

# Mask R-CNN

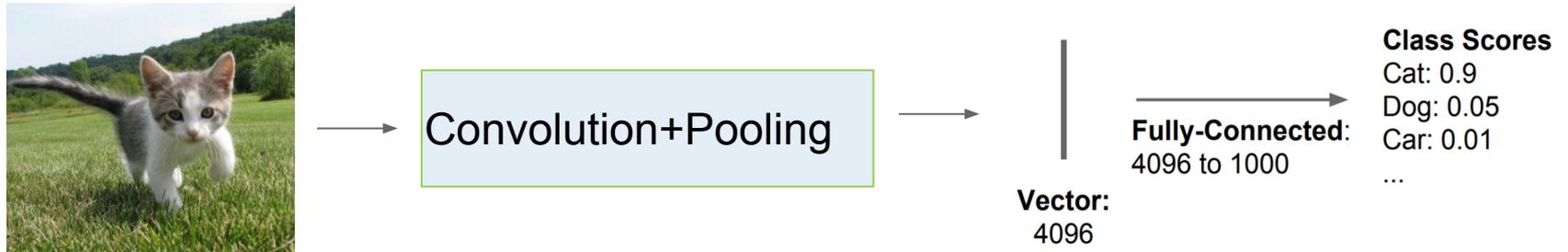
Kaiming He, Georgia, Gkioxari, Piotr Dollar, Ross Girshick

Presenters: Xiaokang Wang, Mengyao Shi

Feb. 13, 2018

# Common computer vision tasks

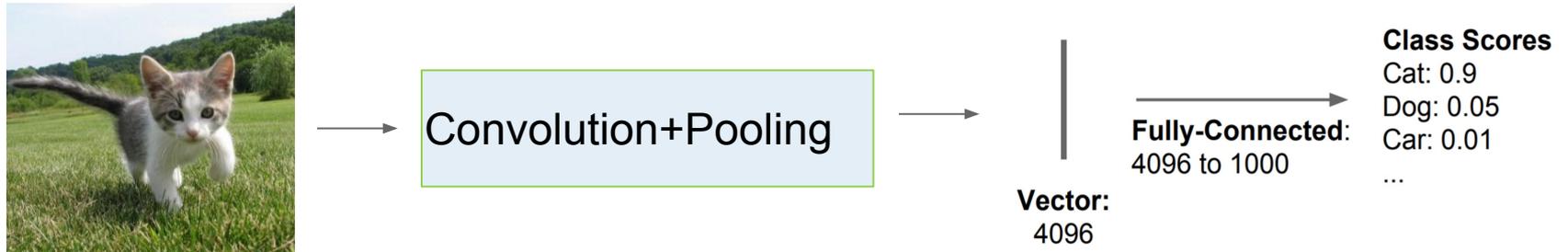
**Image Classification:** one label is generated for an image from the probability from a softmax function.



There is a cat in the image.

# Common computer vision tasks

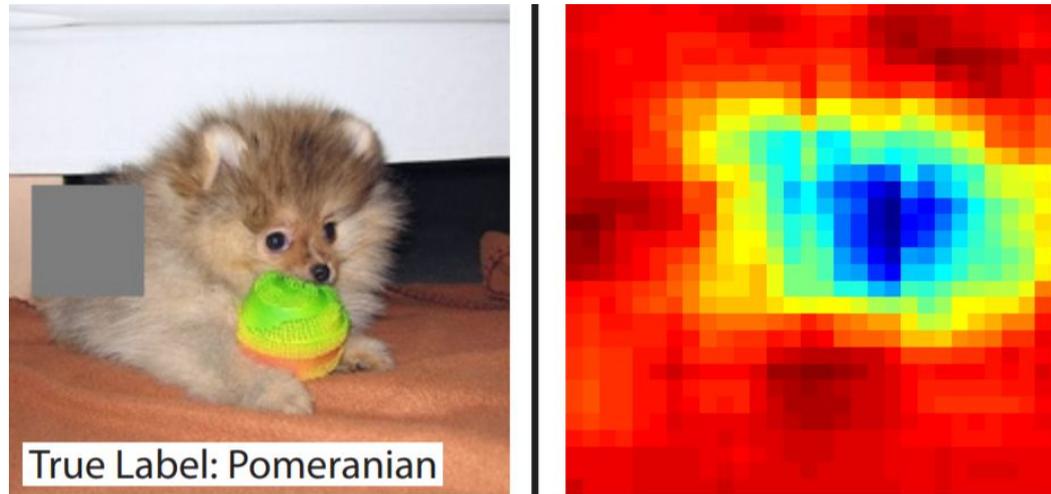
**Image Classification:** one label is generated for an image from the probability from a softmax function.



There is a cat in the image. **Where is the cat?**

# Common computer vision tasks

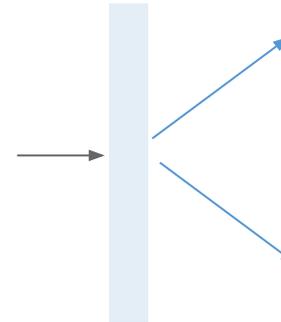
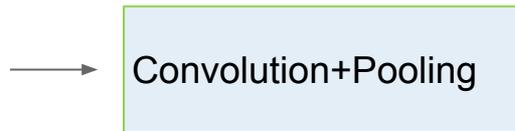
A trained neural network roughly knows where the object is.



Ref: Zeiler, Matthew D., and Rob Fergus. "Visualizing and understanding convolutional networks." *European conference on computer vision*. Springer, Cham, 2014.

# Common computer vision tasks

**Classification+Localization:** one label is generated for an image based on the probability from a softmax function.



**Class Scores**

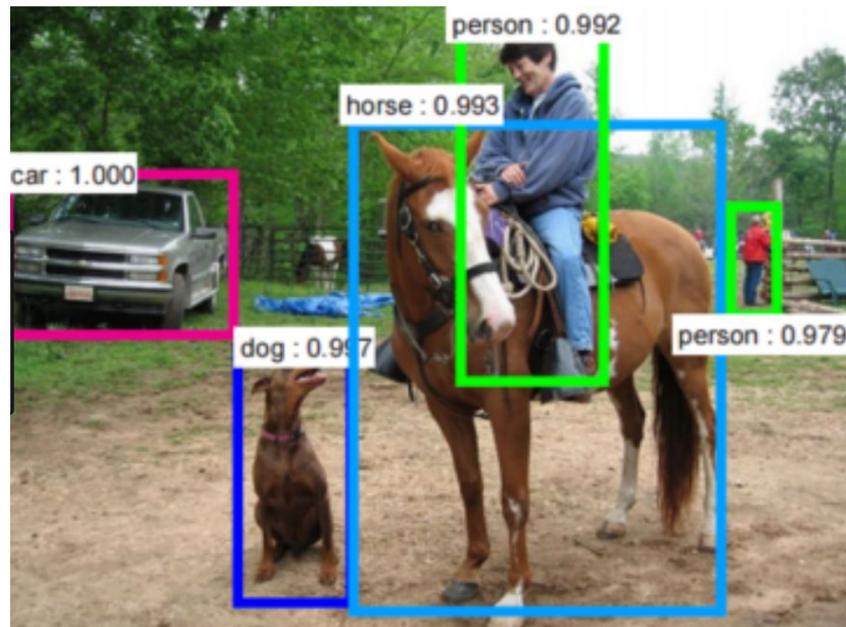
Cat: 0.9  
Dog: 0.05  
Car: 0.01  
...

**Center coordinate + height + weight:**

$x_c$ : 128  
 $y_c$ : 128  
H: 190  
W: 210

# Common computer vision tasks

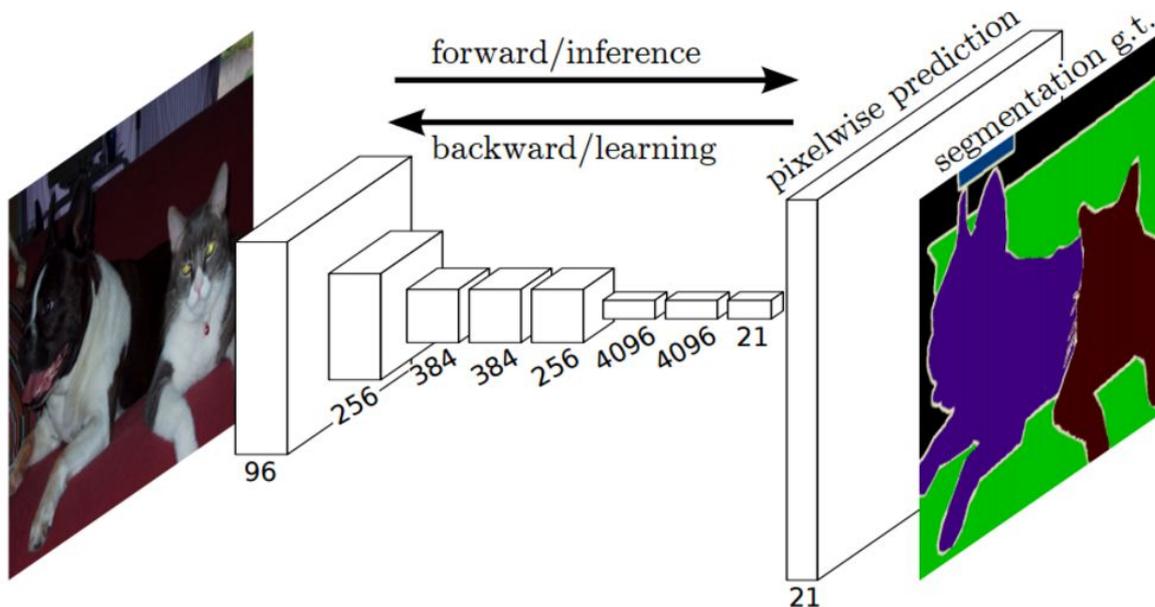
**Object detection:** detect the presence of objects and then **label** and **localize** them.



# Semantic segmentation

Upsampling the features to the same width and height as the input image.

