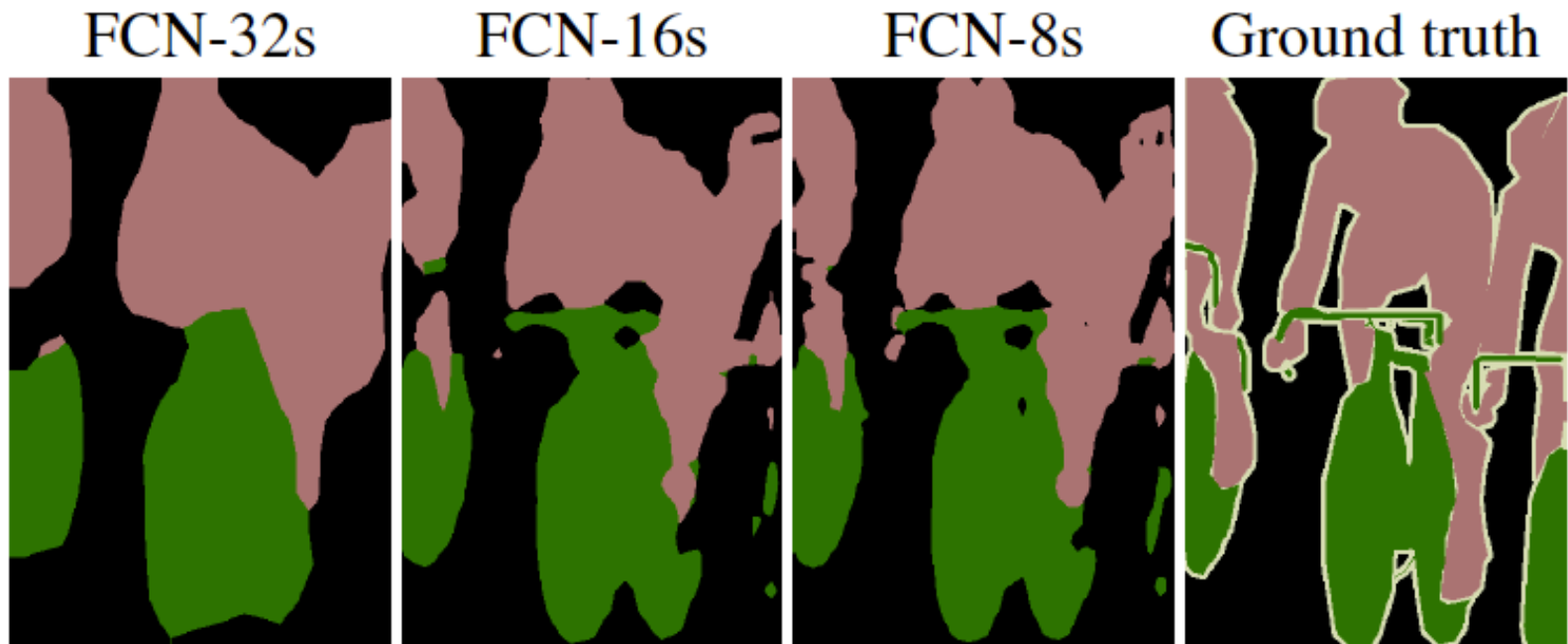Conv7: 2x upsampling

Summing both prediction

16x upsampling of result

# Combining what and where



conv7

Pooling with stride 32 equivalently

C

32x

Pooling with stride 16 equivalently

pool4

C 1x1 filters

C

conv7

Pooling with stride 32 equivalently

C

2x

C

C

C

16x

- Pool4: C 1x1 filters on pool4 -> class prediction
- Conv7: 2x upsampling
- Summing both prediction
- 16x upsampling of result

