

# Deep Image Segmentation

## [Spring 2020 CS-8395 Deep Learning in Medical Image Computing]

Instructor: Yuankai Huo, Ph.D.  
Department of Electrical Engineering and Computer Science  
Vanderbilt University

# Topics



- Review
- Semantic Segmentation
- U-Net
- Instance Segmentation

# Localization & Detection

**Classification**



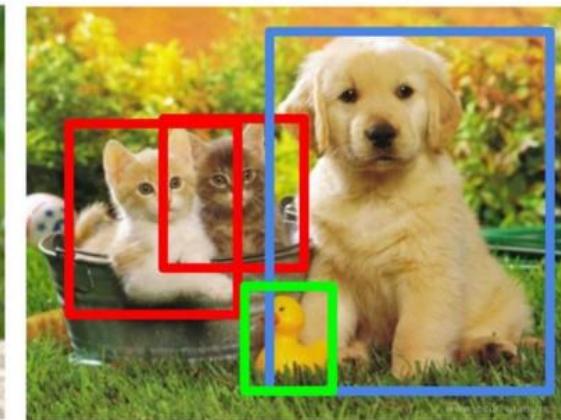
CAT

**Classification + Localization**



CAT

**Object Detection**



CAT, DOG, DUCK

**Instance Segmentation**



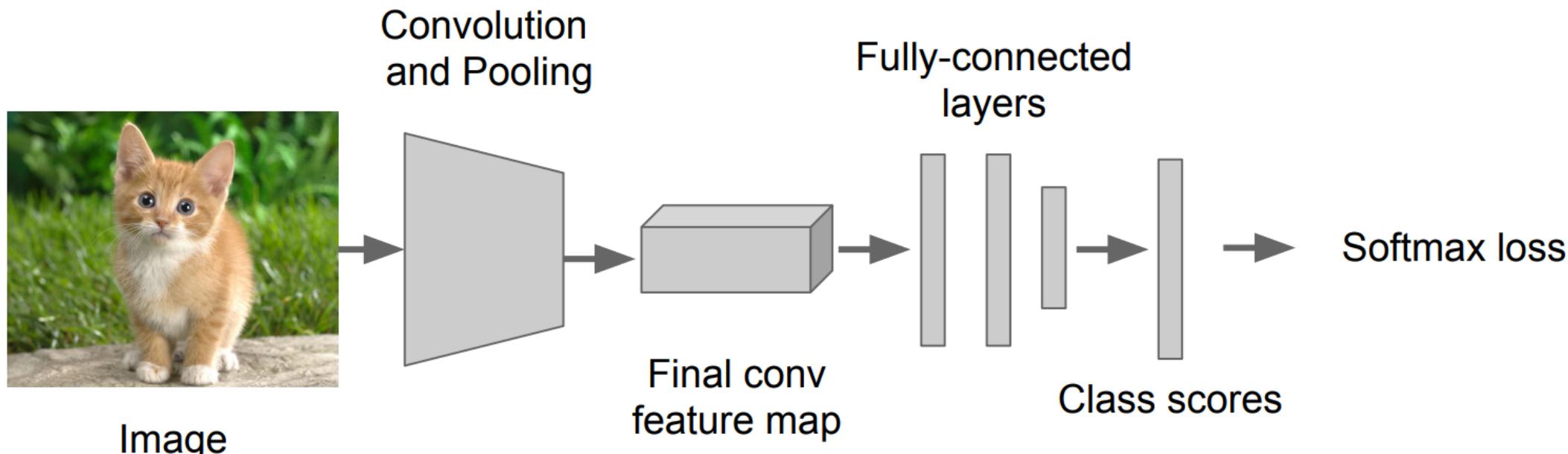
CAT, DOG, DUCK

Single object

Multiple objects

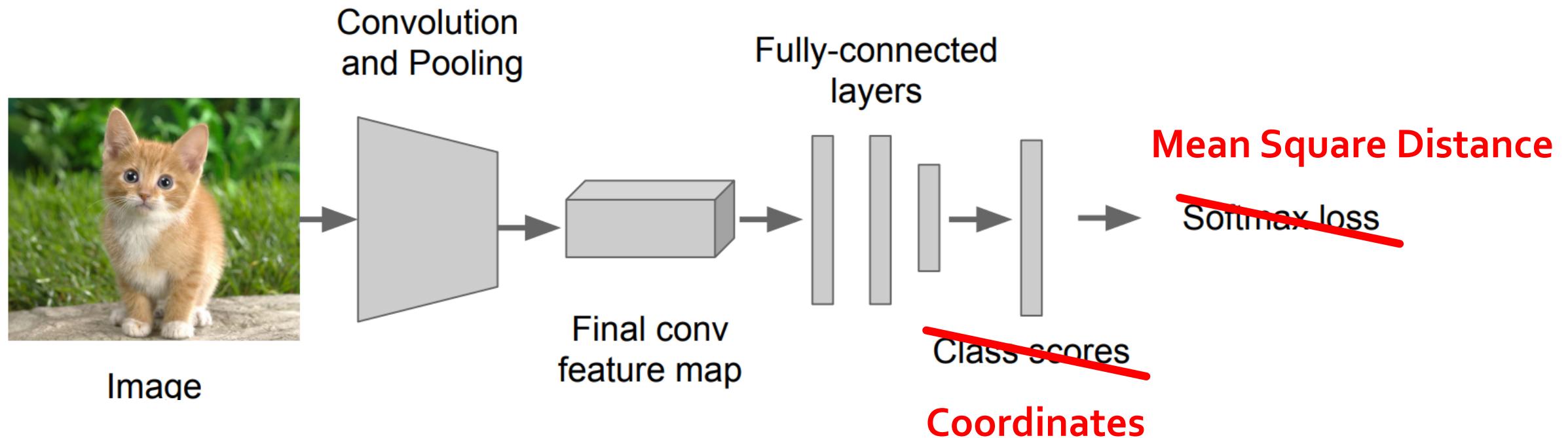
# Simple Recipe for Classification

## “CONV-POOL-FC”



# Simple Recipe for Localization

## “CONV-POOL-FC”



# Single Landmark Detection



Center Point ( $x, y$ )

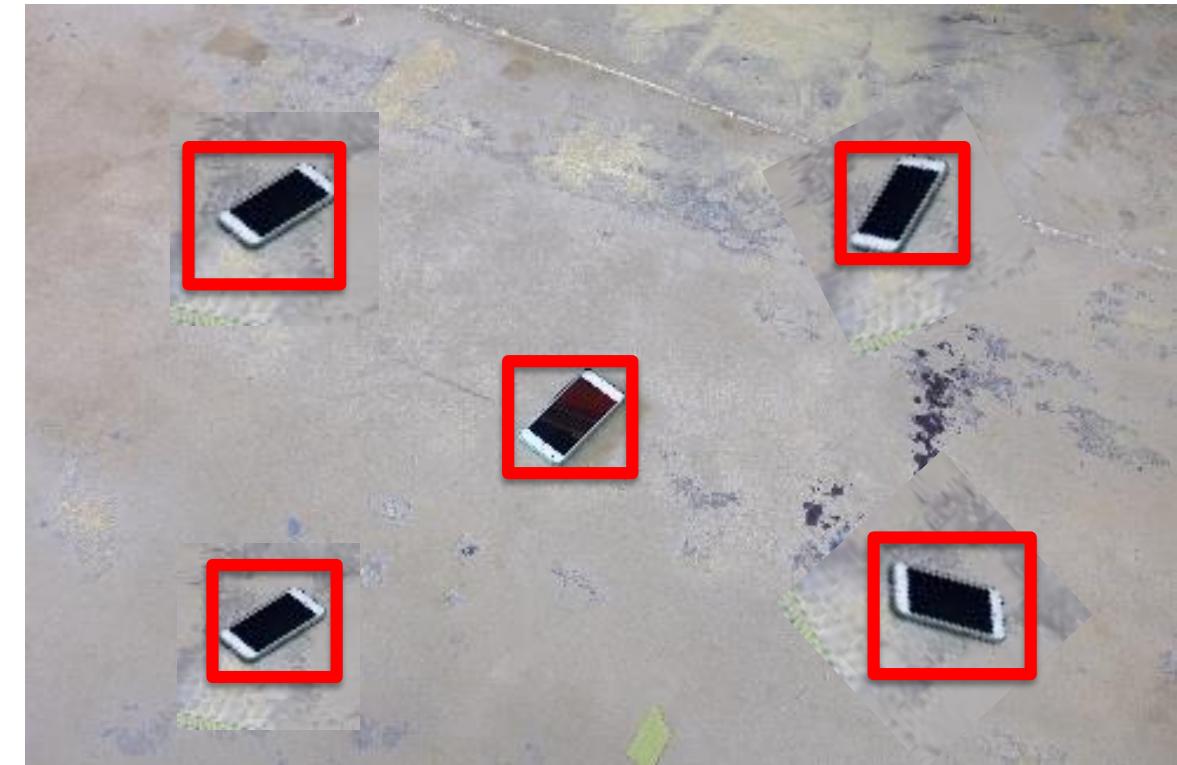


Bounding Box ( $x_1, y_1, x_2, y_2$ )

# Multi Landmark Detection

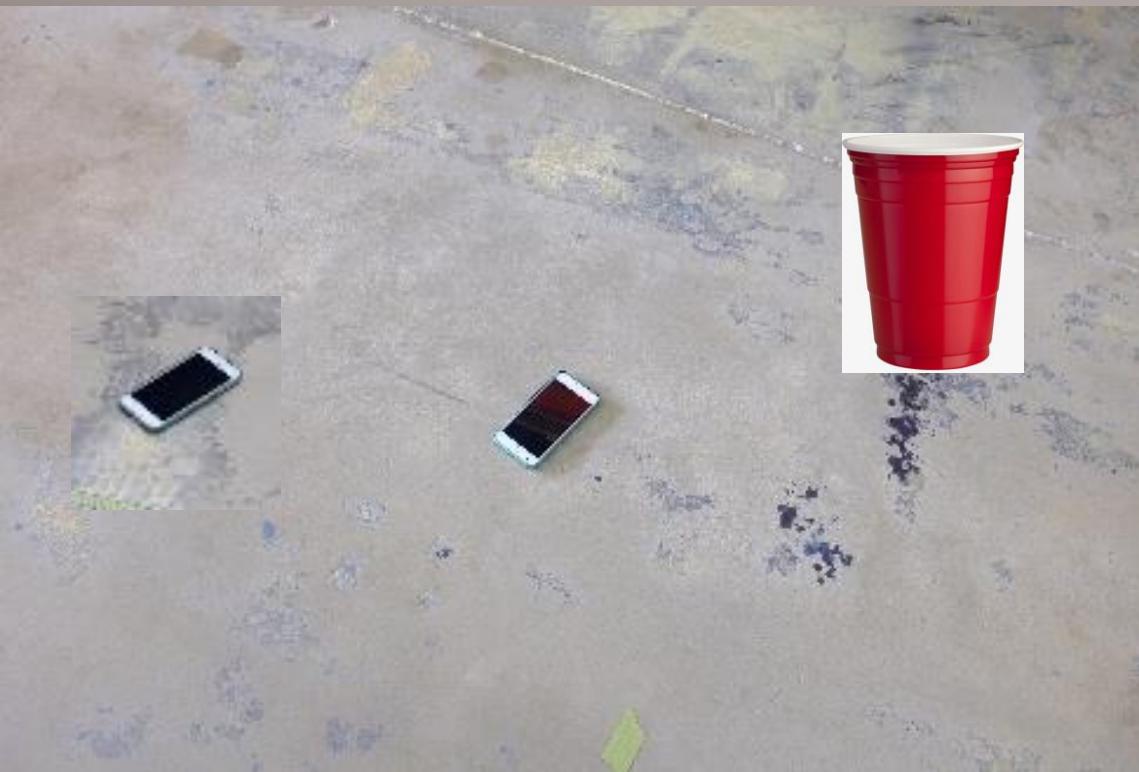


Center Points ( $x_1, y_1, x_2, y_2 \dots$ )



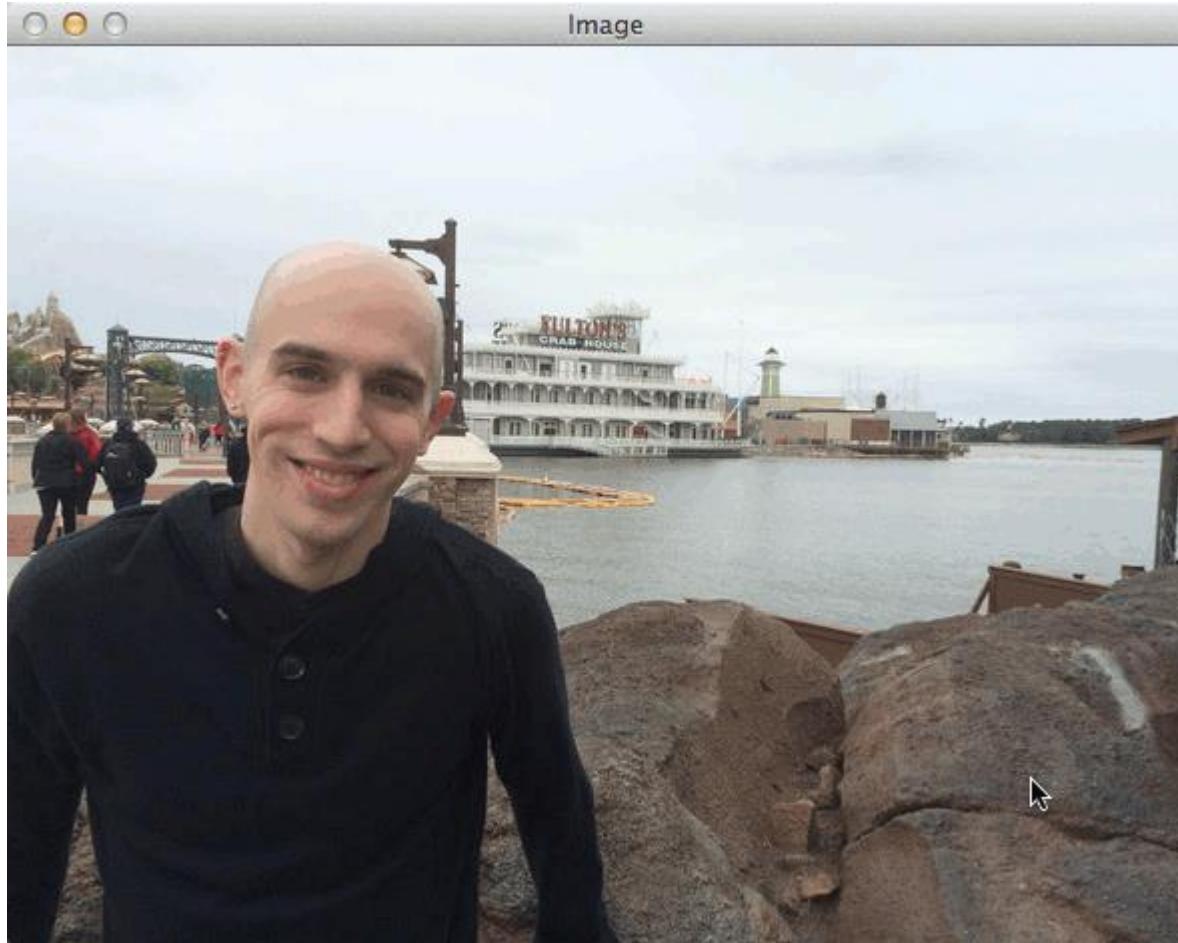
Bounding Boxes ( $x_{11}, y_{11}, x_{12}, y_{12}$   
 $x_{21}, y_{21}, x_{22}, y_{22} \dots$ )

# Now, what to do?



# Idea 1: Sliding Window

Detect Face

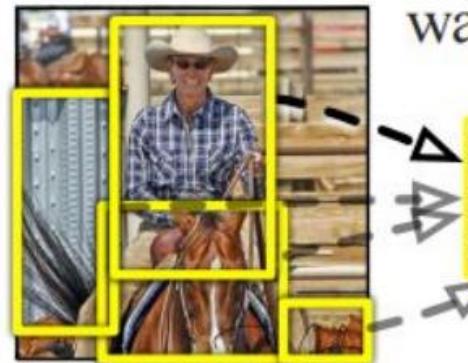


# R-CNN

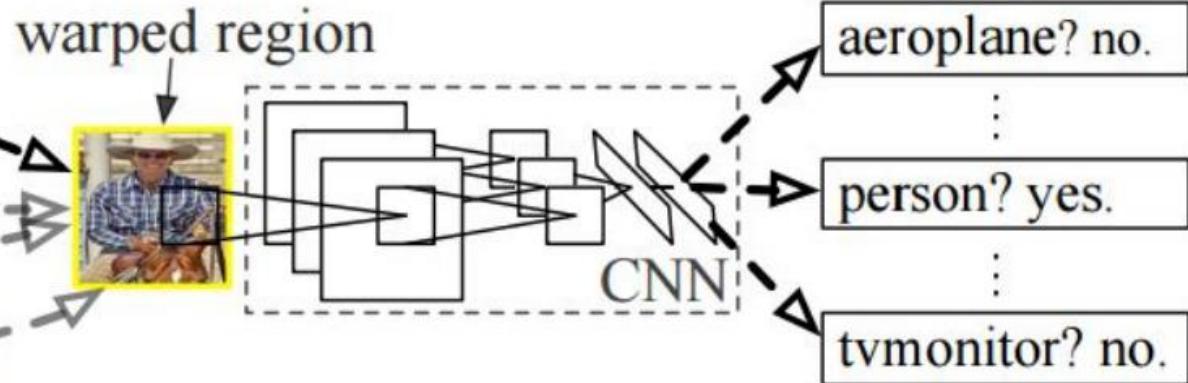
## R-CNN: *Regions with CNN features*



1. Input image



2. Extract region proposals (~2k)



3. Compute CNN features

4. Classify regions

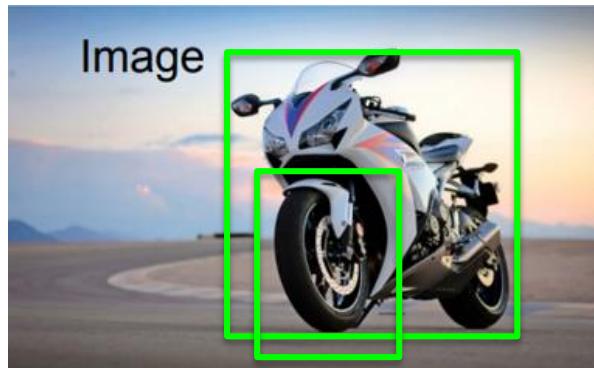
**Stage 1:  
Region Proposal**

**Stage 2:  
Classification**

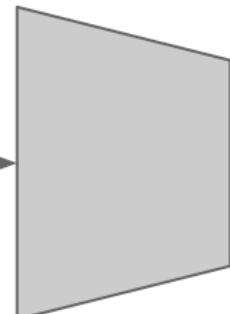
# Detection

## “CONV-POOL-FC”

Region Proposals



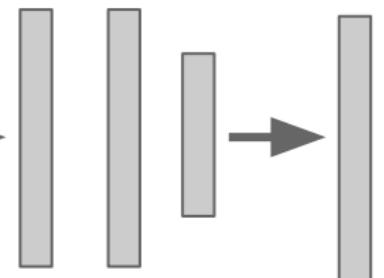
Convolution  
and Pooling



Final conv  
feature map



Fully-connected  
layers



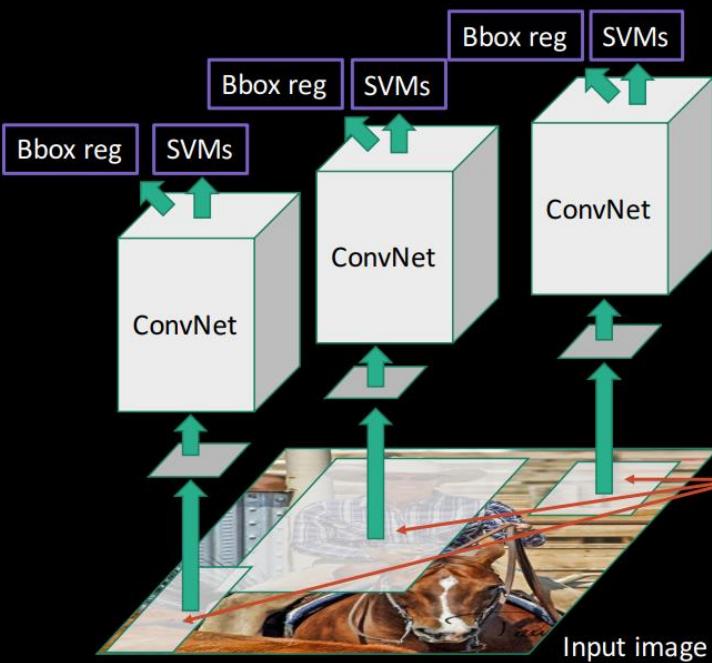
Class scores

Softmax loss

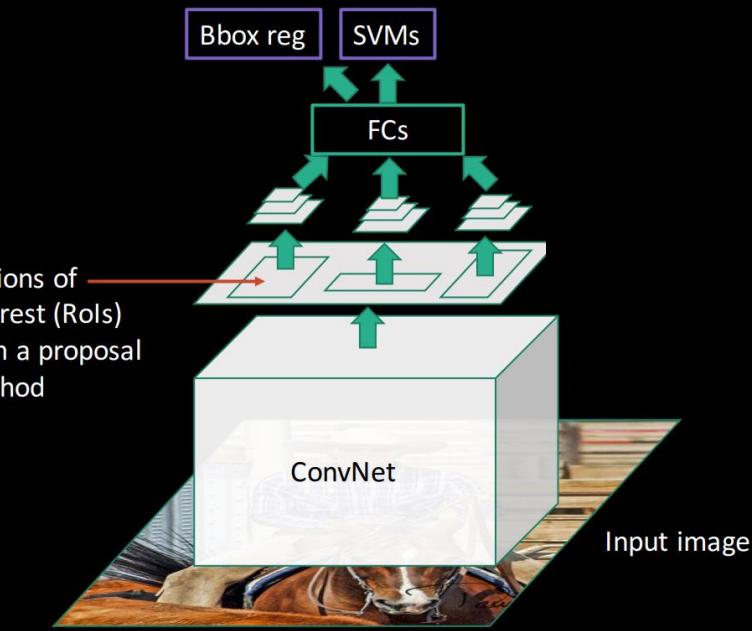
[http://slazebni.cs.illinois.edu/spring17/lec07\\_detection.pdf](http://slazebni.cs.illinois.edu/spring17/lec07_detection.pdf)  
[https://blog.csdn.net/sum\\_nap/article/details/80388110](https://blog.csdn.net/sum_nap/article/details/80388110)  
[http://cs231n.stanford.edu/slides/2016/winter1516\\_lecture8.pdf](http://cs231n.stanford.edu/slides/2016/winter1516_lecture8.pdf)

# RCNN → Fast RCNN

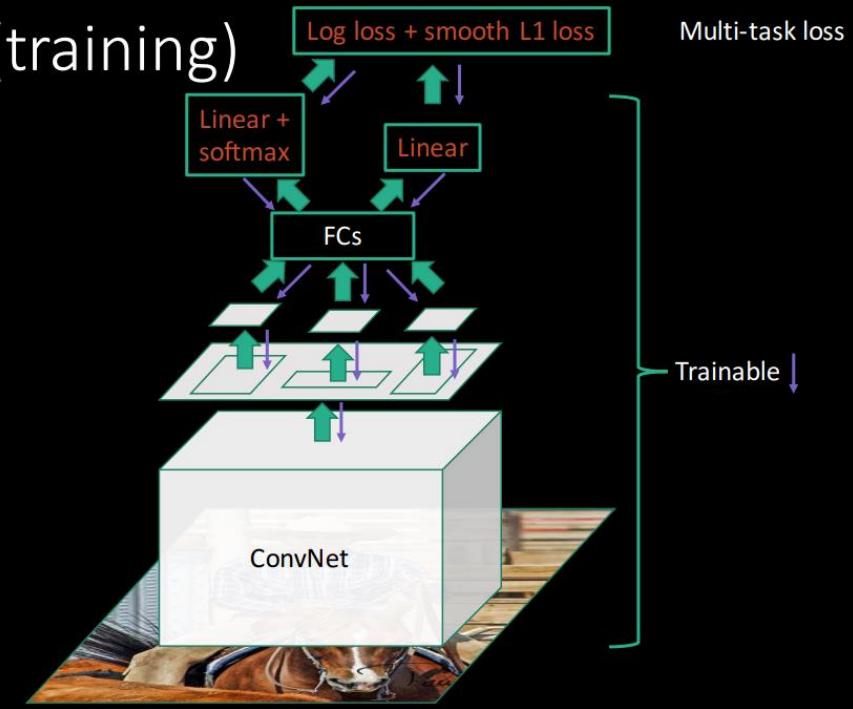
Slow R-CNN



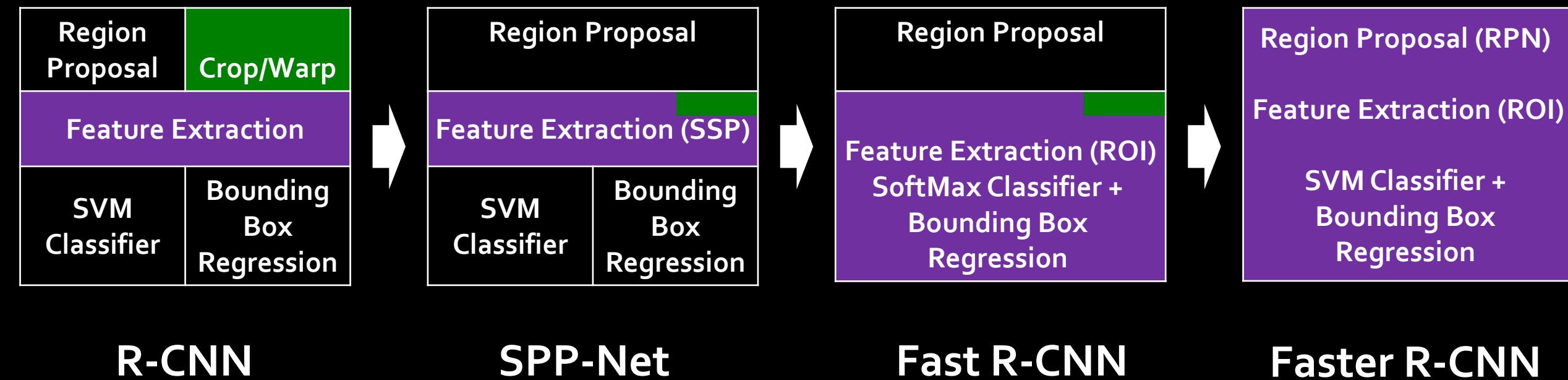
SPP-net



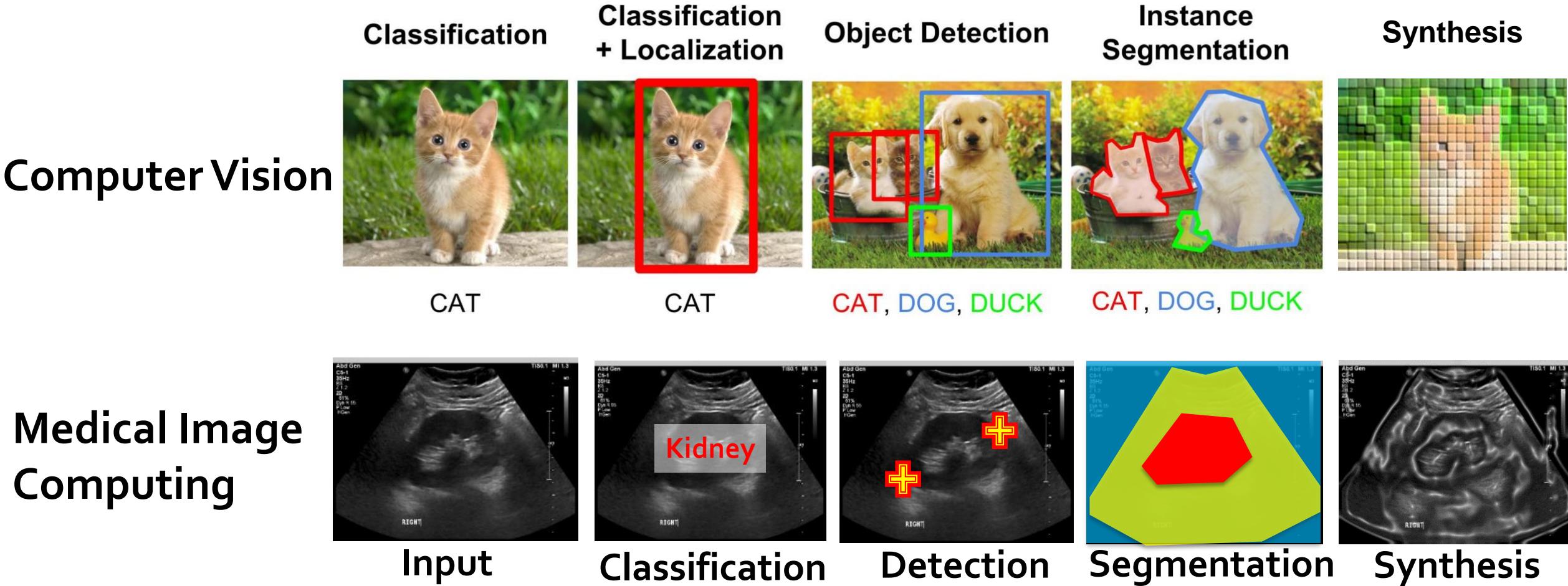
Fast R-CNN  
(training)



