

# mask R-CNN

*Sung Man Cho*

# INDEX

## 1. Related Work

- R-CNN series
- FCN

## 2. Mask R-CNN

- Architecture
- Mask Branch
- RoI Align

# **Related Work**

# Related Work

## R-CNN Series

R-CNN

Fast R-CNN

Faster R-CNN

Mask R-CNN

**Selective Search**

**Selective Search**

**RoI Pooling**

**RPN**

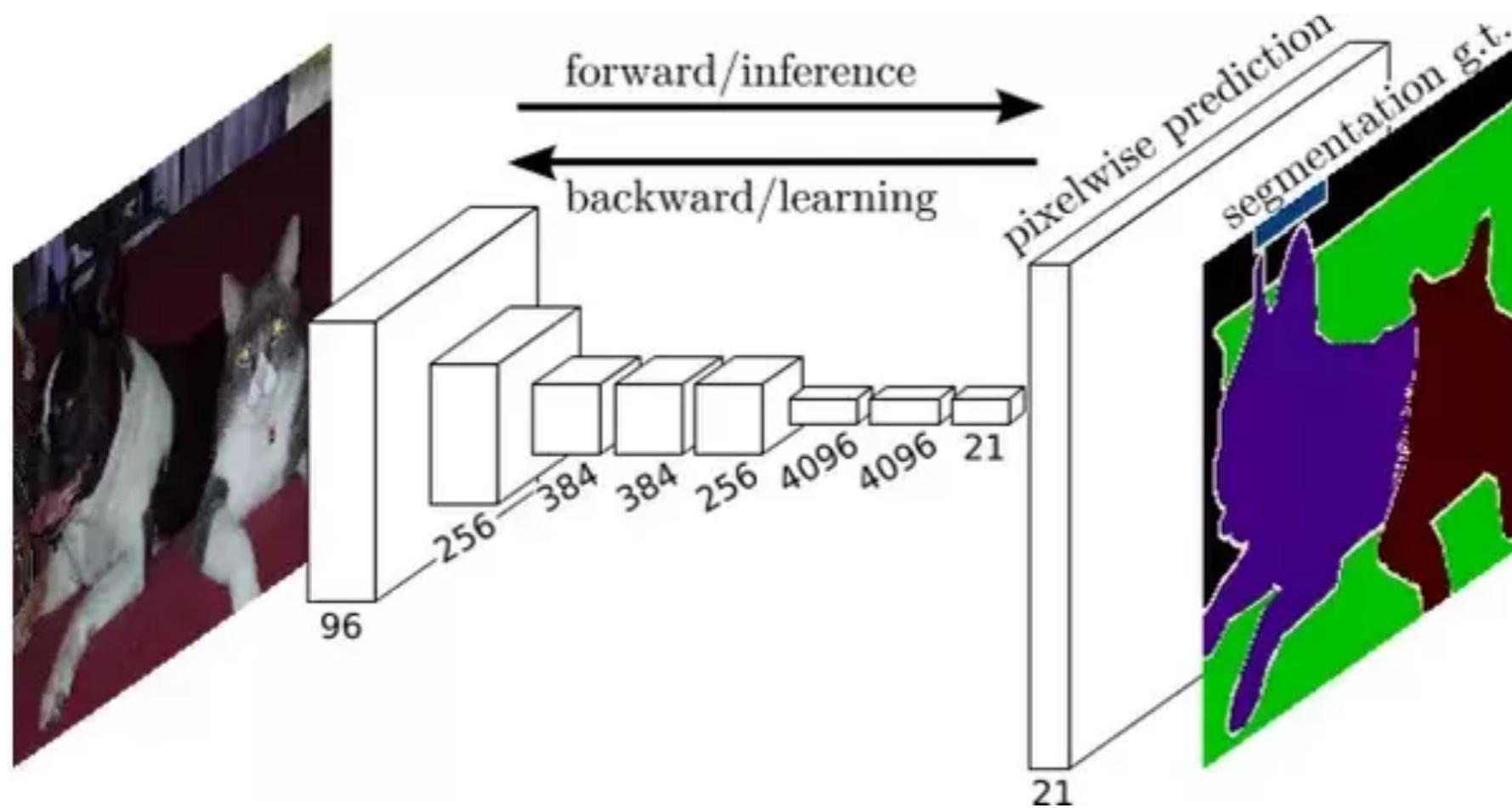
**Masking**

**RoI Align**

# Related Work

## FCN

- Fully Convolutional Layer -> heat map output.

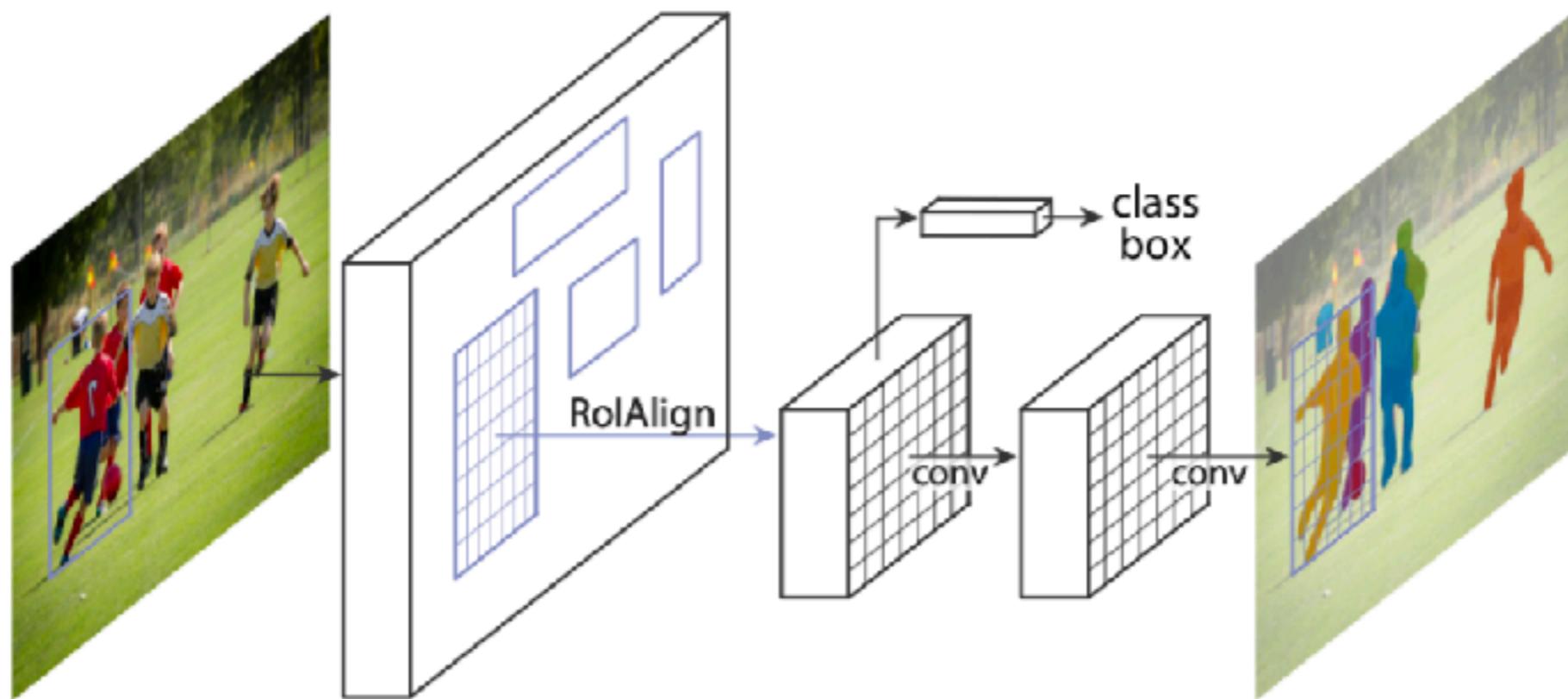


# **Mask R-CNN**

# Mask R-CNN

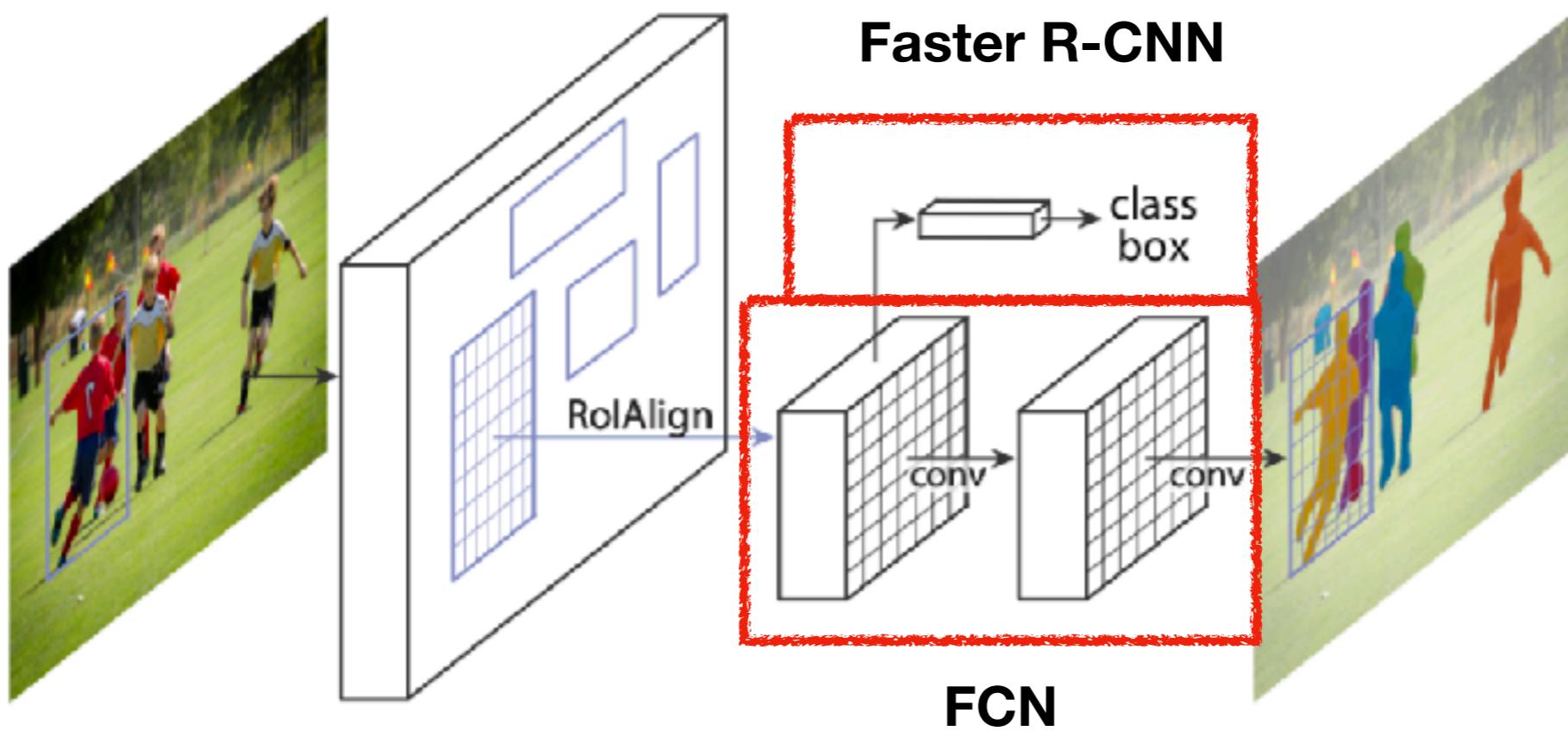
## Architecture

- Goal is enabling framework for **instance segmentation**.



