

# Mask R-CNN

**Kaiming He · Georgia Gkioxari · Piotr Dollar · Ross Girshick**  
Facebook AI Research (FAIR)

Charles Canavaggio

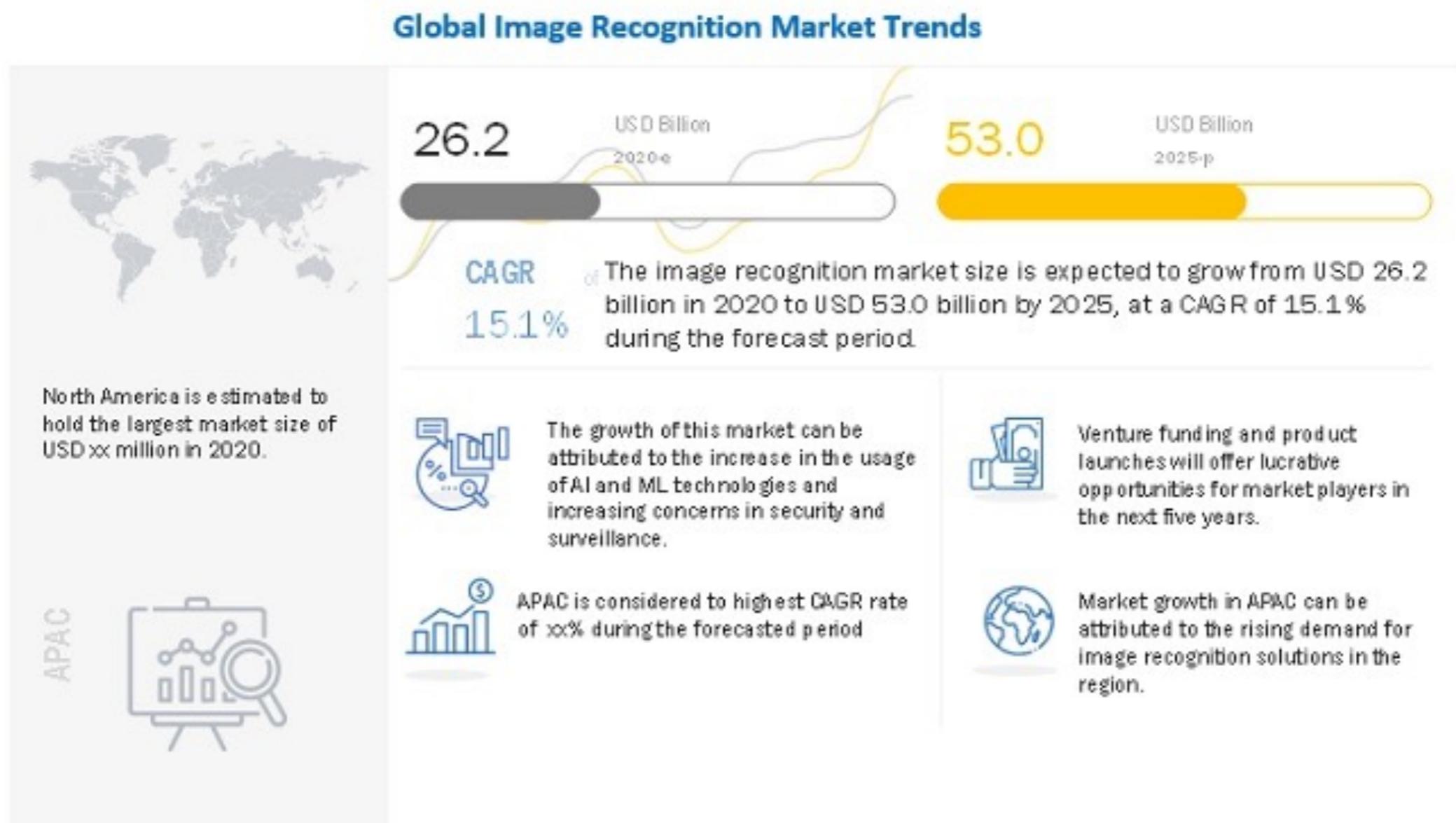


UNIVERSITÉ  
**CÔTE D'AZUR**

Master Data Science M1

2020 – 2021