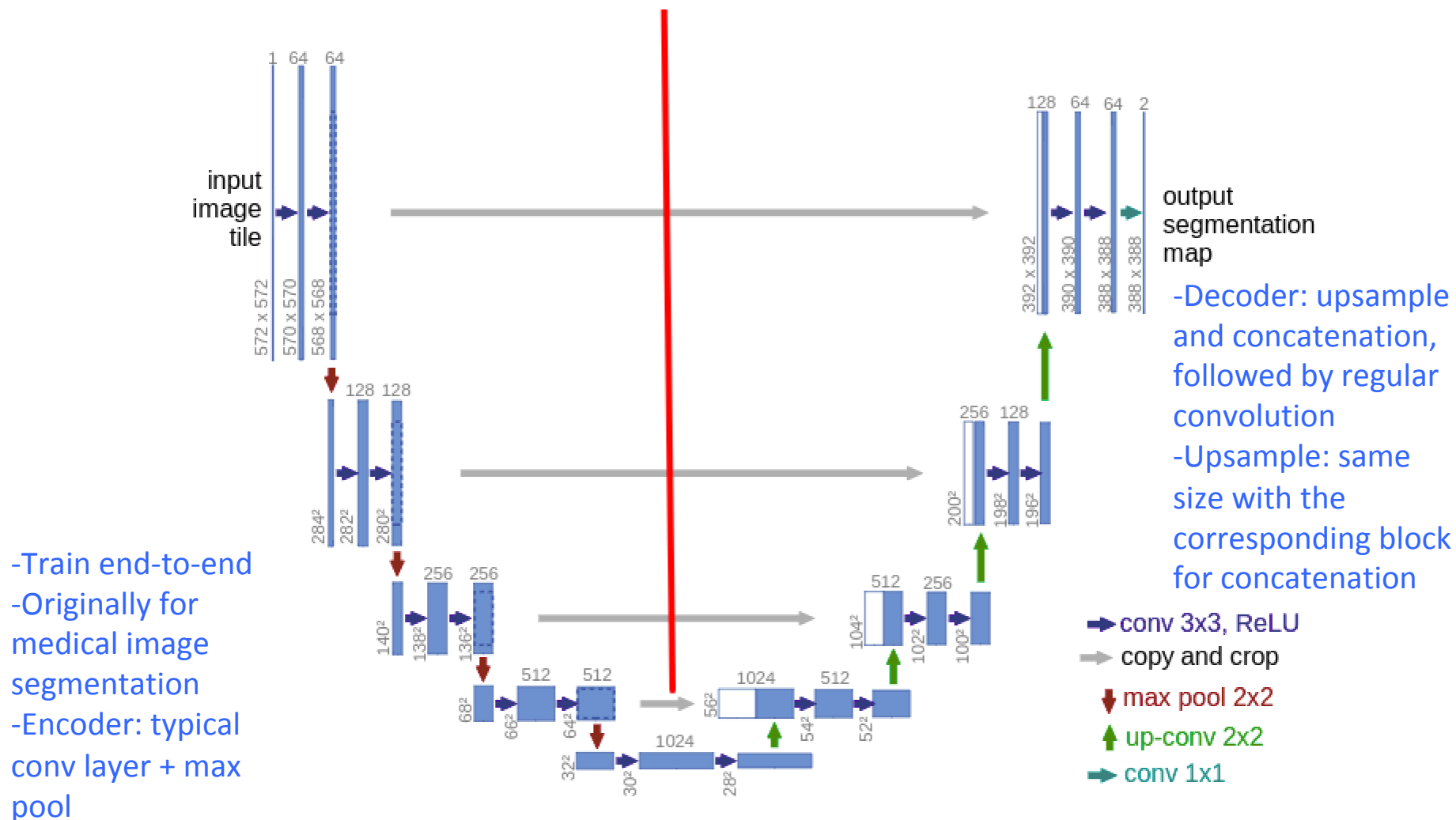(2x to use to fuse with pool3)

14

# Combining what and where



FCN-32s   FCN-16s   FCN-8s   Ground truth

# U-Net: high capacity decoder



input image tile

output segmentation map

-Decoder: upsample and concatenation, followed by regular convolution
-Upsample: same size with the corresponding block for concatenation

-Train end-to-end
-Originally for medical image segmentation
-Encoder: typical conv layer + max pool