# Mask R-CNN

## Architecture



<i>net-depth-features</i>	AP	AP <sub>50</sub>	AP <sub>75</sub>
ResNet-50-C4	30.3	51.2	31.5
ResNet-101-C4	32.7	54.2	34.3
ResNet-50-FPN	33.6	55.2	35.3
ResNet-101-FPN	35.4	57.3	37.5
ResNeXt-101-FPN	<b>36.7</b>	<b>59.5</b>	<b>38.9</b>

# Mask R-CNN

## Architecture



### Faster R-CNN

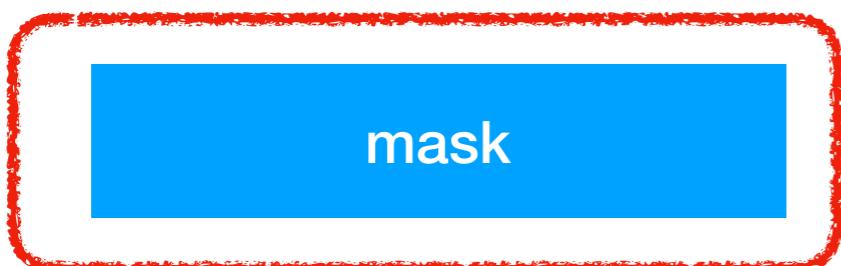
- Classification
- Instance Level RoI

### FCN

- ~~Pixel Level Classification~~
- ~~Per Pixel Softmax (Multinomial)~~
- ~~Multi Instance~~

# Mask R-CNN

## Architecture



### FCN

- ~~Pixel level Classification~~
  - Per Pixel ~~Softmax (Multinomial)~~
  - ~~Multi Instance~~
- 
- Per Pixel Sigmoid (Binary)

	AP	AP <sub>50</sub>	AP <sub>75</sub>
<i>softmax</i>	24.8	44.1	25.1
<i>sigmoid</i>	<b>30.3</b>	<b>51.2</b>	<b>31.5</b>
	+5.5	+7.1	+6.4

# Mask R-CNN

## Architecture

DB  
BBox + Class + Mask

$$L = L_{cls} + L_{box} + L_{mask}$$

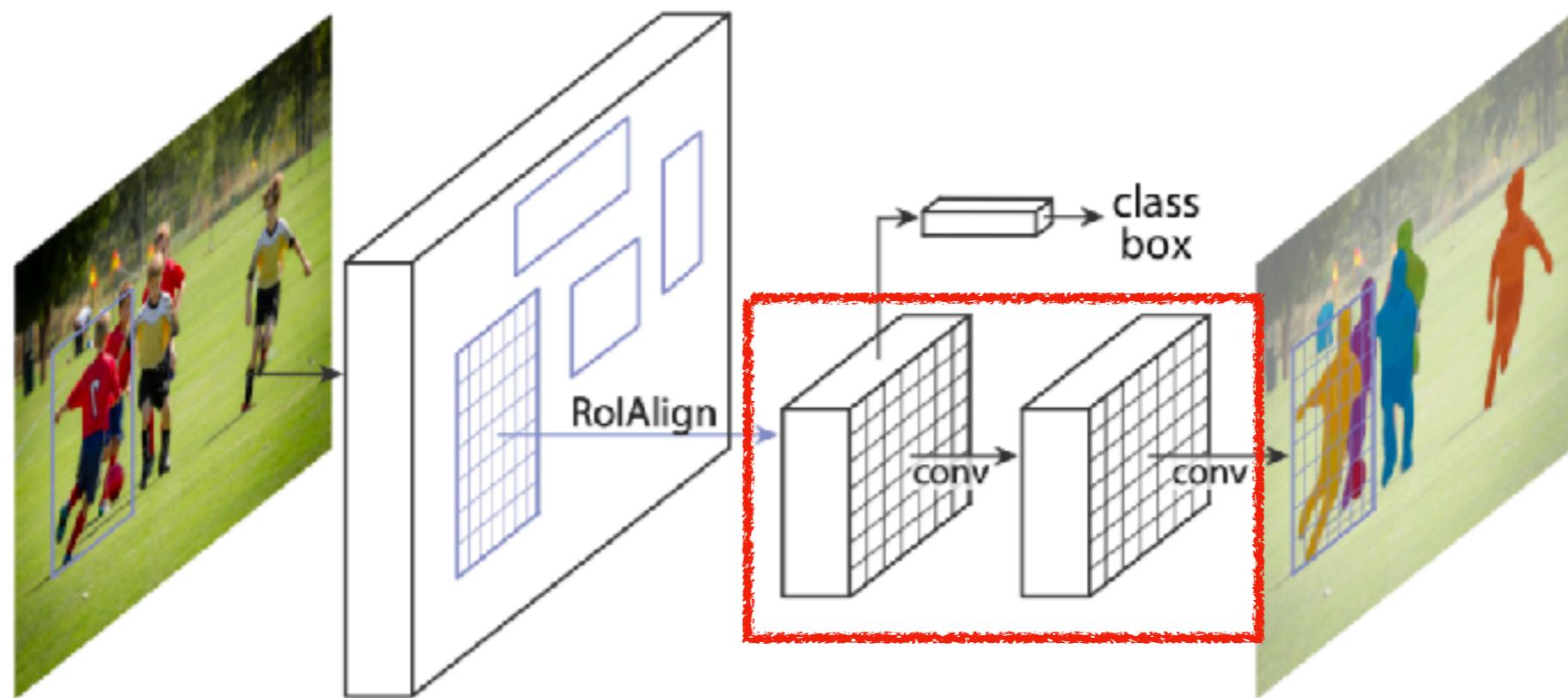
$L_{cls}$ : Softmax Cross Entropy

$L_{box}$ : Regression

$L_{mask}$ : Binary Cross Entropy

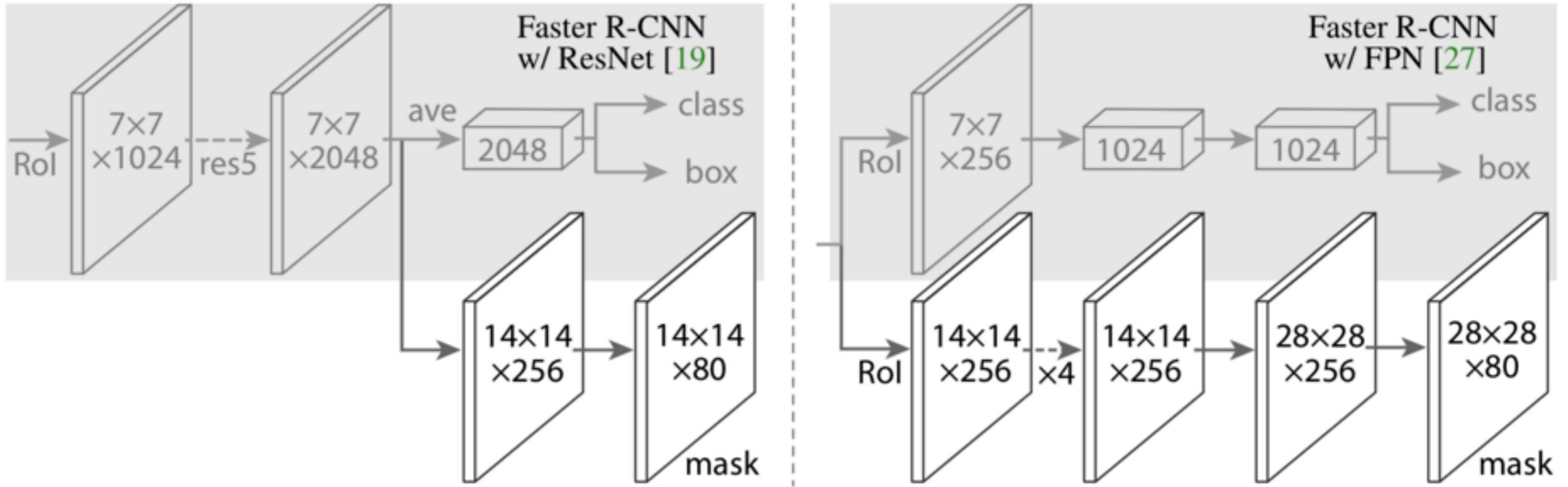
# Mask R-CNN

## Mask Branch



# Mask R-CNN

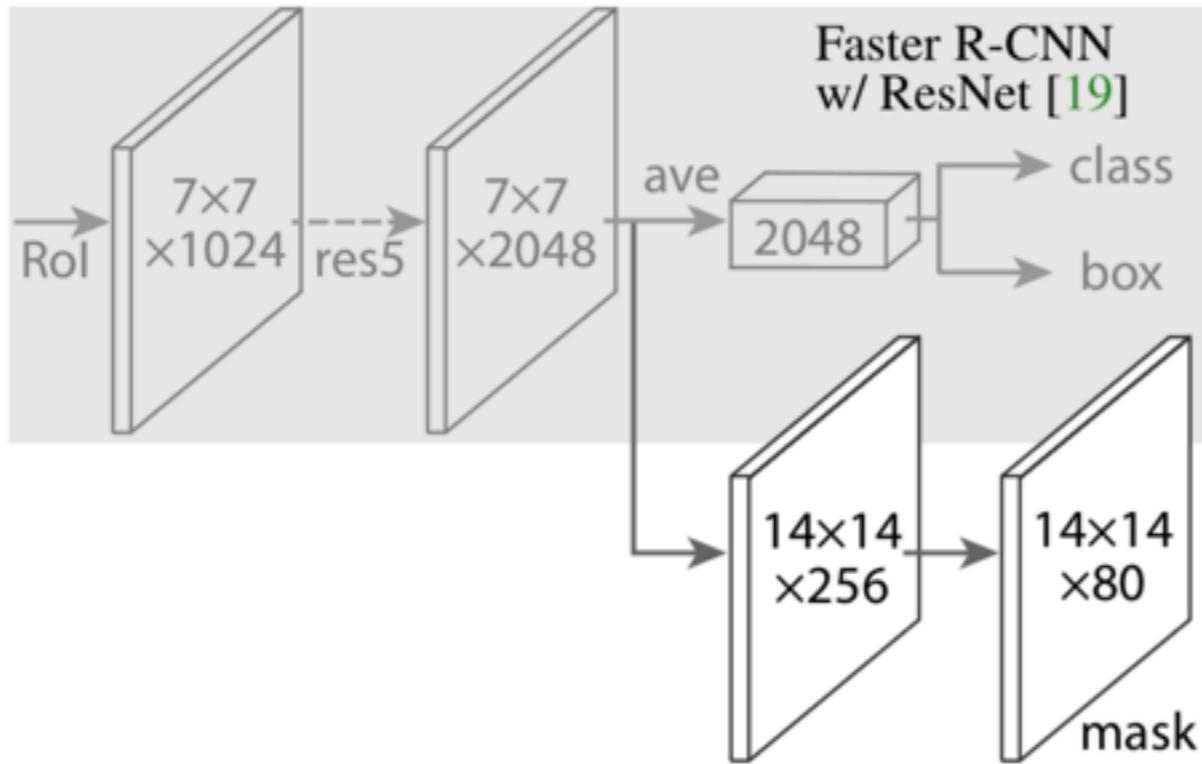
## Mask Branch



while deconv increases it). All convs are  $3 \times 3$ , except the output conv which is  $1 \times 1$ , deconv are  $2 \times 2$  with stride 2, and we use ReLU [31] in hidden layers. *Left*: ‘res5’ denotes ResNet’s fifth stage, which for simplicity we altered so that the first conv operates on a  $7 \times 7$  ROI with stride 1 (instead of  $14 \times 14$  / stride 2 as in [19]). *Right*: ‘ $\times 4$ ’ denotes a stack of four consecutive convs.