# Context



e: estimated, p: projected

# What is Instance Segmentation ?

## Classification



CAT

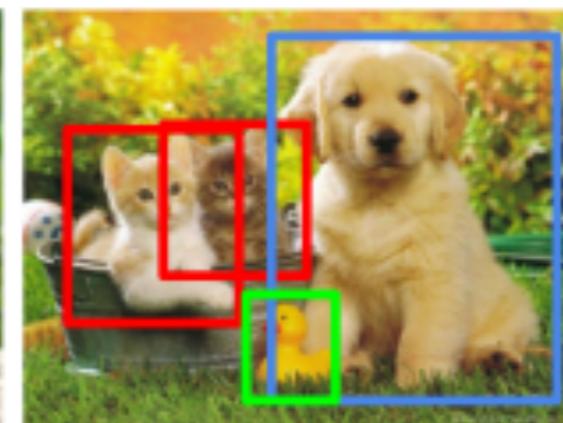
## Single object

## Classification + Localization



CAT

## Object Detection



CAT, DOG, DUCK

## Instance Segmentation

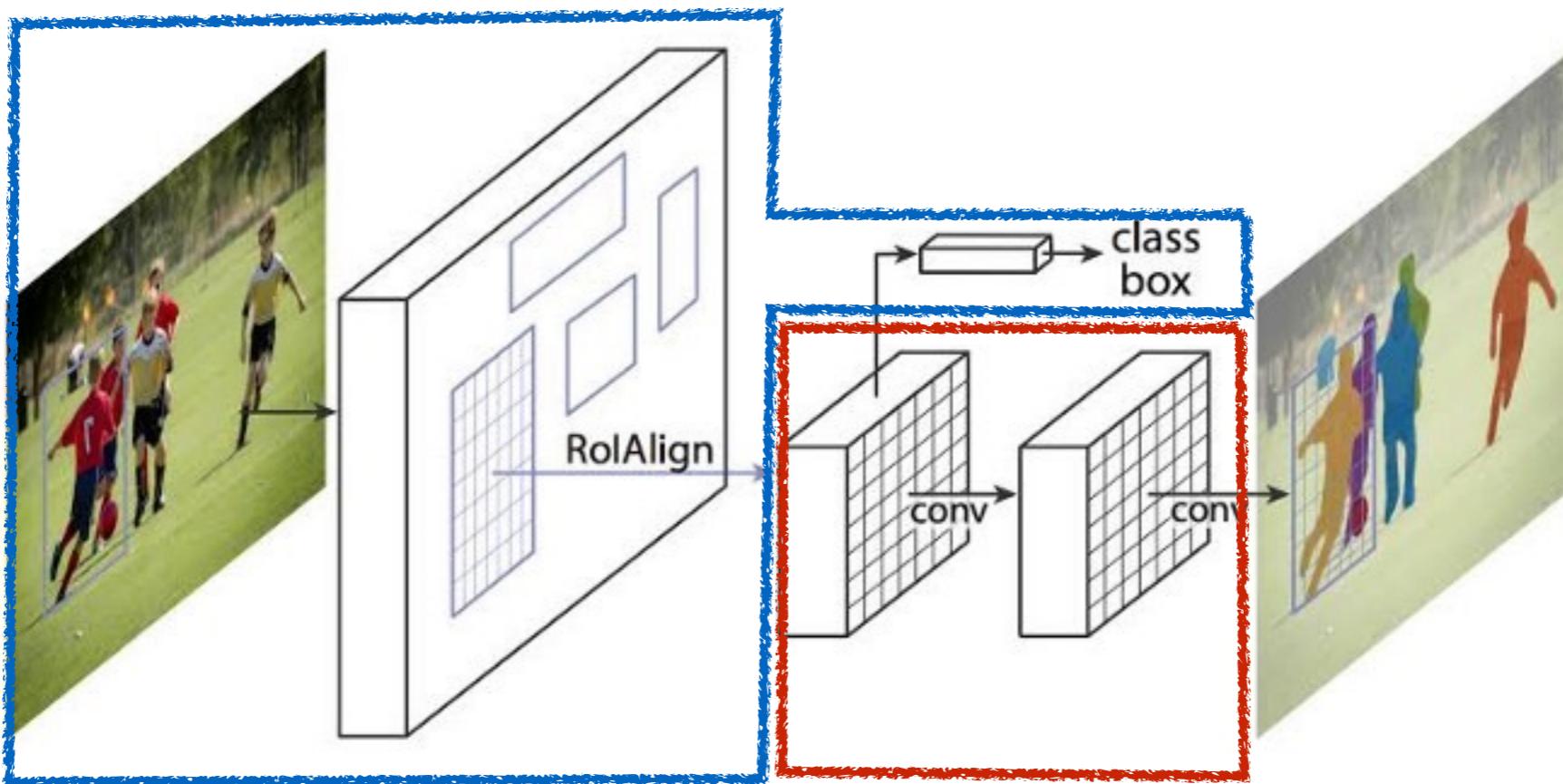


CAT, DOG, DUCK

## Multiple objects

# What is Mask R-CNN ?

## Network Architecture

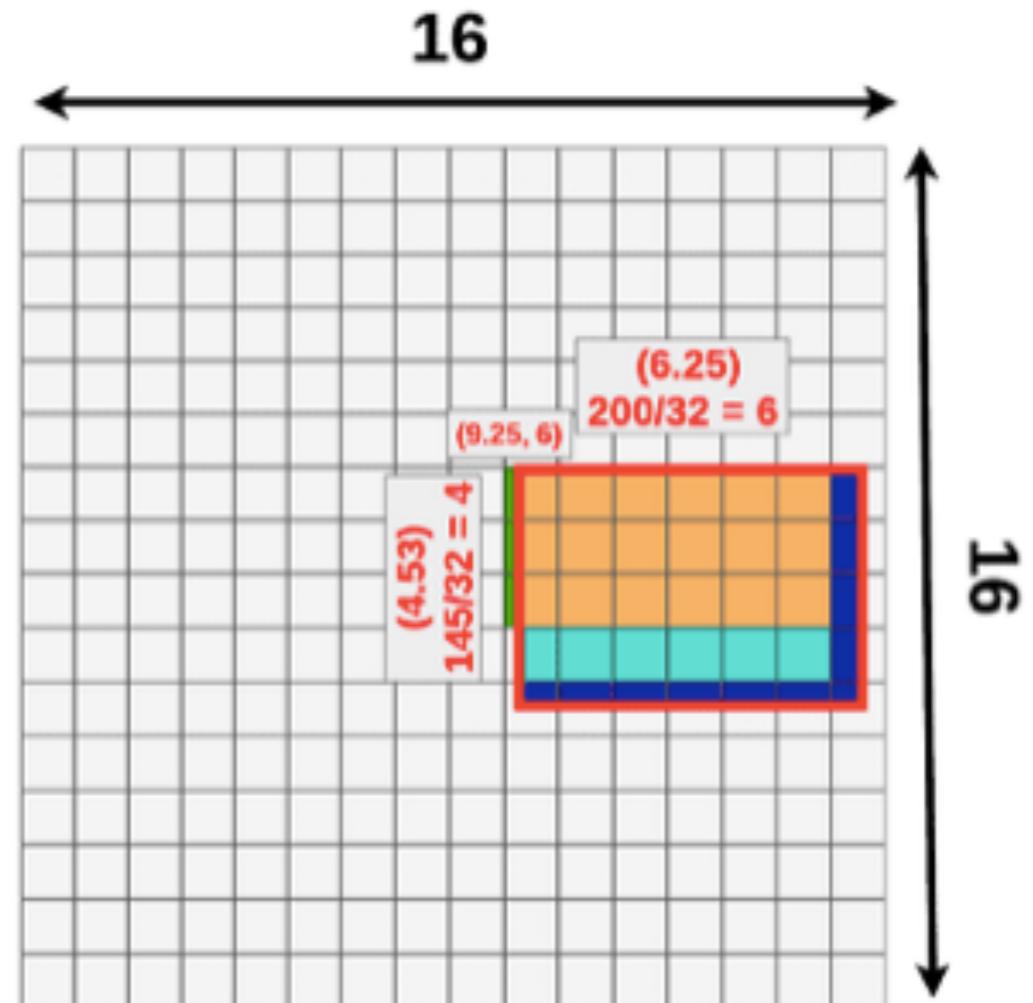
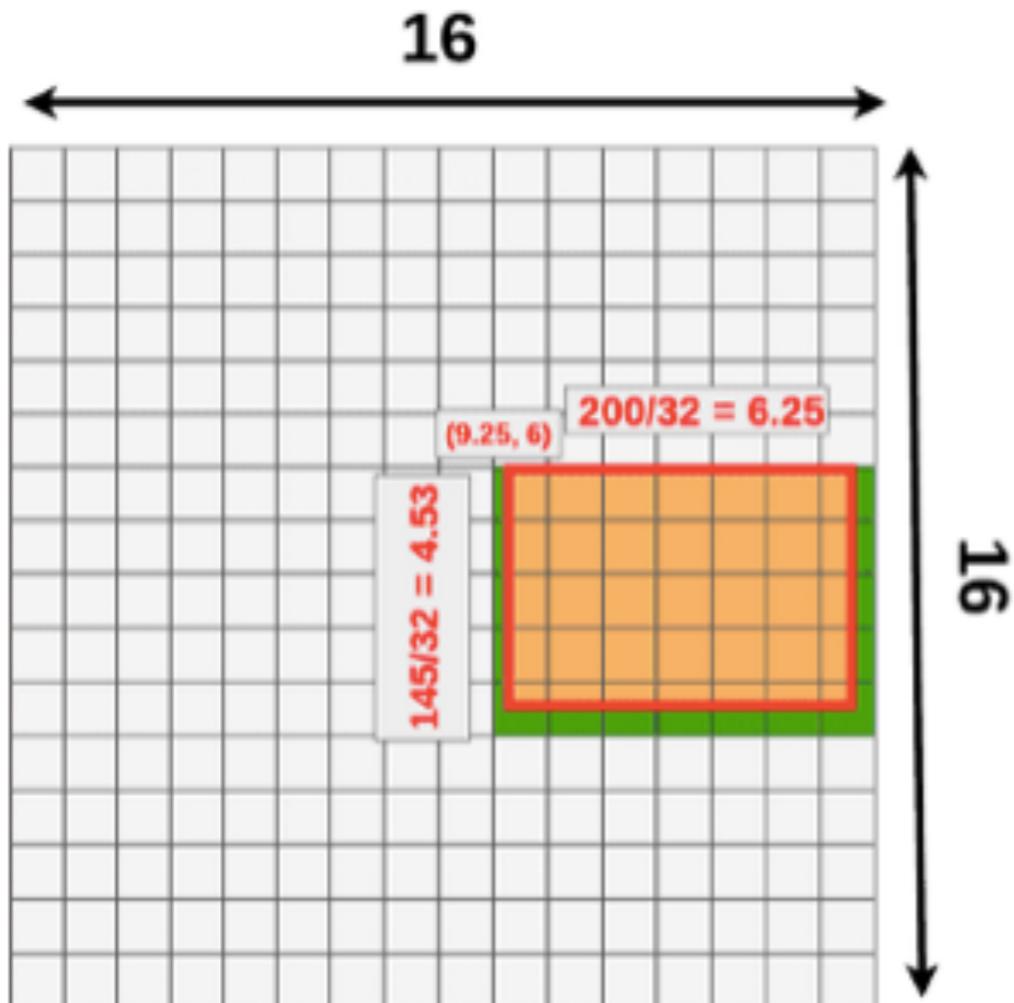