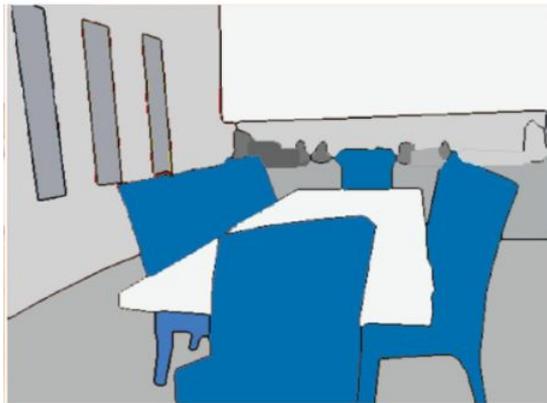
1. **No fully connected layers** in such an architecture
2. **1×1 convolution** is used to collapse all the channels of a feature volume into one channel



# Semantic segmentation vs instance segmentation



Input Image



Semantic Segmentation



Semantic Instance Segmentation

Source;

<https://stackoverflow.com/questions/33947823/what-is-semantic-segmentation-compared-to-segmentation-and-scene-labeling>

# Sematic segmentation vs instance segmentation

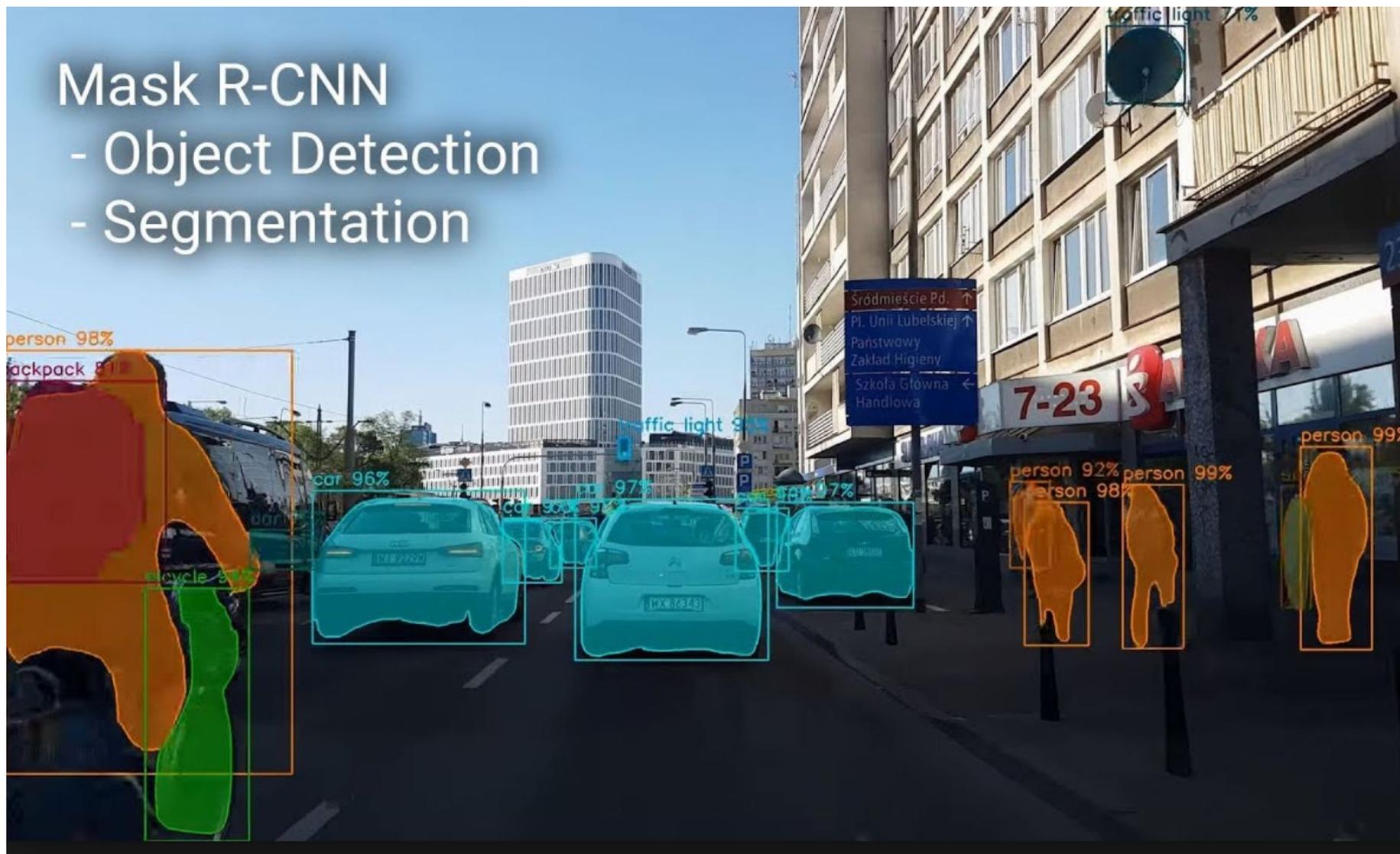


**One pixel** can have **different semantic meaning** in different regions

# What can Mask-RCNN do?

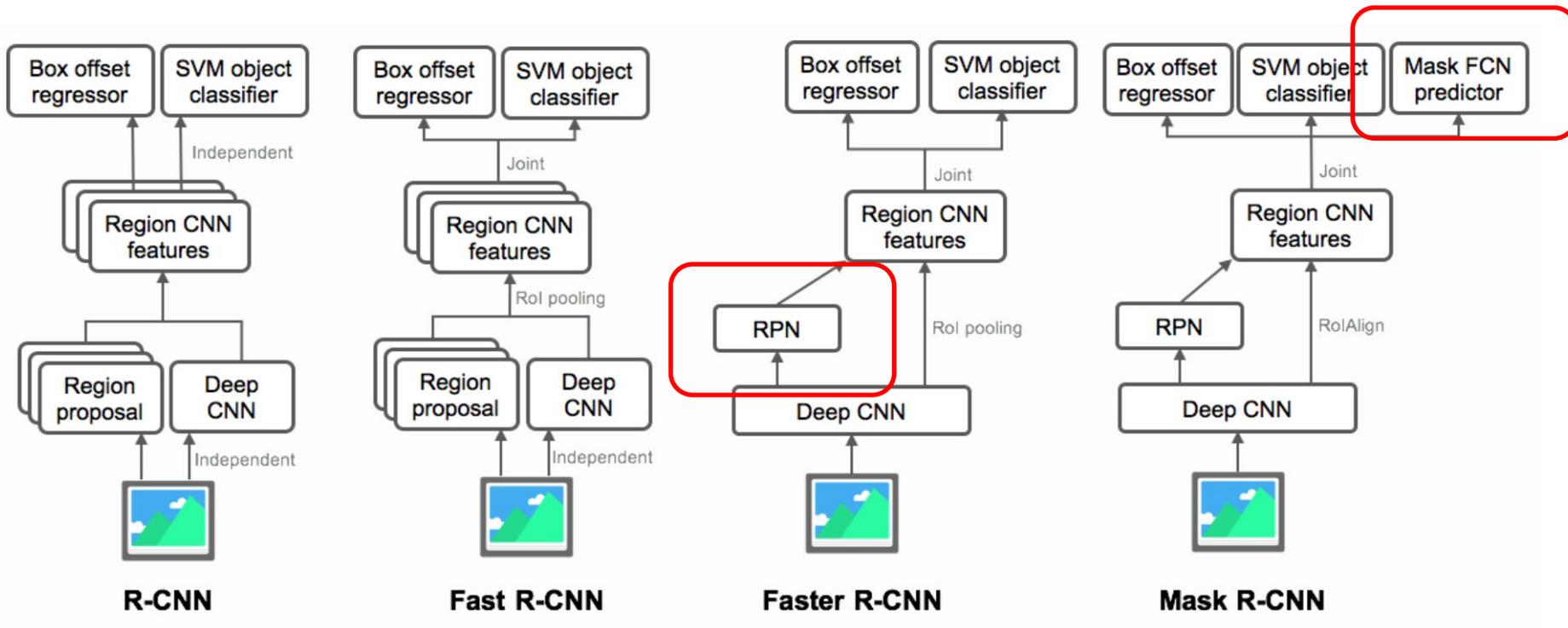
## Mask R-CNN

- Object Detection
- Segmentation



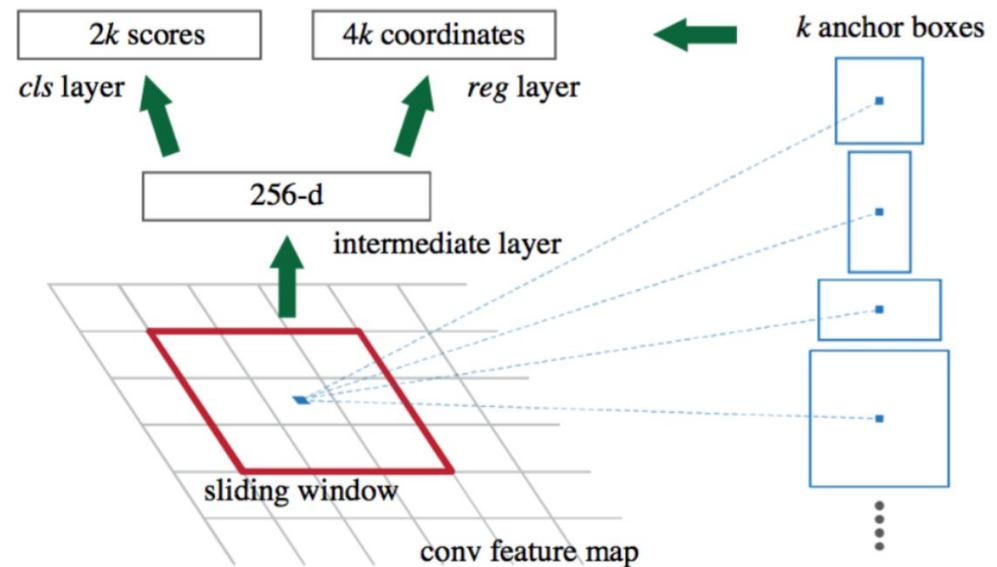
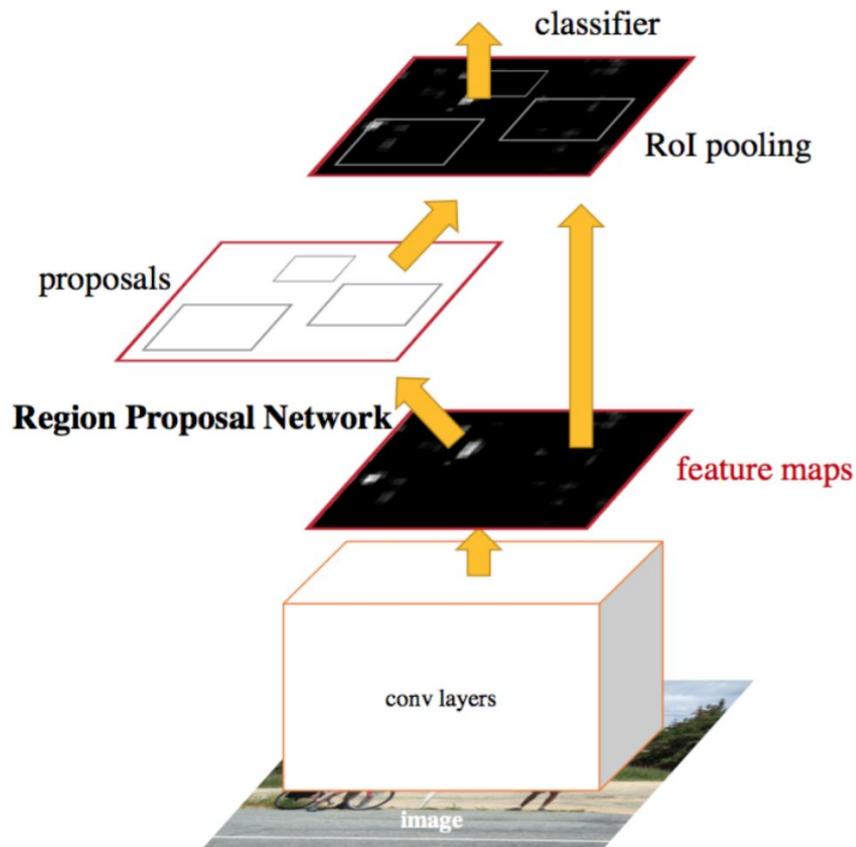
# How does Mask-rcnn work at a high level?