# Mask R-CNN

## Mask Branch



while deconv increases it). All convs are  $3 \times 3$ , except the output conv which is  $1 \times 1$ , deconv are  $2 \times 2$  with stride 2, and we use ReLU [31] in hidden layers. *Left*: ‘res5’ denotes ResNet’s fifth stage, which for simplicity we altered so that the first conv operates on a  $7 \times 7$  RoI with stride 1 (instead of  $14 \times 14$  / stride 2 as in [19]). *Right*: ‘ $\times 4$ ’ denotes a stack of four consecutive convs.

```
def _mask_network(self):
    def f(x):
        target_metadata, output_features, output_proposal_bounding_boxes = x

        mask_features = keras_rcnn.layers.RegionOfInterest(
            extent=(14, 14),
            strides=2
        )([
            target_metadata,
            output_features,
            output_proposal_bounding_boxes
        ])

        mask_features = keras.layers.TimeDistributed(
            keras.layers.Conv2D(
                activation="relu",
                filters=256,
                kernel_size=(3, 3),
                padding="same"
            )
        )(mask_features)

        mask_features = keras.layers.TimeDistributed(
            keras.layers.Conv2D(
                activation="relu",
                filters=256,
                kernel_size=(3, 3),
                padding="same"
            )
        )(mask_features)
```

# Mask R-CNN

## Mask Branch

```
def _mask_network(self):
    def f(x):
        target_metadata, output_features, output_proposal_bounding_boxes = x

        mask_features = keras_rcnn.layers.RegionOfInterest(
            extent=(14, 14),
            strides=2
        )([
            target_metadata,
            output_features,
            output_proposal_bounding_boxes
        ])

        mask_features = keras.layers.TimeDistributed(
            keras.layers.Conv2D(
                activation="relu",
                filters=256,
                kernel_size=(3, 3),
                padding="same"
            )
        )(mask_features)

        mask_features = keras.layers.TimeDistributed(
            keras.layers.Conv2DTranspose(
                activation="relu",
                filters=256,
                kernel_size=(2, 2),
                strides=2
            )
        )(mask_features)

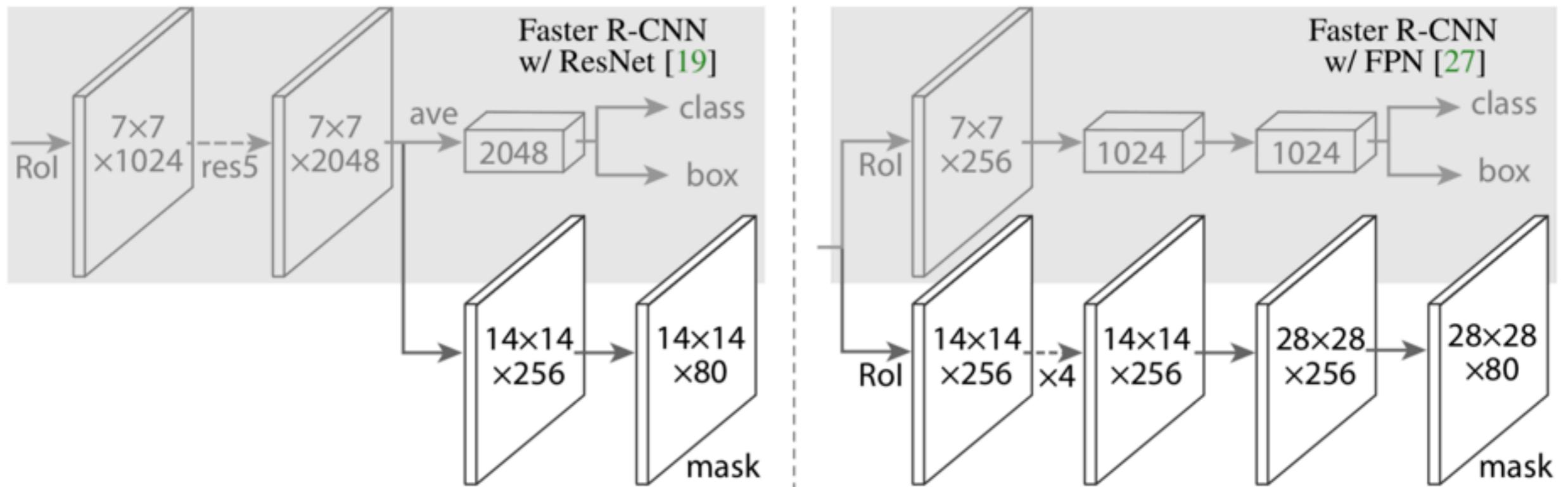
        mask_features = keras.layers.TimeDistributed(
            keras.layers.Conv2D(
                activation="sigmoid",
                filters=self.n_categories,
                kernel_size=(1, 1),
                strides=1
            )
        )(mask_features)

        return mask_features

    return f
```

# Mask R-CNN

## Mask Branch

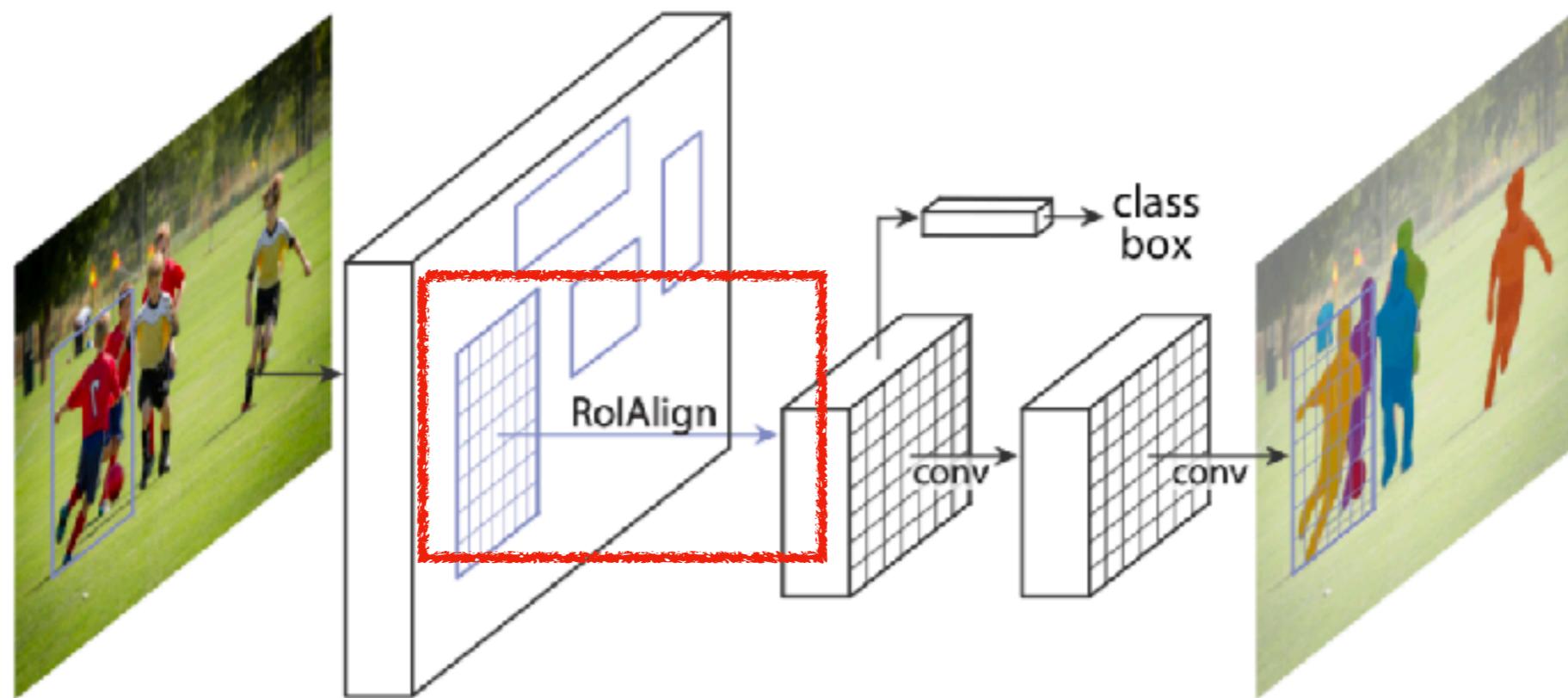


	mask branch	AP	AP <sub>50</sub>	AP <sub>75</sub>
MLP	fc: 1024 → 1024 → 80 · 28 <sup>2</sup>	31.5	53.7	32.8
MLP	fc: 1024 → 1024 → 1024 → 80 · 28 <sup>2</sup>	31.5	54.0	32.6
<b>FCN</b>	conv: 256 → 256 → 256 → 256 → 256 → 80	<b>33.6</b>	<b>55.2</b>	<b>35.3</b>

(e) **Mask Branch** (ResNet-50-FPN): Fully convolutional networks (FCN) vs. multi-layer perceptrons (MLP, fully-connected) for mask prediction. FCNs improve results as they take advantage of explicitly encoding spatial layout.

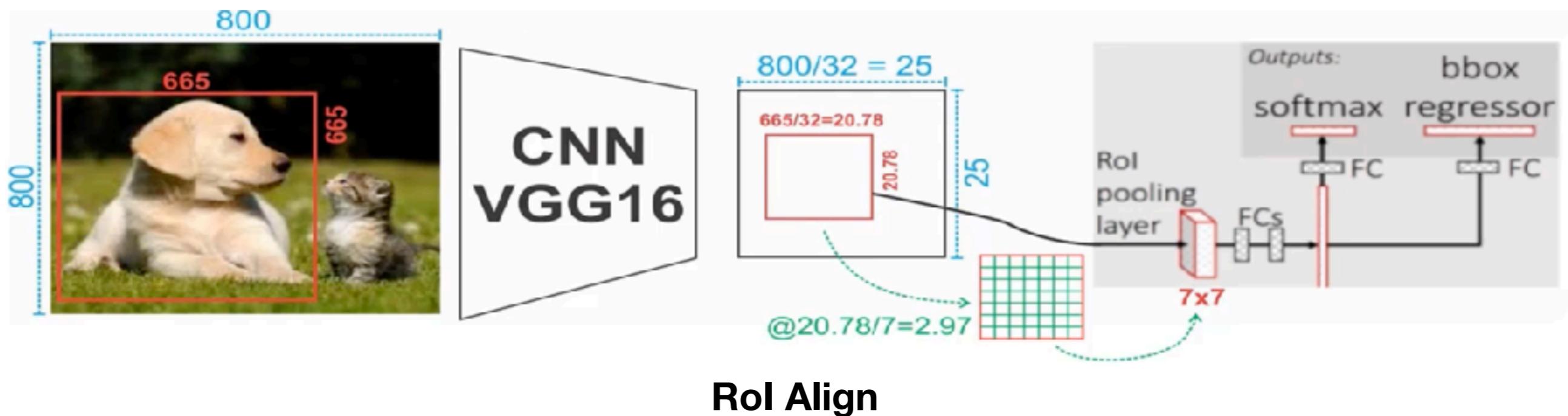
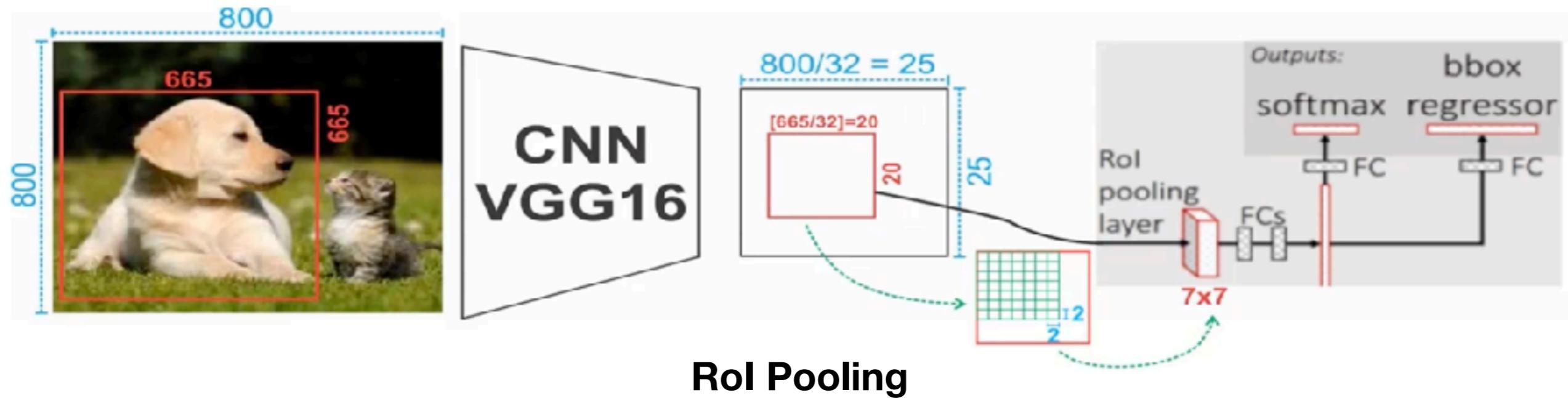
# Mask R-CNN