Instance segmentation



Faster R-CNN

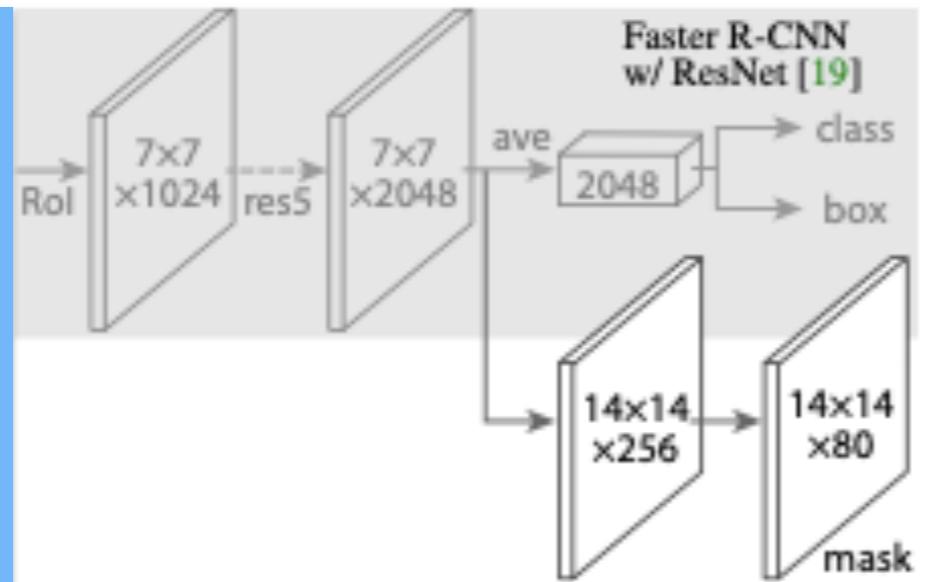
# RoIAlign versus RoIPooling



Comparing RoIAlign(left) and RoIPooling(right) data sources.