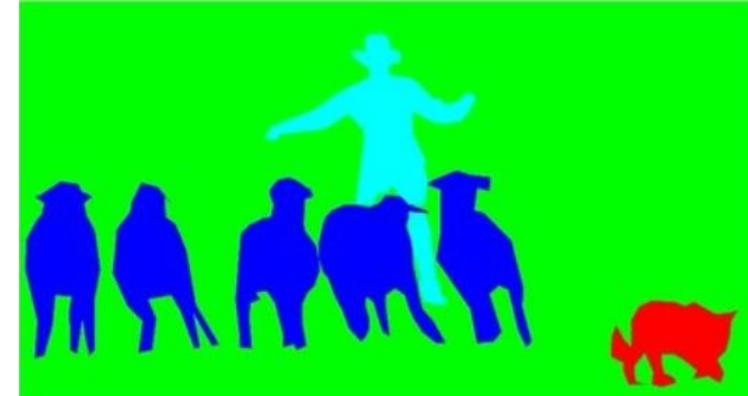
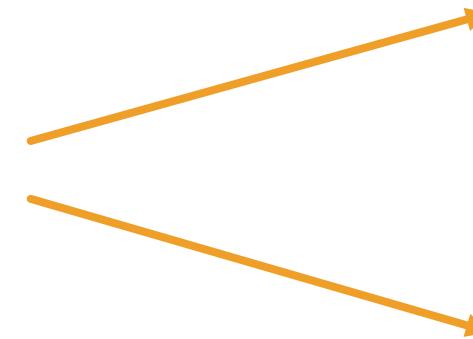
# Summerize



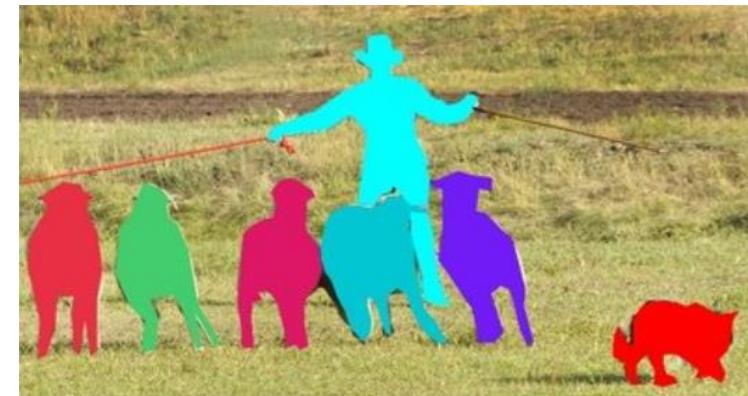
# Compare with Computer Vision



# Segmentation

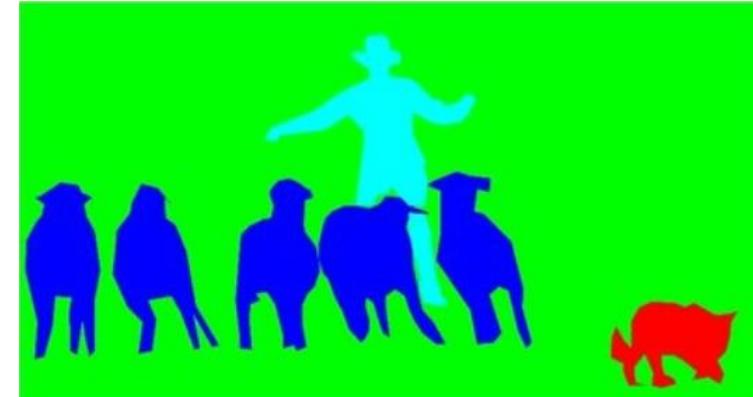


**Semantic Segmentation**

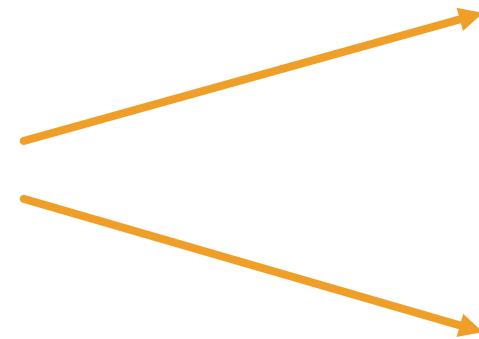


**Instance Segmentation**

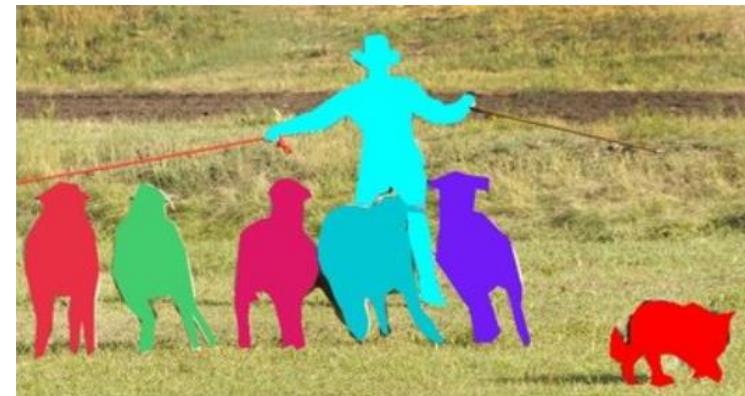
# Segmentation



# Segmentation



**Semantic Segmentation**



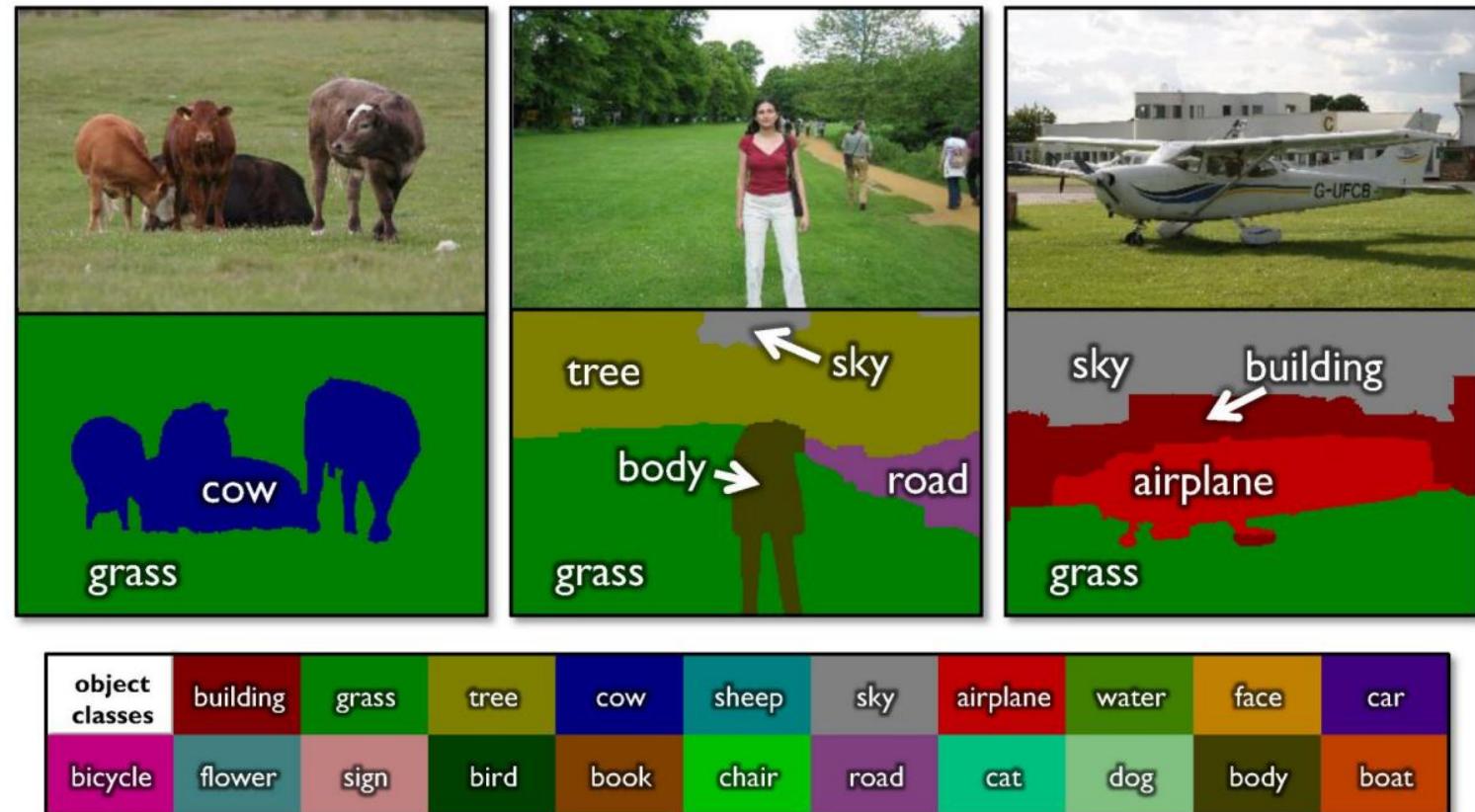
**Instance Segmentation**

# Semantic Segmentation

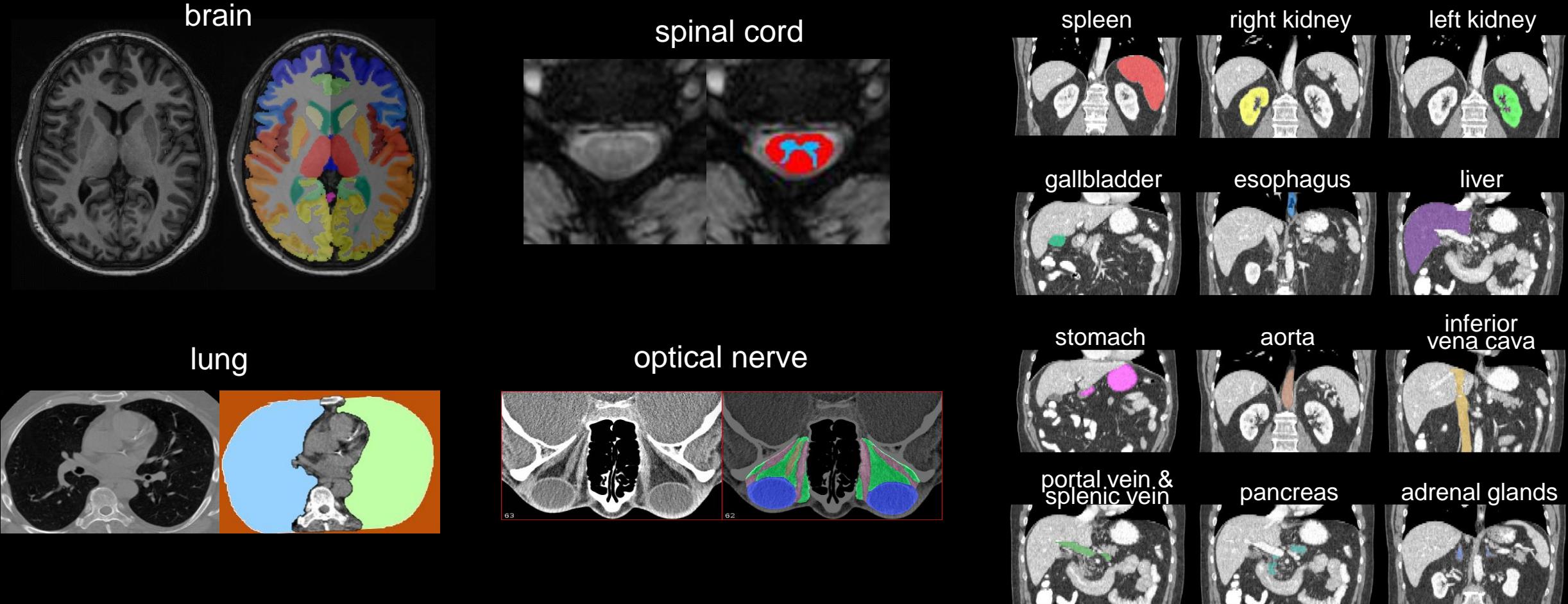
Label every pixel!

Don't differentiate instances (cows)

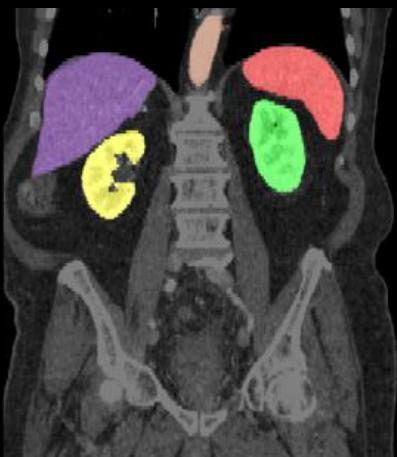
Classic computer vision problem



# Medical Image Segmentation

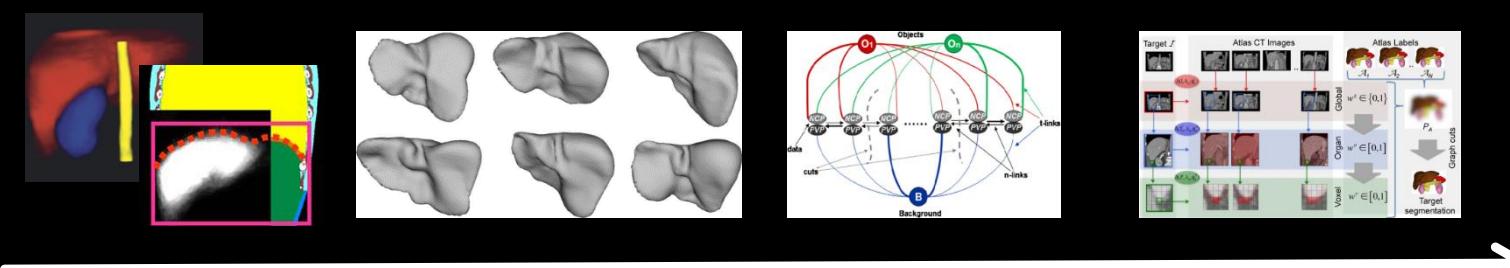


# History



Automatic  
Segmentation?

Probabilistic Atlas + Statistical Shape Model + Graph Cut Multi-Atlas Label Fusion

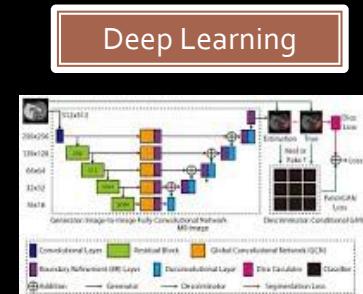


Park et al. (2003)  
Zhou et al. (2005)  
Shimizu et al. (2007)

Okada et al. (2008)  
Heimann et al. (2009)

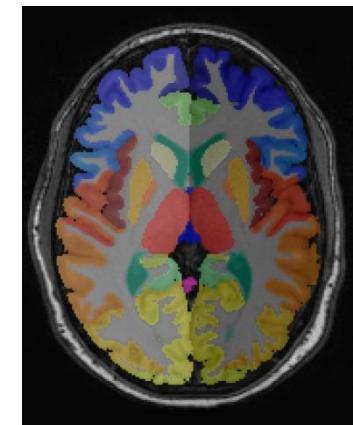
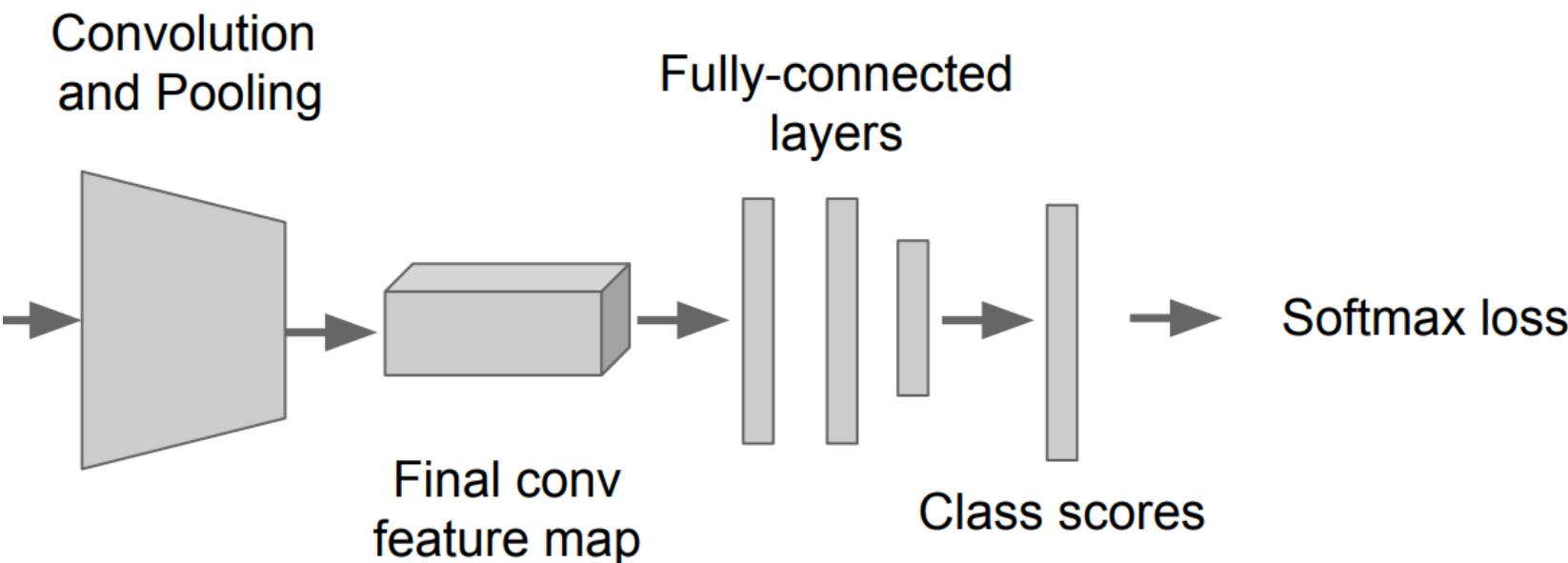
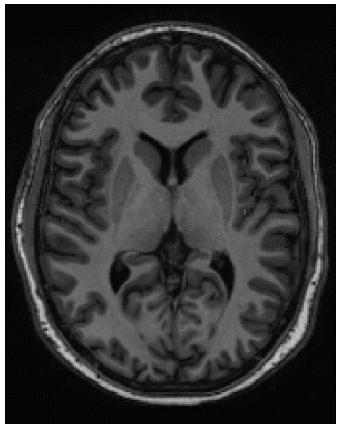
Bagci et al. (2012)  
Chen et al. (2012)  
Linguraru et al. (2012)

Wolz et al. (2013)  
Xu et al. (2015)



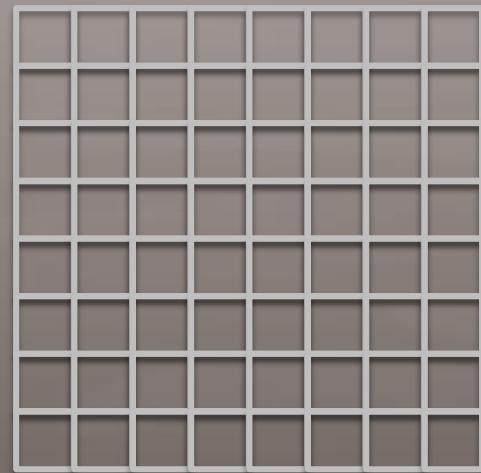
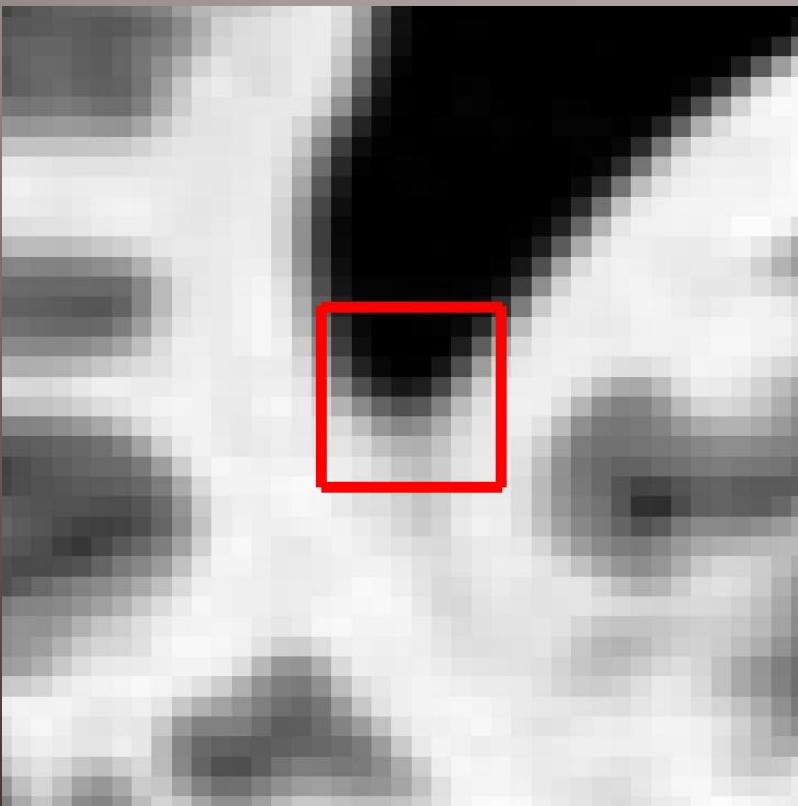
# Segmentation ?