An FCN on Rols is added to **Faster-rcnn**



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

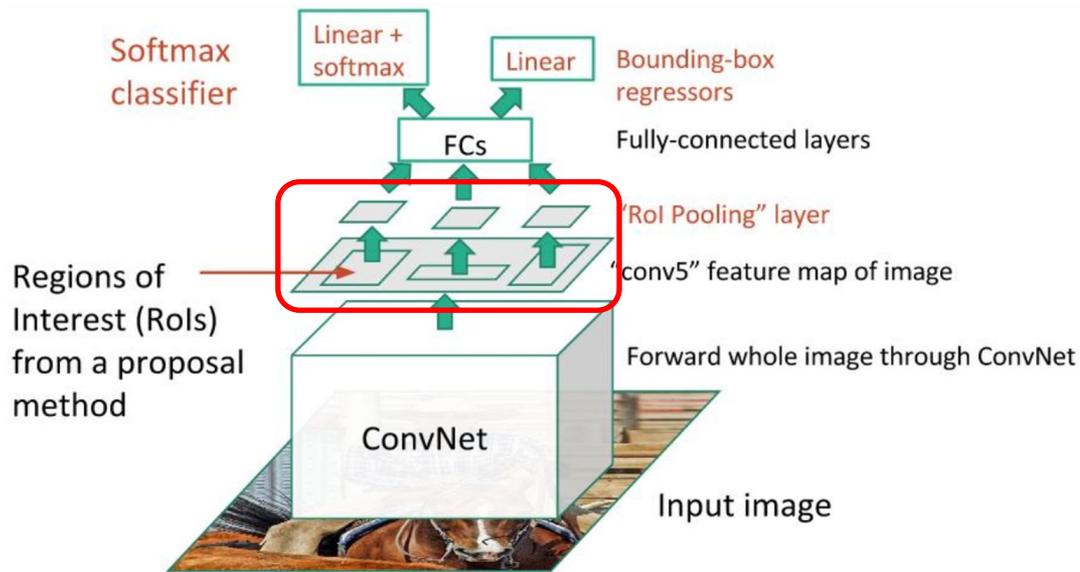
# Faster-rcnn = Fast-rcnn + RPN



Ref: Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. [“Faster R-CNN: Towards real-time object detection with region proposal networks.”](#)

# RoIPool layer in fast-rcnn

**RoI pooling layer** uses max pooling to **covert** the features inside **any valid region of interest into a small feature map with a predefined size.**



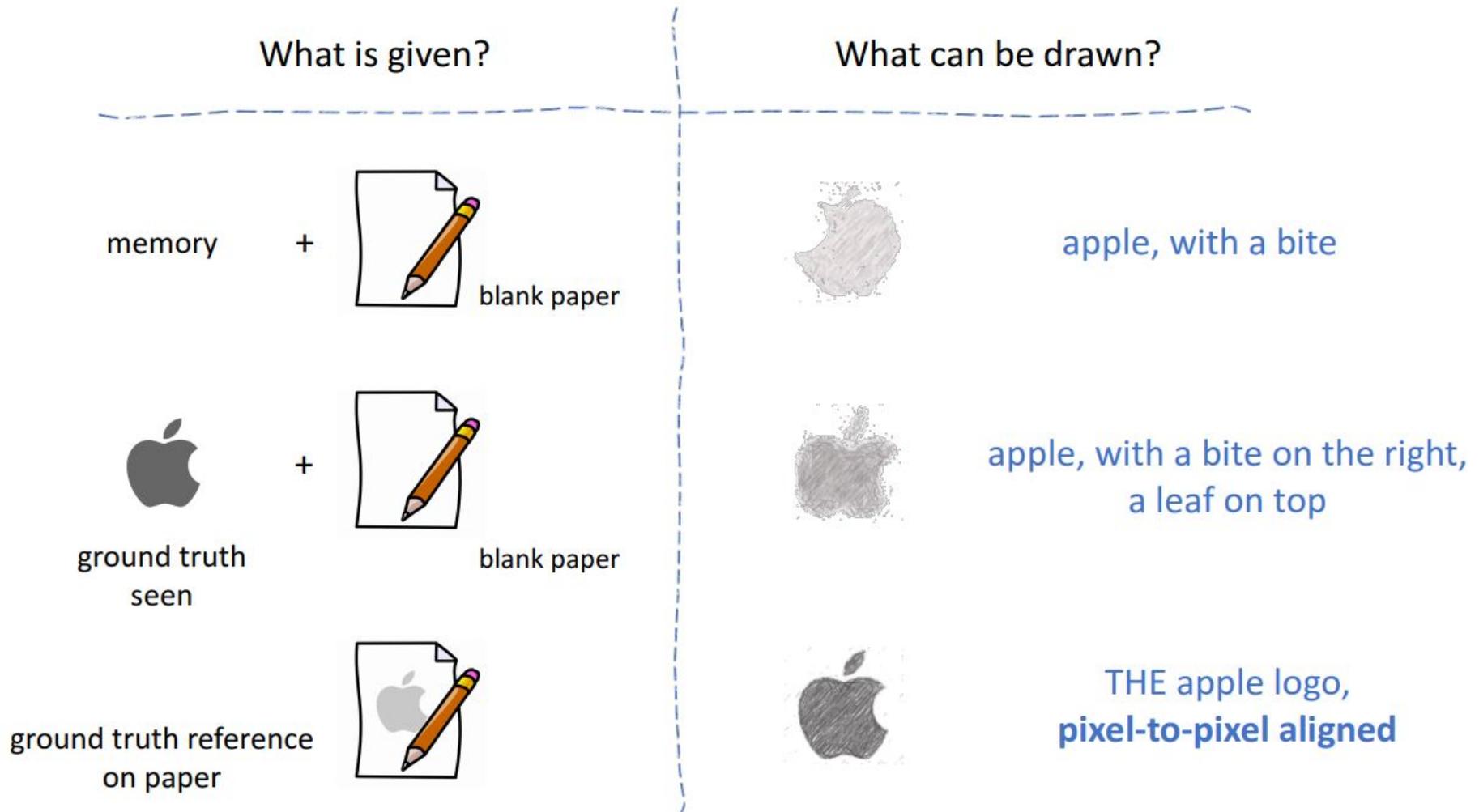
0.88	0.44	0.14	0.16	0.37	0.77	0.96	0.27
0.19	0.45	0.57	0.16	0.63	0.29	0.71	0.70
0.66	0.26	0.82	0.64	0.54	0.73	0.59	0.26
0.85	0.34	0.76	0.84	0.29	0.75	0.62	0.25
0.32	0.74	0.21	0.39	0.34	0.03	0.33	0.48
0.20	0.14	0.16	0.13	0.73	0.65	0.96	0.32
0.19	0.69	0.09	0.86	0.88	0.07	0.01	0.48
0.83	0.24	0.97	0.04	0.24	0.35	0.50	0.91

RoI pooling

0.85	0.84
0.97	0.96

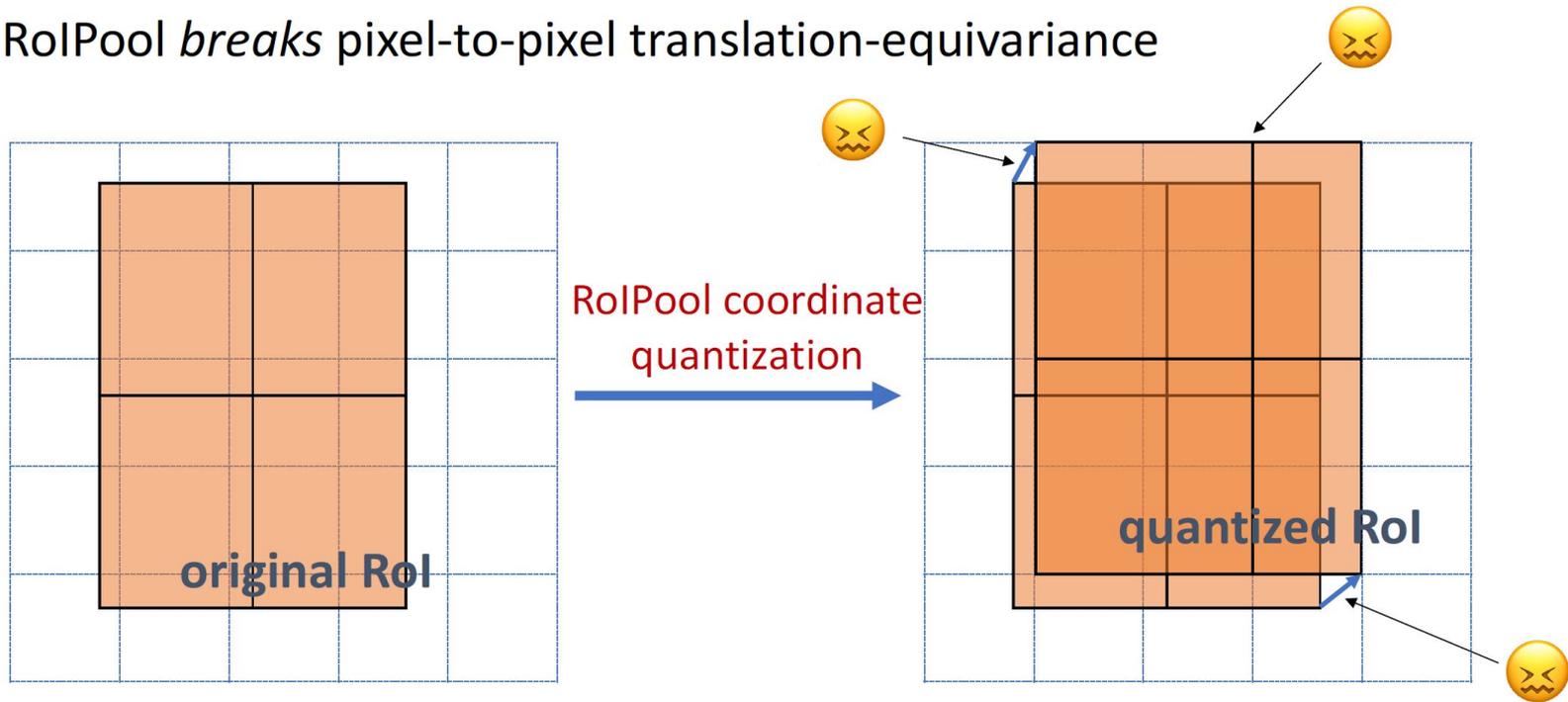
**Divide the  $h \times w$  RoI window into a  $H \times W$  grid of subwindows and then do do max-pooling in each sub-window**