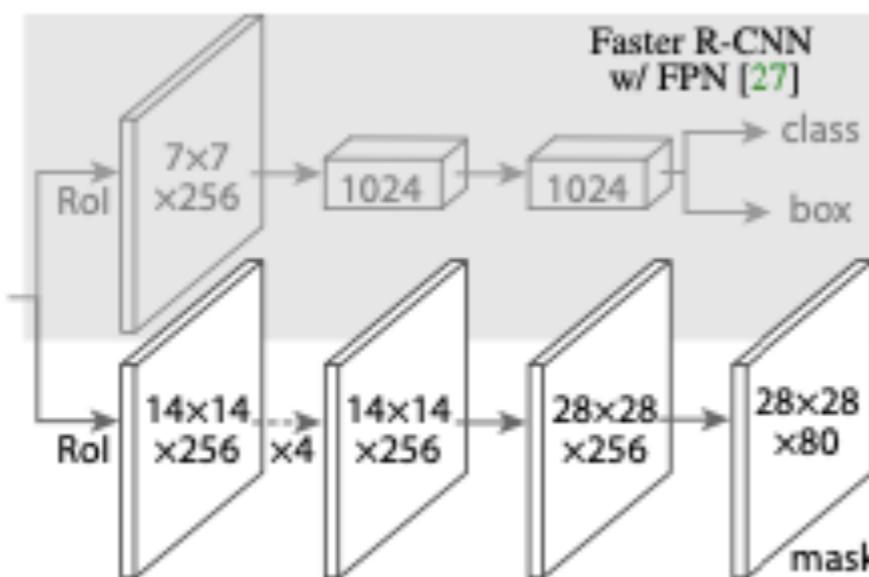
# Network Architecture

ResNet  
or  
ResNext



head

FPN  
(Feature  
Pyramid  
Network)