**“CONV-POOL-FC”**

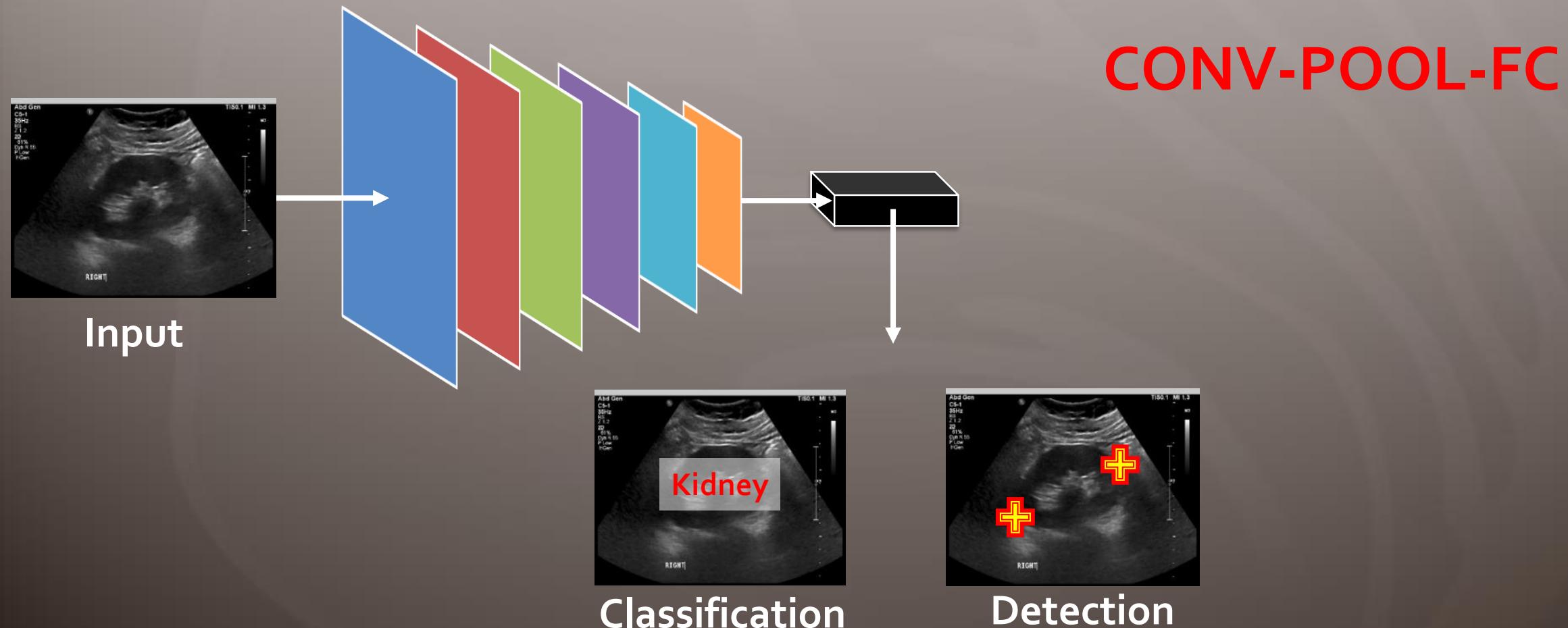


**everything can be solved by classification**

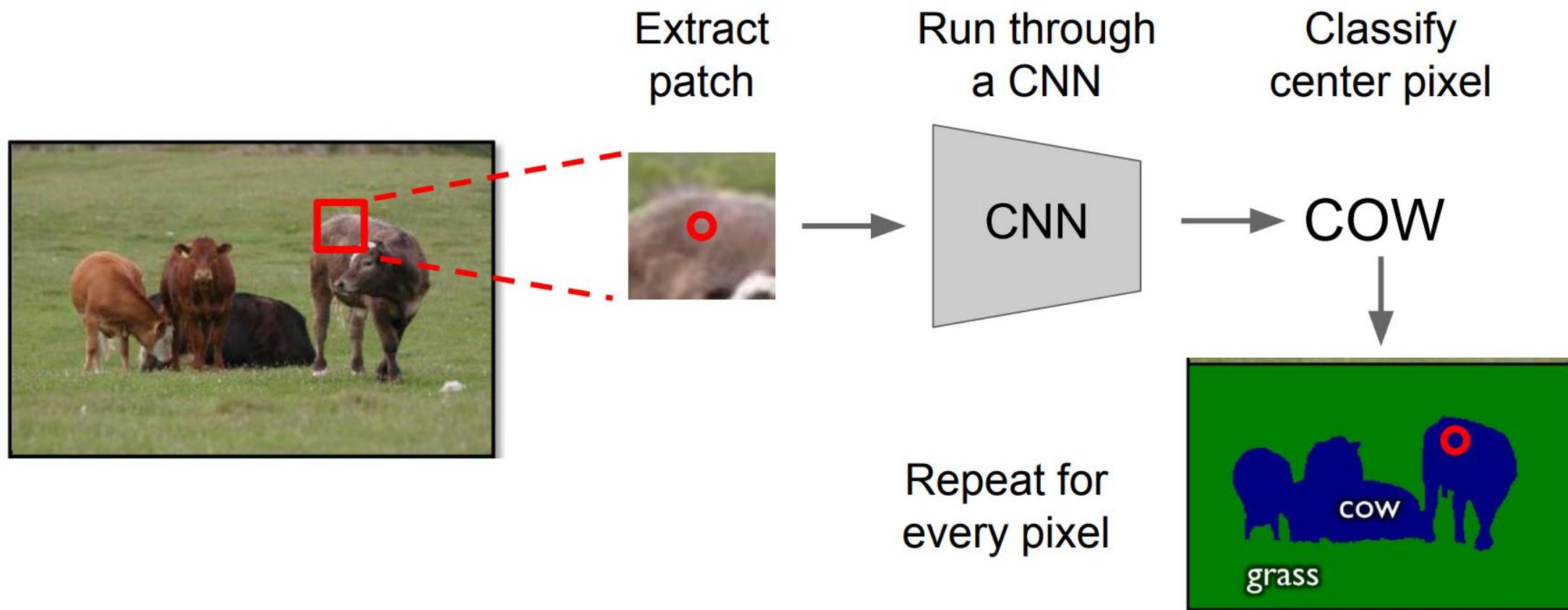
# Segmentation is Classification



# Classification / Regression

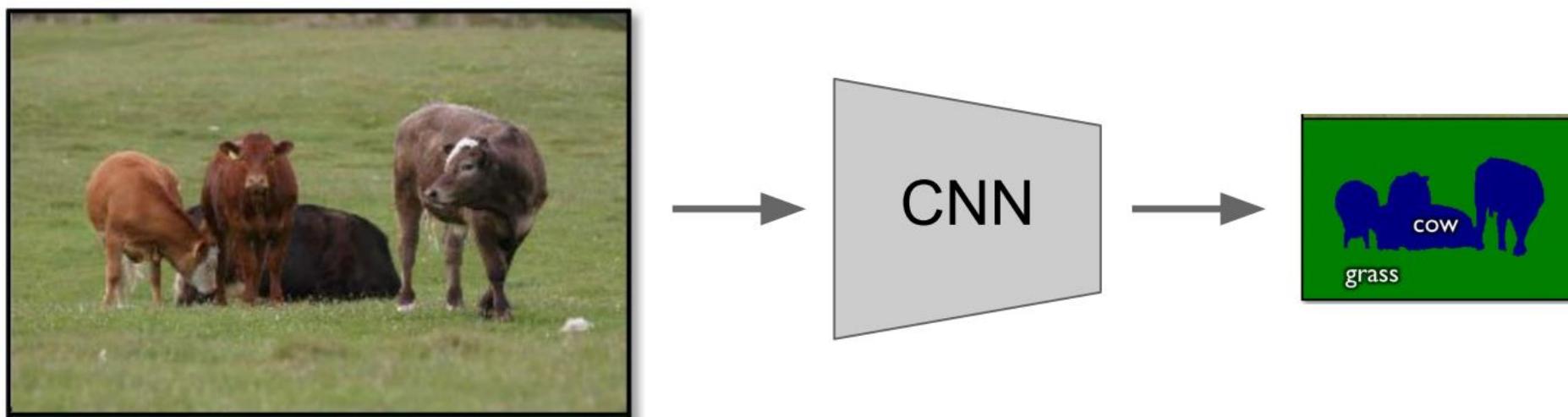


# Classification to Segmentation



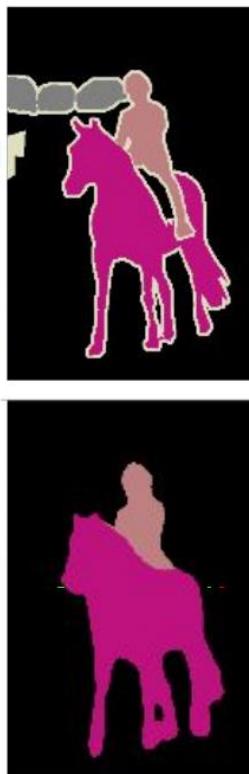
# FCN

Run “fully convolutional” network  
to get all pixels at once

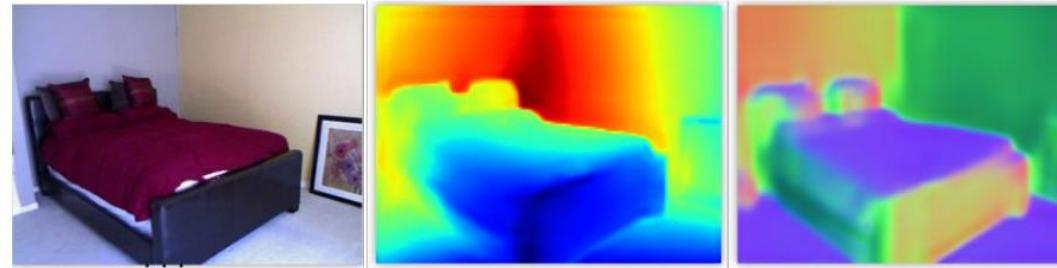


# Pixel In – Pixel Out (Image to Image)

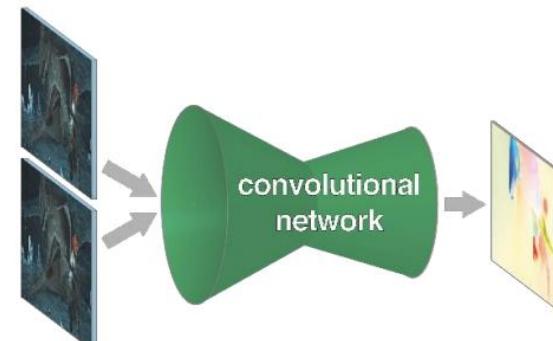
semantic segmentation



monocular depth + normals Eigen & Fergus 2015



colorization  
Zhang et al. 2016



optical flow Fischer et al. 2015

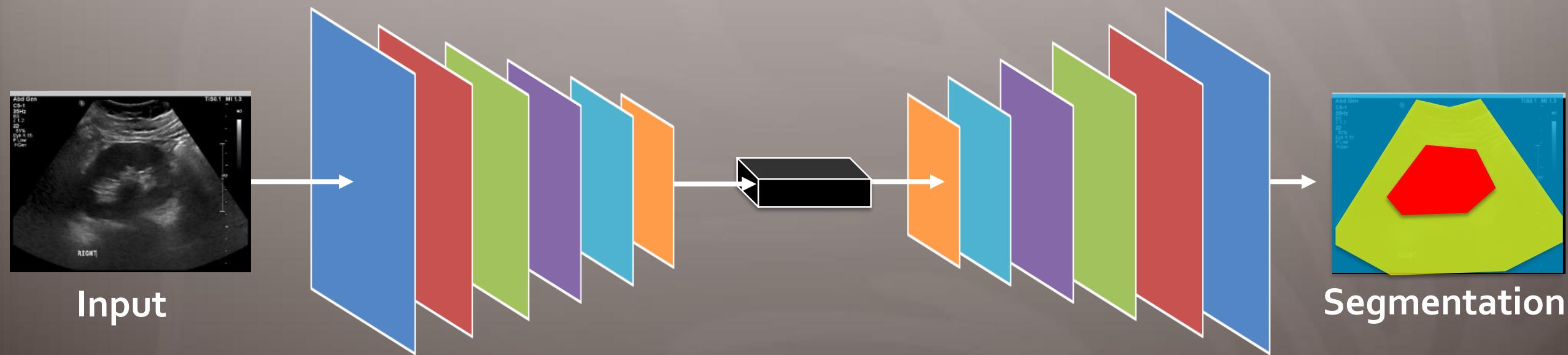


boundary prediction Xie & Tu 2015

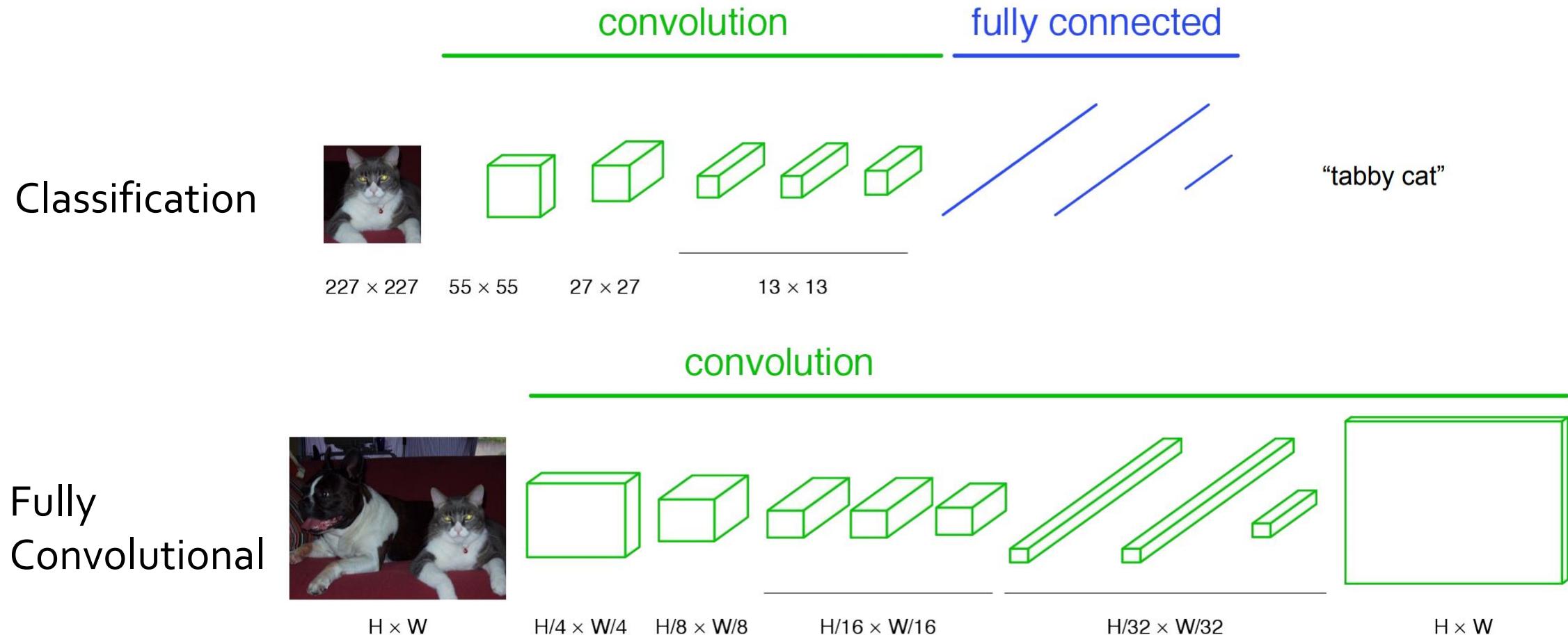


2

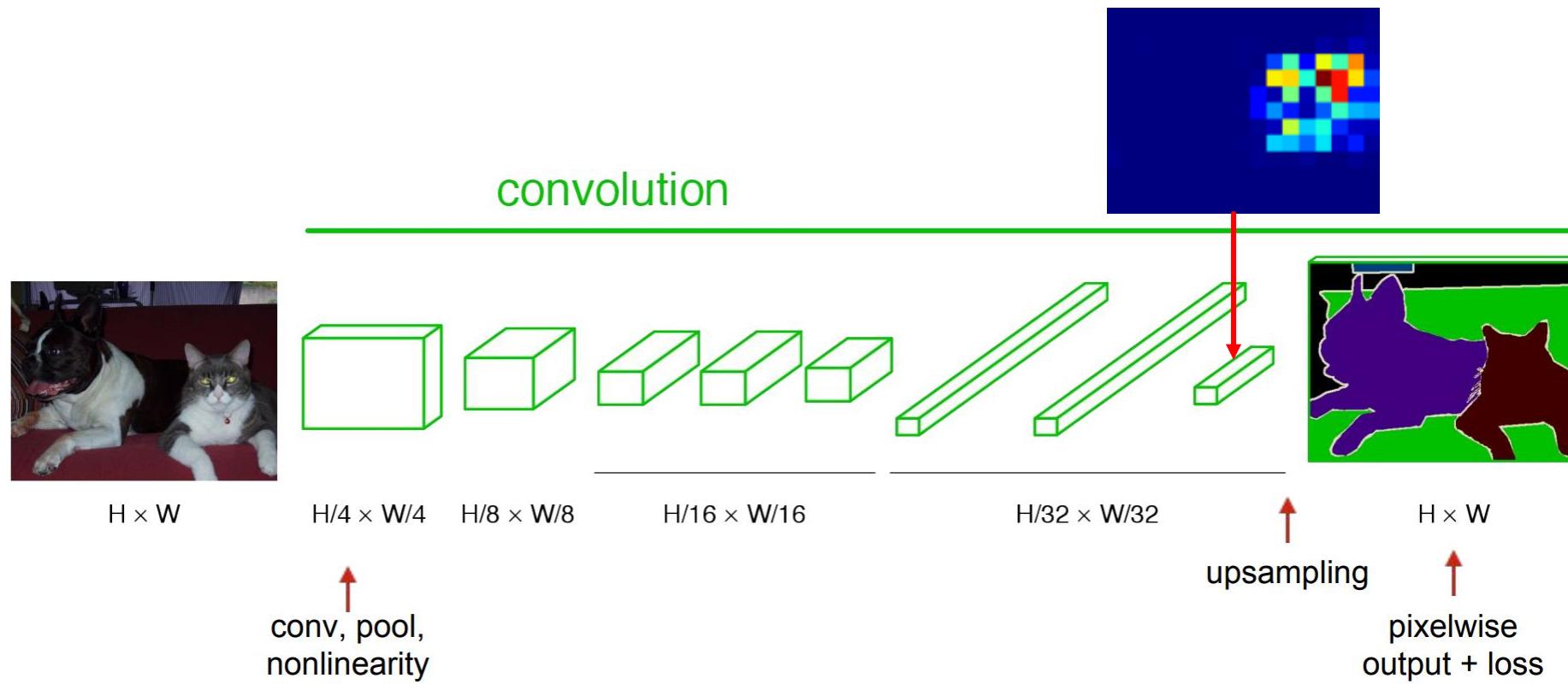
# Segmentation (Image to Image)