## RoIAlign



# Mask R-CNN

## RoIAlign



# Mask R-CNN

## RoIAlign

```
x1, y1, x2, y2 = torch.split(boxes, 1, dim=1)
image_height, image_width = featuremap.size()[2:4]

if self.transform_fpcoor:
    spacing_w = (x2 - x1) / float(self.crop_width)
    spacing_h = (y2 - y1) / float(self.crop_height)

    nx0 = (x1 + spacing_w / 2 - 0.5) / float(image_width - 1)
    ny0 = (y1 + spacing_h / 2 - 0.5) / float(image_height - 1)
    nw = spacing_w * float(self.crop_width - 1) / float(image_width - 1)
    nh = spacing_h * float(self.crop_height - 1) / float(image_height - 1)

    boxes = torch.cat((ny0, nx0, ny0 + nh, nx0 + nw), 1)

else:
    x1 = x1 / float(image_width - 1)
    x2 = x2 / float(image_width - 1)
    y1 = y1 / float(image_height - 1)
    y2 = y2 / float(image_height - 1)
    boxes = torch.cat((y1, x1, y2, x2), 1)

boxes = boxes.detach().contiguous()
box_ind = box_ind.detach()

return CropAndResizeFunction(self.crop_height, self.crop_width, self.extrapolation_value)(featuremap, boxes, box_ind)
```

# Mask R-CNN

## RoIAlign

```
# convert regions from (x, y, w, h) to (x1, y1, x2, y2)
boxes = keras.backend.cast(boxes, keras.backend.floatx())

boxes = boxes / self.stride

x1 = boxes[:, 0]
y1 = boxes[:, 1]
x2 = boxes[:, 2]
y2 = boxes[:, 3]

# normalize the boxes
shape = metadata[0]

h = keras.backend.cast(shape[0], keras.backend.floatx())
w = keras.backend.cast(shape[1], keras.backend.floatx())

x1 /= w
y1 /= h
x2 /= w
y2 /= h

x1 = keras.backend.expand_dims(x1, axis=2)
y1 = keras.backend.expand_dims(y1, axis=2)
x2 = keras.backend.expand_dims(x2, axis=2)
y2 = keras.backend.expand_dims(y2, axis=2)

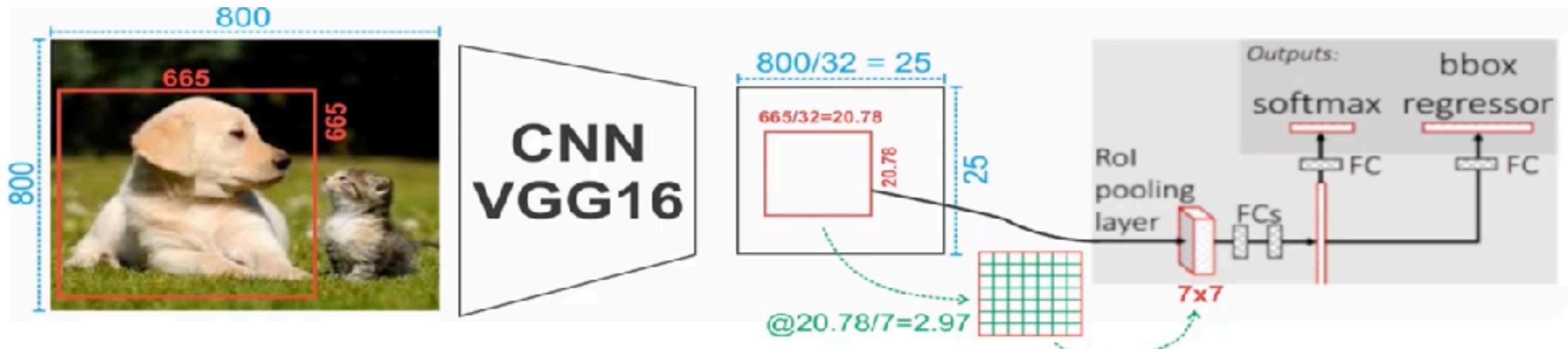
boxes = keras.backend.concatenate([y1, x1, y2, x2], axis=2)
boxes = keras.backend.reshape(boxes, (-1, 4))

slices = keras_rcnn.backend.crop_and_resize(image, boxes, self.extent)

return keras.backend.expand_dims(slices, axis=0)
```

# Mask R-CNN

## RoIAlign



	AP	AP <sub>50</sub>	AP <sub>75</sub>	AP <sup>bb</sup>	AP <sub>50</sub> <sup>bb</sup>	AP <sub>75</sub> <sup>bb</sup>
<i>RoIPool</i>	23.6	46.5	21.6	28.2	52.7	26.9
<b>RoIAlign</b>	<b>30.9</b>	<b>51.8</b>	<b>32.1</b>	<b>34.0</b>	<b>55.3</b>	<b>36.4</b>
	+7.3	+ 5.3	+10.5	+5.8	+2.6	+9.5

(d) **RoIAlign** (ResNet-50-C5, *stride 32*): Mask-level and box-level AP using *large-stride* features. Misalignments are more severe than with stride-16 features (Table 2c), resulting in big accuracy gaps.

# Fully Convolutional Networks

*Jeon Ji Hye*

# Index

- Introduction
- Architecture
- Fully Convolutional Networks
- Up-Sampling

# **Before Presentation**

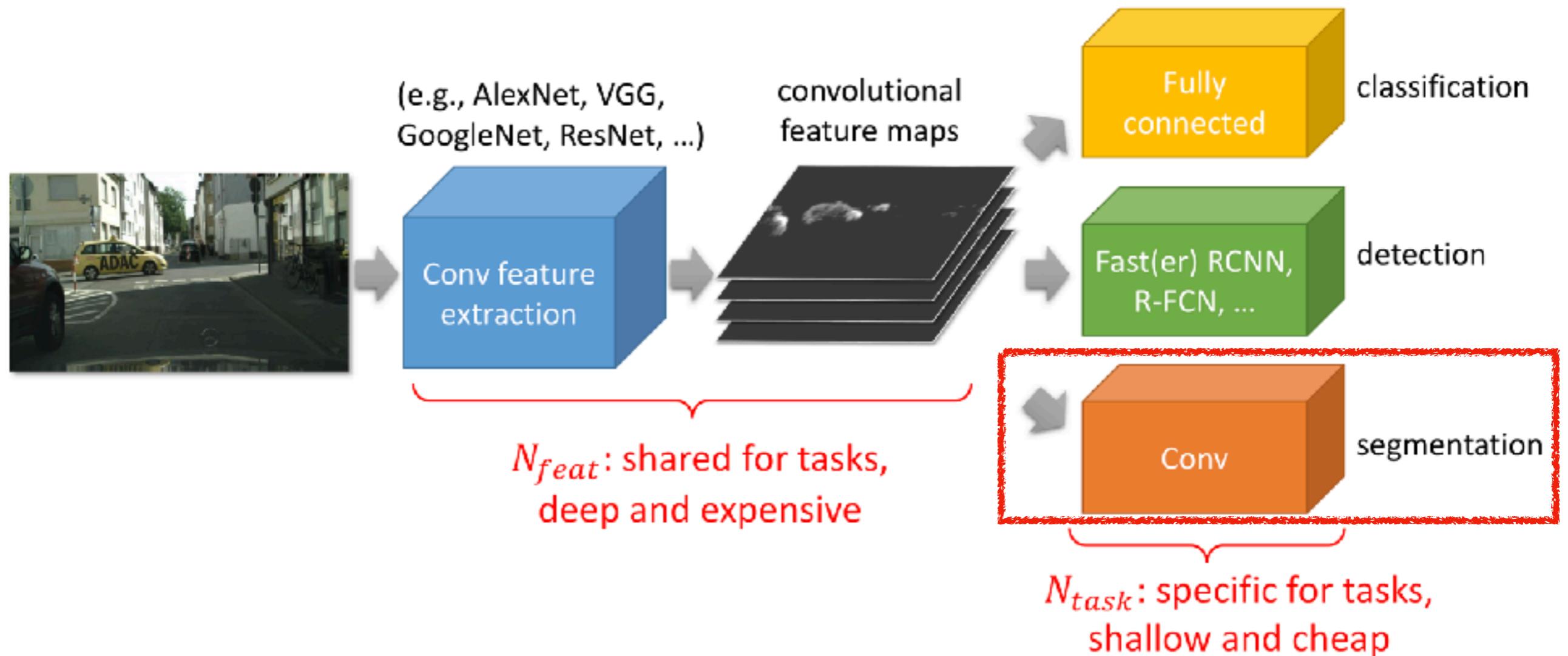
# Two Important Things..