→ conv 3x3, ReLU
→ copy and crop
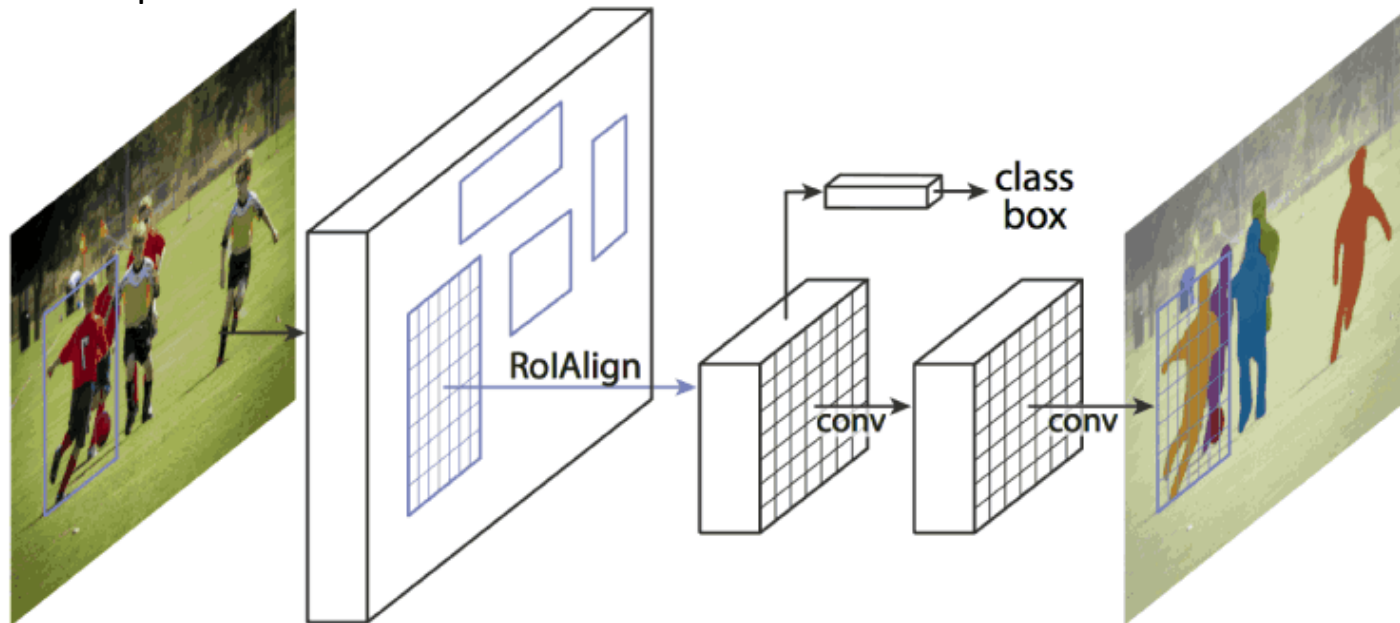↓ max pool 2x2
↑ up-conv 2x2
→ conv 1x1

# U-Net

- Upsampling: not accurate

- Use earlier stages to provide representation for localization: via concatenation
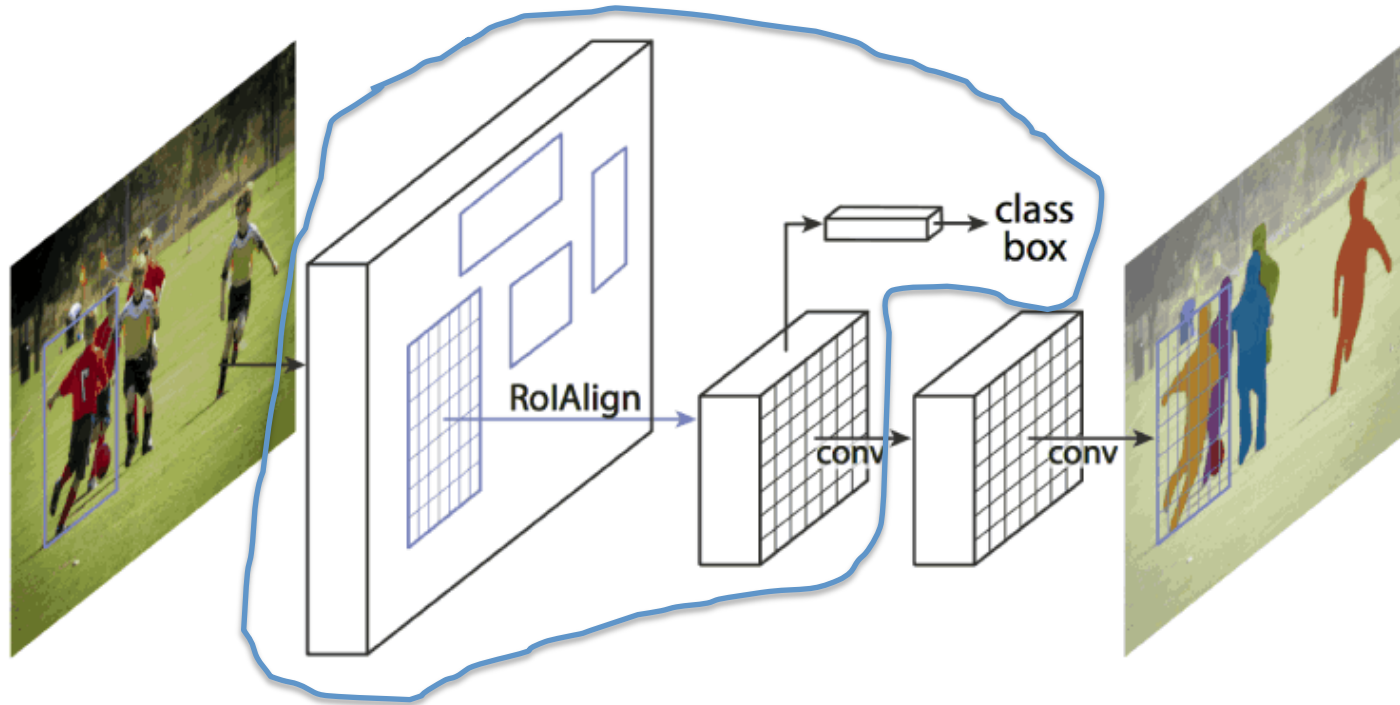
# Mask R-CNN for instance segmentation