# Best way to draw an apple?



# Drawbacks of RoI pooling for segmentation task

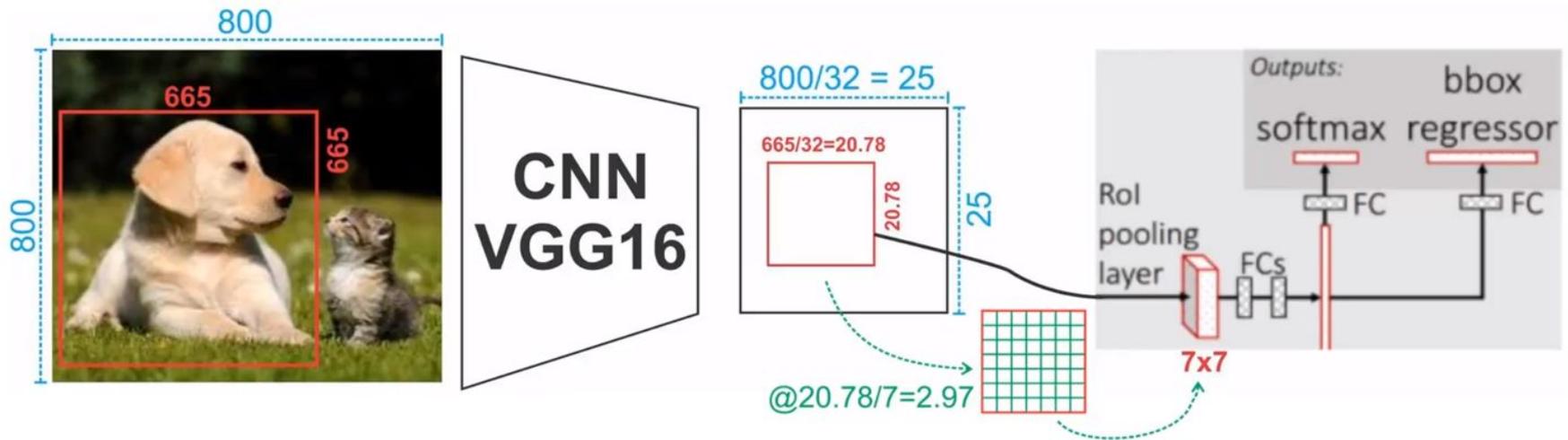
- RoIPool *breaks* pixel-to-pixel translation-equivariance



**Translation equivariance:** shifts in input lead to corresponding changes in the feature map

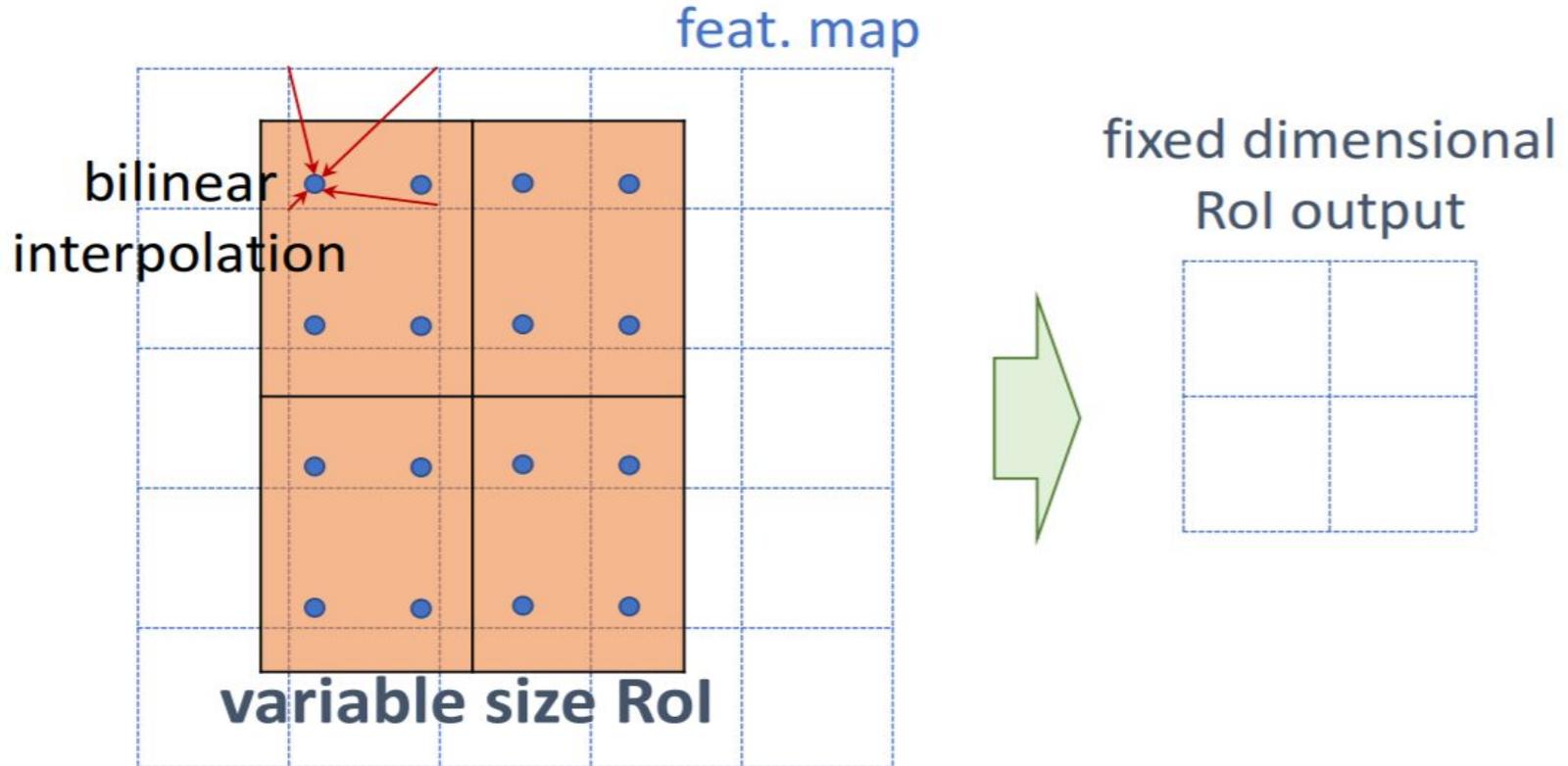
# One pixel in RoI means many pixels in the raw image

The **misalignment** introduced by **Quantization** associated with **RoIPool** matters for a **large stride**



# Solution: RoI align

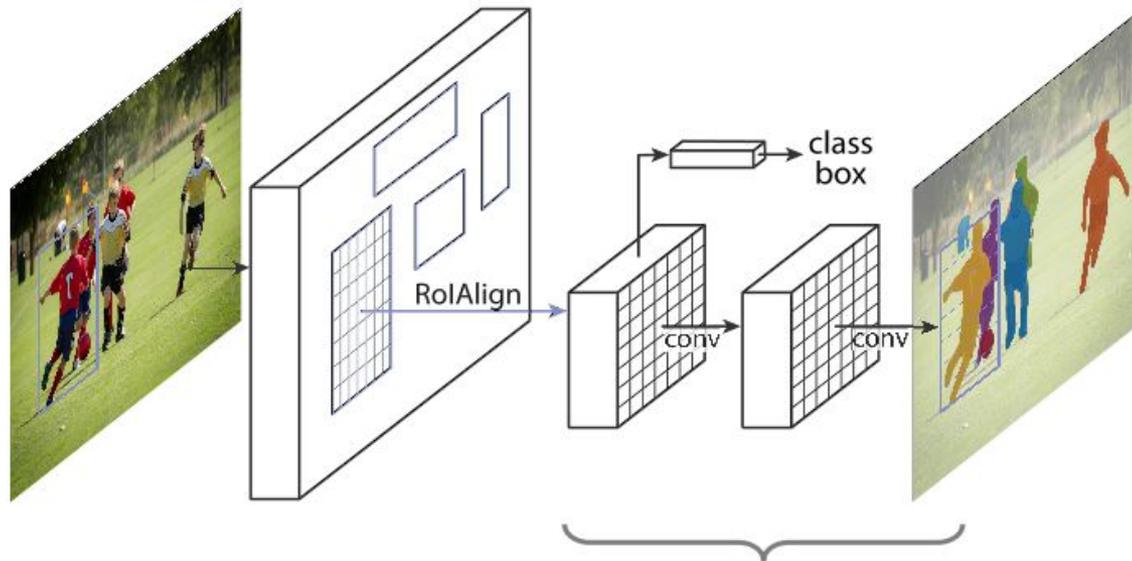
## Property 1: Pixel-to-pixel alignment



[https://en.wikipedia.org/wiki/Bilinear\\_interpolation](https://en.wikipedia.org/wiki/Bilinear_interpolation)