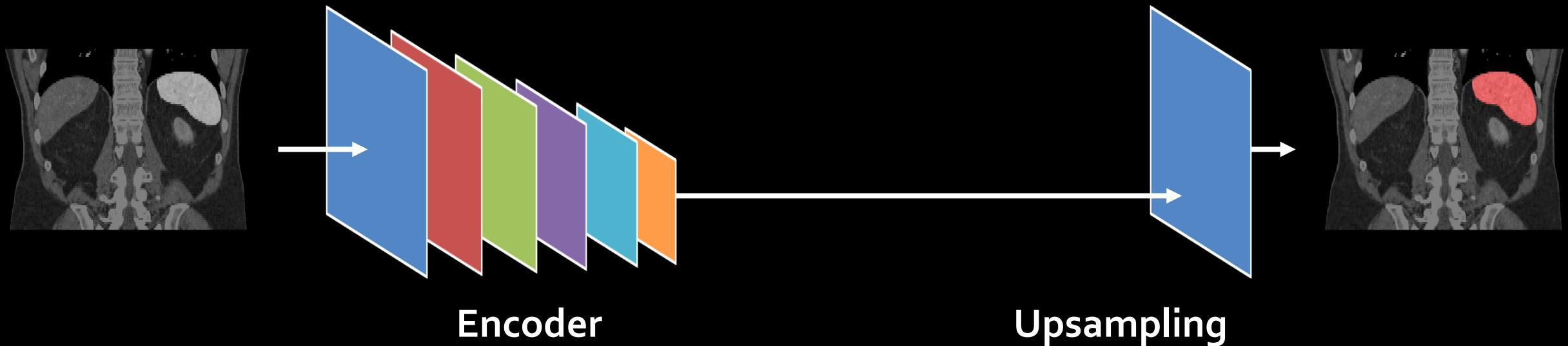
# Change a Classification Network



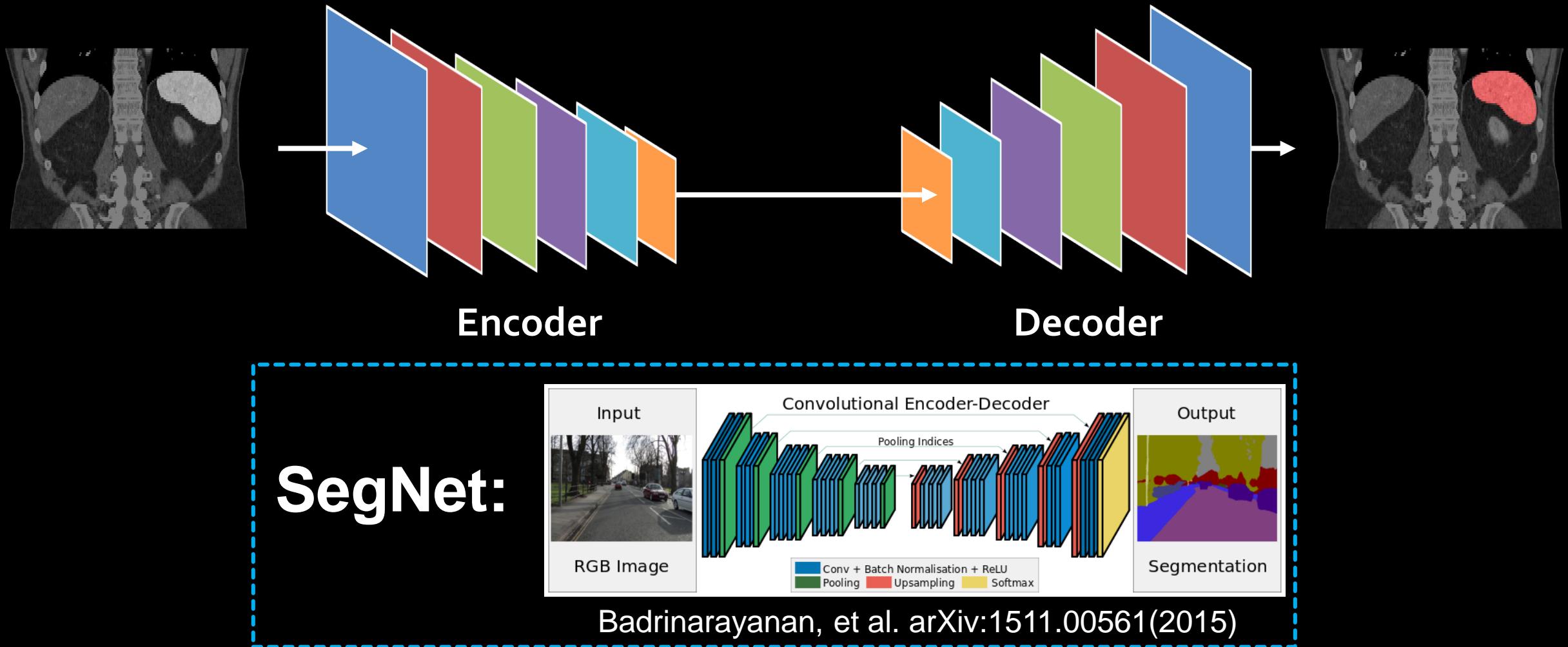
# Naïve Idea



# Fully Convolutional Neural Network

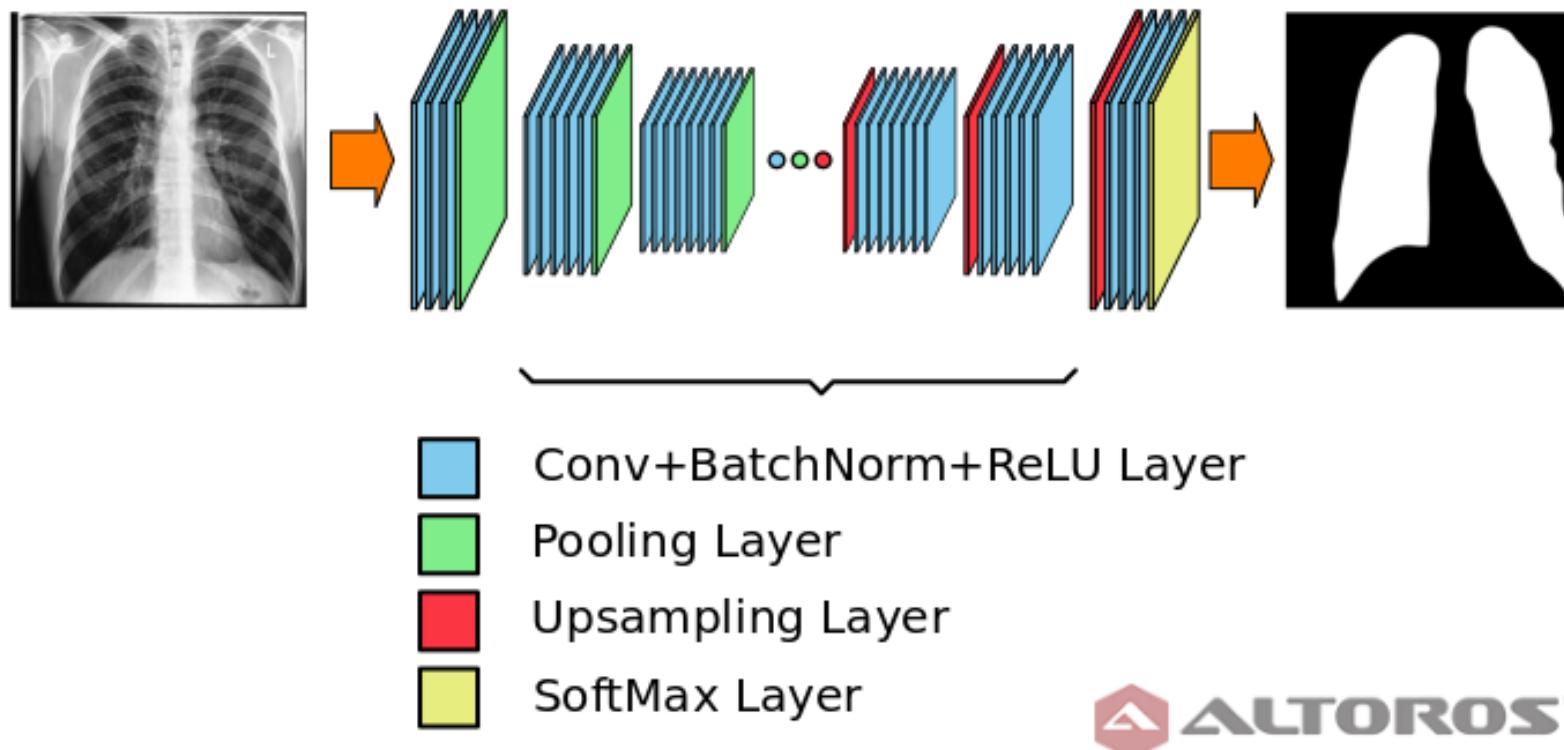


# Encoder – Decoder Strategy

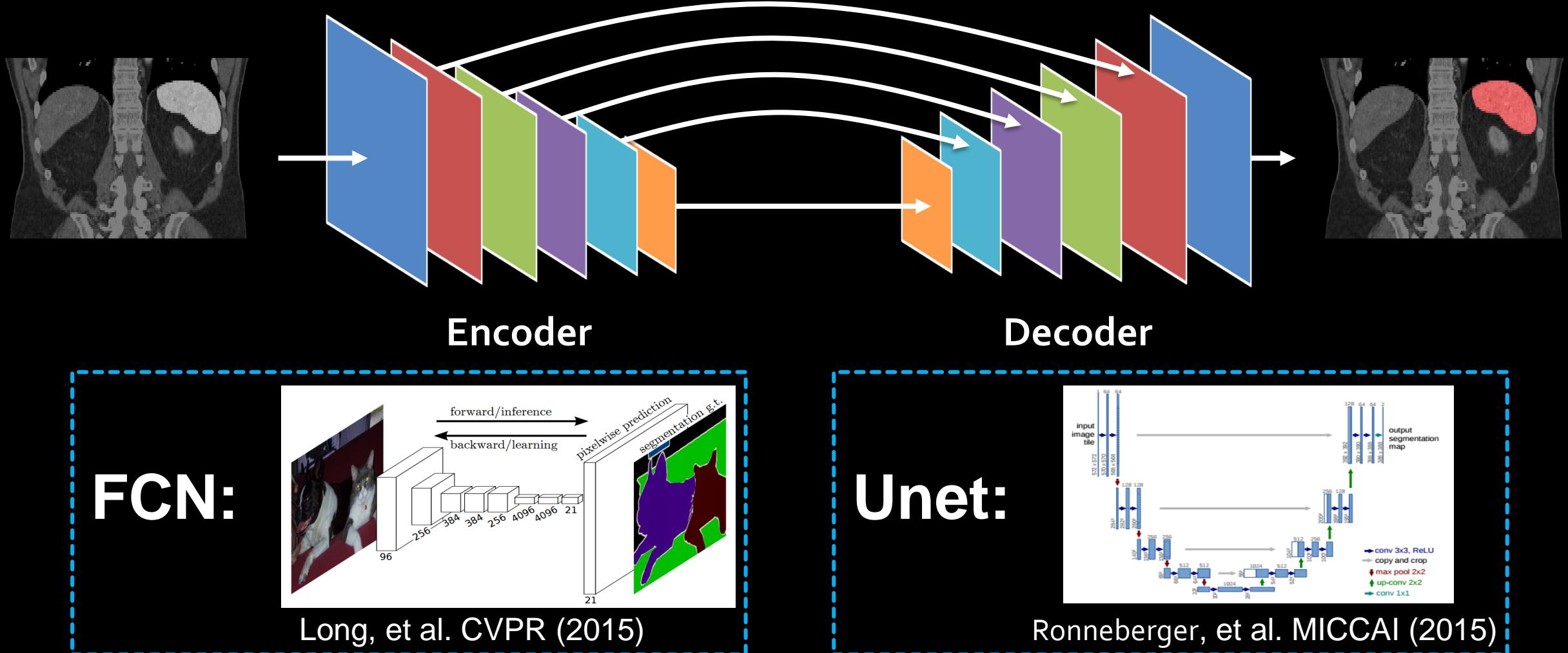


# SegNet

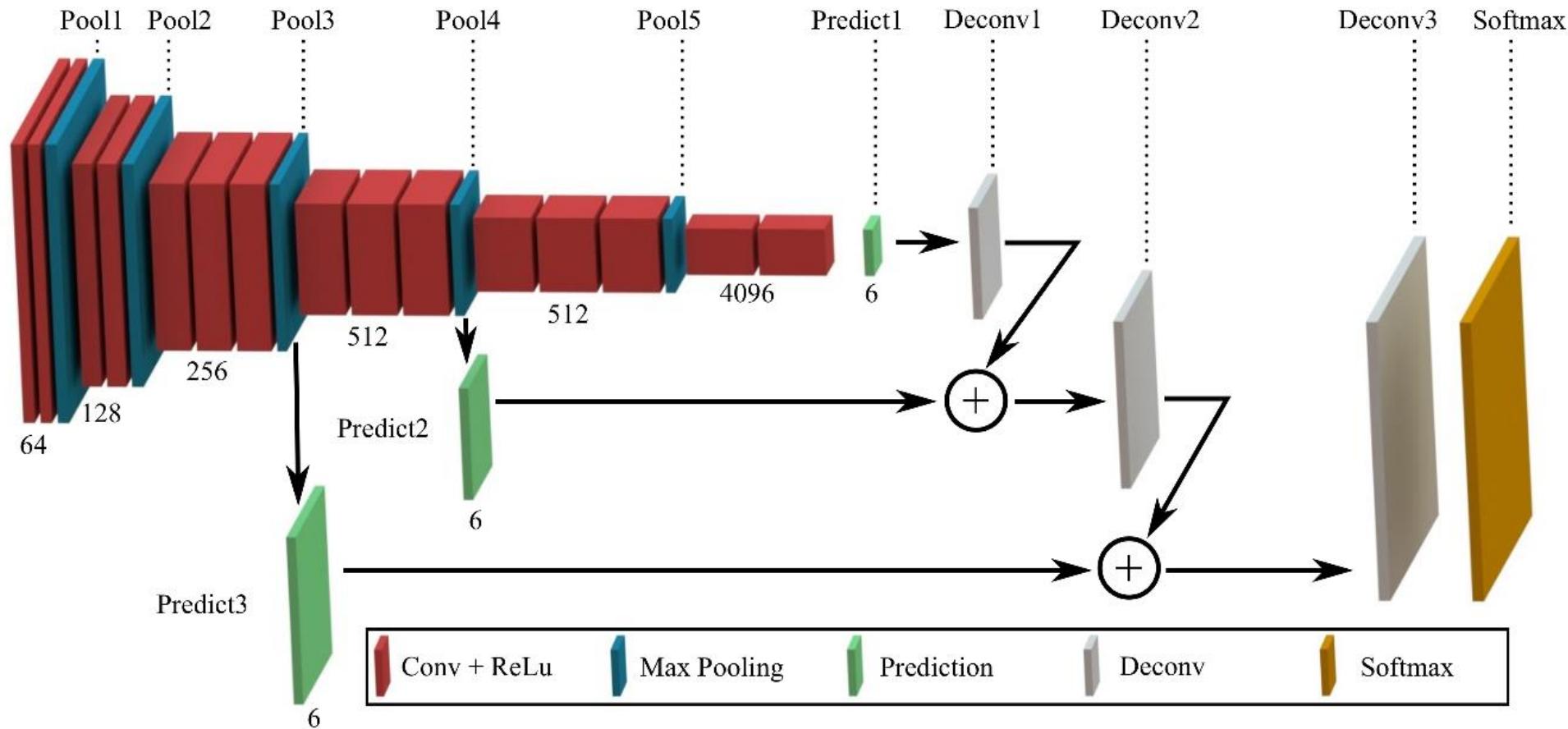
## Encoder-decoder CNN



# Fully Convolutional Network (FCN)

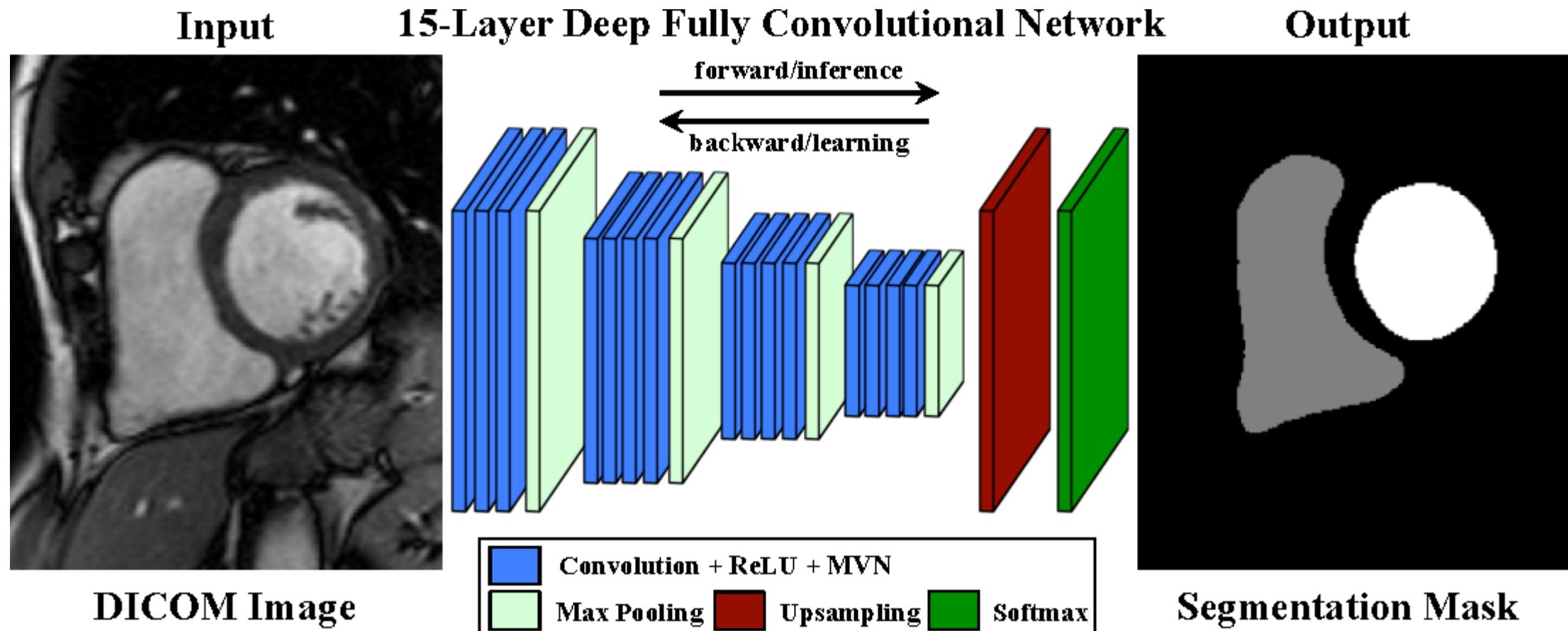


# Fully Convolutional Network

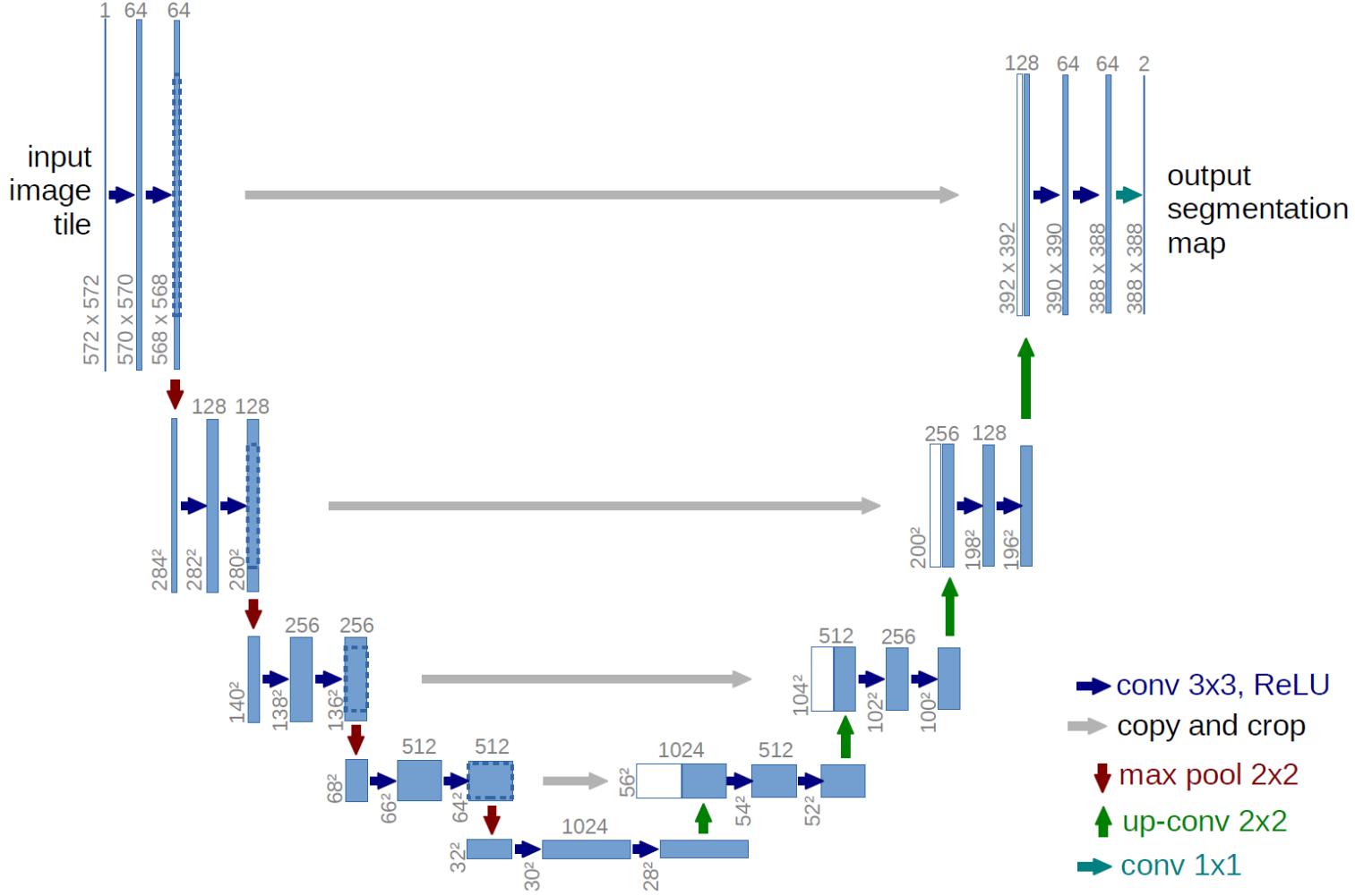


<https://medium.com/@wilburdes/semantic-segmentation-using-fully-convolutional-neural-networks-86e45336f99b>

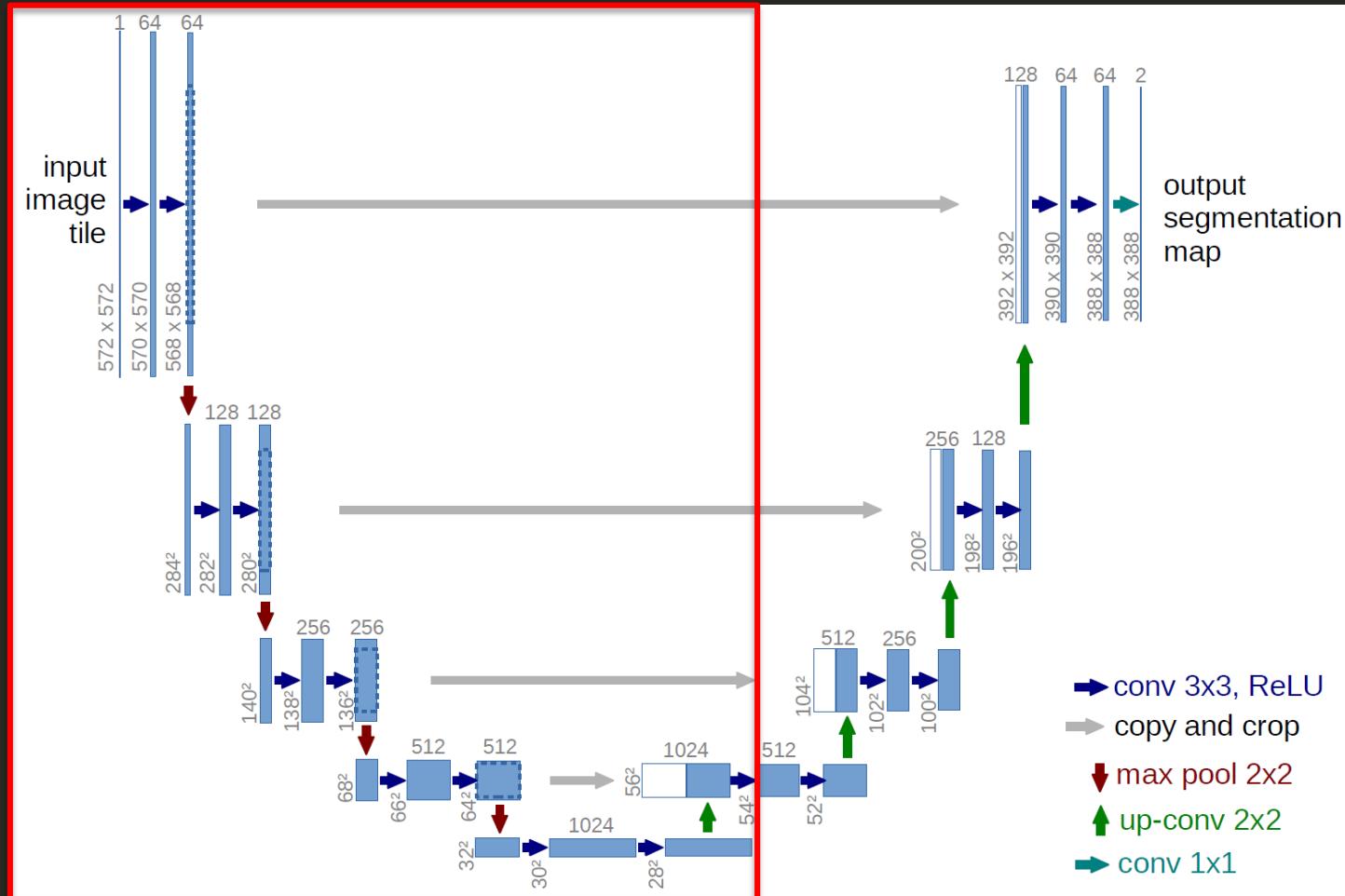
# FCN



# U-Net



# Encoder



```
layer1 = F.relu(self.conv1_input(x))  
layer1 = F.relu(self.conv1(layer1))
```

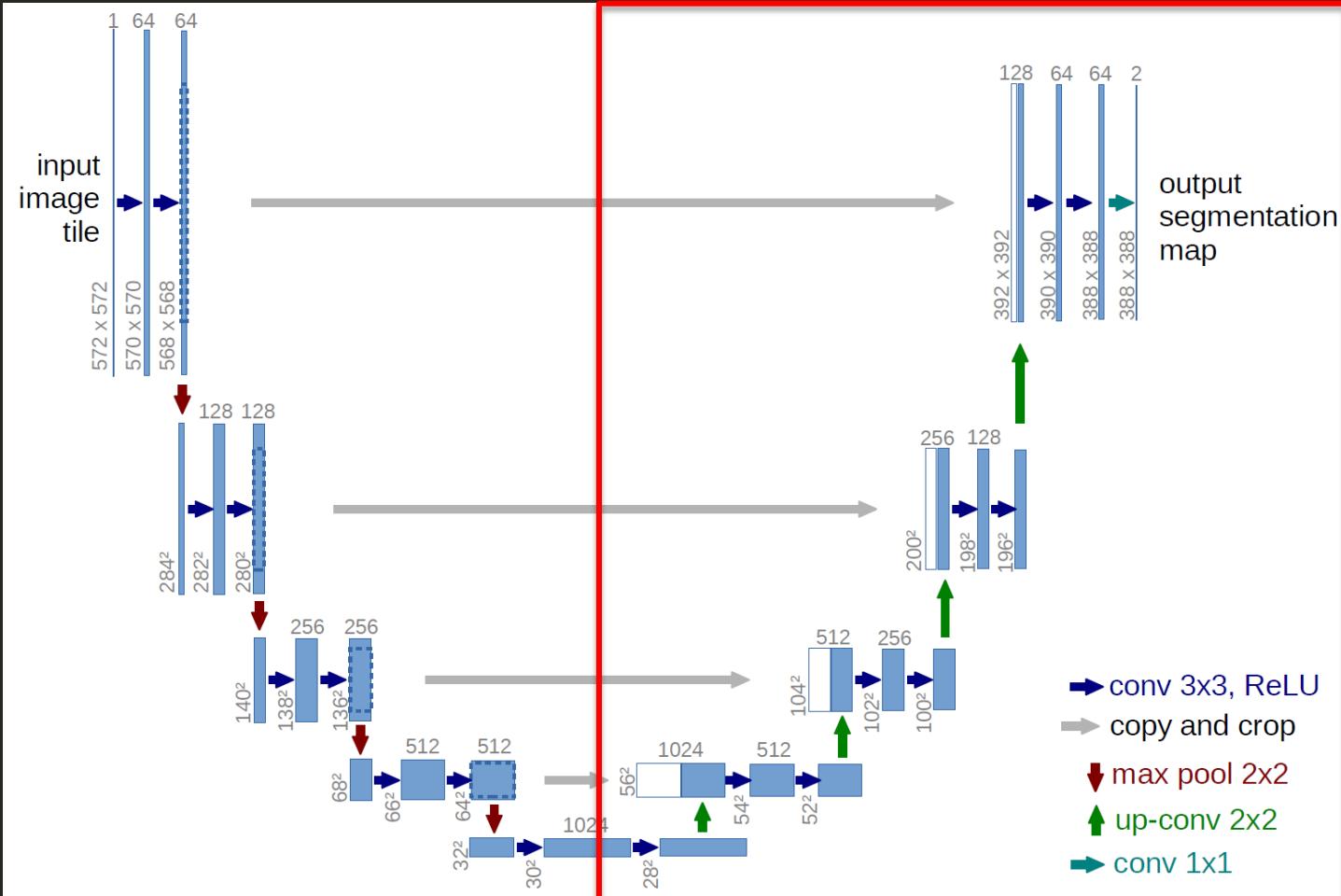
```
layer2 = F.max_pool2d(layer1, 2)  
layer2 = F.relu(self.conv2_input(layer2))  
layer2 = F.relu(self.conv2(layer2))
```

```
layer3 = F.max_pool2d(layer2, 2)  
layer3 = F.relu(self.conv3_input(layer3))  
layer3 = F.relu(self.conv3(layer3))
```

```
layer4 = F.max_pool2d(layer3, 2)  
layer4 = F.relu(self.conv4_input(layer4))  
layer4 = F.relu(self.conv4(layer4))
```

```
layer5 = F.max_pool2d(layer4, 2)  
layer5 = F.relu(self.conv5_input(layer5))  
layer5 = F.relu(self.conv5(layer5))
```

# Decoder



```

layer6 = F.relu(self.conv6_up(layer5))
layer6 = torch.cat((layer4, layer6), 1)
layer6 = F.relu(self.conv6_input(layer6))
layer6 = F.relu(self.conv6(layer6))

layer7 = F.relu(self.conv7_up(layer6))
layer7 = torch.cat((layer3, layer7), 1)
layer7 = F.relu(self.conv7_input(layer7))
layer7 = F.relu(self.conv7(layer7))

layer8 = F.relu(self.conv8_up(layer7))
layer8 = torch.cat((layer2, layer8), 1)
layer8 = F.relu(self.conv8_input(layer8))
layer8 = F.relu(self.conv8(layer8))

layer9 = F.relu(self.conv9_up(layer8))
layer9 = torch.cat((layer1, layer9), 1)
layer9 = F.relu(self.conv9_input(layer9))
layer9 = F.relu(self.conv9(layer9))
layer9 = self.final(self.conv9_output(layer9))

return layer9

```

# All Code for U-Net

## Encoder

```
layer1 = F.relu(self.conv1_input(x))
layer1 = F.relu(self.conv1(layer1))

layer2 = F.max_pool2d(layer1, 2)
layer2 = F.relu(self.conv2_input(layer2))
layer2 = F.relu(self.conv2(layer2))

layer3 = F.max_pool2d(layer2, 2)
layer3 = F.relu(self.conv3_input(layer3))
layer3 = F.relu(self.conv3(layer3))

layer4 = F.max_pool2d(layer3, 2)
layer4 = F.relu(self.conv4_input(layer4))
layer4 = F.relu(self.conv4(layer4))

layer5 = F.max_pool2d(layer4, 2)
layer5 = F.relu(self.conv5_input(layer5))
layer5 = F.relu(self.conv5(layer5))
```

## Decoder

```
layer6 = F.relu(self.conv6_up(layer5))
layer6 = torch.cat((layer4, layer6), 1)
layer6 = F.relu(self.conv6_input(layer6))
layer6 = F.relu(self.conv6(layer6))

layer7 = F.relu(self.conv7_up(layer6))
layer7 = torch.cat((layer3, layer7), 1)
layer7 = F.relu(self.conv7_input(layer7))
layer7 = F.relu(self.conv7(layer7))

layer8 = F.relu(self.conv8_up(layer7))
layer8 = torch.cat((layer2, layer8), 1)
layer8 = F.relu(self.conv8_input(layer8))
layer8 = F.relu(self.conv8(layer8))

layer9 = F.relu(self.conv9_up(layer8))
layer9 = torch.cat((layer1, layer9), 1)
layer9 = F.relu(self.conv9_input(layer9))
layer9 = F.relu(self.conv9(layer9))
layer9 = self.final(self.conv9_output(layer9))

return layer9
```

## Parameters

self.conv1_input =	nn.Conv2d(1, 64, 3, padding=1)
self.conv1 =	nn.Conv2d(64, 64, 3, padding=1)
self.conv2_input =	nn.Conv2d(64, 128, 3, padding=1)
self.conv2 =	nn.Conv2d(128, 128, 3, padding=1)
self.conv3_input =	nn.Conv2d(128, 256, 3, padding=1)
self.conv3 =	nn.Conv2d(256, 256, 3, padding=1)
self.conv4_input =	nn.Conv2d(256, 512, 3, padding=1)
self.conv4 =	nn.Conv2d(512, 512, 3, padding=1)
self.conv5_input =	nn.Conv2d(512, 1024, 3, padding=1)
self.conv5 =	nn.Conv2d(1024, 1024, 3, padding=1)
self.conv6_up =	nn.ConvTranspose2d(1024, 512, 2, 2)
self.conv6_input =	nn.Conv2d(1024, 512, 3, padding=1)
self.conv6 =	nn.Conv2d(512, 512, 3, padding=1)
self.conv7_up =	nn.ConvTranspose2d(512, 256, 2, 2)
self.conv7_input =	nn.Conv2d(512, 256, 3, padding=1)
self.conv7 =	nn.Conv2d(256, 256, 3, padding=1)
self.conv8_up =	nn.ConvTranspose2d(256, 128, 2, 2)
self.conv8_input =	nn.Conv2d(256, 128, 3, padding=1)
self.conv8 =	nn.Conv2d(128, 128, 3, padding=1)
self.conv9_up =	nn.ConvTranspose2d(128, 64, 2, 2)
self.conv9_input =	nn.Conv2d(128, 64, 3, padding=1)
self.conv9 =	nn.Conv2d(64, 64, 3, padding=1)
self.conv9_output =	nn.Conv2d(64, 2, 1)

# Convolution 2D

`torch.nn.Conv2d(in_channels, out_channels, kernel_size, stride=1, padding=0, dilation=1, groups=1, bias=True)`

```
layer1 = F.relu(self.conv1_input(x))
layer1 = F.relu(self.conv1(layer1))

layer2 = F.max_pool2d(layer1, 2)
layer2 = F.relu(self.conv2_input(layer2))
layer2 = F.relu(self.conv2(layer2))

layer3 = F.max_pool2d(layer2, 2)
layer3 = F.relu(self.conv3_input(layer3))
layer3 = F.relu(self.conv3(layer3))

layer4 = F.max_pool2d(layer3, 2)
layer4 = F.relu(self.conv4_input(layer4))
layer4 = F.relu(self.conv4(layer4))

layer5 = F.max_pool2d(layer4, 2)
layer5 = F.relu(self.conv5_input(layer5))
layer5 = F.relu(self.conv5(layer5))
```

- **in\_channels** (*int*) – Number of channels in the input image
- **out\_channels** (*int*) – Number of channels produced by the convolution
- **kernel\_size** (*int or tuple*) – Size of the convolving kernel
- **stride** (*int or tuple, optional*) – Stride of the convolution. Default: 1
- **padding** (*int or tuple, optional*) – Zero-padding added to both sides of the input.  
Default: 0
- **dilation** (*int or tuple, optional*) – Spacing between kernel elements. Default: 1
- **groups** (*int, optional*) – Number of blocked connections from input channels to output channels. Default: 1
- **bias** (*bool, optional*) – If `True`, adds a learnable bias to the output. Default: `True`

<https://pytorch.org/docs/stable/nn.html>

# Convolution

`torch.nn.Conv2d(in_channels, out_channels, kernel_size, stride=1, padding=0, dilation=1, groups=1, bias=True)`

1x1	1x0	1x1	0	0
0x0	1x1	1x0	1	0
0x1	0x0	1x1	1	1
0	0	1	1	0
0	1	1	0	0

Input x Filter

4		

Feature Map

# Convolution

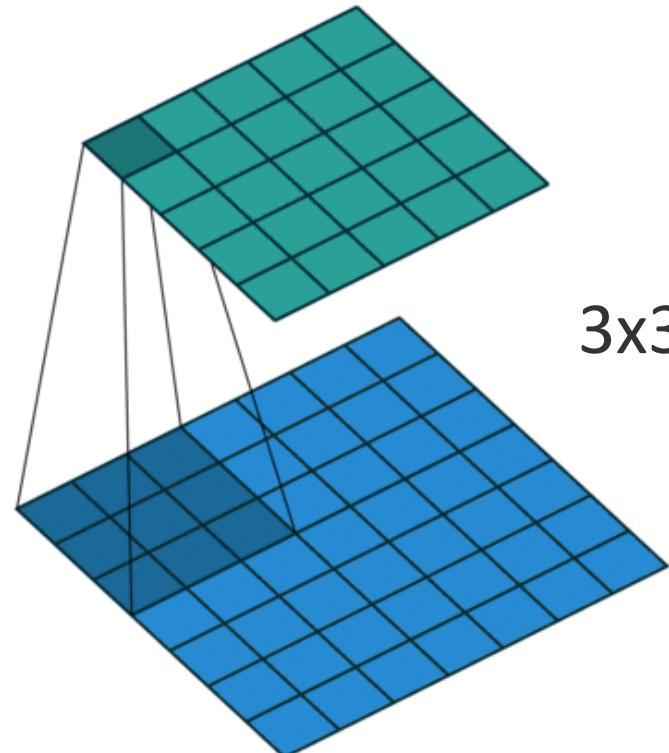
`torch.nn.Conv2d(in_channels, out_channels, kernel_size, stride=1, padding=0, dilation=1, groups=1, bias=True)`

1x1	1x0	1x1	0	0
0x0	1x1	1x0	1	0
0x1	0x0	1x1	1	1
0	0	1	1	0
0	1	1	0	0

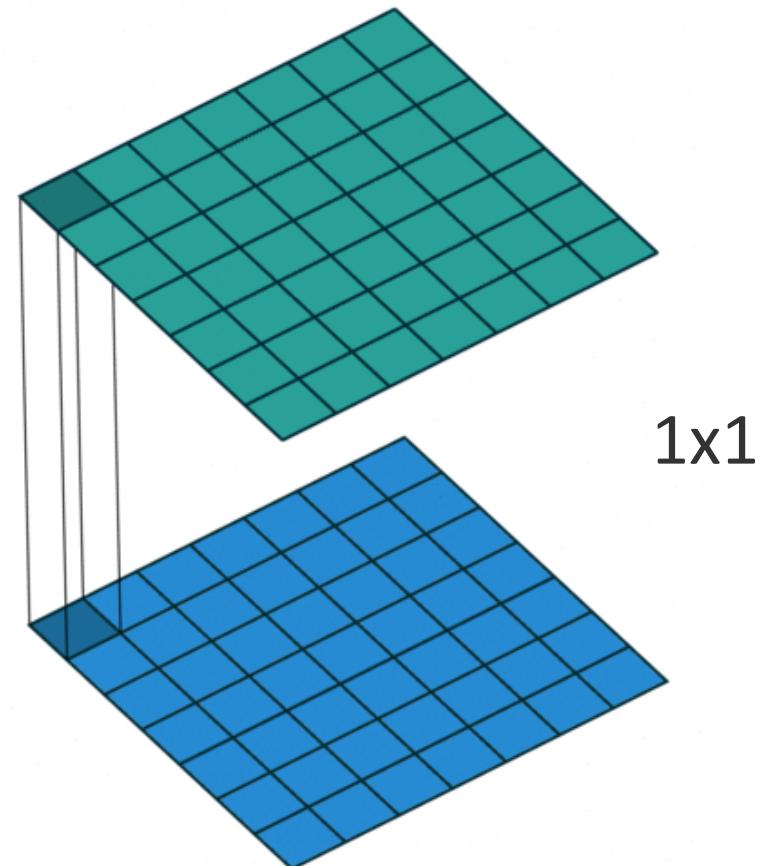
4		

# Kernel Size

`torch.nn.Conv2d(in_channels, out_channels, kernel_size, stride=1, padding=0, dilation=1, groups=1, bias=True)`