Fully Convolutional Networks  
: 1x1 Convolution Layer

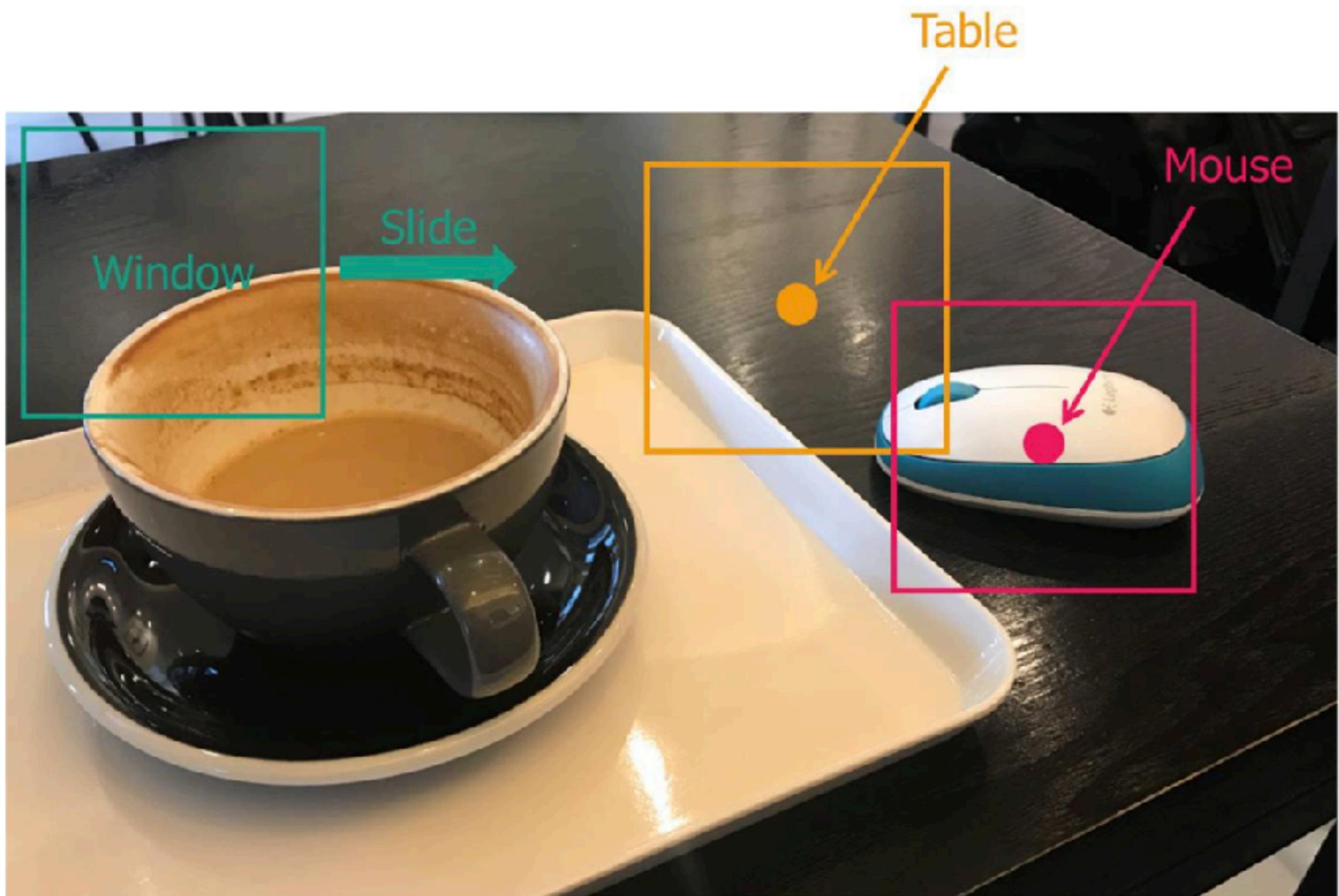
Skip Connection

# Introduction

# Introduction

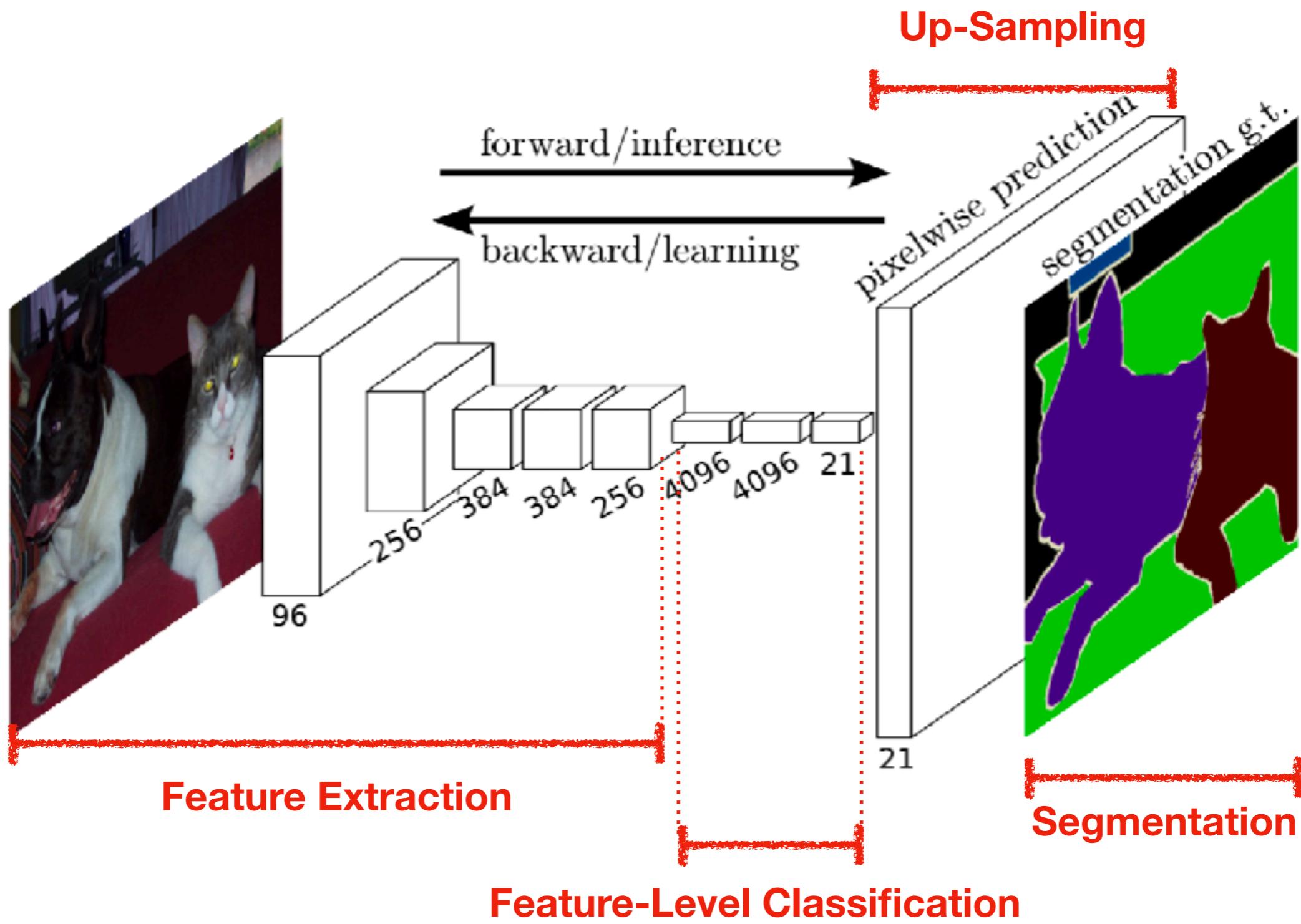


# Previous Work

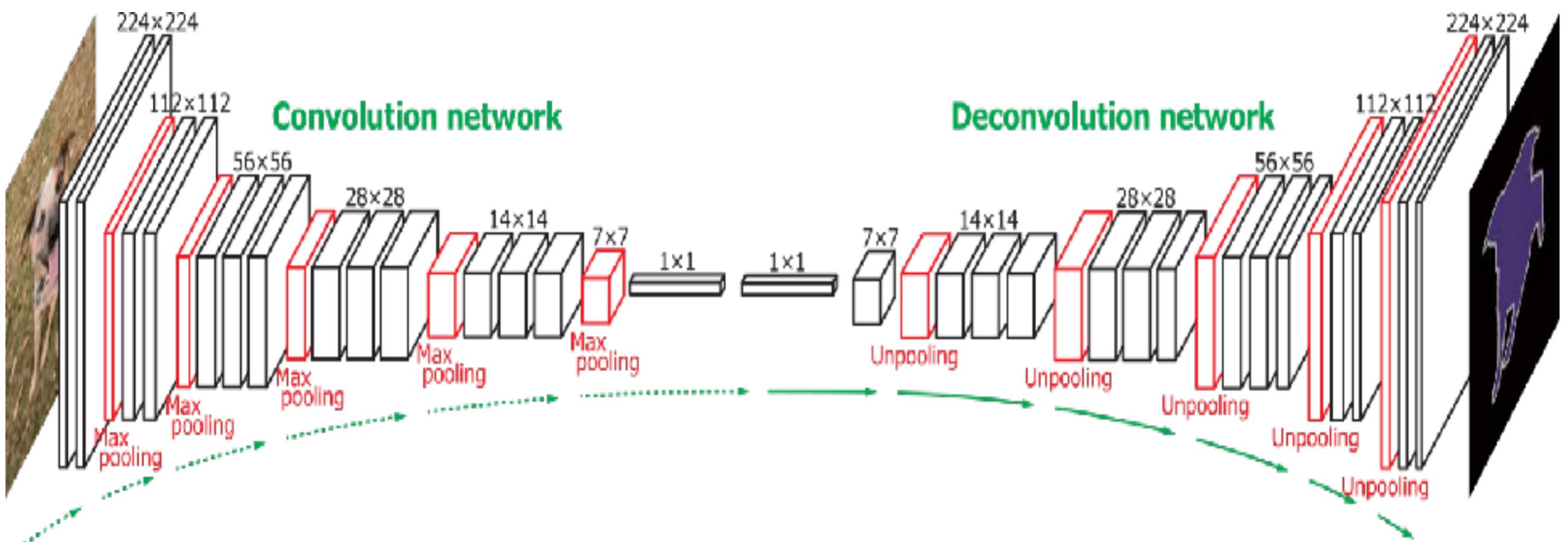


# Architecture

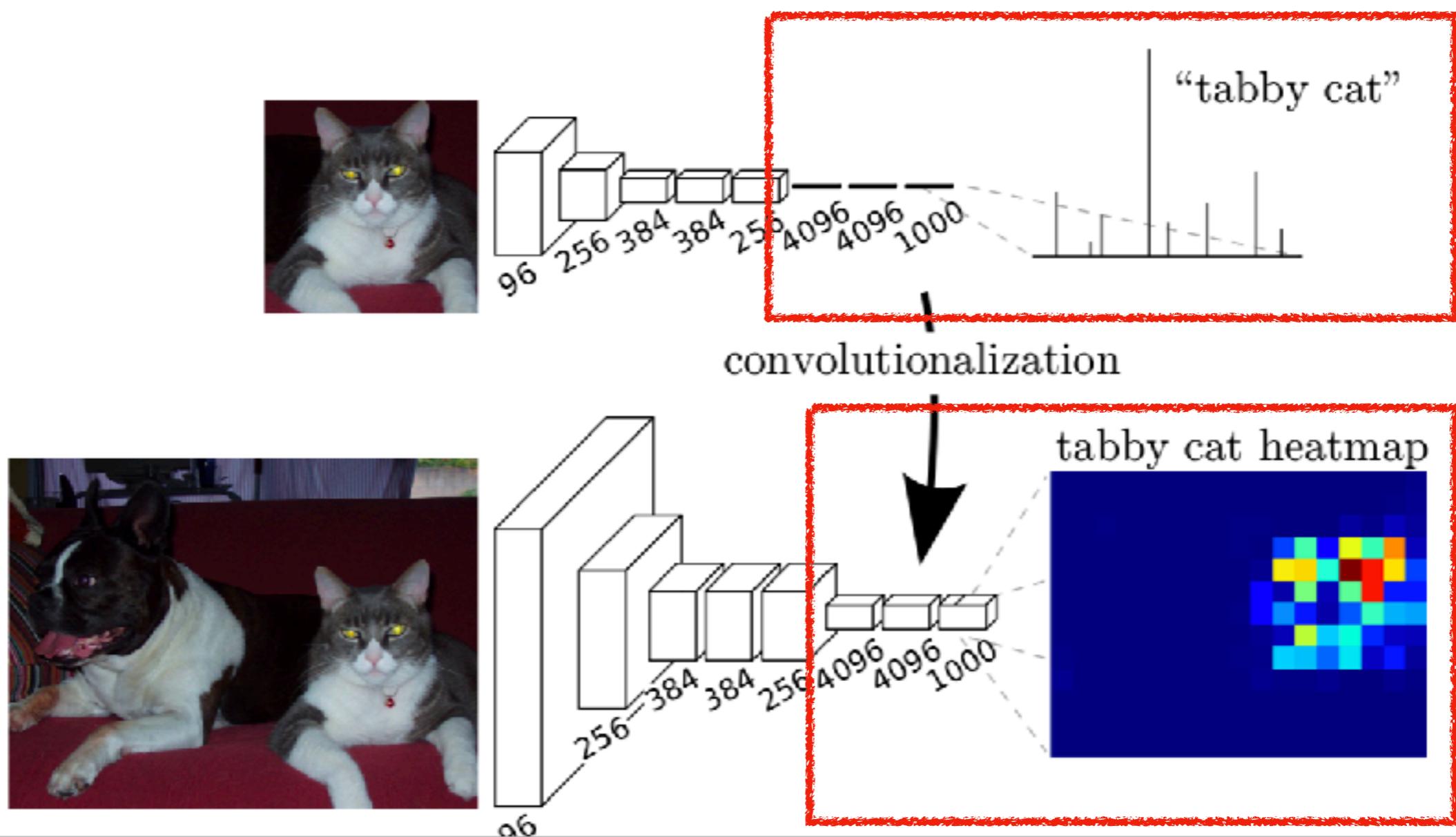
# Architecture



# Architecture

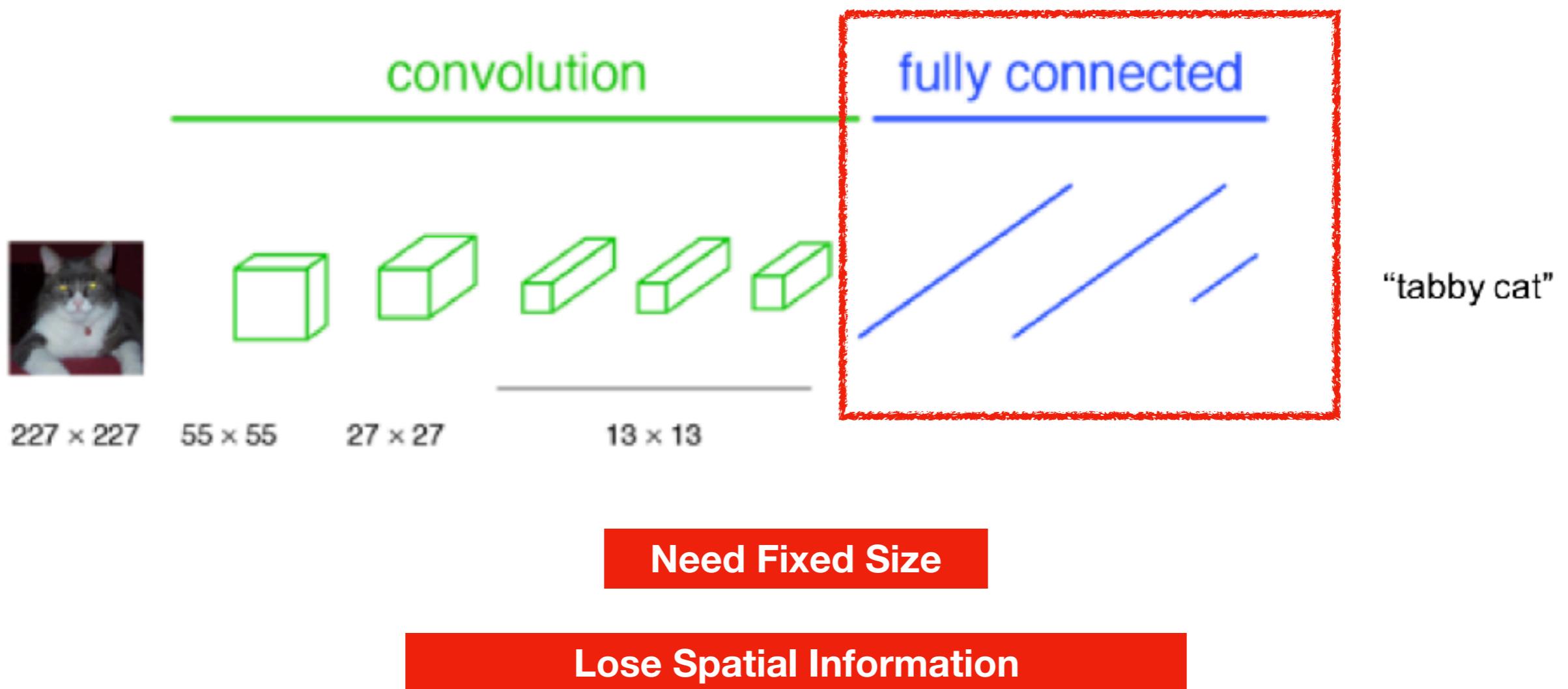


# Architecture



# **Fully Convolutional Networks**

# Fully Connected Problem



# FC Layer - Too many Param

INPUT: [224x224x3] memory:  $224 \times 224 \times 3 = 150K$  params: 0 (not counting biases)

CONV3-64: [224x224x64] memory:  $224 \times 224 \times 64 = 3.2M$  params:  $(3 \times 3 \times 3) \times 64 = 1,728$

CONV3-64: [224x224x64] memory:  $224 \times 224 \times 64 = 3.2M$  params:  $(3 \times 3 \times 64) \times 64 = 36,864$

POOL2: [112x112x64] memory:  $112 \times 112 \times 64 = 800K$  params: 0

CONV3-128: [112x112x128] memory:  $112 \times 112 \times 128 = 1.6M$  params:  $(3 \times 3 \times 64) \times 128 = 73,728$

CONV3-128: [112x112x128] memory:  $112 \times 112 \times 128 = 1.6M$  params:  $(3 \times 3 \times 128) \times 128 = 147,456$

POOL2: [56x56x128] memory:  $56 \times 56 \times 128 = 400K$  params: 0

CONV3-256: [56x56x256] memory:  $56 \times 56 \times 256 = 800K$  params:  $(3 \times 3 \times 128) \times 256 = 294,912$