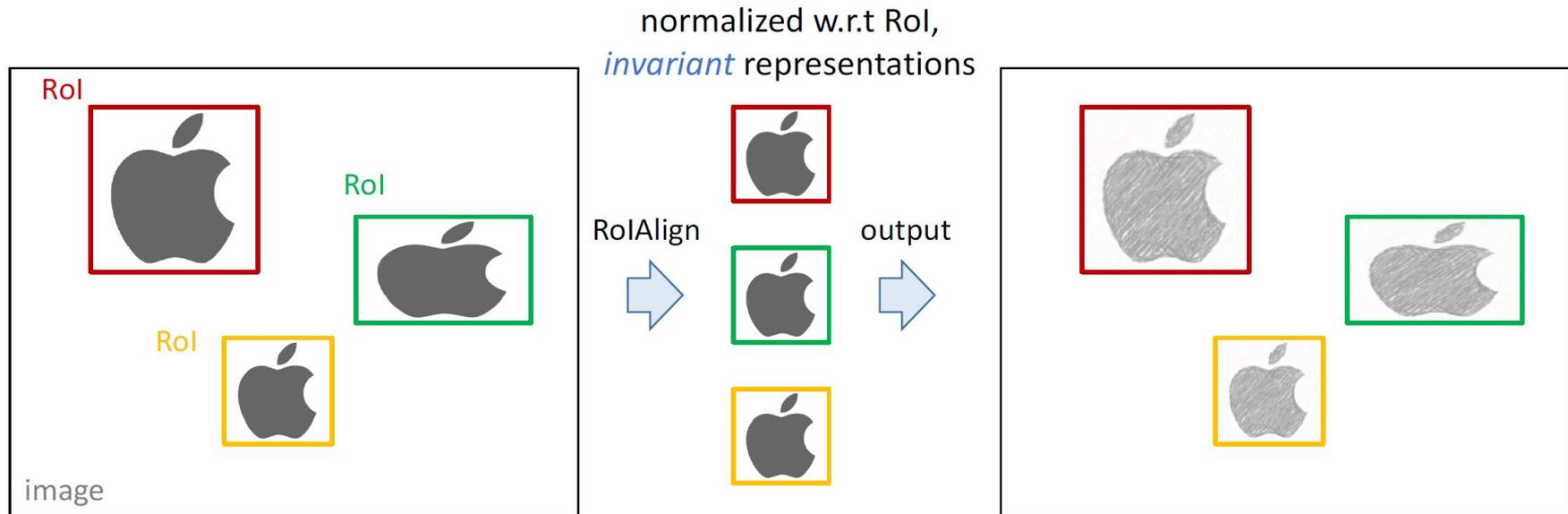
## Property 2: equivariance

### Equivariance in Mask R-CNN



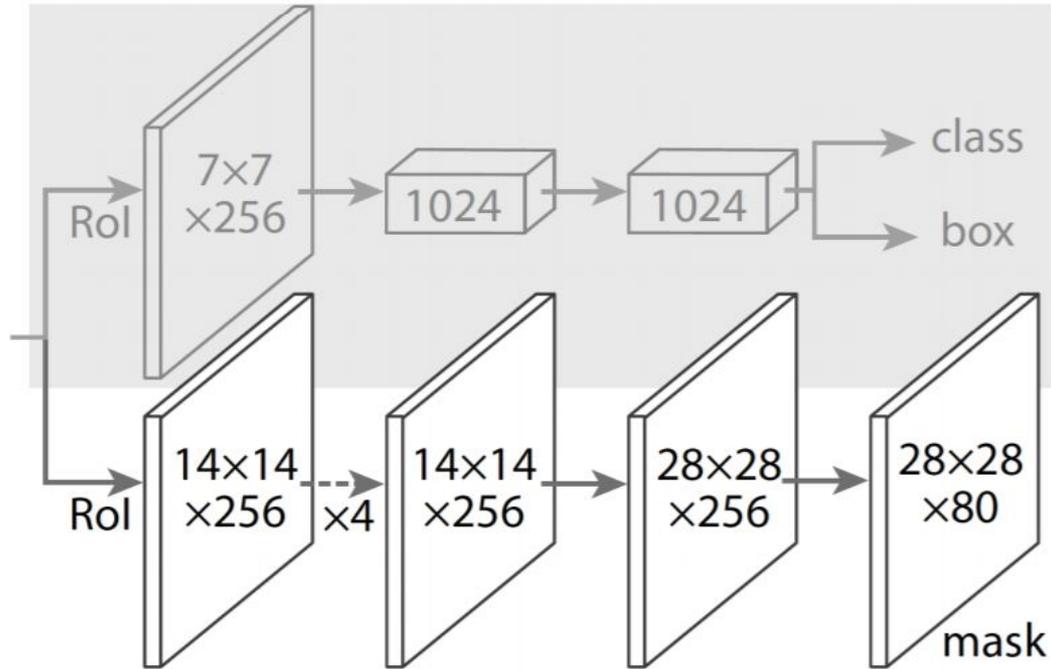
2. Fully-Conv on RoI:  
equivariant to translation within RoI

## Property 2: Scale-invariant



- RoIAlign creates *scale-invariant* representations
- RoIAlign + “output pasted back” provides *scale-equivariance*

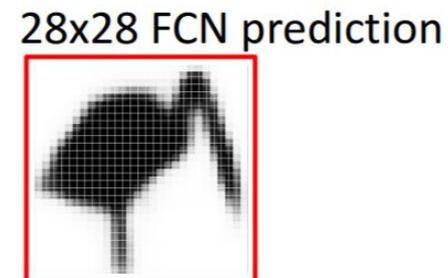
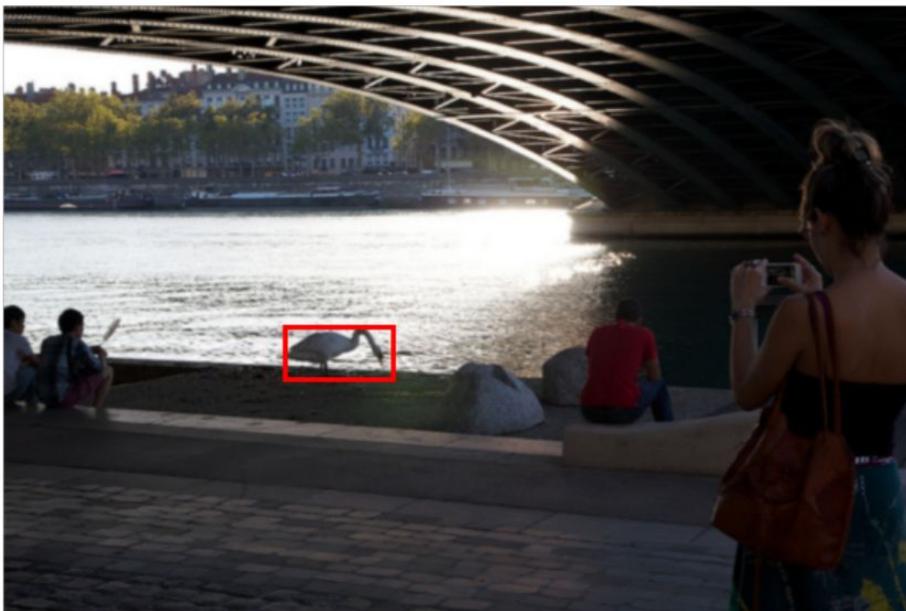
# FCN Mask head



Per class **binary masks** are generated rather than a **multinomial mask**. No class competition.