BackBone

# Main Results : Instance Segmentation

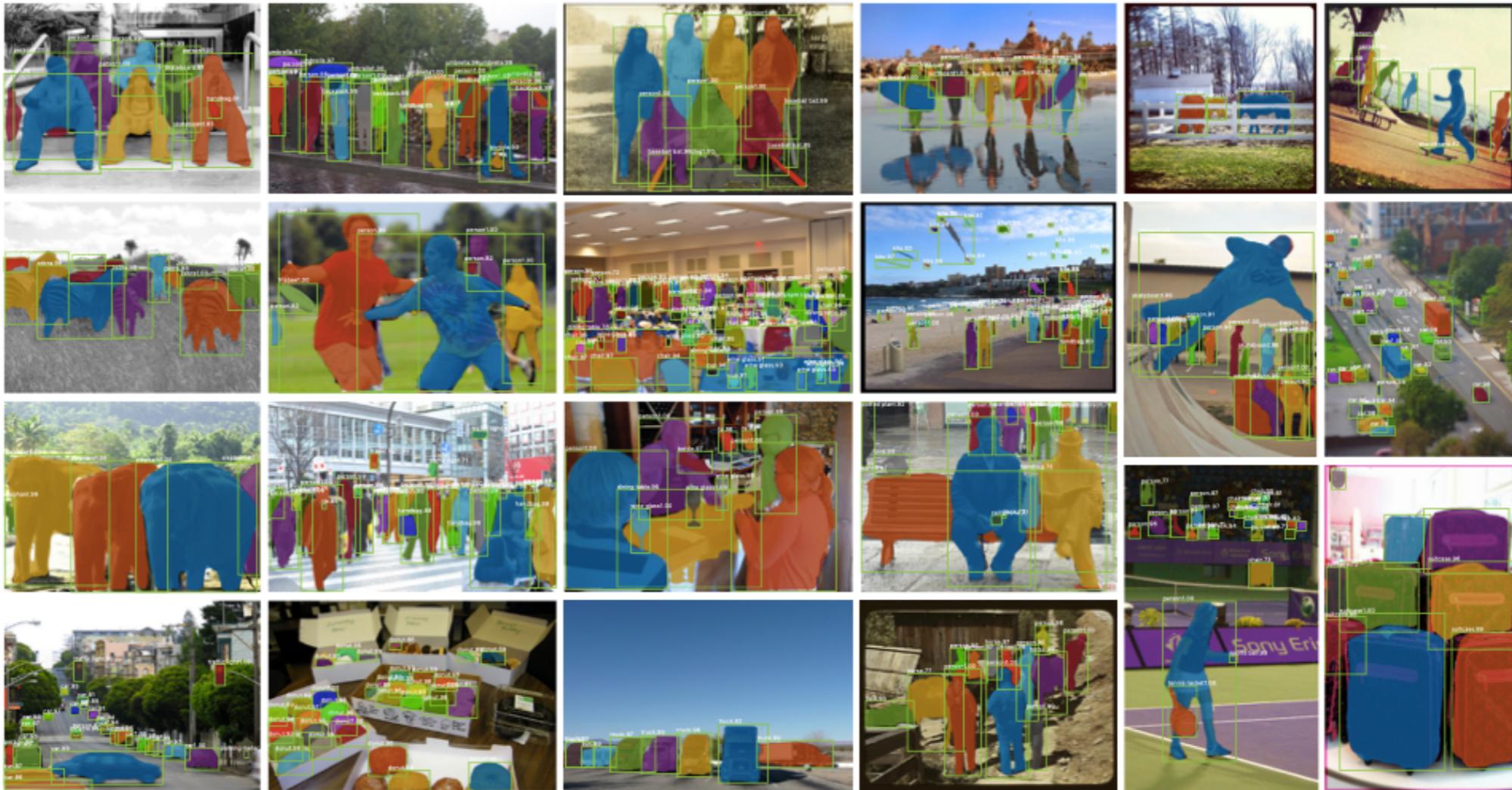


Figure 5. More results of **Mask R-CNN** on COCO test images, using ResNet-101-FPN and running at 5 fps, with 35.7 mask AP (Table 1).

	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>

Table 1. **Instance segmentation** mask AP on COCO test-dev. MNC [10] and FCIS [26] are the winners of the COCO 2015 and 2016 segmentation challenges, respectively. Without bells and whistles, Mask R-CNN outperforms the more complex FCIS++, which includes multi-scale train/test, horizontal flip test, and OHEM [38]. All entries are *single-model* results.

# Ablation Experiments : Instance Segmentation

<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.

	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>

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

	<th>bilinear?</th> <th>agg.</th> <th>AP</th> <th>AP<sub>50</sub></th> <th>AP<sub>75</sub></th>	bilinear?	agg.	AP	AP <sub>50</sub>	AP <sub>75</sub>
<i>RoIPool</i> [12]			max	26.9	48.8	26.4
<i>RoIWarp</i> [10]		✓	max	27.2	49.2	27.1
<i>RoIAlign</i>	✓	✓	max	<b>30.2</b>	<b>51.0</b>	<b>31.8</b>
	✓	✓	ave	<b>30.3</b>	<b>51.2</b>	<b>31.5</b>

(c) **RoIAlign** (ResNet-50-C4): Mask results with various RoI layers. Our *RoIAlign* layer improves AP by ~3 points and AP<sub>75</sub> by ~5 points. Using proper alignment is the only factor that contributes to the large gap between RoI layers.

	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
<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	+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.

	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
FCN	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.

Table 2. **Ablations**. We train on `trainval35k`, test on `minival`, and report *mask* AP unless otherwise noted.

# Ablation Experiments : Bounding Box

	backbone	AP <sup>bb</sup>	AP <sub>50</sub> <sup>bb</sup>	AP <sub>75</sub> <sup>bb</sup>	AP <sub>S</sub> <sup>bb</sup>	AP <sub>M</sub> <sup>bb</sup>	AP <sub>L</sub> <sup>bb</sup>
Faster R-CNN+++ [19]	ResNet-101-C4	34.9	55.7	37.4	15.6	38.7	50.9
Faster R-CNN w FPN [27]	ResNet-101-FPN	36.2	59.1	39.0	18.2	39.0	48.2
Faster R-CNN by G-RMI [21]	Inception-ResNet-v2 [41]	34.7	55.5	36.7	13.5	38.1	52.0
Faster R-CNN w TDM [39]	Inception-ResNet-v2-TDM	36.8	57.7	39.2	16.2	39.8	<b>52.1</b>
Faster R-CNN, RoIAlign	ResNet-101-FPN	37.3	59.6	40.3	19.8	40.2	48.8
<b>Mask R-CNN</b>	ResNet-101-FPN	38.2	60.3	41.7	20.1	41.1	50.2
<b>Mask R-CNN</b>	ResNeXt-101-FPN	<b>39.8</b>	<b>62.3</b>	<b>43.4</b>	<b>22.1</b>	<b>43.2</b>	51.2