CONV3-256: [56x56x256] memory:  $56 \times 56 \times 256 = 800K$  params:  $(3 \times 3 \times 256) \times 256 = 589,824$

CONV3-256: [56x56x256] memory:  $56 \times 56 \times 256 = 800K$  params:  $(3 \times 3 \times 256) \times 256 = 589,824$

POOL2: [28x28x256] memory:  $28 \times 28 \times 256 = 200K$  params: 0

CONV3-512: [28x28x512] memory:  $28 \times 28 \times 512 = 400K$  params:  $(3 \times 3 \times 256) \times 512 = 1,179,648$

CONV3-512: [28x28x512] memory:  $28 \times 28 \times 512 = 400K$  params:  $(3 \times 3 \times 512) \times 512 = 2,359,296$

CONV3-512: [28x28x512] memory:  $28 \times 28 \times 512 = 400K$  params:  $(3 \times 3 \times 512) \times 512 = 2,359,296$

POOL2: [14x14x512] memory:  $14 \times 14 \times 512 = 100K$  params: 0

CONV3-512: [14x14x512] memory:  $14 \times 14 \times 512 = 100K$  params:  $(3 \times 3 \times 512) \times 512 = 2,359,296$

CONV3-512: [14x14x512] memory:  $14 \times 14 \times 512 = 100K$  params:  $(3 \times 3 \times 512) \times 512 = 2,359,296$

CONV3-512: [14x14x512] memory:  $14 \times 14 \times 512 = 100K$  params:  $(3 \times 3 \times 512) \times 512 = 2,359,296$

POOL2: [7x7x512] memory:  $7 \times 7 \times 512 = 25K$  params: 0

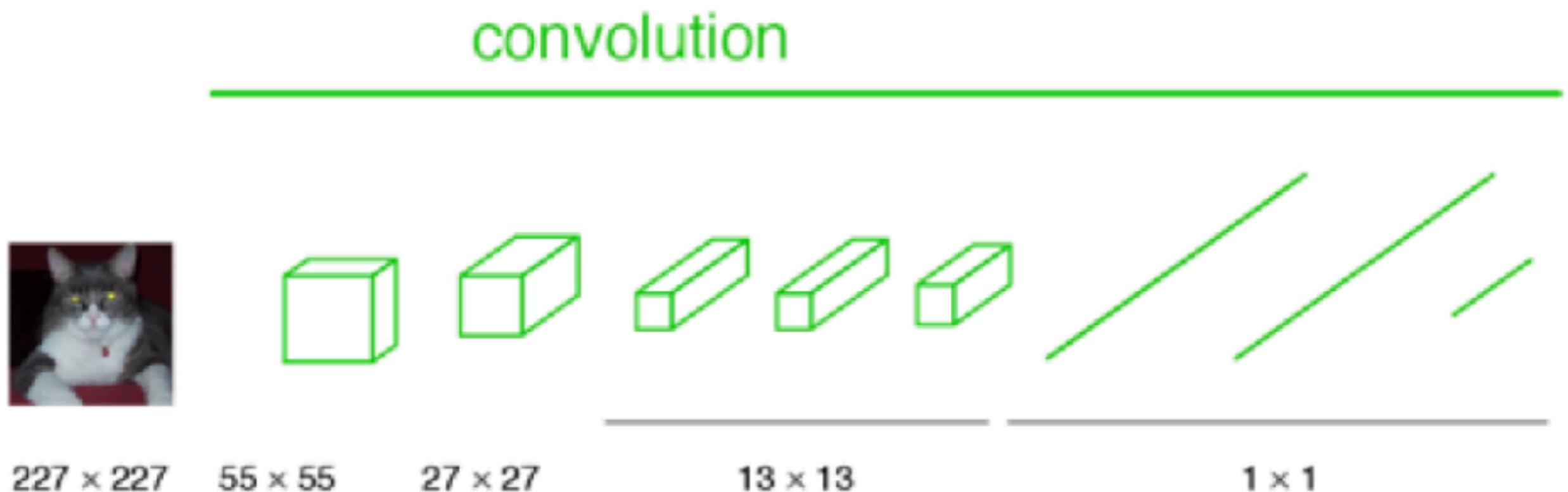
FC: [1x1x4096] memory: 4096 params:  $7 \times 7 \times 512 \times 4096 = 102,760,448$