Mask R-CNN = Faster R-CNN + FCN
Add mask prediction in addition to class and box



The Mask R-CNN framework for instance segmentation
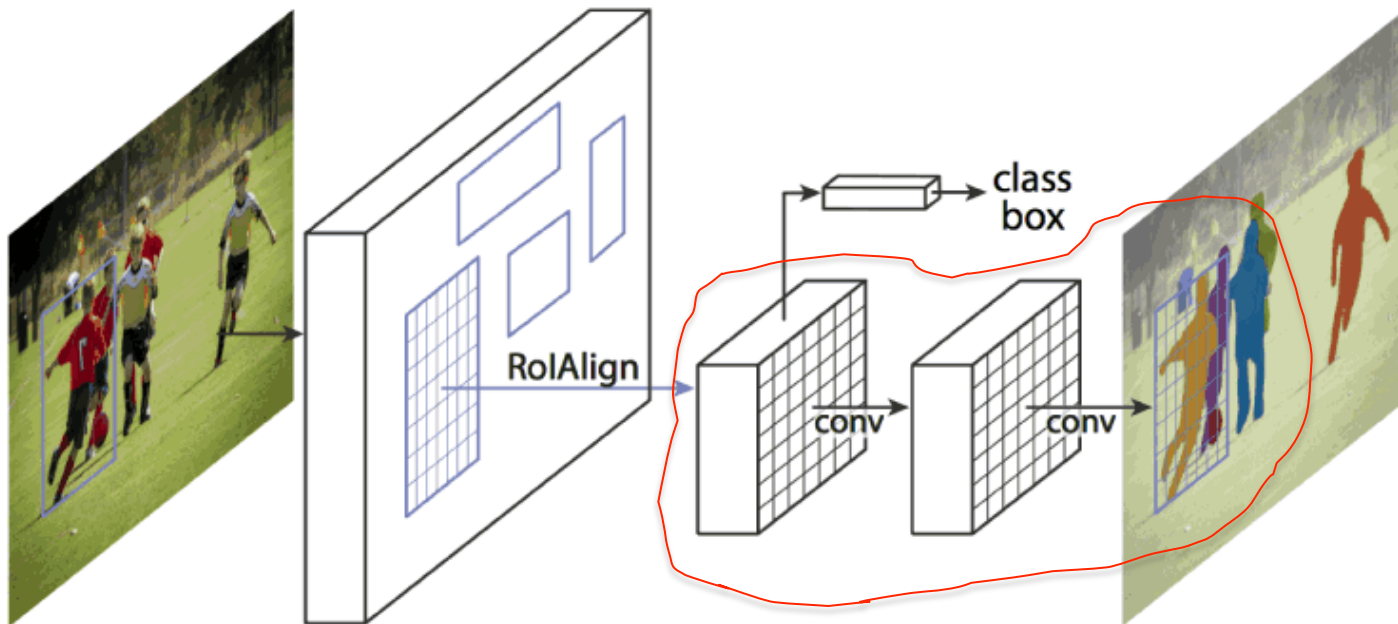
# Mask R-CNN for instance segmentation



The Mask R-CNN framework for instance segmentation

# Mask R-CNN for instance segmentation

Predict mxm mask from each ROI
Mask: encode input object's spatial layout



The Mask R-CNN framework for instance segmentation