Table 3. **Object detection single-model** results (bounding box AP), vs. state-of-the-art on test-dev. Mask R-CNN using ResNet-101-FPN outperforms the base variants of all previous state-of-the-art models (the mask output is ignored in these experiments). The gains of Mask R-CNN over [27] come from using RoIAlign (+1.1 AP<sup>bb</sup>), multitask training (+0.9 AP<sup>bb</sup>), and ResNeXt-101 (+1.6 AP<sup>bb</sup>).

# Mask R-CNN for Human Pose Estimation



Figure 7. Keypoint detection results on COCO test using Mask R-CNN (ResNet-50-FPN), with person segmentation masks predicted from the same model. This model has a keypoint AP of 63.1 and runs at 5 fps.

	AP <sup>kp</sup>	AP <sup>kp</sup> <sub>50</sub>	AP <sup>kp</sup> <sub>75</sub>	AP <sup>kp</sup> <sub>M</sub>	AP <sup>kp</sup> <sub>L</sub>
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
Mask R-CNN, keypoint-only	62.7	87.0	68.4	57.4	71.1
Mask R-CNN, keypoint & mask	<b>63.1</b>	<b>87.3</b>	<b>68.7</b>	57.8	<b>71.4</b>

# Conclusion

Mask R-CNN provide stronger results with :

- a deeper network

# Conclusion

Mask R-CNN provide stronger results with :

- a deeper network
- a sigmoid activation function for the mask branch

# Conclusion

Mask R-CNN provide stronger results with :

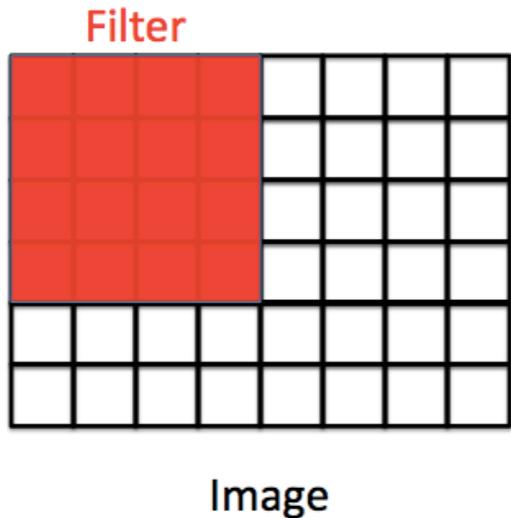
- a deeper network
- a sigmoid activation function for the mask branch
- uses of RoIAlign

# Conclusion

Mask R-CNN provide stronger results with :

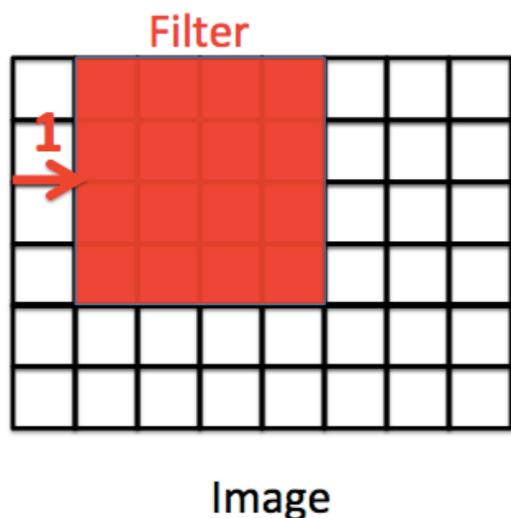
- a deeper network
- a sigmoid activation function for the mask branch
- uses of RoIAlign
- uses of FCN for mask prediction

# Appendix



stride = **Stride** is a parameter of the **neural network's** filter that modifies the amount of movement over the image or video.

If stride = 1, the filter will move one pixel.