# How is a mask generated from FCN head?

- Pixel-to-pixel aligned



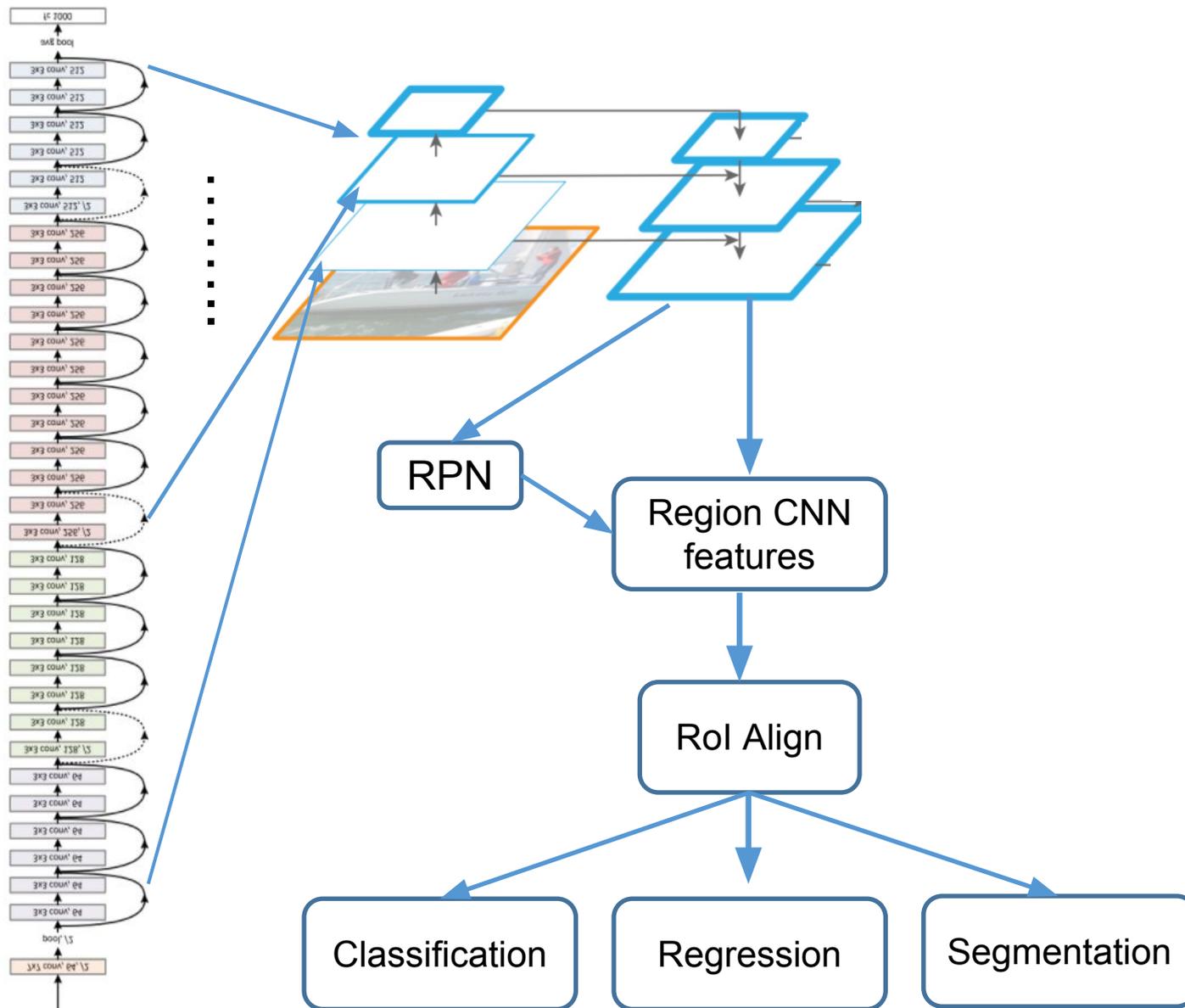
resized soft prediction



final mask

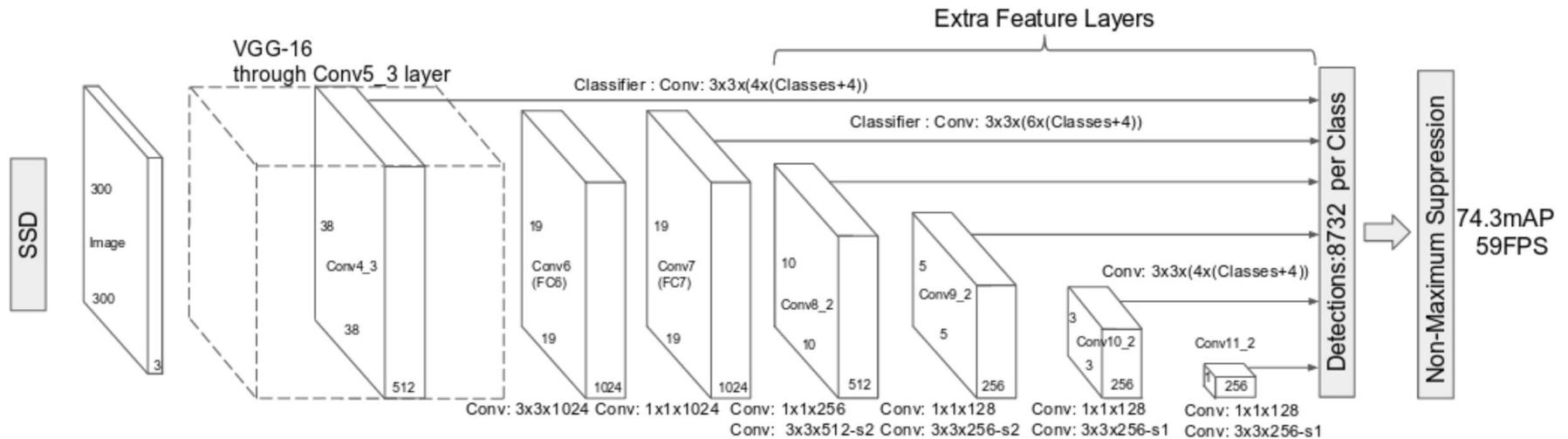


# Backbone of Mask-rcnn: FPN

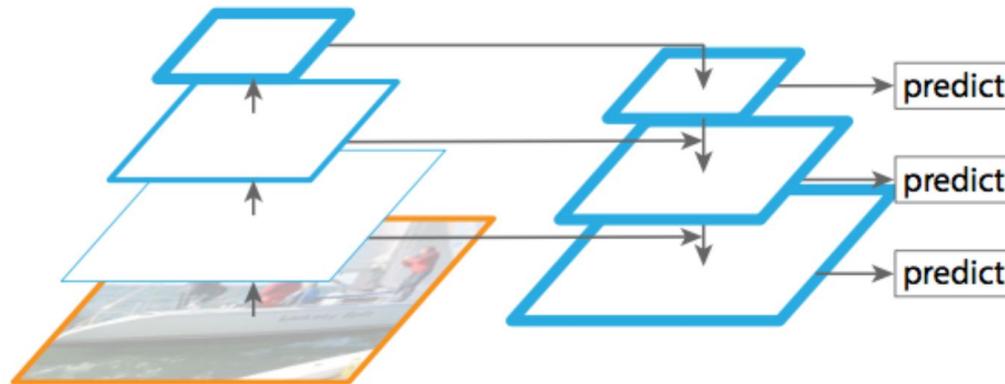


# Backbone of Mask-rcnn: FPN

## SSD: skip connection



## FPN: lateral connection

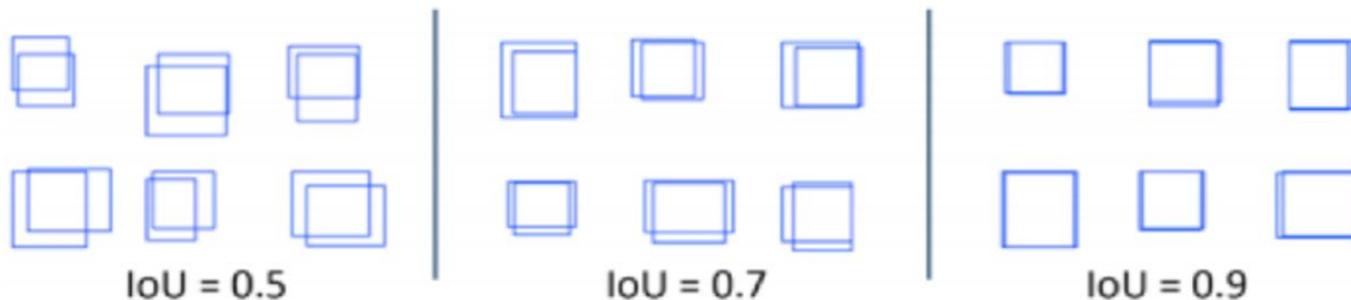


# Experiments

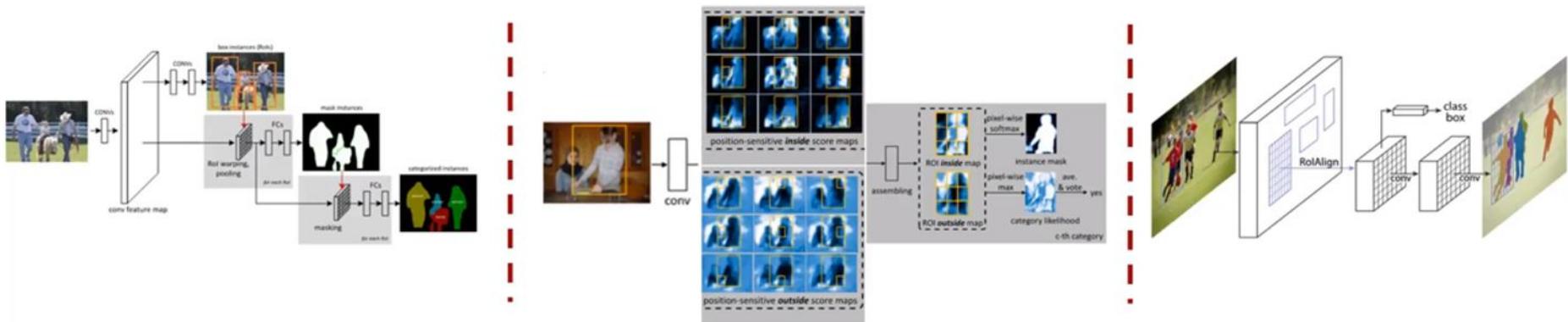
- Main dataset: MS COCO
  - 80 classes
  - 80k train image
  - 35k sub- set of val images
  - 5k images for ablation experiments
- Metric

## Challenges Score: AP

- AP is averaged over multiple IoU values between 0.5 and 0.95.
- More comprehensive metric than the traditional AP at a fixed IoU value (0.5 for PASCAL).



# State-of-art in 2016: FCIS



	MNC	FCIS	Mask R-CNN
Mask	All classes are in one mask (translation-invar.)	Position sensitive mask (translation-var.) [1]	Binary mask for each object
Feature map	RoI-Warping	?	RoIAlign
Input classification	Depend mask	Join mask pred. & class.	Parallel-independent

# Main results: Mask R-CNN vs MNC and FCIS

	backbone	AP	AP <sub>50</sub>	AP <sub>75</sub>	AP <sub>S</sub>	AP <sub>M</sub>	AP <sub>L</sub>
MNC [10]	ResNet-101-C4	24.6	44.3	24.8	4.7	25.9	43.6
FCIS [26] +OHEM	ResNet-101-C5-dilated	29.2	49.5	-	7.1	31.3	50.0
FCIS+++ [26] +OHEM	ResNet-101-C5-dilated	33.6	54.5	-	-	-	-
<b>Mask R-CNN</b>	ResNet-101-C4	33.1	54.9	34.8	12.1	35.6	51.1
<b>Mask R-CNN</b>	ResNet-101-FPN	35.7	58.0	37.8	15.5	38.1	52.4
<b>Mask R-CNN</b>	ResNeXt-101-FPN	<b>37.1</b>	<b>60.0</b>	<b>39.4</b>	<b>16.9</b>	<b>39.9</b>	<b>53.5</b>

# Ablation: different backbone 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>

(a) **Backbone Architecture:** Better backbones bring expected gains: deeper networks do better, FPN outperforms C4 features, and ResNeXt improves on ResNet.

More efficient encoder and decoder help generally.

# Ablation: Binary vs Multinomial Loss

	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

(b) **Multinomial vs. Independent Masks**  
(ResNet-50-C4): *Decoupling* via per-class binary masks (sigmoid) gives large gains over multinomial masks (softmax).

# Ablation: RoIAlign vs RoPool

baseline: ResNet-50-Conv5 backbone, **stride=32**

	mask AP			box AP		
	AP	AP <sub>50</sub>	AP <sub>75</sub>	AP <sup>bb</sup>	AP <sup>bb</sup> <sub>50</sub>	AP <sup>bb</sup> <sub>75</sub>
<i>RoIPool</i>	23.6	46.5	21.6	28.2	52.7	26.9
<i>RoIAlign</i>	<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	<b>+10.5</b>	+5.8	+2.6	+9.5

- huge gain at high IoU, in case of big stride (32)

# Bounding box detection results are improved by RoIAlign

ResNeXt-101-FPN on COCO.

	backbone	AP <sup>bb</sup>	AP <sub>50</sub> <sup>bb</sup>
Faster R-CNN+++ [15]	ResNet-101-C4	34.9	55.7
Faster R-CNN w FPN [22]	ResNet-101-FPN	36.2	59.1
Faster R-CNN by G-RMI [17]	Inception-ResNet-v2 [32]	34.7	55.5
Faster R-CNN w TDM [31]	Inception-ResNet-v2-TDM	36.8	57.7
Faster R-CNN, RoIAlign	ResNet-101-FPN	37.3	59.6

bbox improved by:

- RoIAlign

# Bounding box detection benefits from mult-task learning

Bounding box detection benefits from the segmentation task

Backbone: ResNeXt-101-FPN on COCO.

	backbone	AP <sup>bb</sup>	AP <sub>50</sub> <sup>bb</sup>
Faster R-CNN+++ [15]	ResNet-101-C4	34.9	55.7
Faster R-CNN w FPN [22]	ResNet-101-FPN	36.2	59.1
Faster R-CNN by G-RMI [17]	Inception-ResNet-v2 [32]	34.7	55.5
Faster R-CNN w TDM [31]	Inception-ResNet-v2-TDM	36.8	57.7
Faster R-CNN, RoIAlign	ResNet-101-FPN	37.3	59.6
<b>Mask R-CNN</b>	ResNet-101-FPN	38.2	60.3

bbox improved by:

- RoIAlign
- Multi-task training w/ mask

# Qualitative Results on Instance Segmentation

<https://www.youtube.com/watch?v=g7z4mkfRjI4&t=639s>

# Generality of Mask-rcnn: Human keypoint detection

- keypoint = 1-hot mask
- human pose = 17 masks

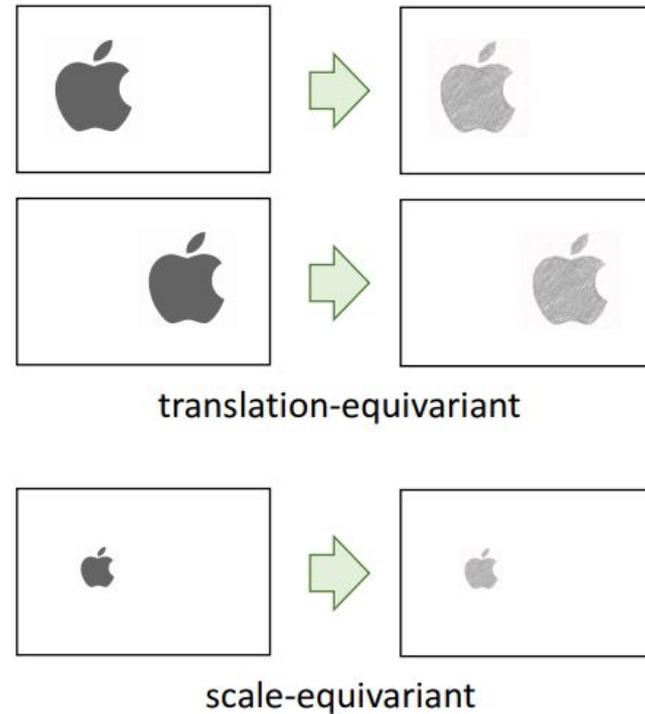
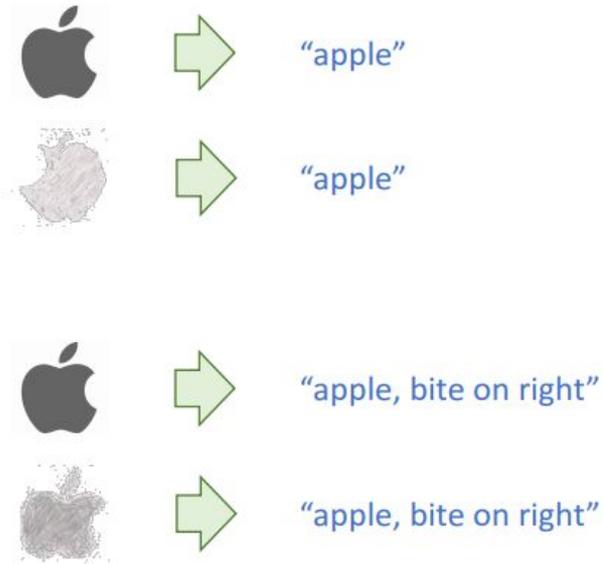


	$AP^{kp}$	$AP_{50}^{kp}$	$AP_{75}^{kp}$	$AP_M^{kp}$	$AP_L^{kp}$
CMU-Pose+++ [6]	61.8	84.9	67.5	57.1	68.2
G-RMI [32] <sup>†</sup>	62.4	84.0	68.5	<b>59.1</b>	68.1
<b>Mask R-CNN, keypoint-only</b>	62.7	87.0	68.4	57.4	71.1
<b>Mask R-CNN, keypoint &amp; mask</b>	<b>63.1</b>	<b>87.3</b>	<b>68.7</b>	57.8	<b>71.4</b>

**Thank you!**

Supplementary slides follow

# Invariance vs. Equivariance



see also “What is wrong with convolutional neural nets?”, Geoffrey Hinton, 2017

# Invariance vs equivariance

**Translation** means that each point/pixel in the image has been moved the same amount in the same direction



Convolution operation is translation equivariant  
Maxpooling is translation invariance

# Method for upsampling in segmentation task

1. **Uppooling**
2. **Interpolation**
3. **deconvolution**

0.1	0.5	<b>1.2</b>	-0.7
<b>0.8</b>	-0.2	-0.5	0.3
0.4	<b>0.9</b>	-0.1	-0.2
-0.6	0.1	<b>0.5</b>	0.3

max-pooling

0.8	1.2
0.9	0.5

		x	
x			
	x		
		x	

max locations

0	0	<b>0.5</b>	0
<b>1.3</b>	0	0	0
0	<b>0.4</b>	0	0
0	0	<b>0.1</b>	0

unpooling

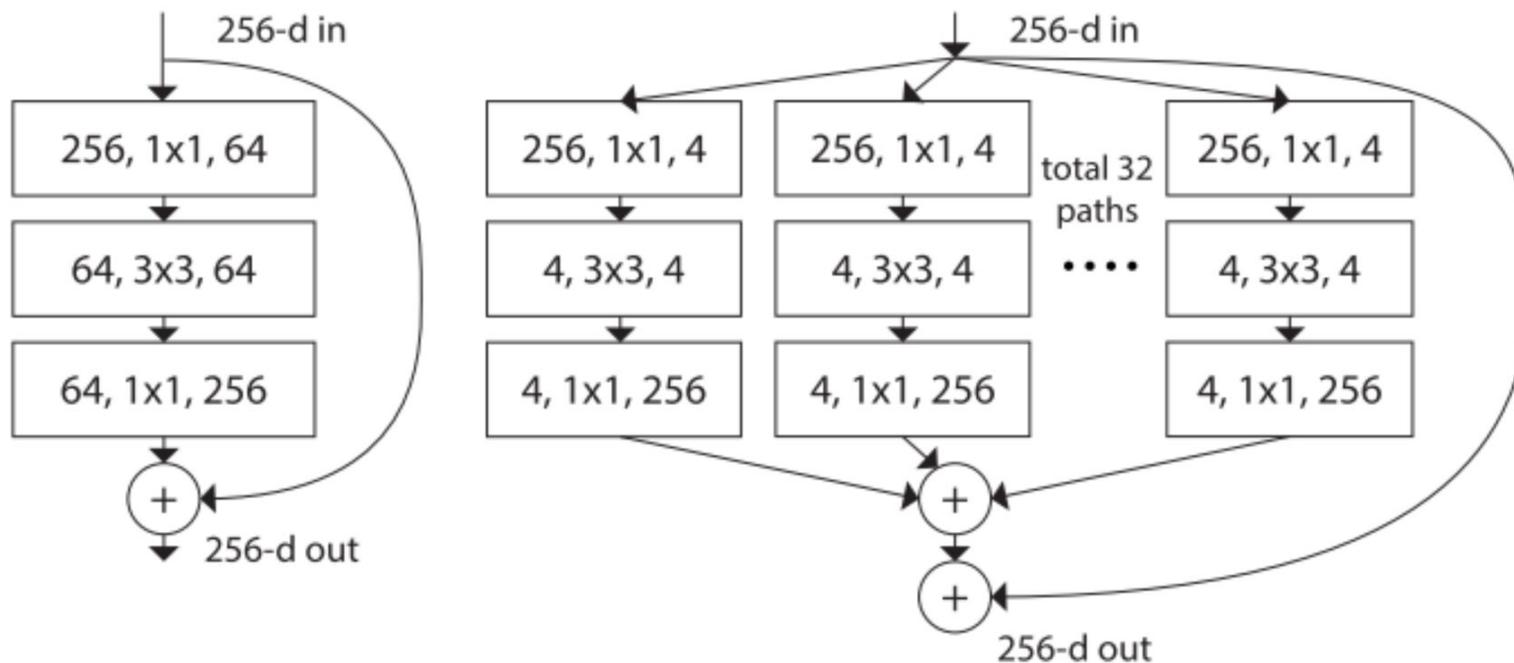
1.3	0.5
0.4	0.1

4 days ago

Ref:

# ResNeXt: what is new compared to Resnet?

Similar to Inception net but **share the same topology** among the multiple paths and use **summation rather than concatenation** to merge



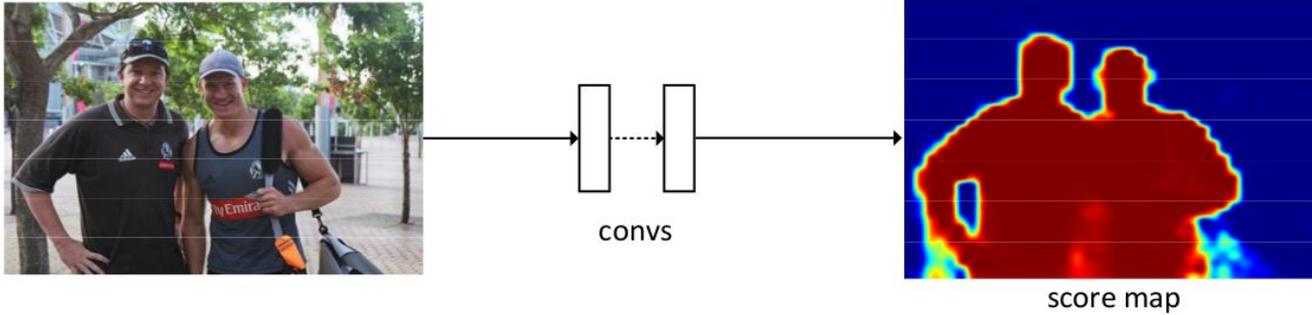
**Cardinality**—the number of independent paths

**Width:** the number of channels in output.

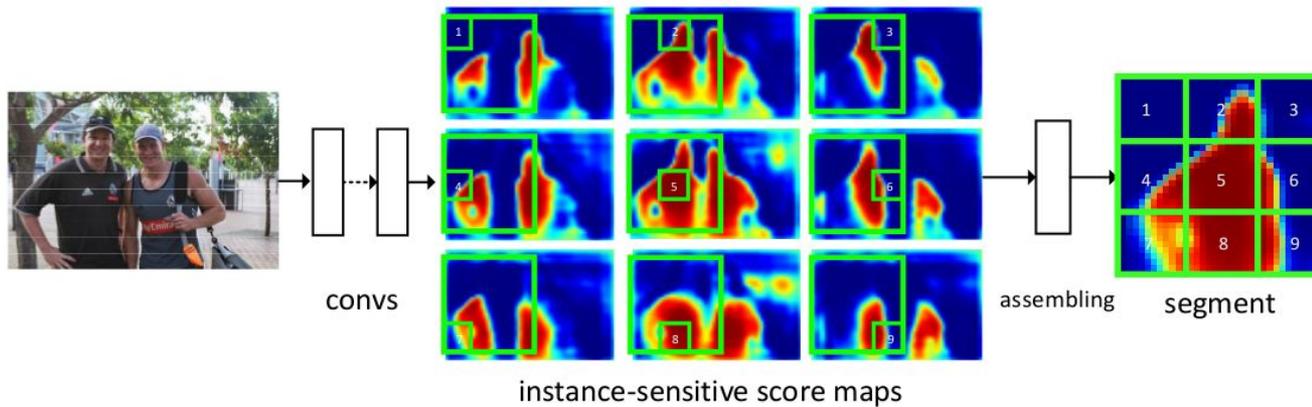
Note: convolution is 3d computation.