FC: [1x1x4096] memory: 4096 params:  $4096 \times 4096 = 16,777,216$

FC: [1x1x1000] memory: 1000 params:  $4096 \times 1000 = 4,096,000$

VGG 16

# $1 \times 1$ Conv-Layer



Remain Spatial Information

per patch -> whole image  
: reduce calculation

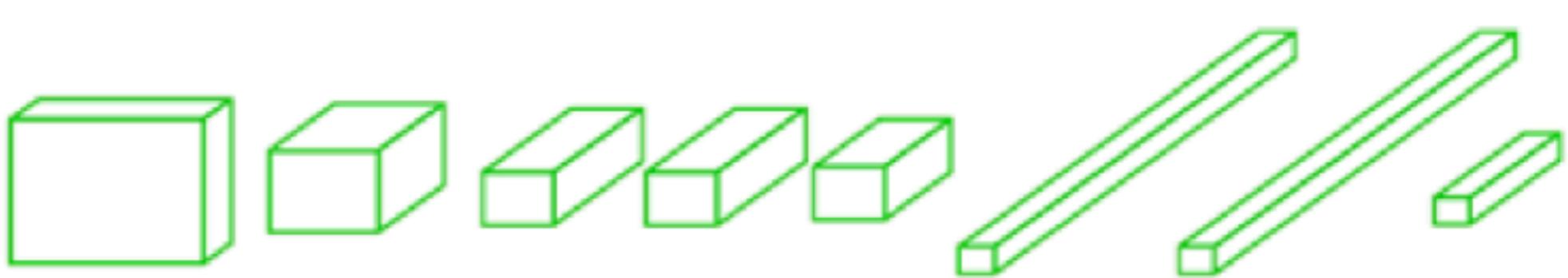
# 1x1 Conv-Layer

convolution

---



$H \times W$



$H/4 \times W/4$

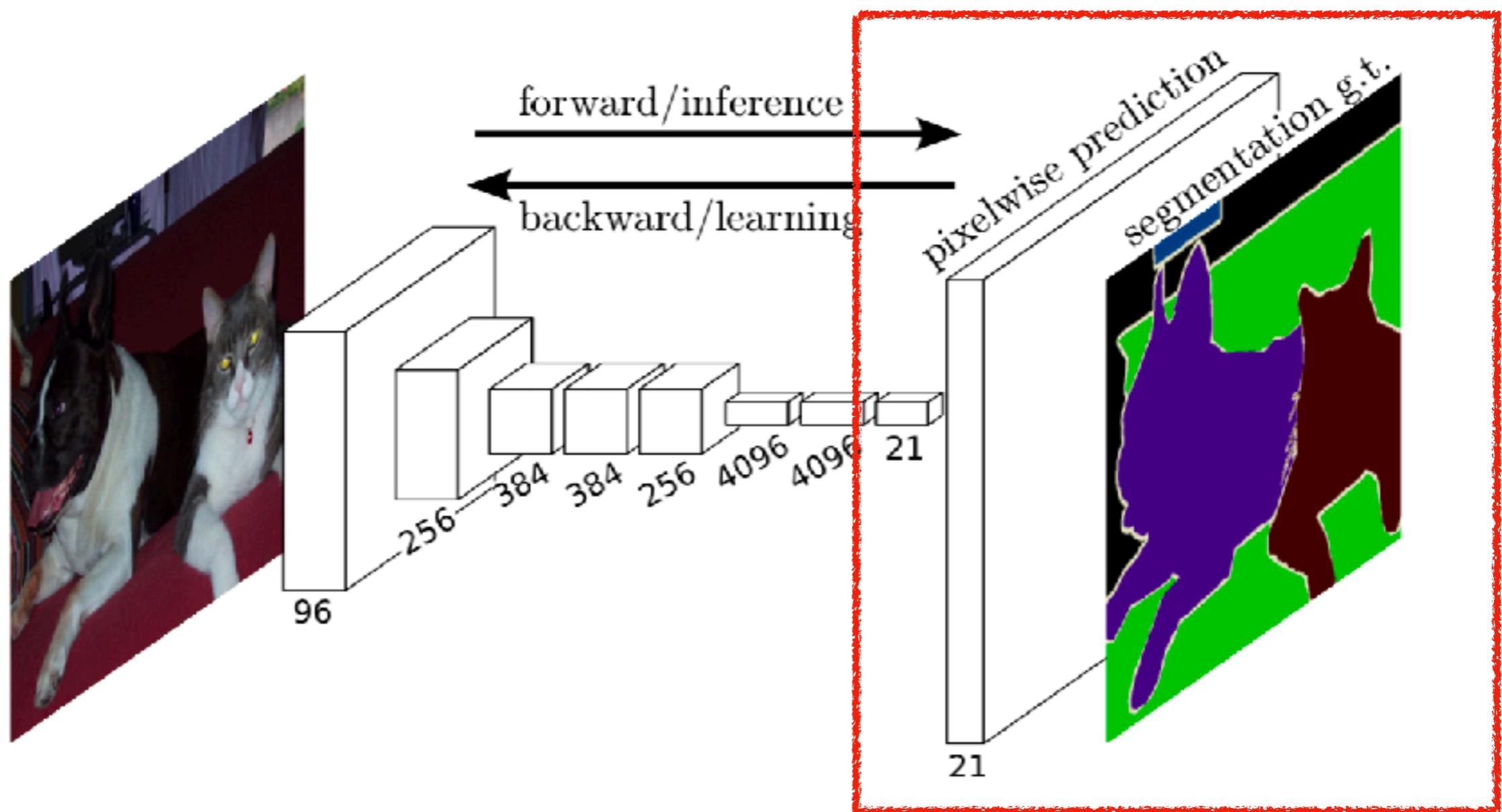
$H/8 \times W/8$

$H/16 \times W/16$

$H/32 \times W/32$

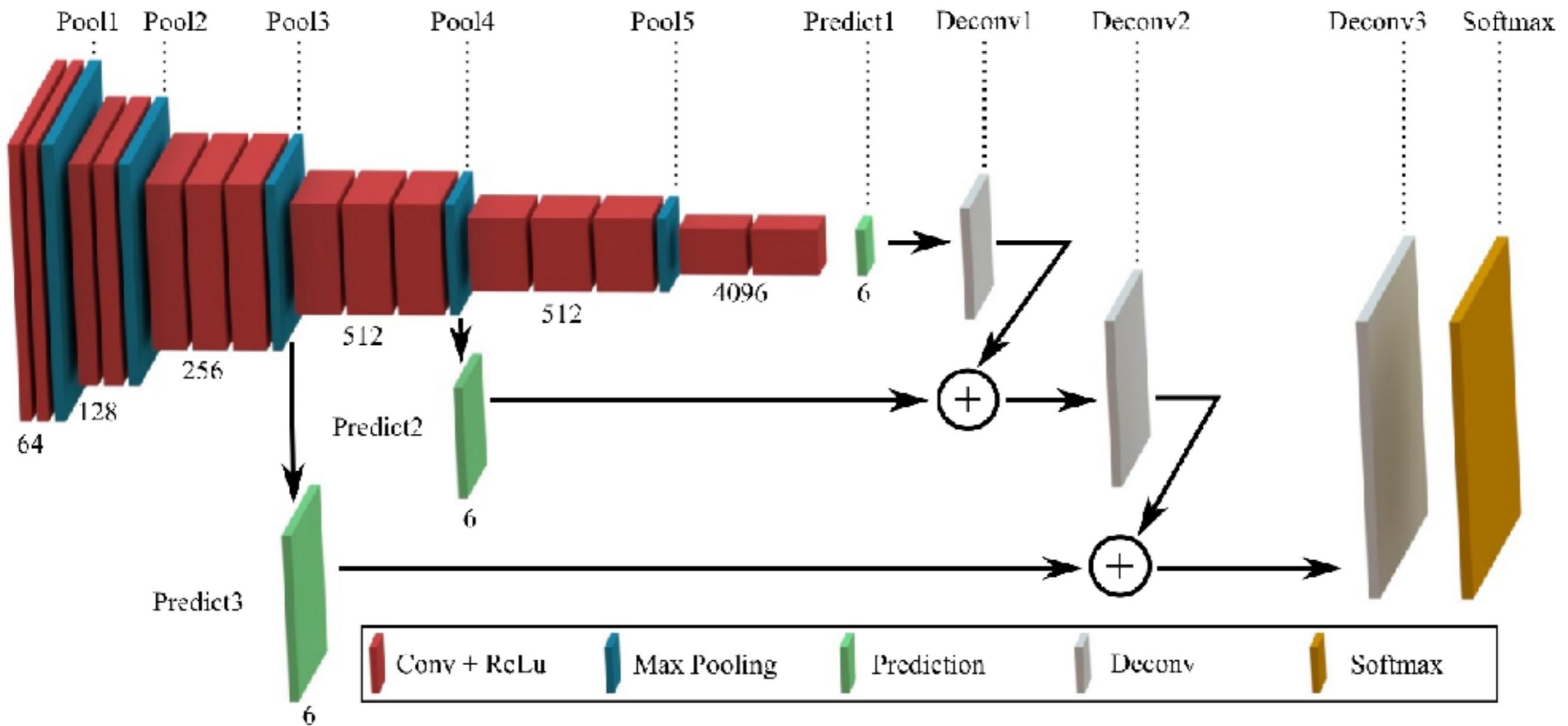
# **Up-Sampling**

# Coarse To Dense



**Do not use shift-and-stitch**

# Coarse To Dense



# Up-Sampling Problem

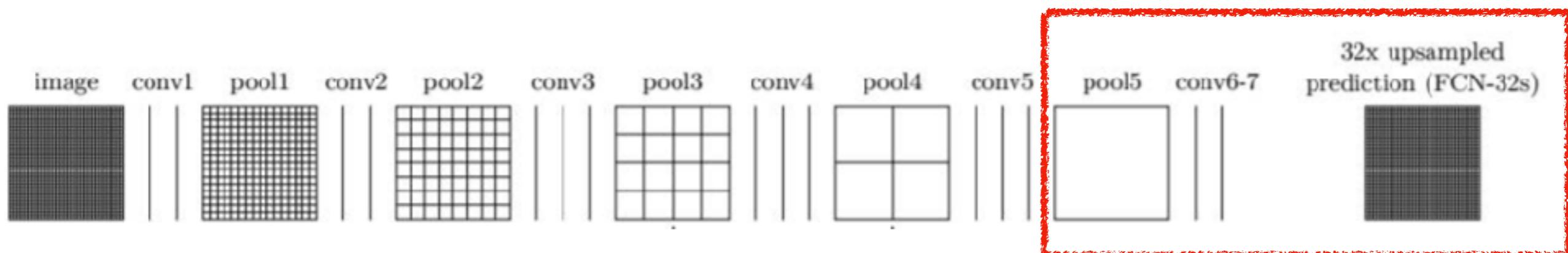
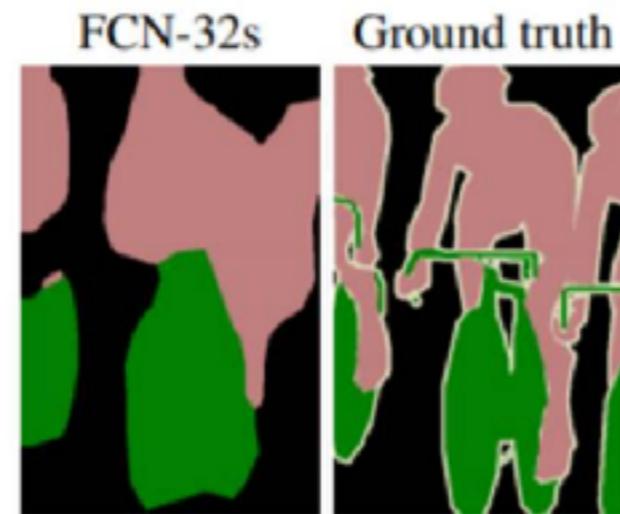


Figure 6. (Conventional) Fully Convolutional Networks : FCN-32s

**x32 Up-Sampling**



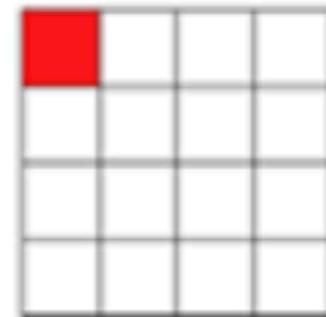
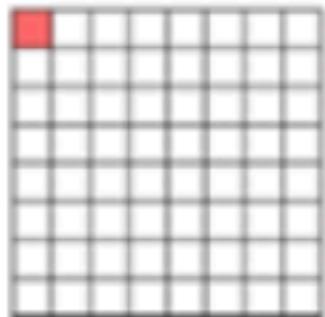
**Lose Detail**