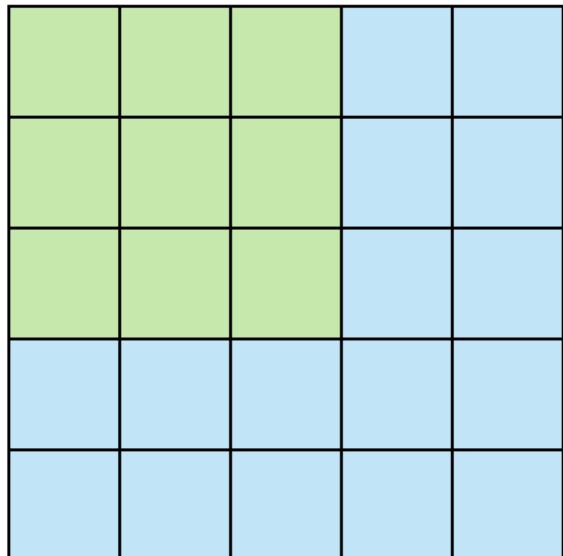
3x3



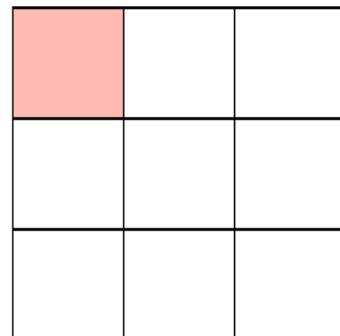
1x1

# Stride

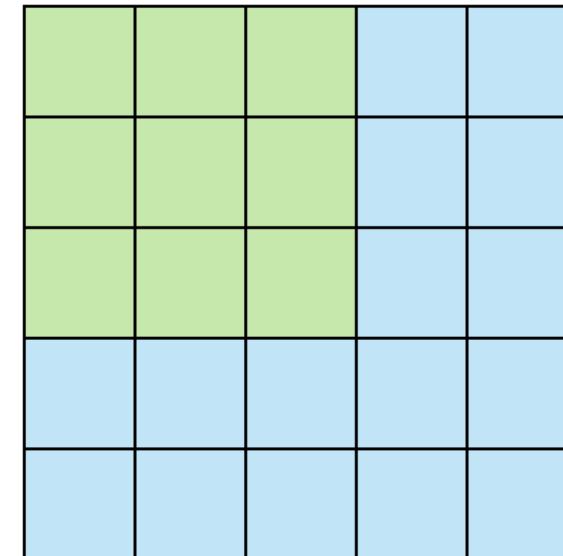
`torch.nn.Conv2d(in_channels, out_channels, kernel_size, stride=1, padding=0, dilation=1, groups=1, bias=True)`



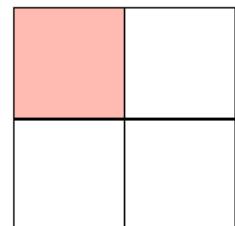
Stride 1



Feature Map



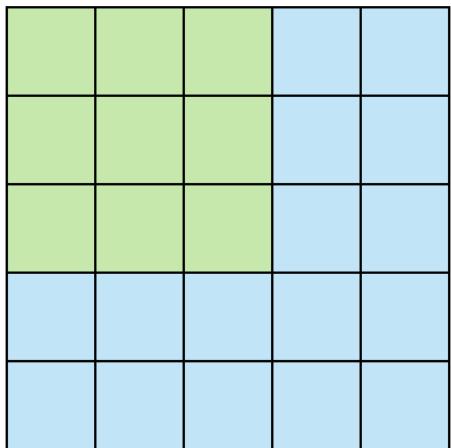
Stride 2



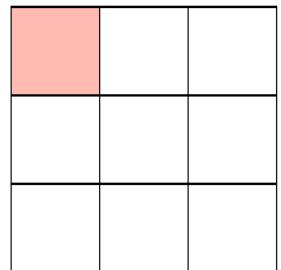
Feature Map

# Padding

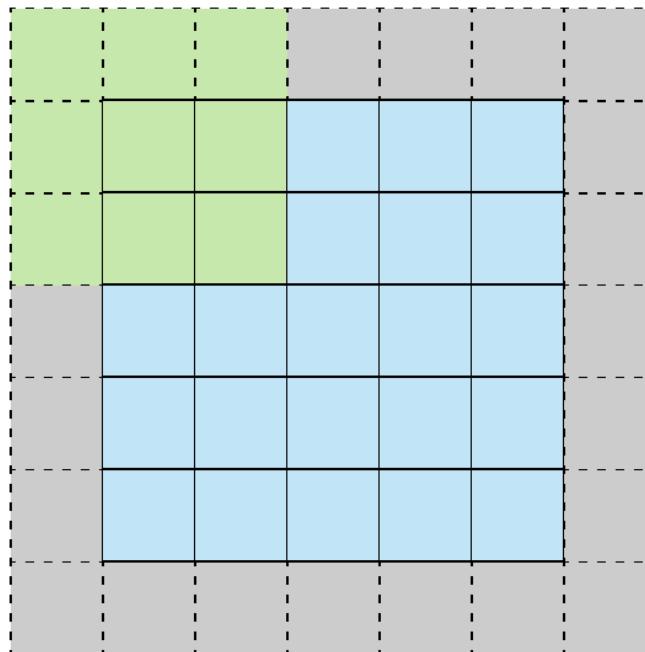
`torch.nn.Conv2d(in_channels, out_channels, kernel_size, stride=1, padding=0, dilation=1, groups=1, bias=True)`



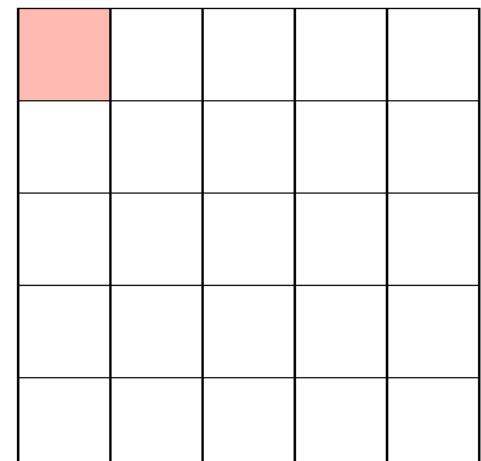
Stride 1



Feature Map



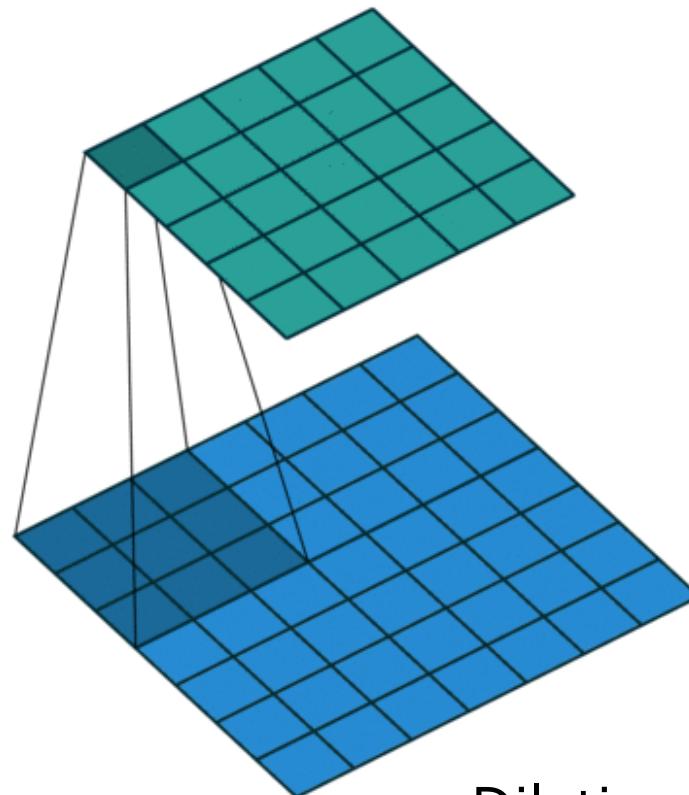
Stride 1 with Padding



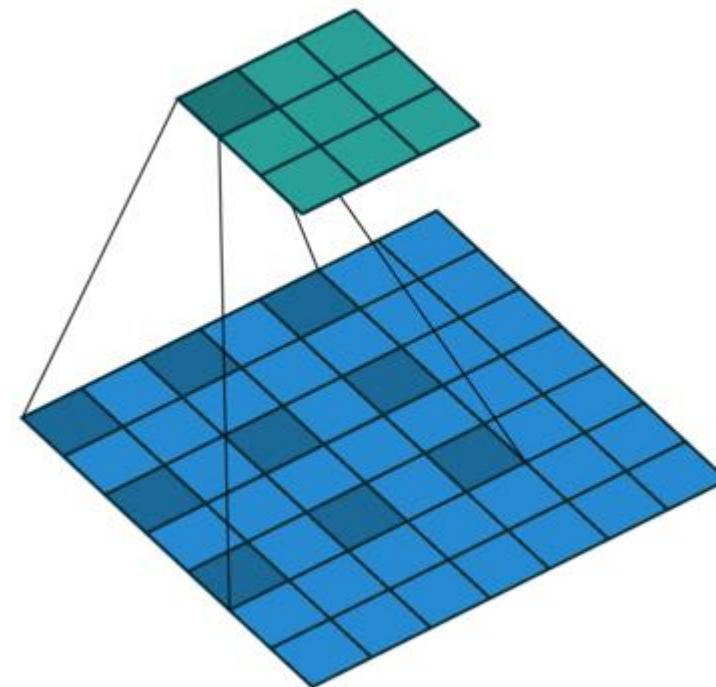
Feature Map

# Dilation (typically default)

`torch.nn.Conv2d(in_channels, out_channels, kernel_size, stride=1, padding=0, dilation=1, groups=1, bias=True)`



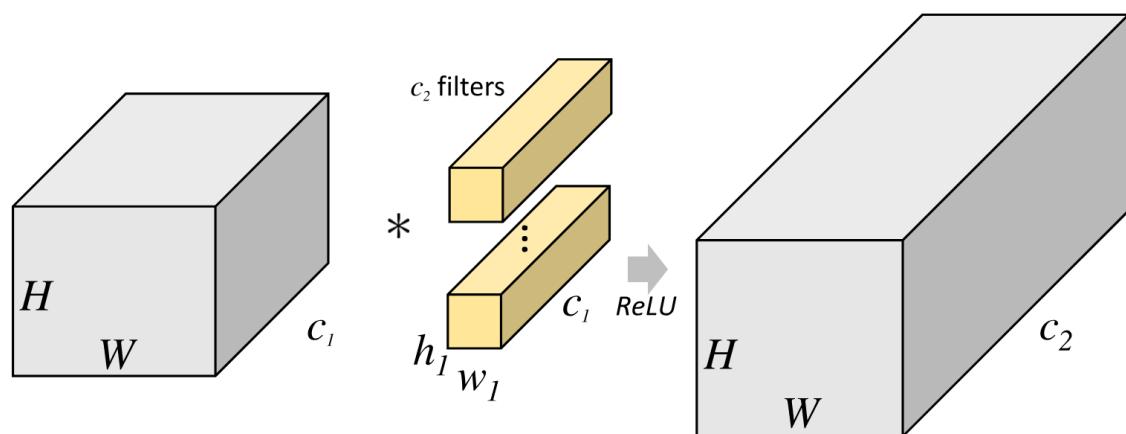
Dilation = 1



Dilation = 2

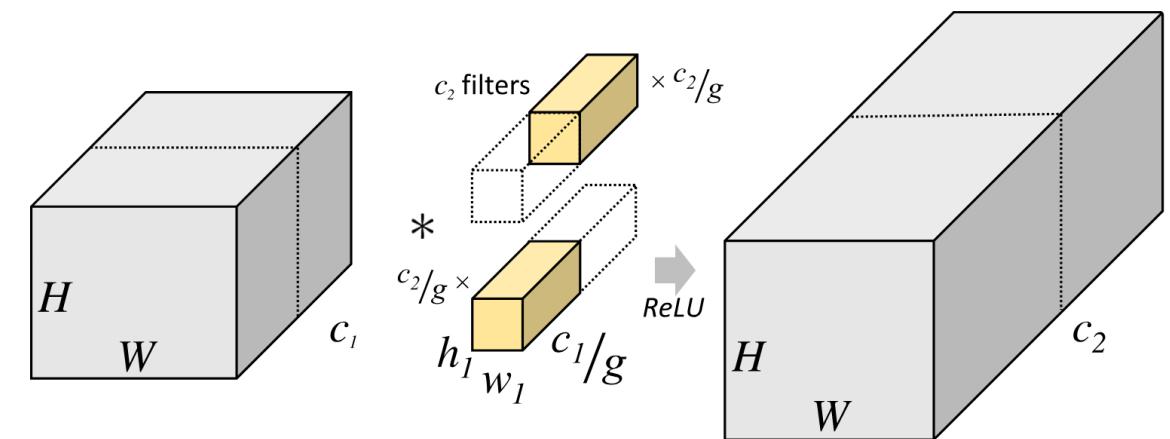
# Groups (typically default)

`torch.nn.Conv2d(in_channels, out_channels, kernel_size, stride=1, padding=0, dilation=1, groups=1, bias=True)`



A normal convolutional layer. Yellow blocks represent learned parameters, gray blocks represent feature maps/input images (working memory).

Groups = 1



A convolutional layer with 2 filter groups. Note that each of the filters in the grouped convolutional layer is now exactly half the depth, i.e. half the parameters and half the compute as the original filter.

Groups = 2

# Input Channels

`torch.nn.Conv2d(in_channels, out_channels, kernel_size, stride=1, padding=0, dilation=1, groups=1, bias=True)`

0	0	0	0	0	0	0	...
0	156	155	156	158	158	158	...
0	153	154	157	159	159	159	...
0	149	151	155	158	159	159	...
0	146	146	149	153	158	158	...
0	145	143	143	148	158	158	...
...	...	...	...	...	...	...	...

Input Channel #1 (Red)

0	0	0	0	0	0	0	...
0	167	166	167	169	169	169	...
0	164	165	168	170	170	170	...
0	160	162	166	169	170	170	...
0	156	156	159	163	168	168	...
0	155	153	153	158	168	168	...
...	...	...	...	...	...	...	...

Input Channel #2 (Green)

0	0	0	0	0	0	0	...
0	163	162	163	165	165	165	...
0	160	161	164	166	166	166	...
0	156	158	162	165	166	166	...
0	155	155	158	162	167	167	...
0	154	152	152	157	167	167	...
...	...	...	...	...	...	...	...

Input Channel #3 (Blue)

-1	-1	1
0	1	-1
0	1	1

Kernel Channel #1



308

1	0	0
1	-1	-1
1	0	-1

Kernel Channel #2



-498

0	1	1
0	1	0
1	-1	1

Kernel Channel #3



164

-25				...
				...
				...
				...
...	...	...	...	...

Bias = 1

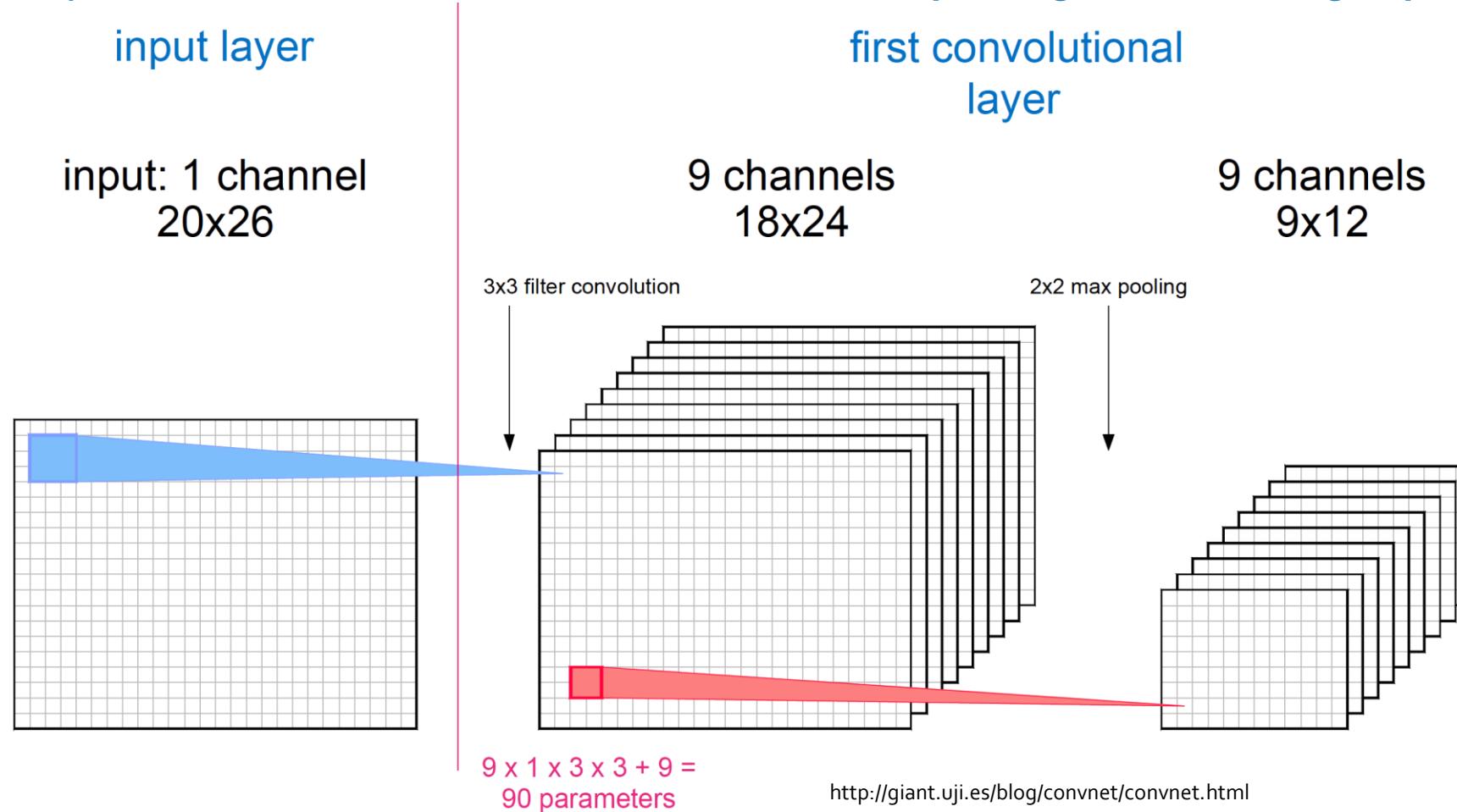


+ 1 = -25

[http://machinelearningguru.com/computer\\_vision/basics/convolution/convolution\\_layer.html](http://machinelearningguru.com/computer_vision/basics/convolution/convolution_layer.html)

# Output Channels

`torch.nn.Conv2d(in_channels, out_channels, kernel_size, stride=1, padding=0, dilation=1, groups=1, bias=True)`

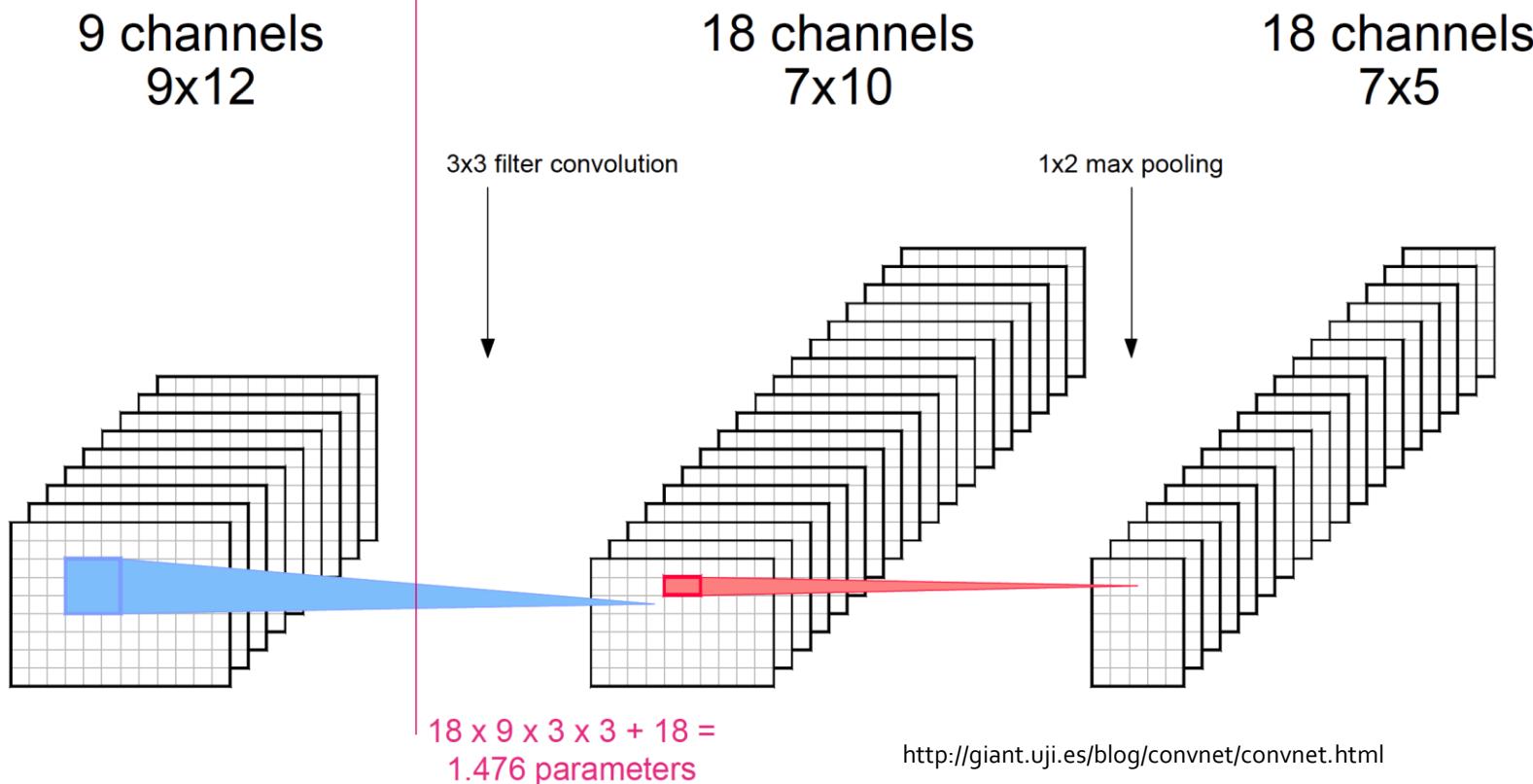


<http://giant.uji.es/blog/convnet/convnet.html>

# Multiple Input and Output Channels

`torch.nn.Conv2d(in_channels, out_channels, kernel_size, stride=1, padding=0, dilation=1, groups=1, bias=True)`

second convolutional  
layer



# Bias

`torch.nn.Conv2d(in_channels, out_channels, kernel_size, stride=1, padding=0, dilation=1, groups=1, bias=True)`

Input Volume (+pad 1) (7x7x3)	Filter W0 (3x3x3)
x[:, :, 0]	w0[:, :, 0]
0 0 0 0 0 0 0 0 2 1 1 1 0 0 0 1 1 1 0 2 0 0 1 1 0 1 1 0 0 2 1 0 2 1 0 0 0 2 2 0 2 0 0 0 0 0 0 0 0	-1 -1 0 1 -1 -1 0 0 1 -1 0 0 -1 0 0 1 0 0 1 0 -1
x[:, :, 1]	w0[:, :, 1]
0 0 0 0 0 0 0 0 1 2 2 1 0 0 0 0 2 1 0 0 0 0 0 0 1 2 0 0 0 2 1 1 2 1 0 0 2 2 1 2 2 0 0 0 0 0 0 0 0	0 1 1 0 1 -1 1 0 1 0 1 1 0 -1 -1 0 1 1 1 1 0
x[:, :, 2]	w0[:, :, 2]
0 0 0 0 0 0 0 0 2 0 1 0 1 0 0 2 1 0 1 0 0 0 1 0 1 2 2 0 0 0 1 2 1 2 0 0 1 0 0 0 2 0 0 0 0 0 0 0 0	1 0

Filter W1 (3x3x3)

w1[:, :, 0]	o[:, :, 0]
1 -1 0	-4 3 1
-1 0 0	1 -4 0
-1 1 0	-3 4 -4
w1[:, :, 1]	o[:, :, 1]
1 0 0	3 0 2
-1 0 0	0 4 0
0 1 -1	-2 -5 1
w1[:, :, 2]	

Bias b0 (1x1x1)

b0[:, :, 0]
0

toggle movement

Input Volume (+pad 1) (7x7x3) Filter W0 (3x3x3) Filter W1 (3x3x3) Output Volume (3x3x2)

x[:, :, 0]	w0[:, :, 0]	w1[:, :, 0]	o[:, :, 0]
0 0 0 0 0 0 0	-1 -1 0	1 -1 0	-4 3 1
0 2 1 1 1 0 0	1 -1 -1	-1 0 0	1 -4 0
0 1 1 1 0 2 0	0 0 1	-1 1 0	-3 4 -4
0 1 1 0 1 1 0	w0[:, :, 1]	w1[:, :, 1]	o[:, :, 1]
0 2 1 0 2 1 0	-1 0 0	1 0 0	3 0 2
0 0 2 2 0 2 0	1 0 0	-1 0 0	0 4 0
0 0 0 0 0 0 0	1 0 -1	0 1 -1	-2 -5 1

x[:, :, 1]	w0[:, :, 1]	w1[:, :, 1]	o[:, :, 1]
0 0 0 0 0 0 0	0 1 1	1 0 0	3 0 2
0 1 2 2 1 0 0	0 1 -1	-1 0 0	0 4 0
0 0 2 1 0 0 0	1 0 0	-1 0 0	-2 -5 1
0 0 0 1 2 0 0	1 0 -1	0 1 -1	
0 2 1 1 2 1 0	w0[:, :, 2]	w1[:, :, 2]	o[:, :, 2]
0 2 2 1 2 2 0	0 1 1	0 -1 -1	
0 0 0 0 0 0 0	0 0 2	0 1 1	

x[:, :, 2]	w0[:, :, 2]	w1[:, :, 2]	o[:, :, 2]
0 0 0 0 0 0 0	0 0 0	0 -1 -1	
0 2 0 1 0 1 0	0 1 -1	0 1 1	
0 2 1 0 1 0 0	-1 0 1	1 1 0	
0 0 0 1 2 0 0	0 0 0	0 0 0	
0 2 1 1 2 1 0	Bias b0 (1x1x1)	Bias b0 (1x1x1)	1
0 2 2 1 2 2 0	b0[:, :, 0]	b1[:, :, 0]	0
0 0 0 0 0 0 0	0 0 0	0 0 0	

toggle movement

# ReLU

`torch.nn.ReLU(inplace=False)`

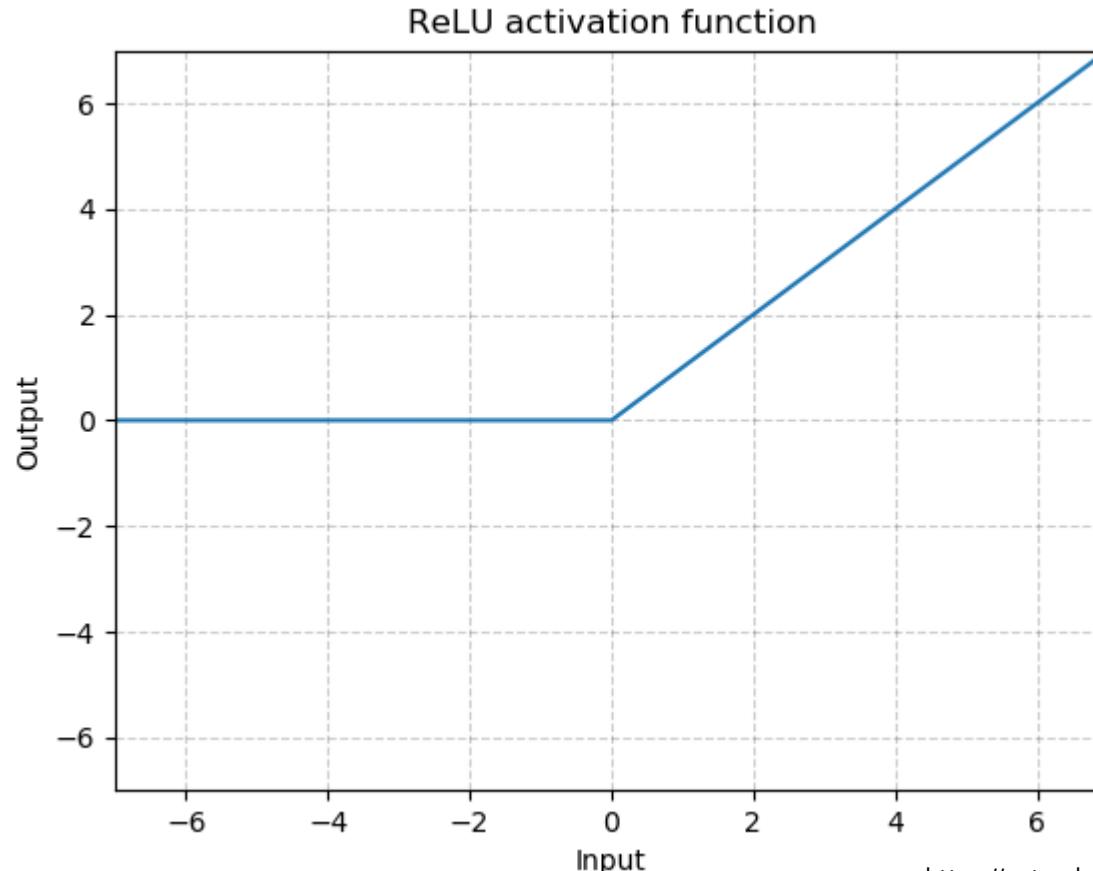
```
layer1 = F.relu(self.conv1_input(x))
layer1 = F.relu(self.conv1(layer1))

layer2 = F.max_pool2d(layer1, 2)
layer2 = F.relu(self.conv2_input(layer2))
layer2 = F.relu(self.conv2(layer2))

layer3 = F.max_pool2d(layer2, 2)
layer3 = F.relu(self.conv3_input(layer3))
layer3 = F.relu(self.conv3(layer3))

layer4 = F.max_pool2d(layer3, 2)
layer4 = F.relu(self.conv4_input(layer4))
layer4 = F.relu(self.conv4(layer4))

layer5 = F.max_pool2d(layer4, 2)
layer5 = F.relu(self.conv5_input(layer5))
layer5 = F.relu(self.conv5(layer5))
```