# Position sensitive score map

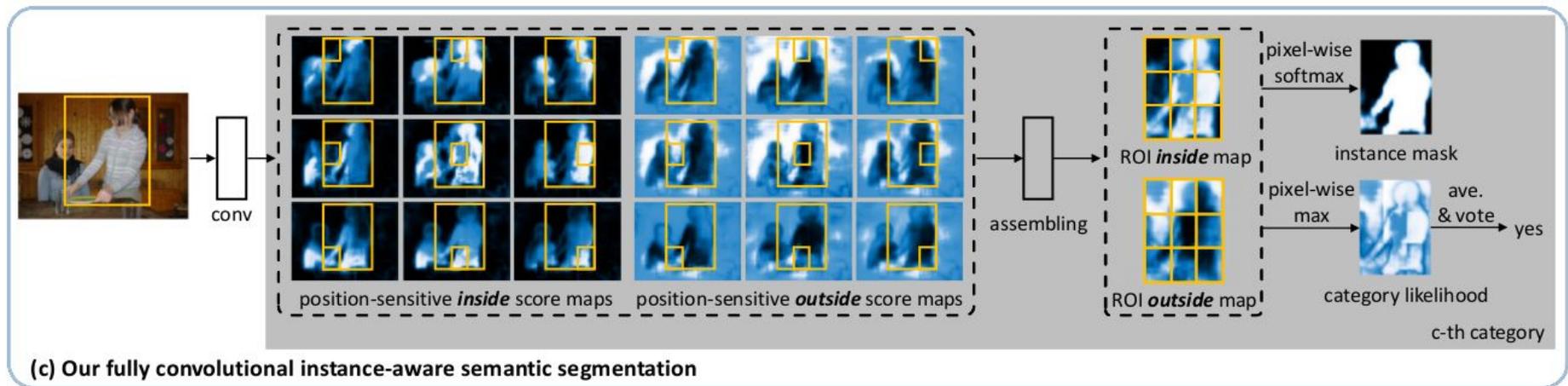
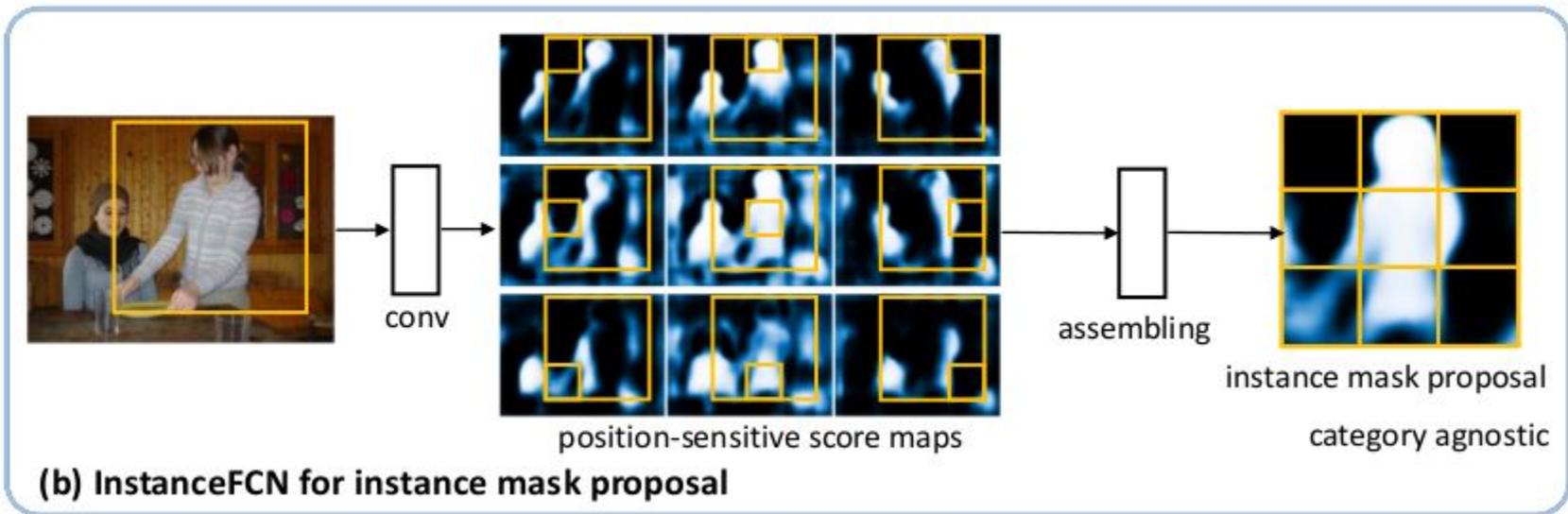
## FCN for semantic segmentation



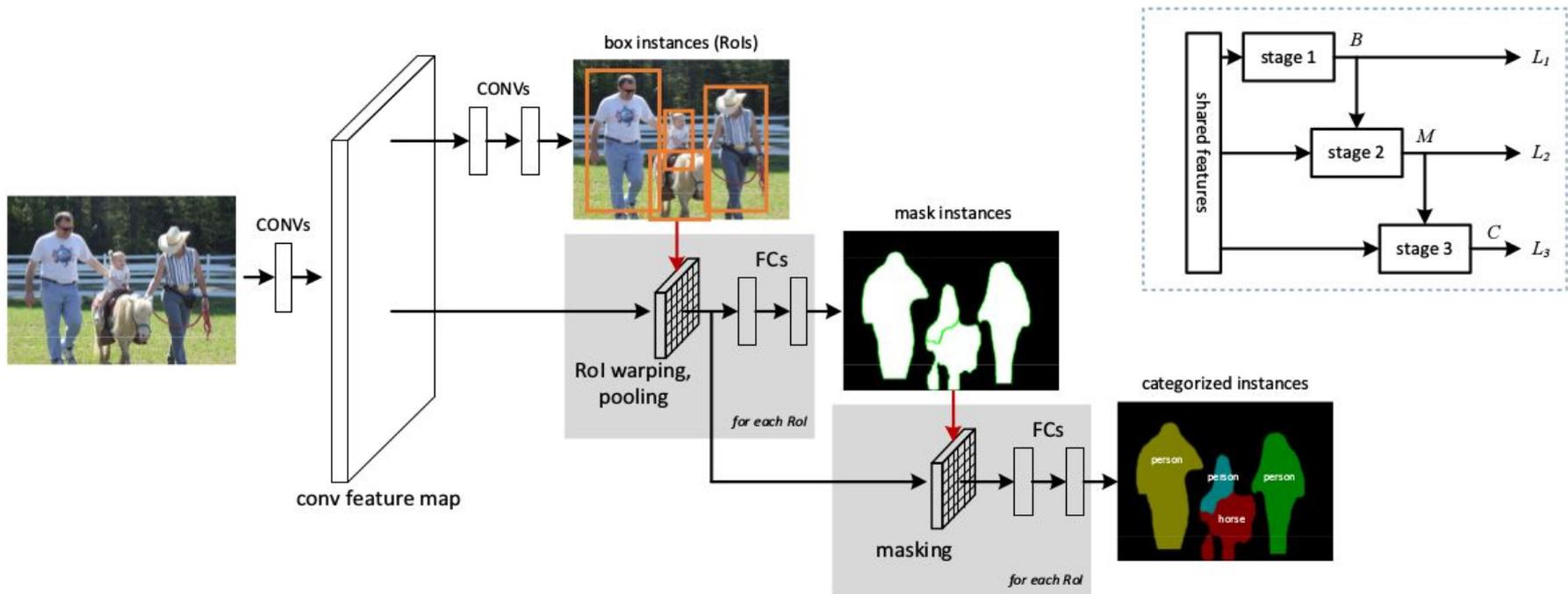
## InstanceFCN for instance segment proposal



# Instance FCN vs FCIS

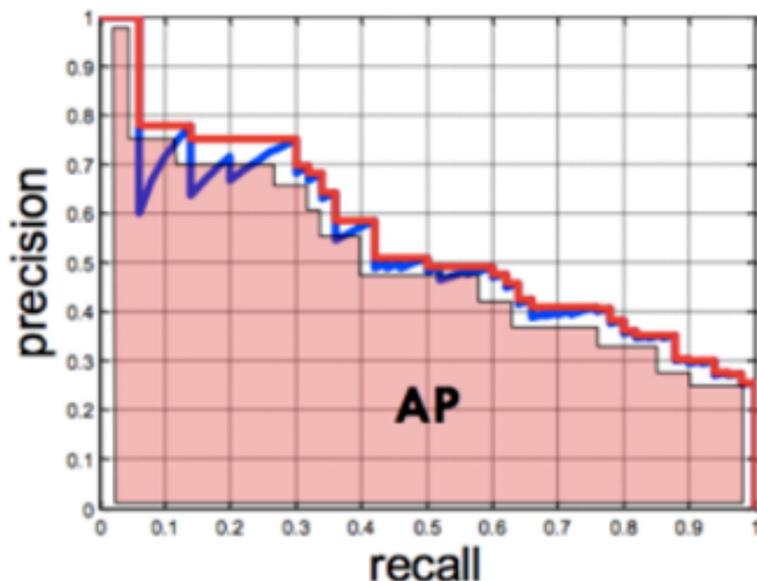


# Multi-task Network Cascades



Ref: Yi Li, et al. Instance-aware Semantic Segmentation via Multi-task Network Cascades

# AP metric



Commonly, IoU > 0.5 means that it was a hit, otherwise it was a fail. For each class, one can calculate the

- True Positive ( $TP(c)$ ): a proposal was made for class  $c$  and there actually was an object of class  $c$
- False Positive ( $FP(c)$ ): a proposal was made for class  $c$ , but there is no object of class  $c$
- Average Precision for class  $c$ :  $\frac{\#TP(c)}{\#TP(c)+\#FP(c)}$

The mAP (mean average precision) =  $\frac{1}{|\text{classes}|} \sum_{c \in \text{classes}} \frac{\#TP(c)}{\#TP(c)+\#FP(c)}$

Fig source: <http://darkpgmr.tistory.com/162>

The RoI warping layer crops a feature map region and warps it into a target size by interpolation.

Ref: Jifeng Dai, kaiming He, Jian Sun Instance-aware Semantic Segmentation via Multi-task Network Cascades

# Timing

**ResNet-50-FPN on COCO trainval135k**

**8-GPU**

**0.72s per 16-image mini-batch**