# Deep jet

image

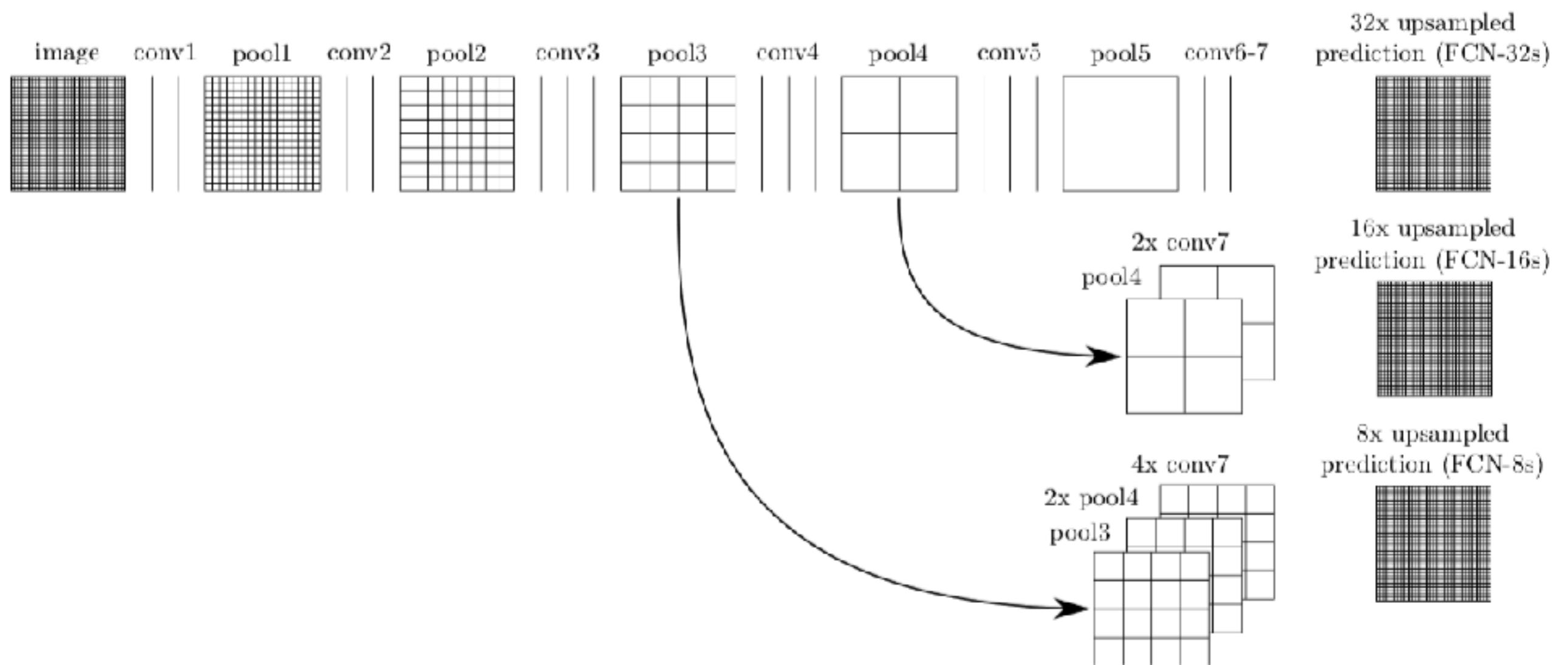


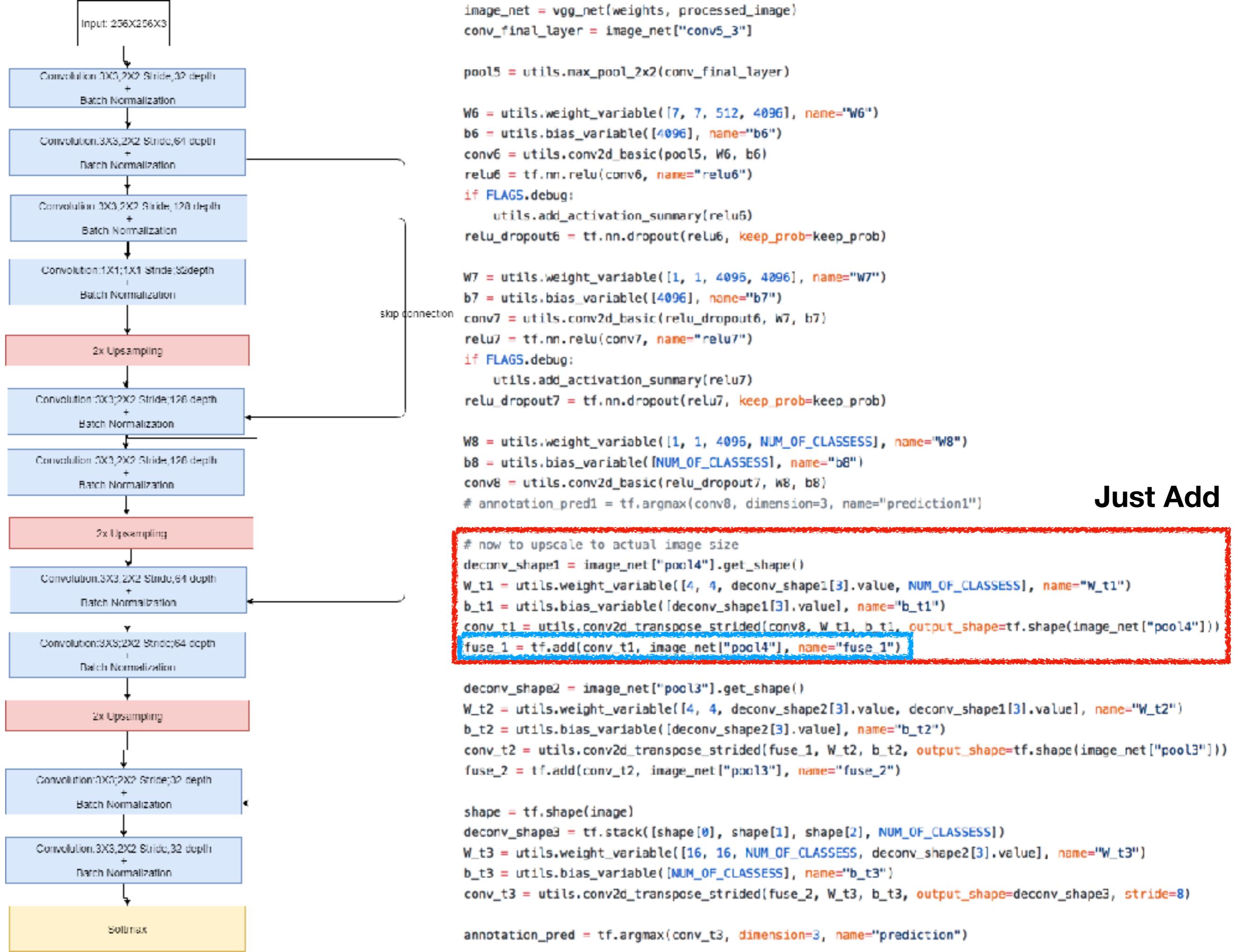
intermediate layers



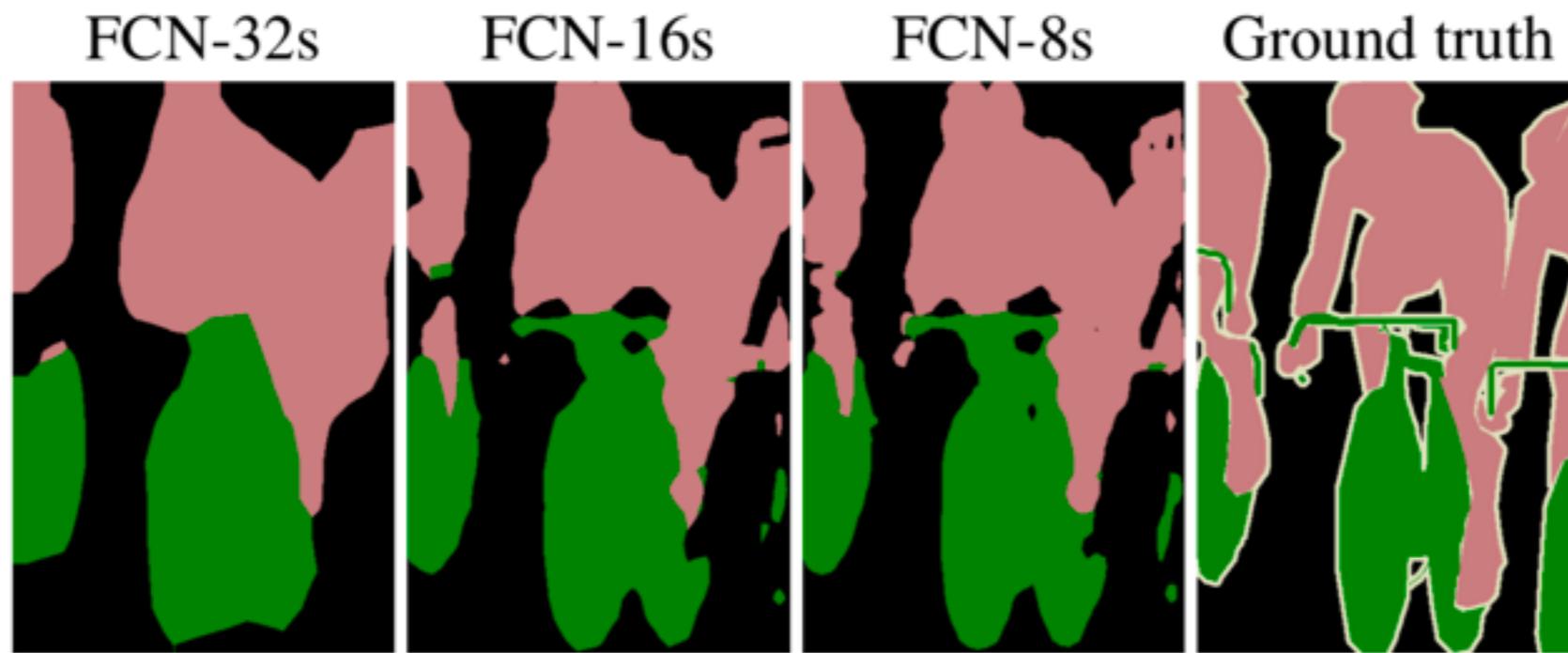
fuse features into **deep jet**

# Skip Layer





# Skip Layer



Pascal VOC 2011 Results

	pixel acc.	mean acc.	mean IU	f.w. IU
FCN-32s-fixed	83.0	59.7	45.4	72.0
FCN-32s	89.1	73.3	59.4	81.4
FCN-16s	90.0	75.7	62.4	83.0
FCN-8s	<b>90.3</b>	<b>75.9</b>	<b>62.7</b>	<b>83.2</b>

# **Results**

# Results

	FCN- AlexNet	FCN- VGG16	FCN- GoogLeNet <sup>4</sup>
mean IU	39.8	<b>56.0</b>	42.5
forward time	50 ms	210 ms	59 ms
conv. layers	8	16	22
parameters	57M	134M	6M
rf size	355	404	907
max stride	32	32	32

PASCAL VOC 2011

# Results