<https://pytorch.org/docs/stable/nn.html>

# Max Pooling

`torch.nn.MaxPool2d(kernel_size, stride=None, padding=0, dilation=1, return_indices=False, ceil_mode=False)`

```
layer1 = F.relu(self.conv1_input(x))
layer1 = F.relu(self.conv1(layer1))

layer2 = F.max_pool2d(layer1, 2)
layer2 = F.relu(self.conv2_input(layer2))
layer2 = F.relu(self.conv2(layer2))

layer3 = F.max_pool2d(layer2, 2)
layer3 = F.relu(self.conv3_input(layer3))
layer3 = F.relu(self.conv3(layer3))

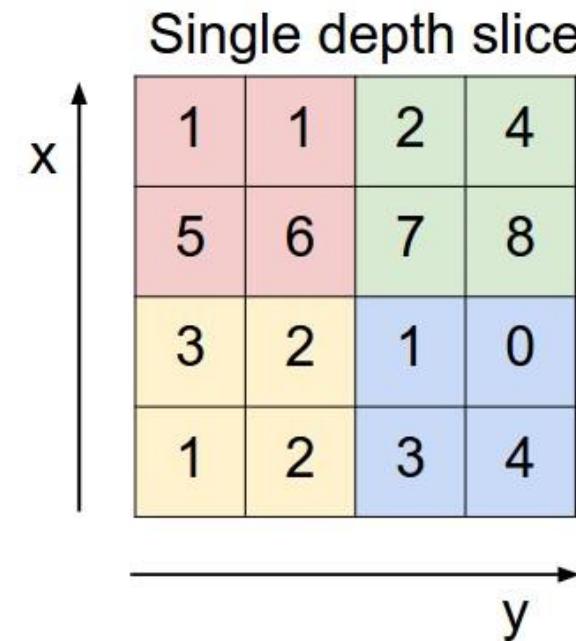
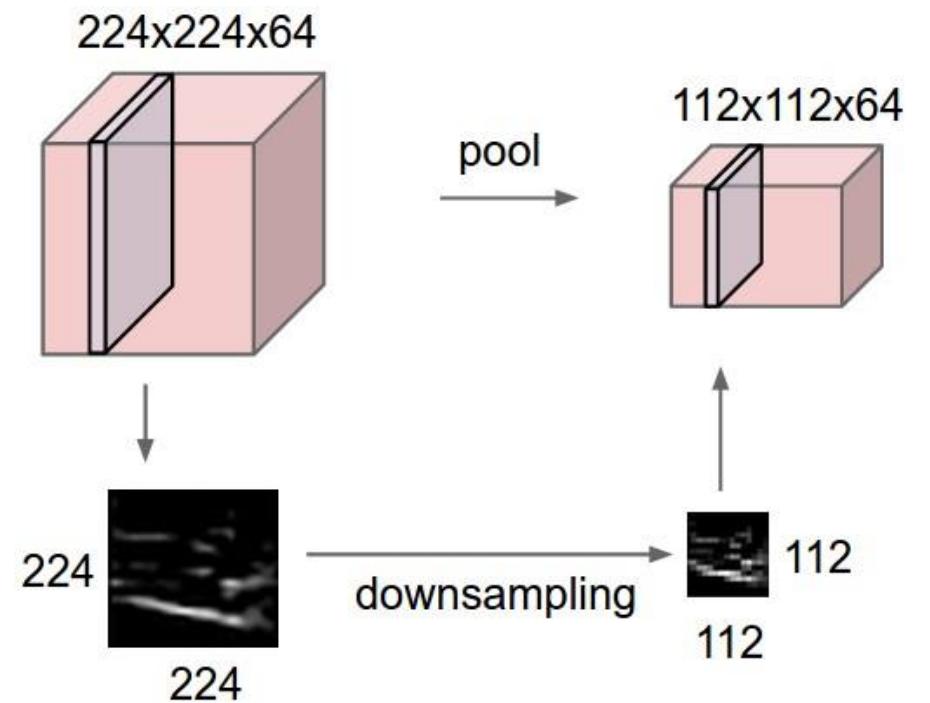
layer4 = F.max_pool2d(layer3, 2)
layer4 = F.relu(self.conv4_input(layer4))
layer4 = F.relu(self.conv4(layer4))

layer5 = F.max_pool2d(layer4, 2)
layer5 = F.relu(self.conv5_input(layer5))
layer5 = F.relu(self.conv5(layer5))
```

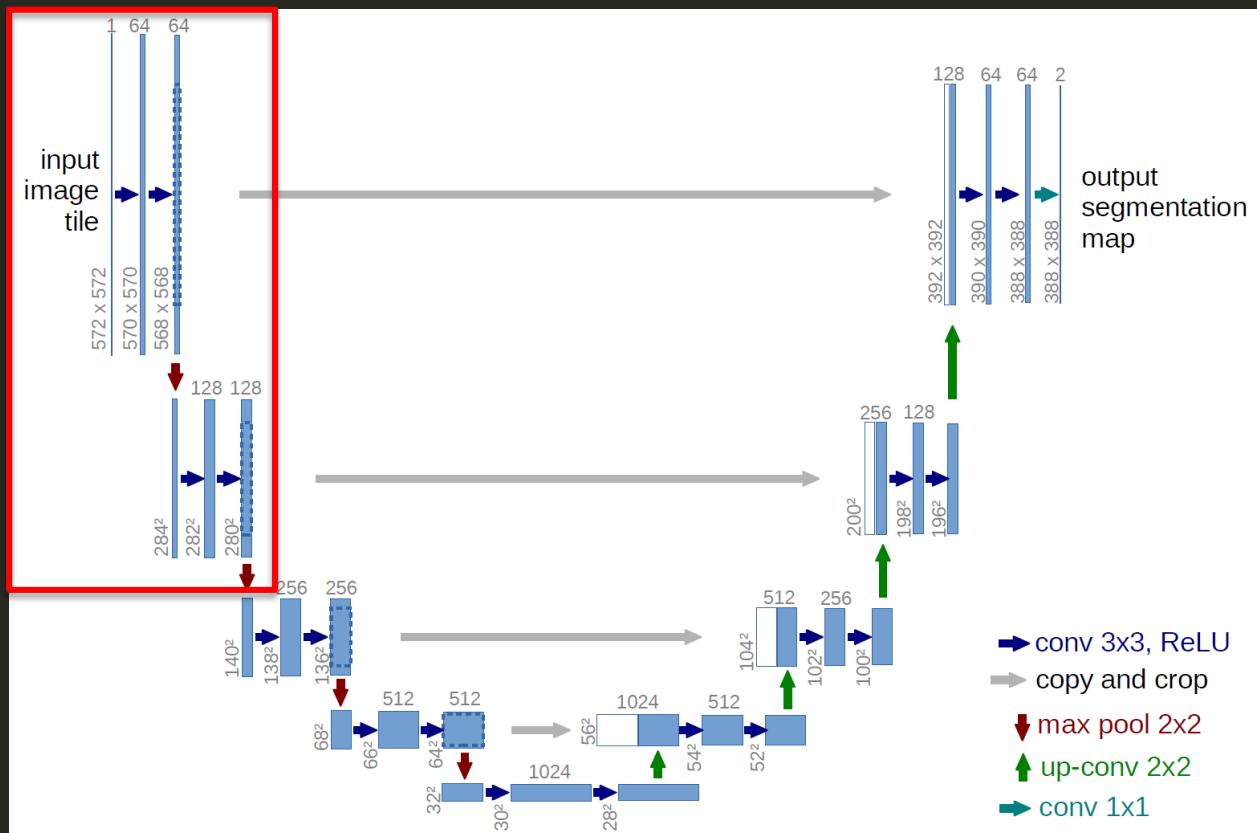
- **kernel\_size** – the size of the window to take a max over
- **stride** – the stride of the window. Default value is `kernel_size`
- **padding** – implicit zero padding to be added on both sides
- **dilation** – a parameter that controls the stride of elements in the window
- **return\_indices** – if `True`, will return the max indices along with the outputs.  
    Useful when Unpooling later
- **ceil\_mode** – when True, will use `ceil` instead of `floor` to compute the output shape

# Max Pooling (kernel size = 2)

`torch.nn.MaxPool2d(kernel_size, stride=None, padding=0, dilation=1, return_indices=False, ceil_mode=False)`



# Encoder



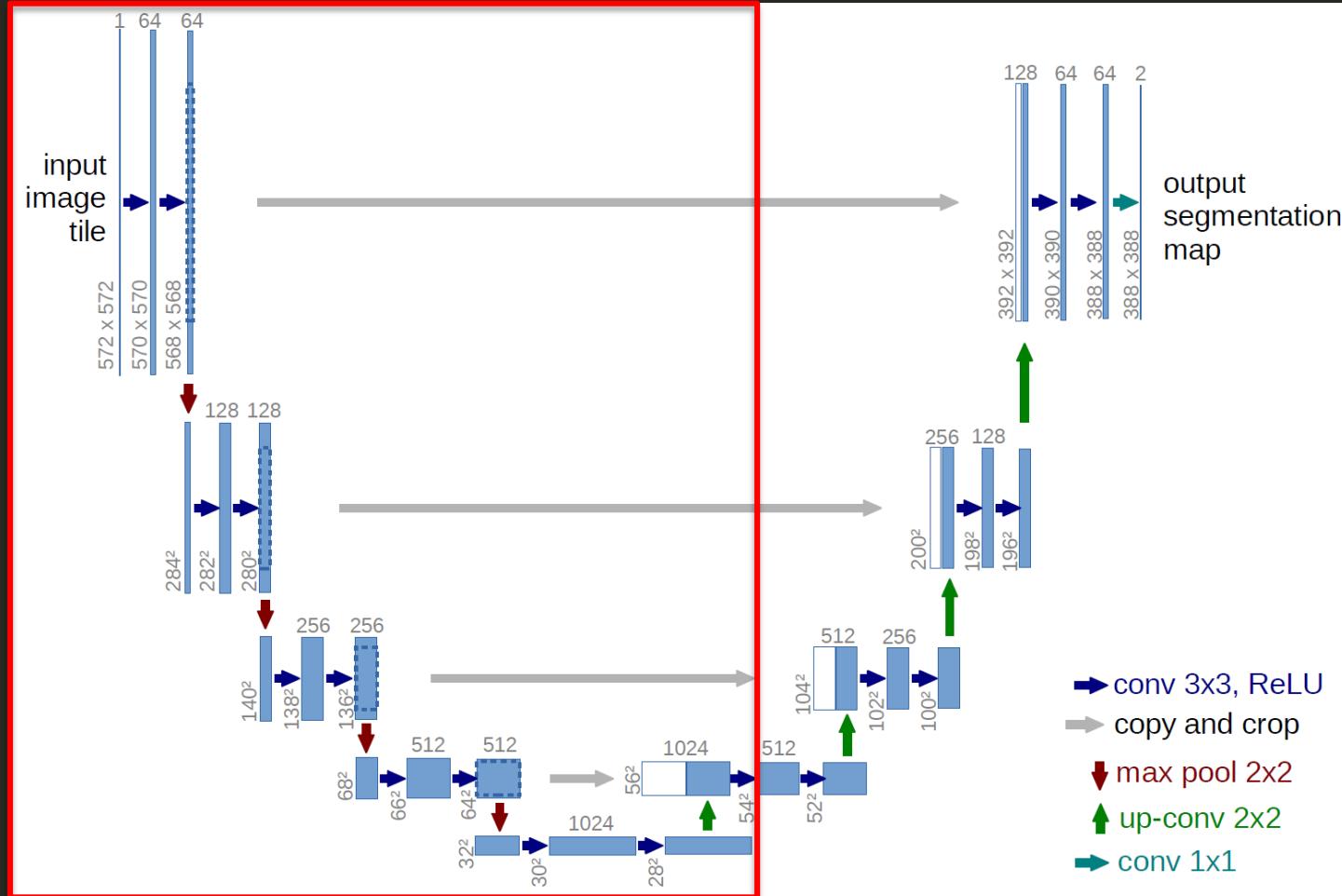
```
layer1 = F.relu(self.conv1_input(x))  
layer1 = F.relu(self.conv1(layer1))
```

```
layer2 = F.max_pool2d(layer1, 2)  
layer2 = F.relu(self.conv2_input(layer2))  
layer2 = F.relu(self.conv2(layer2))
```

```
layer3 = F.max_pool2d(layer2, 2)
```

*self.conv1\_input* = `nn.Conv2d(1, 64, 3, padding=1)`  
*self.conv1* = `nn.Conv2d(64, 64, 3, padding=1)`  
*self.conv2\_input* = `nn.Conv2d(64, 128, 3, padding=1)`  
*self.conv2* = `nn.Conv2d(128, 128, 3, padding=1)`

# Encoder



```
layer1 = F.relu(self.conv1_input(x))  
layer1 = F.relu(self.conv1(layer1))
```

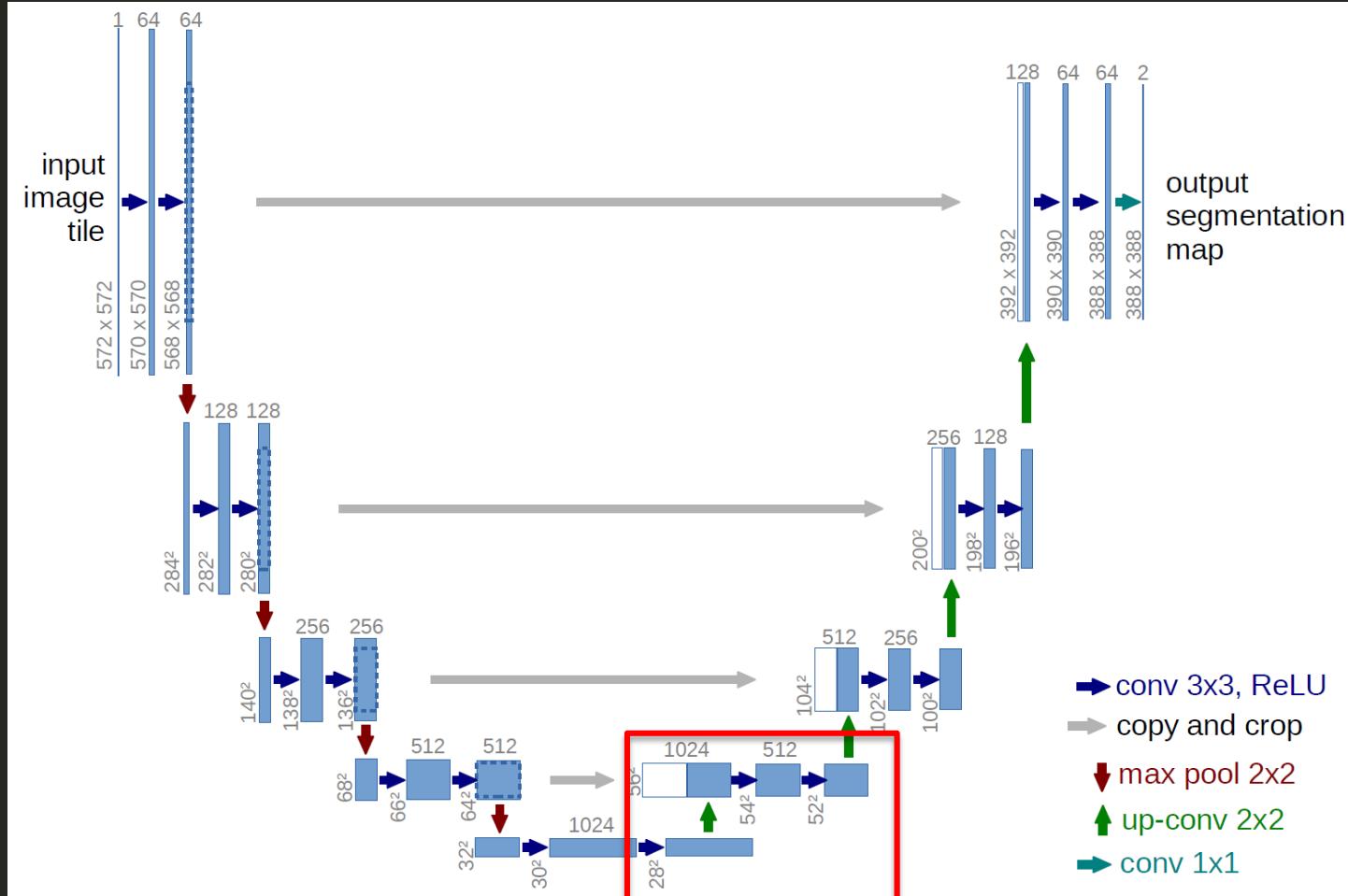
```
layer2 = F.max_pool2d(layer1, 2)  
layer2 = F.relu(self.conv2_input(layer2))  
layer2 = F.relu(self.conv2(layer2))
```

```
layer3 = F.max_pool2d(layer2, 2)  
layer3 = F.relu(self.conv3_input(layer3))  
layer3 = F.relu(self.conv3(layer3))
```

```
layer4 = F.max_pool2d(layer3, 2)  
layer4 = F.relu(self.conv4_input(layer4))  
layer4 = F.relu(self.conv4(layer4))
```

```
layer5 = F.max_pool2d(layer4, 2)  
layer5 = F.relu(self.conv5_input(layer5))  
layer5 = F.relu(self.conv5(layer5))
```

# Decoder



```

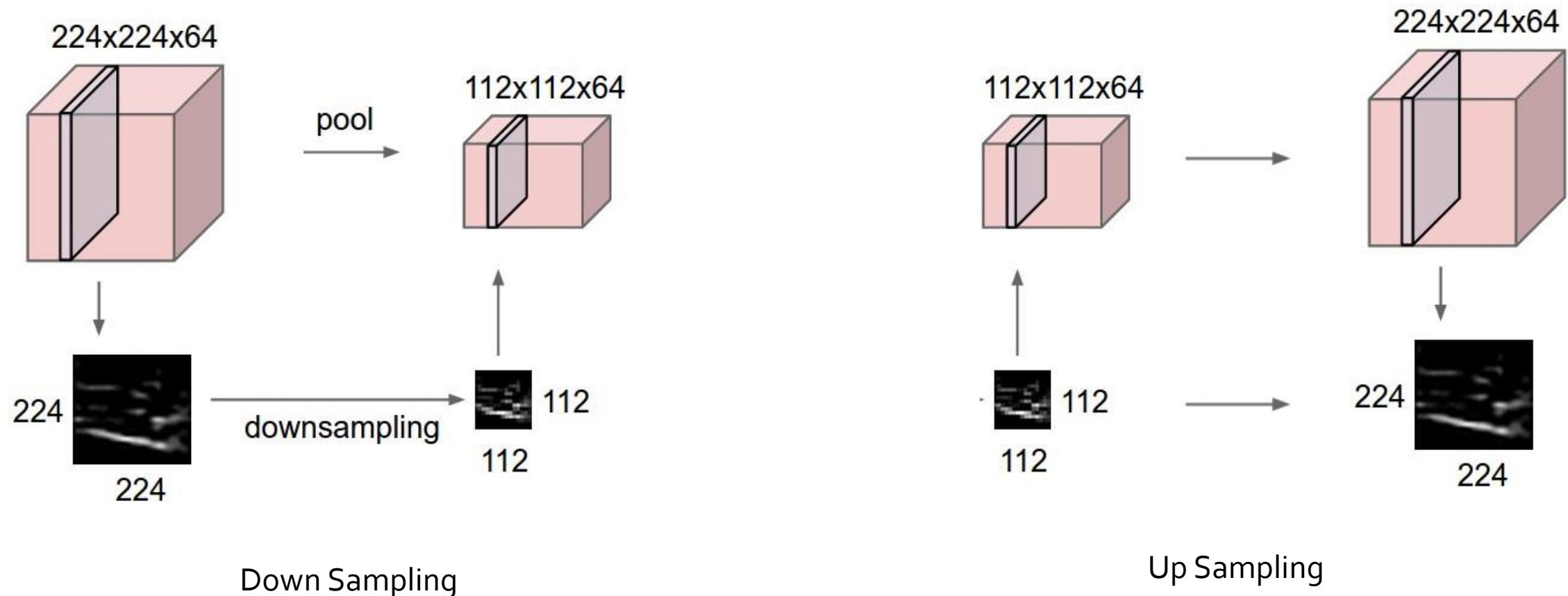
layer6 = F.relu(self.conv6_up(layer5))
layer6 = torch.cat((layer4, layer6), 1)
layer6 = F.relu(self.conv6_input(layer6))
layer6 = F.relu(self.conv6(layer6))

```

<i>self.conv6_up</i> =	<code>nn.ConvTranspose2d(1024, 512, 2, 2)</code>
<i>self.conv6_input</i> =	<code>nn.Conv2d(1024, 512, 3, padding=1)</code>
<i>self.conv6</i> =	<code>nn.Conv2d(512, 512, 3, padding=1)</code>

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

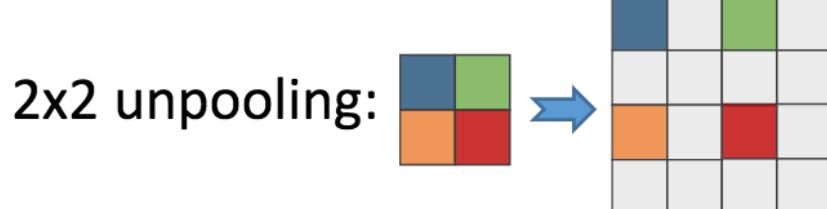
# Upsampling



<https://github.com/cs231n/cs231n.github.io/blob/master/convolutional-networks.md>

# Unpooling

`torch.nn.MaxUnpool2d(kernel_size, stride=None, padding=0)`

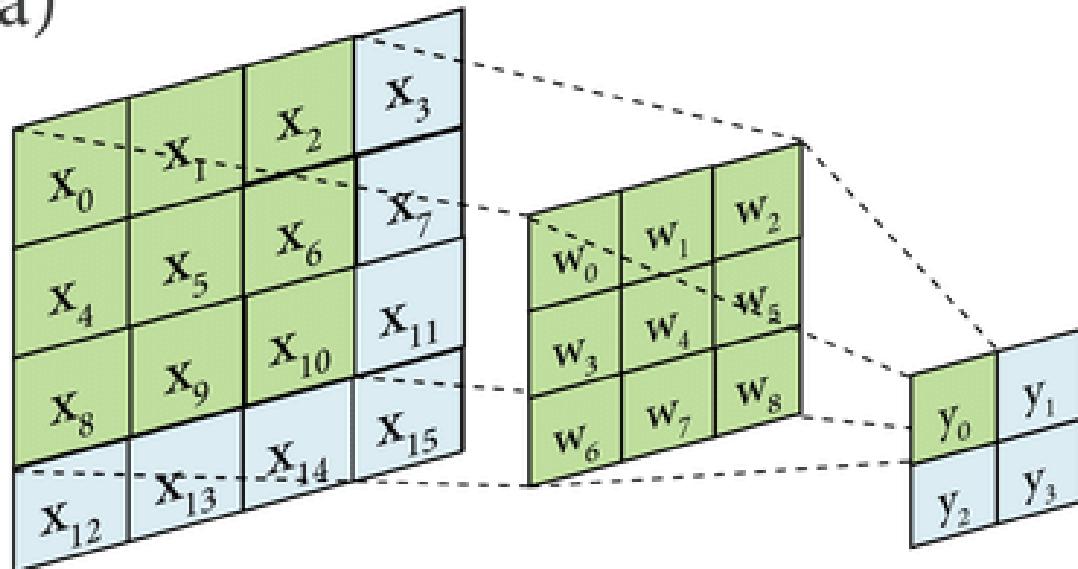


Input	Max-pooling [3, 3]	Inverse																																								
<table border="1"><tr><td>0.8</td><td>0</td><td>3.3</td><td>0</td><td>4</td><td>1</td></tr><tr><td>1.2</td><td>1.2</td><td>0</td><td>2.1</td><td>0</td><td>0</td></tr><tr><td>0</td><td>3</td><td>0.2</td><td>3.1</td><td>5.2</td><td>0</td></tr></table>	0.8	0	3.3	0	4	1	1.2	1.2	0	2.1	0	0	0	3	0.2	3.1	5.2	0	Value: <table border="1"><tr><td>3.3</td><td>5.2</td></tr></table>	3.3	5.2	<table border="1"><tr><td>3.3</td><td>0</td><td>0</td><td>5.2</td><td>0</td><td>0</td></tr><tr><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td></tr><tr><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td></tr></table>	3.3	0	0	5.2	0	0	0	0	0	0	0	0	0	0	0	0	0	0		
0.8	0	3.3	0	4	1																																					
1.2	1.2	0	2.1	0	0																																					
0	3	0.2	3.1	5.2	0																																					
3.3	5.2																																									
3.3	0	0	5.2	0	0																																					
0	0	0	0	0	0																																					
0	0	0	0	0	0																																					
Input	Max-pooling [3, 3]	Unpooling																																								
<table border="1"><tr><td>0.8</td><td>0</td><td>3.3</td><td>0</td><td>4</td><td>1</td></tr><tr><td>1.2</td><td>1.2</td><td>0</td><td>2.1</td><td>0</td><td>0</td></tr><tr><td>0</td><td>3</td><td>0.2</td><td>3.1</td><td>5.2</td><td>0</td></tr></table>	0.8	0	3.3	0	4	1	1.2	1.2	0	2.1	0	0	0	3	0.2	3.1	5.2	0	Value: <table border="1"><tr><td>3.3</td><td>5.2</td></tr></table> Position: <table border="1"><tr><td>(0,2)</td><td>(2,1)</td></tr></table>	3.3	5.2	(0,2)	(2,1)	<table border="1"><tr><td>0</td><td>0</td><td>3.3</td><td>0</td><td>0</td><td>0</td></tr><tr><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td></tr><tr><td>0</td><td>0</td><td>0</td><td>0</td><td>5.2</td><td>0</td></tr></table>	0	0	3.3	0	0	0	0	0	0	0	0	0	0	0	0	0	5.2	0
0.8	0	3.3	0	4	1																																					
1.2	1.2	0	2.1	0	0																																					
0	3	0.2	3.1	5.2	0																																					
3.3	5.2																																									
(0,2)	(2,1)																																									
0	0	3.3	0	0	0																																					
0	0	0	0	0	0																																					
0	0	0	0	5.2	0																																					

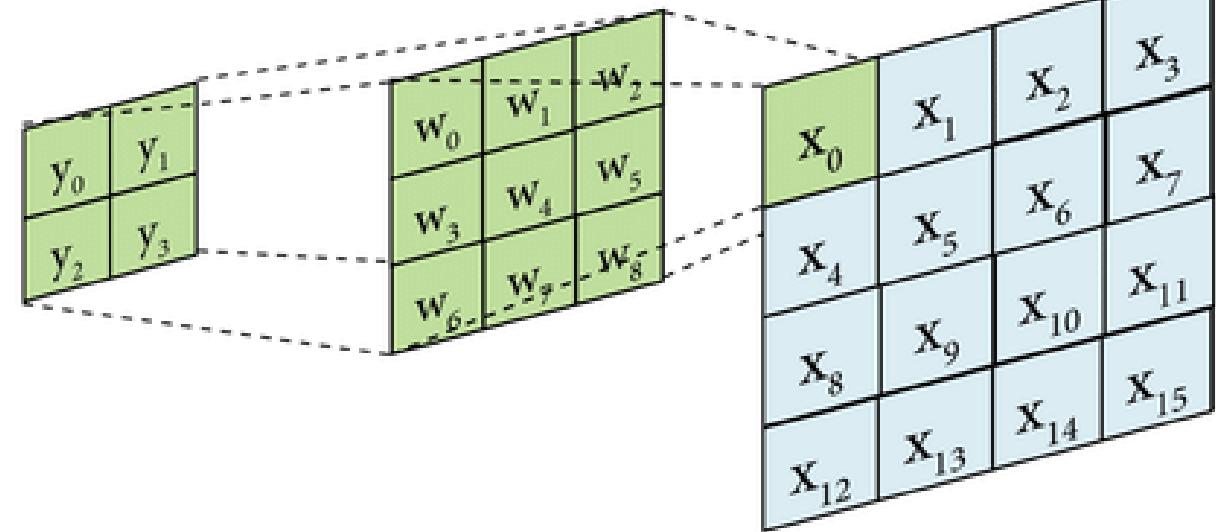
<https://stackoverflow.com/questions/36548736/tensorflow-unpooling>

# Decoder

a)



b)



# Deconvolution (Transpose Convolution)

`torch.nn.ConvTranspose2d(in_channels, out_channels, kernel_size, stride=1, padding=0, output_padding=0, groups=1, bias=True, dilation=1)`

```
layer6 = F.relu(self.conv6_up(layer5))
layer6 = torch.cat((layer4, layer6), 1)
layer6 = F.relu(self.conv6_input(layer6))
layer6 = F.relu(self.conv6(layer6))

layer7 = F.relu(self.conv7_up(layer6))
layer7 = torch.cat((layer3, layer7), 1)
layer7 = F.relu(self.conv7_input(layer7))
layer7 = F.relu(self.conv7(layer7))

layer8 = F.relu(self.conv8_up(layer7))
layer8 = torch.cat((layer2, layer8), 1)
layer8 = F.relu(self.conv8_input(layer8))
layer8 = F.relu(self.conv8(layer8))

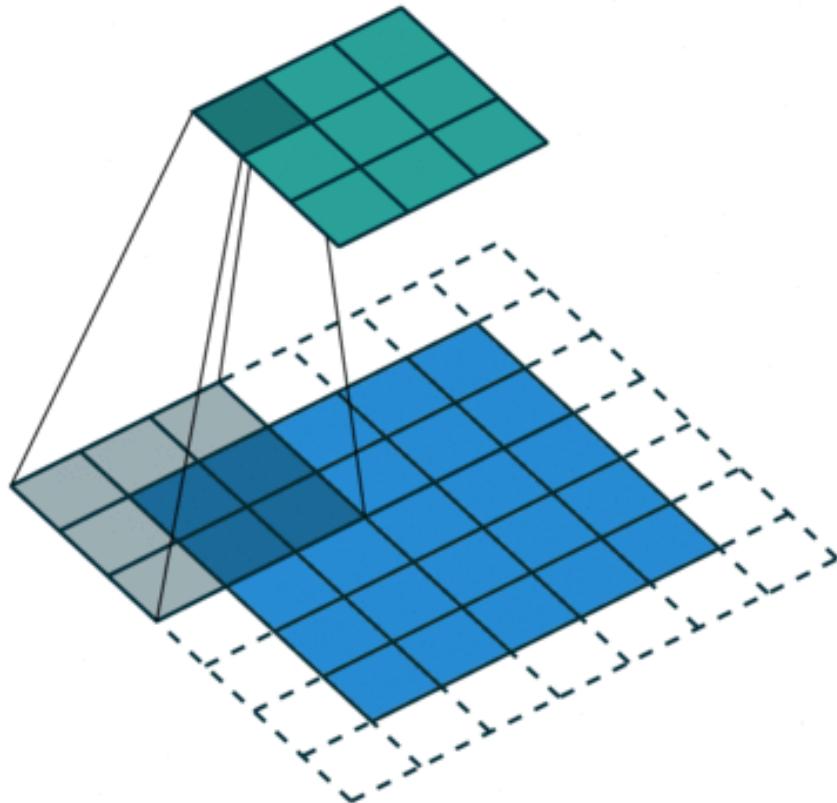
layer9 = F.relu(self.conv9_up(layer8))
layer9 = torch.cat((layer1, layer9), 1)
layer9 = F.relu(self.conv9_input(layer9))
layer9 = F.relu(self.conv9(layer9))
layer9 = self.final(self.conv9_output(layer9))

return layer9
```

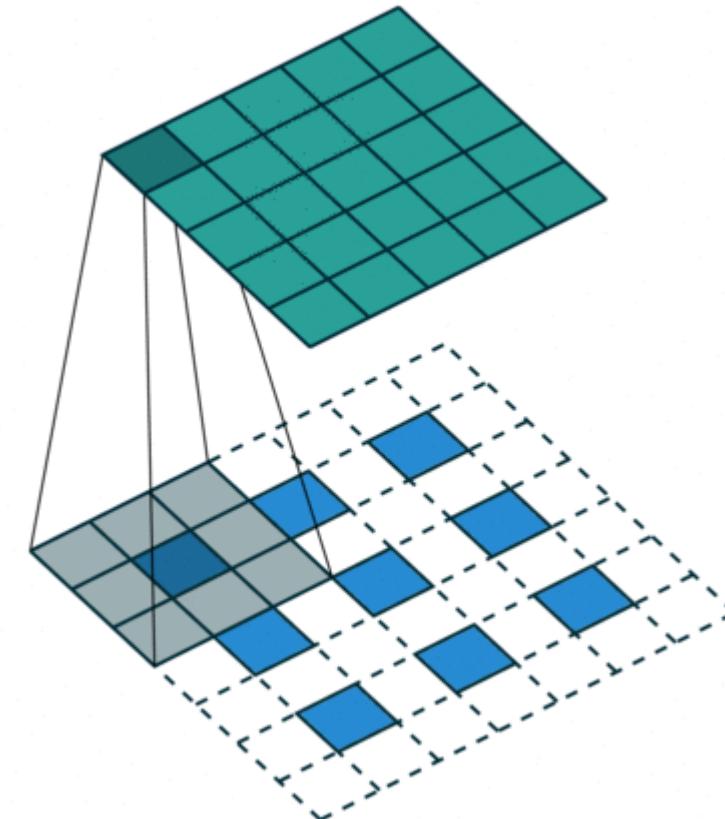
- ***in\_channels* (*int*)** – Number of channels in the input image
- ***out\_channels* (*int*)** – Number of channels produced by the convolution
- ***kernel\_size* (*int* or *tuple*)** – Size of the convolving kernel
- ***stride* (*int* or *tuple*, optional)** – Stride of the convolution. Default: 1
- ***padding* (*int* or *tuple*, optional)** –  $\text{kernel\_size} - 1 - \text{padding}$  zero-padding will be added to both sides of each dimension in the input. Default: 0
- ***output\_padding* (*int* or *tuple*, optional)** – Additional size added to one side of each dimension in the output shape. Default: 0
- ***groups* (*int*, optional)** – Number of blocked connections from input channels to output channels. Default: 1
- ***bias* (*bool*, optional)** – If `True`, adds a learnable bias to the output. Default: `True`
- ***dilation* (*int* or *tuple*, optional)** – Spacing between kernel elements. Default: 1

<https://pytorch.org/docs/stable/nn.html>

# Transpose Convolution



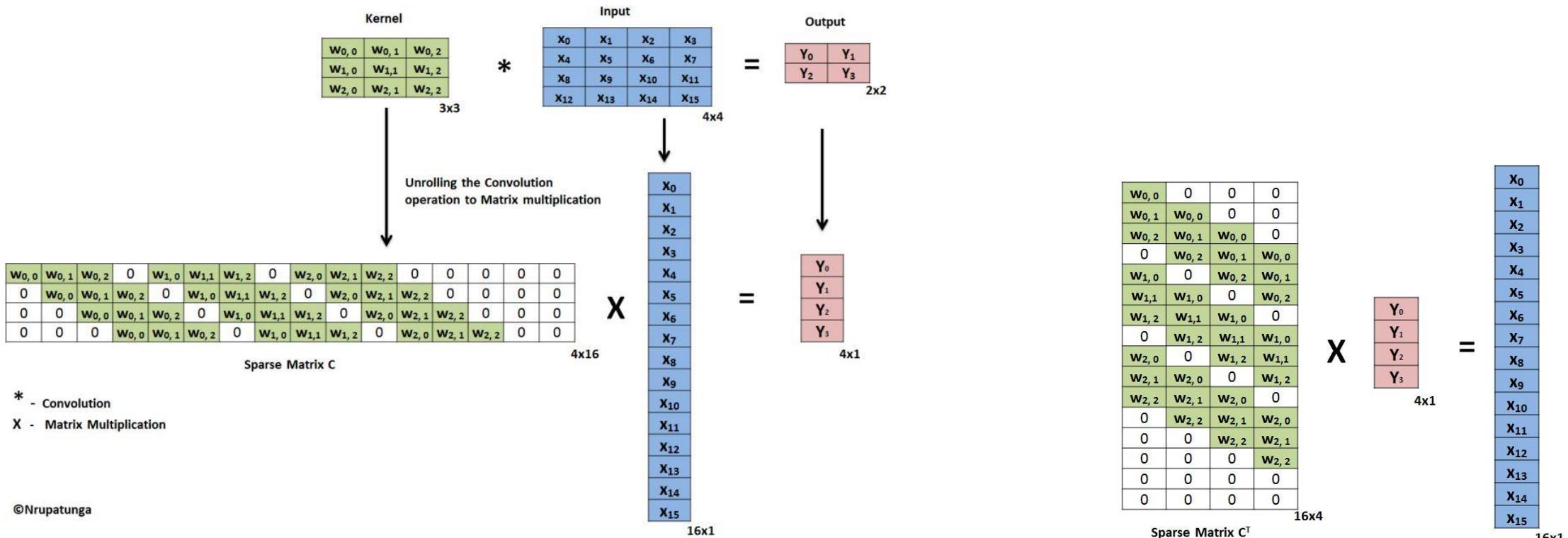
Convolution



Transpose Convolution

[https://github.com/vdumoulin/conv\\_arithmetic/blob/master/README.md](https://github.com/vdumoulin/conv_arithmetic/blob/master/README.md)

# Transpose Convolution



Convolution

<https://nrupatunga.github.io/convolution-2/>

Transpose Convolution

©Nrupatunga

# Concatenation

`torch.cat(seq, dim=0, out=None)`

```
layer6 = F.relu(self.conv6_up(layer5))
layer6 = torch.cat((layer4, layer6), 1)
layer6 = F.relu(self.conv6_input(layer6))
layer6 = F.relu(self.conv6(layer6))

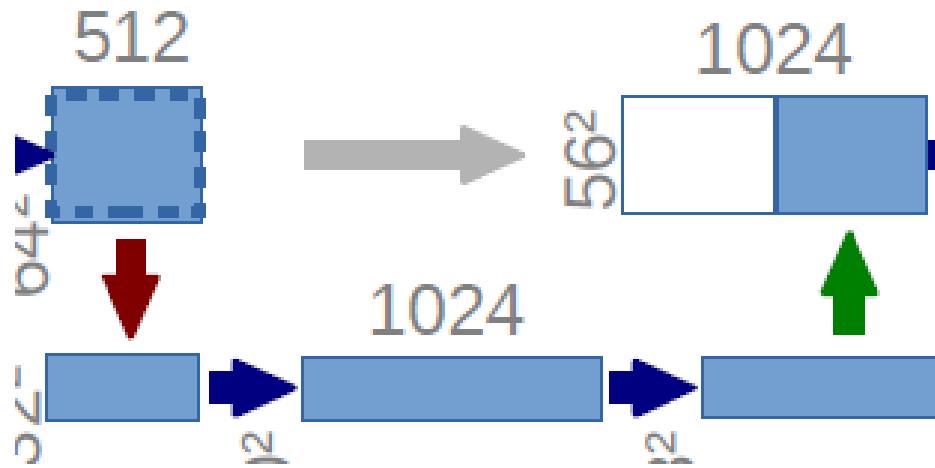
layer7 = F.relu(self.conv7_up(layer6))
layer7 = torch.cat((layer3, layer7), 1)
layer7 = F.relu(self.conv7_input(layer7))
layer7 = F.relu(self.conv7(layer7))

layer8 = F.relu(self.conv8_up(layer7))
layer8 = torch.cat((layer2, layer8), 1)
layer8 = F.relu(self.conv8_input(layer8))
layer8 = F.relu(self.conv8(layer8))

layer9 = F.relu(self.conv9_up(layer8))
layer9 = torch.cat((layer1, layer9), 1)
layer9 = F.relu(self.conv9_input(layer9))
layer9 = F.relu(self.conv9(layer9))
layer9 = self.final(self.conv9_output(layer9))

return layer9
```

- **seq** (*sequence of Tensors*) – any python sequence of tensors of the same type. Non-empty tensors provided must have the same shape, except in the cat dimension.
- **dim** (*int, optional*) – the dimension over which the tensors are concatenated
- **out** (*Tensor, optional*) – the output tensor



<https://pytorch.org/docs/stable/nn.html>

# All Code for U-Net

## Encoder

```
layer1 = F.relu(self.conv1_input(x))
layer1 = F.relu(self.conv1(layer1))

layer2 = F.max_pool2d(layer1, 2)
layer2 = F.relu(self.conv2_input(layer2))
layer2 = F.relu(self.conv2(layer2))

layer3 = F.max_pool2d(layer2, 2)
layer3 = F.relu(self.conv3_input(layer3))
layer3 = F.relu(self.conv3(layer3))

layer4 = F.max_pool2d(layer3, 2)
layer4 = F.relu(self.conv4_input(layer4))
layer4 = F.relu(self.conv4(layer4))

layer5 = F.max_pool2d(layer4, 2)
layer5 = F.relu(self.conv5_input(layer5))
layer5 = F.relu(self.conv5(layer5))
```

## Decoder

```
layer6 = F.relu(self.conv6_up(layer5))
layer6 = torch.cat((layer4, layer6), 1)
layer6 = F.relu(self.conv6_input(layer6))
layer6 = F.relu(self.conv6(layer6))

layer7 = F.relu(self.conv7_up(layer6))
layer7 = torch.cat((layer3, layer7), 1)
layer7 = F.relu(self.conv7_input(layer7))
layer7 = F.relu(self.conv7(layer7))

layer8 = F.relu(self.conv8_up(layer7))
layer8 = torch.cat((layer2, layer8), 1)
layer8 = F.relu(self.conv8_input(layer8))
layer8 = F.relu(self.conv8(layer8))

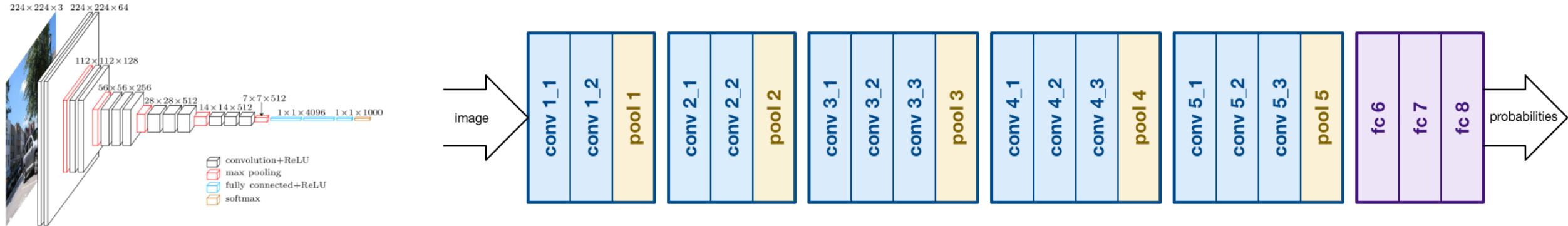
layer9 = F.relu(self.conv9_up(layer8))
layer9 = torch.cat((layer1, layer9), 1)
layer9 = F.relu(self.conv9_input(layer9))
layer9 = F.relu(self.conv9(layer9))
layer9 = self.final(self.conv9_output(layer9))

return layer9
```

## Parameters

self.conv1_input =	nn.Conv2d(1, 64, 3, padding=1)
self.conv1 =	nn.Conv2d(64, 64, 3, padding=1)
self.conv2_input =	nn.Conv2d(64, 128, 3, padding=1)
self.conv2 =	nn.Conv2d(128, 128, 3, padding=1)
self.conv3_input =	nn.Conv2d(128, 256, 3, padding=1)
self.conv3 =	nn.Conv2d(256, 256, 3, padding=1)
self.conv4_input =	nn.Conv2d(256, 512, 3, padding=1)
self.conv4 =	nn.Conv2d(512, 512, 3, padding=1)
self.conv5_input =	nn.Conv2d(512, 1024, 3, padding=1)
self.conv5 =	nn.Conv2d(1024, 1024, 3, padding=1)
self.conv6_up =	nn.ConvTranspose2d(1024, 512, 2, 2)
self.conv6_input =	nn.Conv2d(1024, 512, 3, padding=1)
self.conv6 =	nn.Conv2d(512, 512, 3, padding=1)
self.conv7_up =	nn.ConvTranspose2d(512, 256, 2, 2)
self.conv7_input =	nn.Conv2d(512, 256, 3, padding=1)
self.conv7 =	nn.Conv2d(256, 256, 3, padding=1)
self.conv8_up =	nn.ConvTranspose2d(256, 128, 2, 2)
self.conv8_input =	nn.Conv2d(256, 128, 3, padding=1)
self.conv8 =	nn.Conv2d(128, 128, 3, padding=1)
self.conv9_up =	nn.ConvTranspose2d(128, 64, 2, 2)
self.conv9_input =	nn.Conv2d(128, 64, 3, padding=1)
self.conv9 =	nn.Conv2d(64, 64, 3, padding=1)
self.conv9_output =	nn.Conv2d(64, 2, 1)

# Vgg 16 Classification



```
def make_layers(cfg, batch_norm=False):
    layers = []
    in_channels = 3
    for v in cfg:
        if v == 'M':
            layers += [nn.MaxPool2d(kernel_size=2, stride=2)]
        else:

            conv2d = nn.Conv2d(in_channels, v, 3, padding=1)
            if batch_norm:
                layers += [conv2d, nn.BatchNorm2d(v), nn.ReLU(inplace=True)]
            else:
                layers += [conv2d, nn.ReLU(inplace=True)]
            in_channels = v
    return nn.Sequential(*layers)
```

```
self.classifier = nn.Sequential(
    #fc6
    nn.Linear(512*7*7, 4096),
    nn.ReLU(),
    nn.Dropout(),

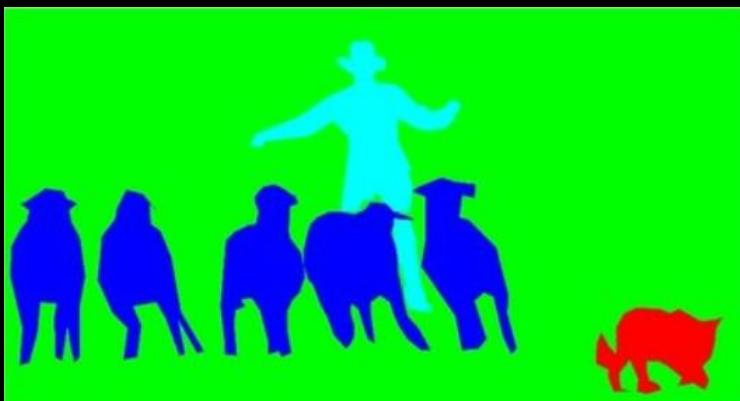
    #fc7
    nn.Linear(4096, 4096),
    nn.ReLU(),
    nn.Dropout(),

    #fc8
    nn.Linear(4096, num_classes))
```

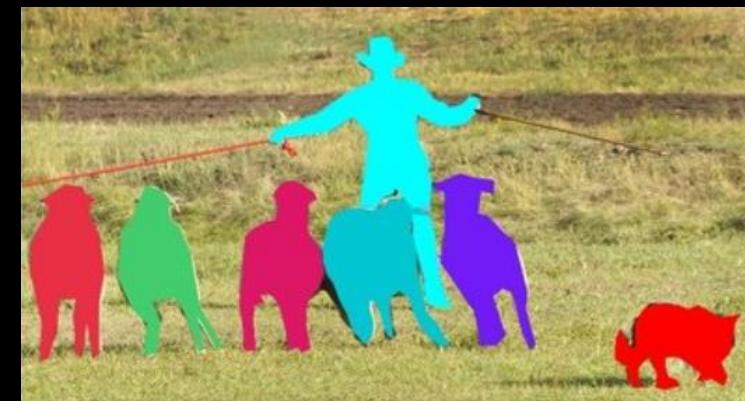
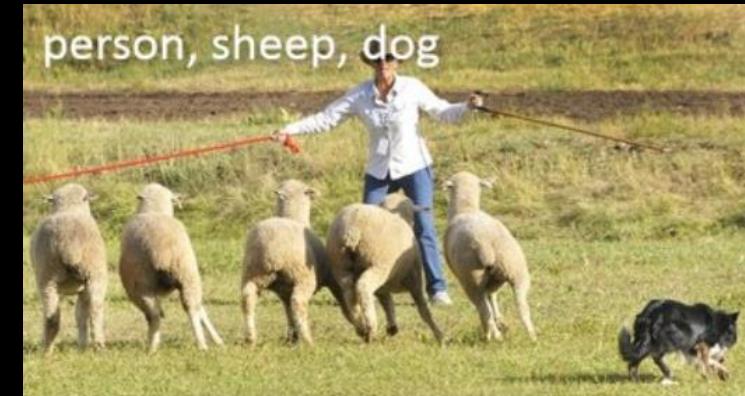
```
def forward(self, x):
    x = self.features(x)
    x = x.view(x.size(0), -1)
    x = self.classifier(x)
    return x
```

# Instance Segmentation

Semantic Segmentation



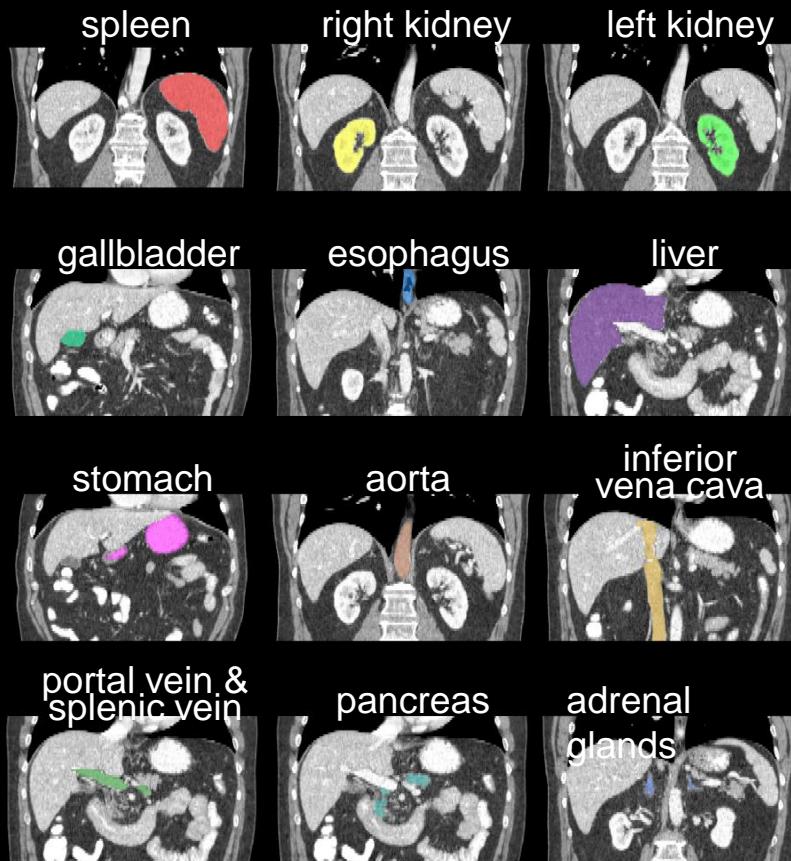
Instance Segmentation



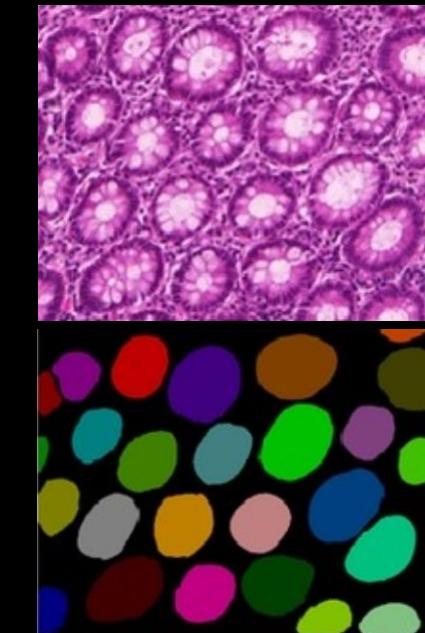
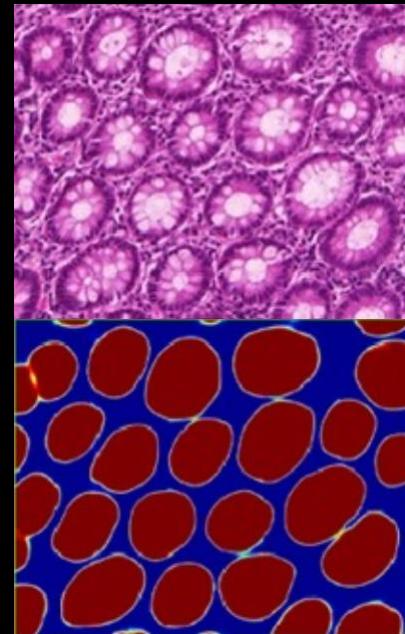
<https://arxiv.org/abs/1405.0312>

# Instance Segmentation

Semantic Segmentation



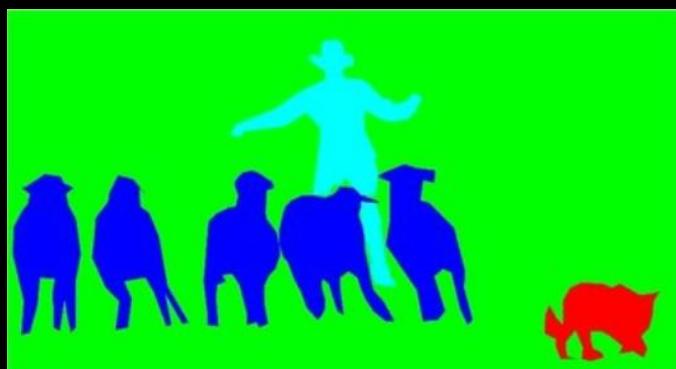
Instance Segmentation



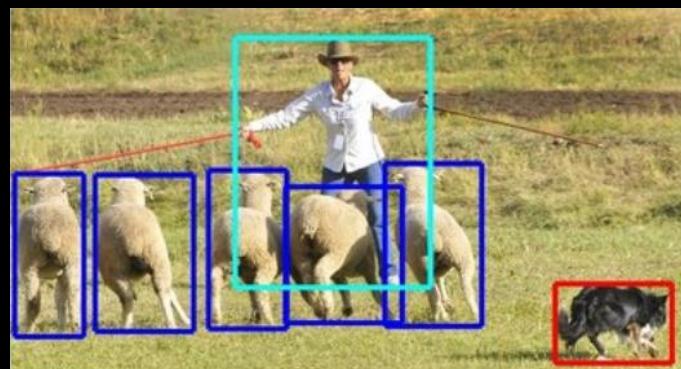
<https://www.sciencedirect.com/science/article/pii/S1361841516302043>

# How To Perform Instance Segmentation

Semantic Segmentation



Detection

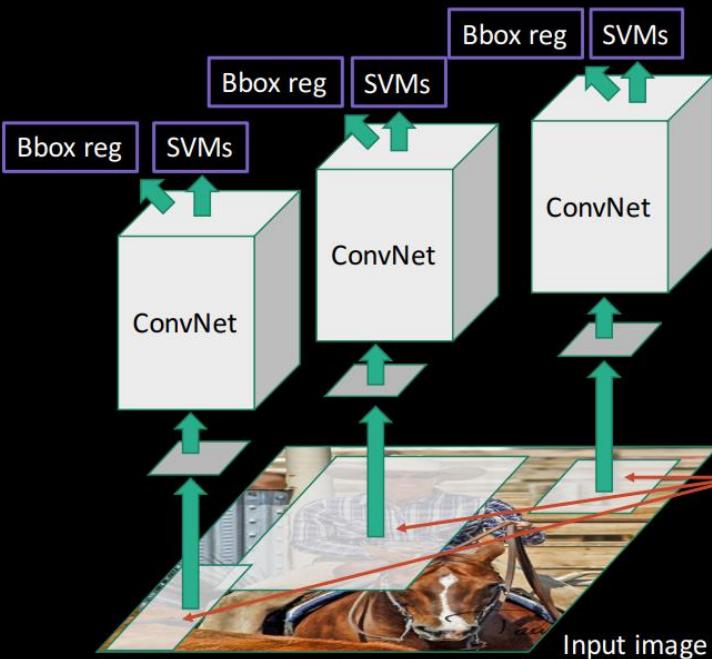


Instance Segmentation

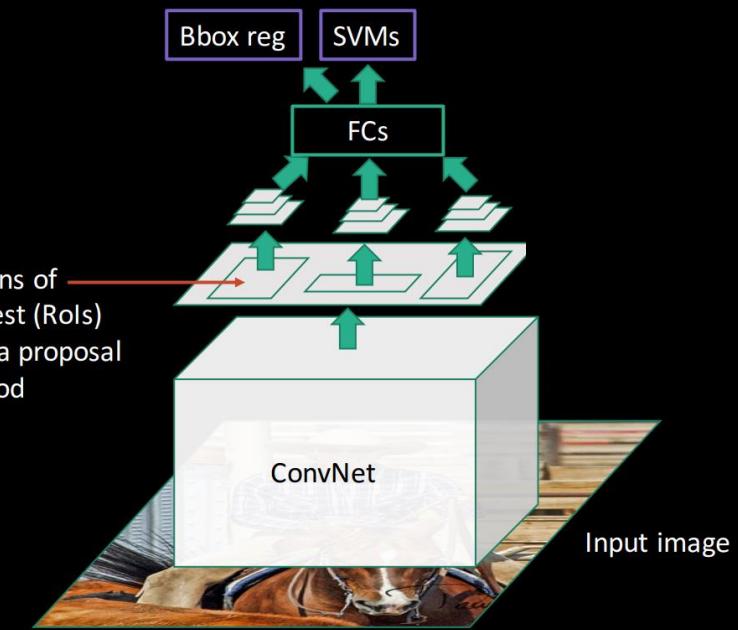


# Detection using RCNN Family

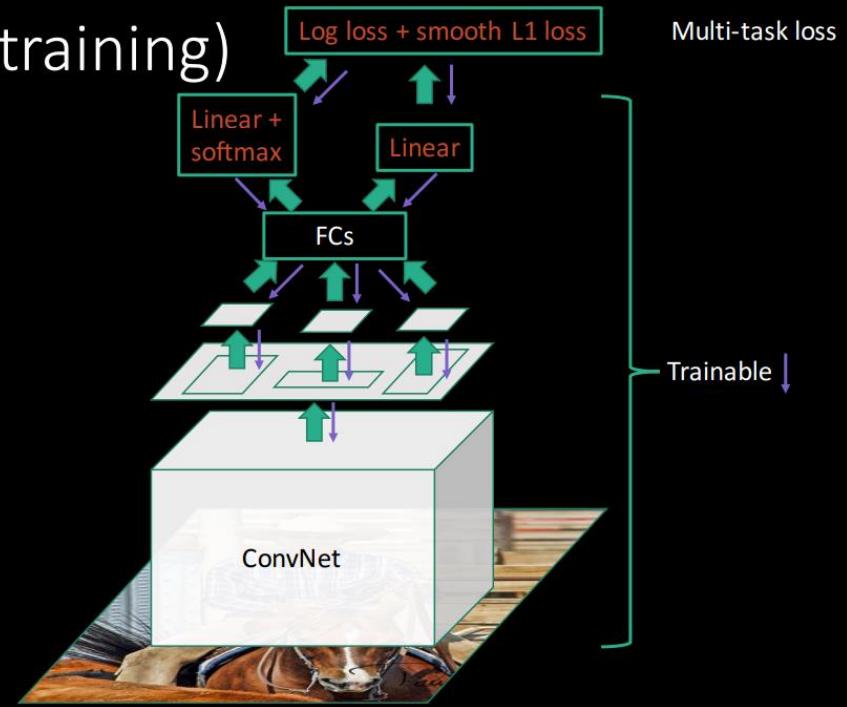
Slow R-CNN



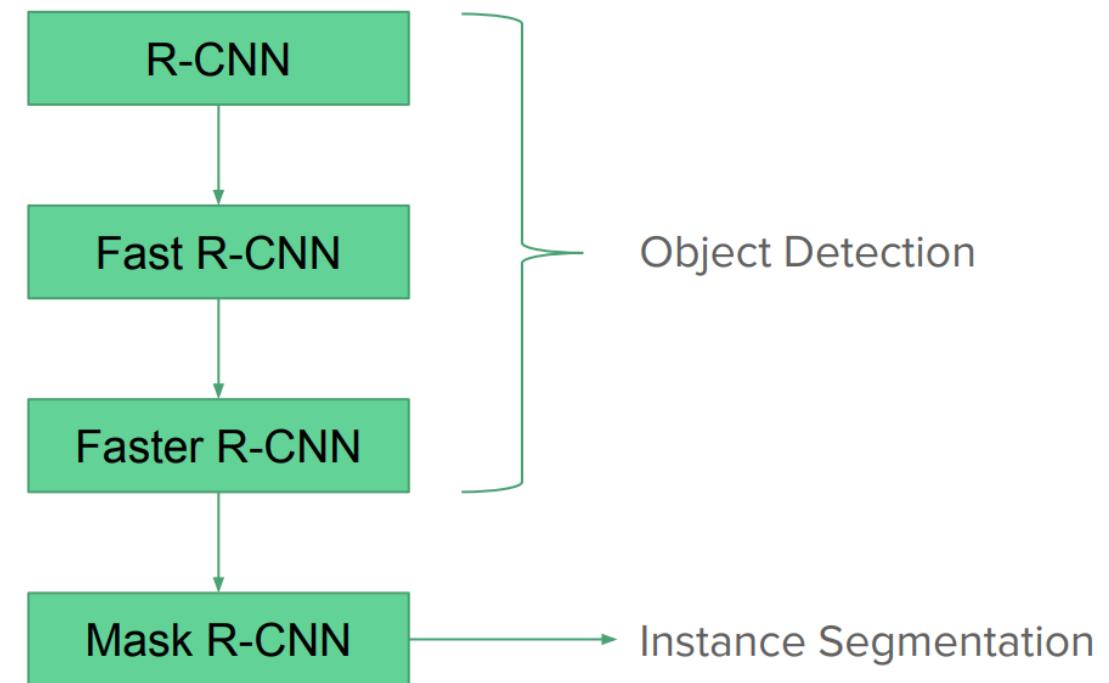
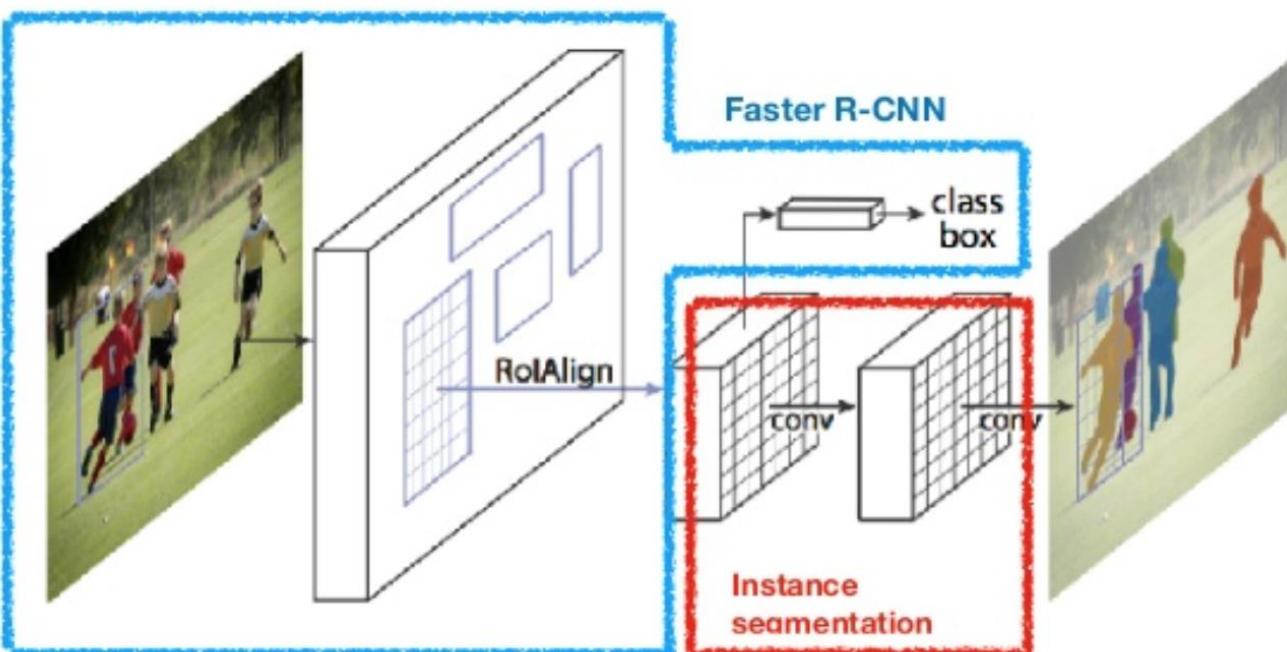
SPP-net



Fast R-CNN  
(training)

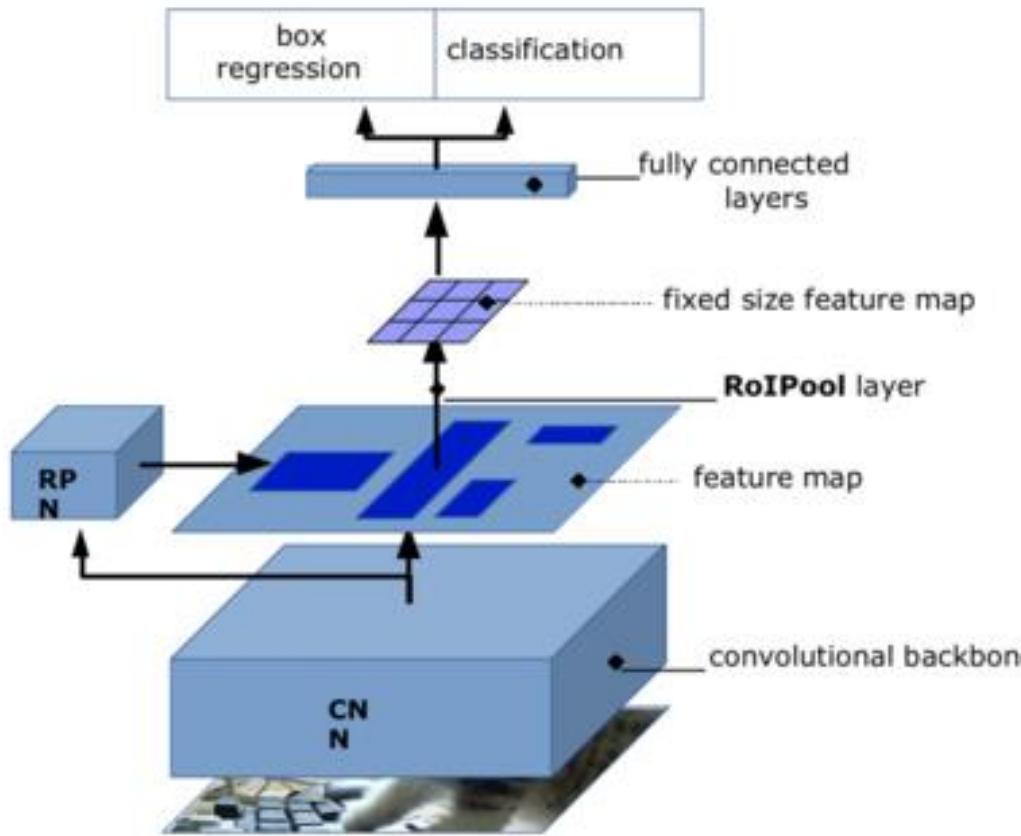


# Mask R-CNN



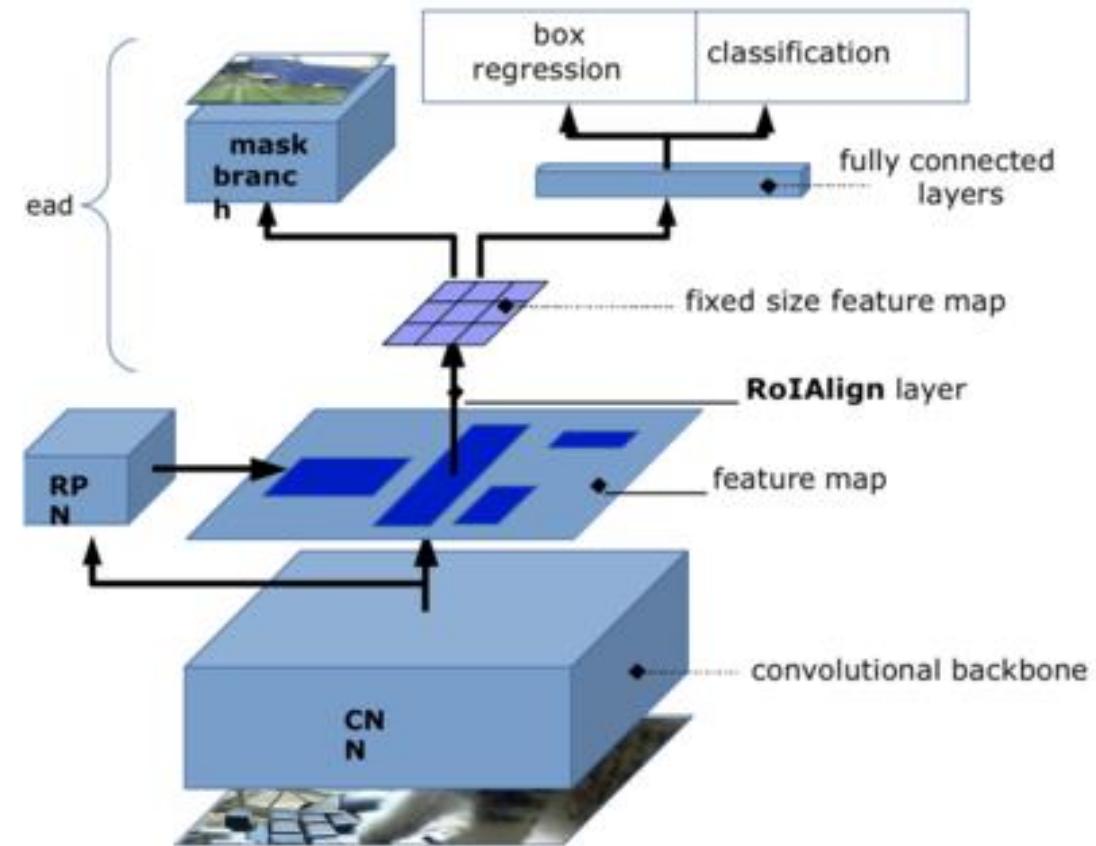
<https://lilianweng.github.io/lil-log/2017/12/31/object-recognition-for-dummies-part-3.html>

# Faster R-CNN vs. Mask R-CNN



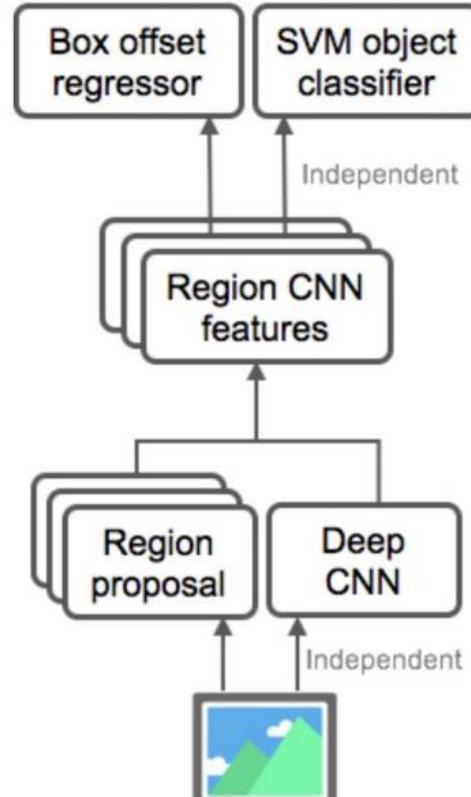
a)

[https://wiki.ubc.ca/CNNs\\_in\\_Image\\_Segmentation](https://wiki.ubc.ca/CNNs_in_Image_Segmentation)

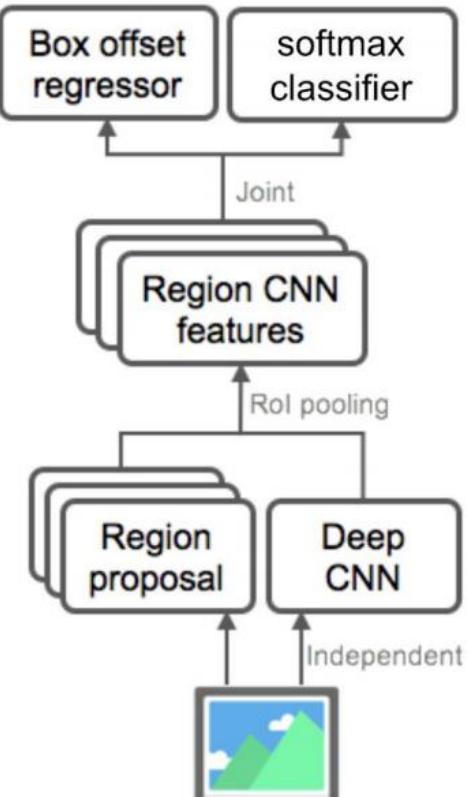


b)

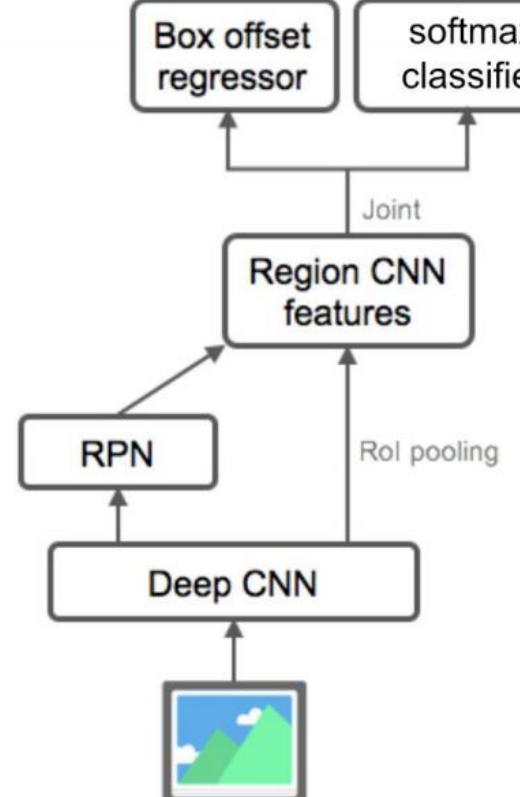
# R-CNN Family



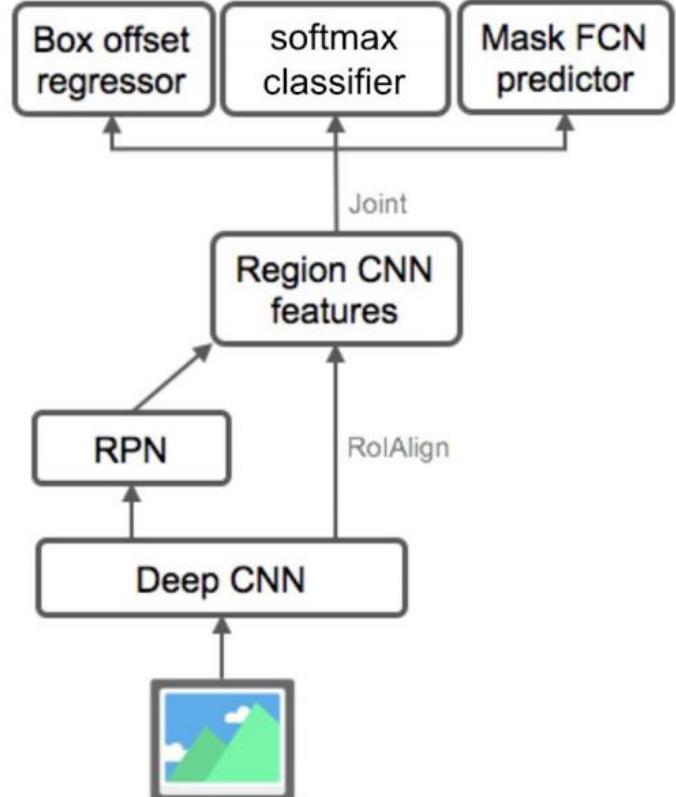
**R-CNN**



**Fast R-CNN**

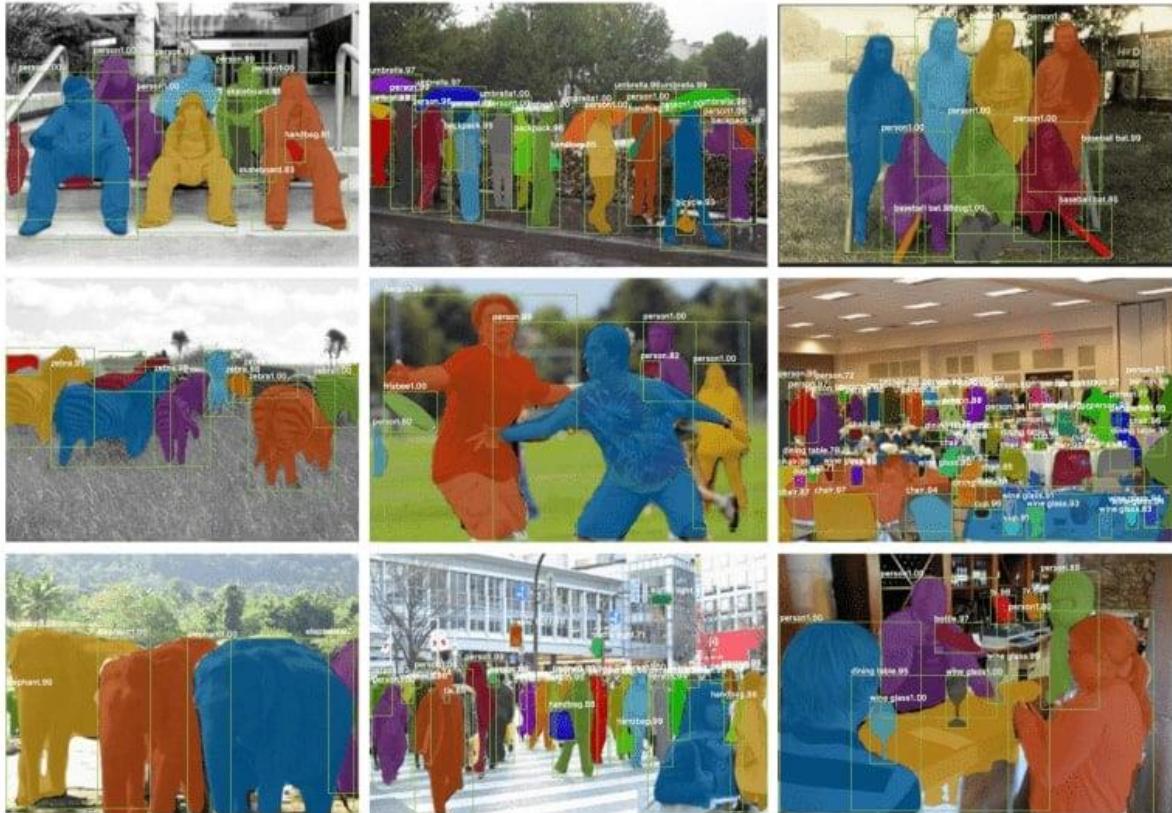


**Faster R-CNN**

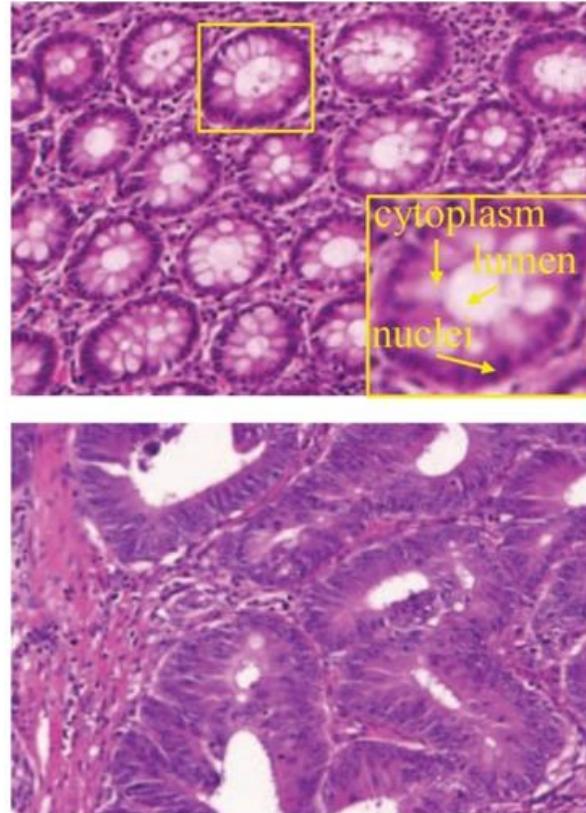


**Mask R-CNN**

# Instance Segmentation



<https://sigmoidal.io/dl-computer-vision-beyond-classification/>



<https://www.ocutri.com/medical/?language=en>