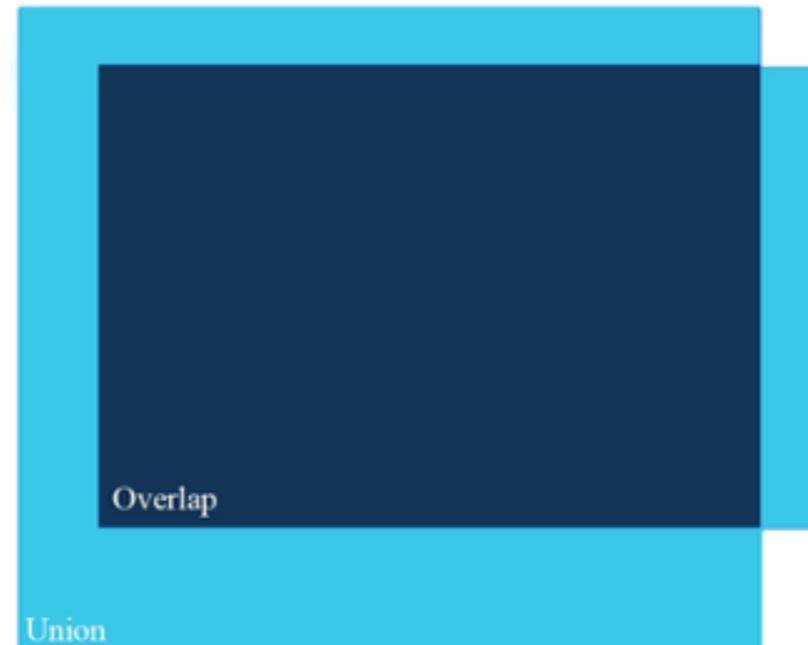
# Appendix

IoU = Intersection over Union,  
IoU measures the overlap  
between 2 boundaries

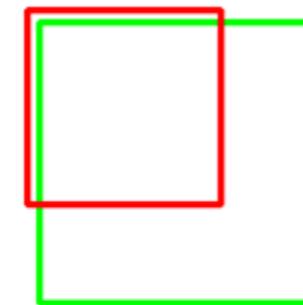


Ground truth  
Prediction

$$IoU = \frac{\text{area of overlap}}{\text{area of union}}$$

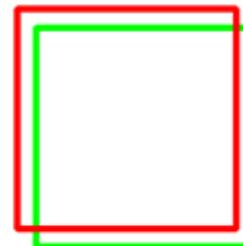


IoU: 0.4034



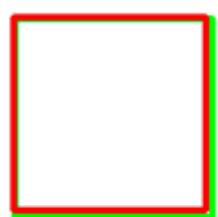
Poor

IoU: 0.7330



Good

IoU: 0.9264



Excellent

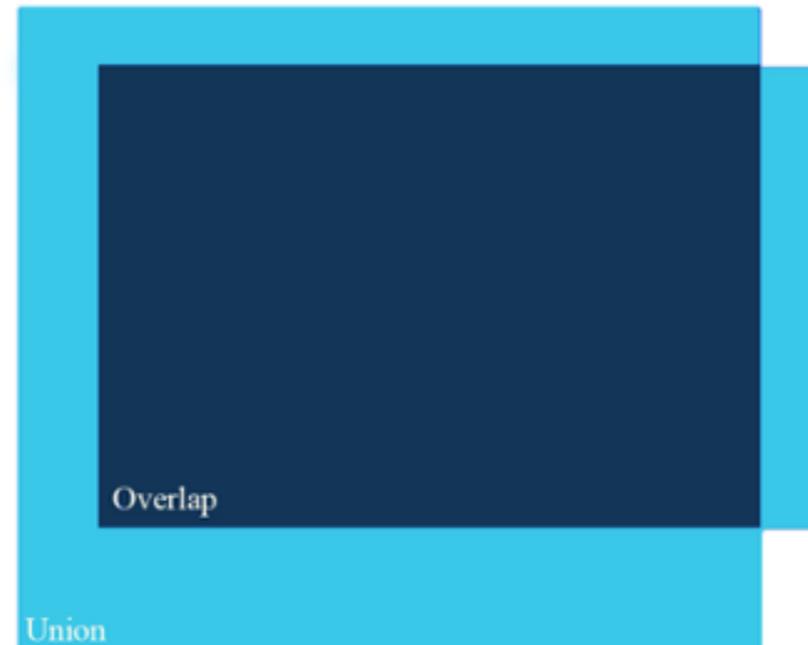
# Appendix

IoU = Intersection over Union,  
IoU measures the overlap  
between 2 boundaries

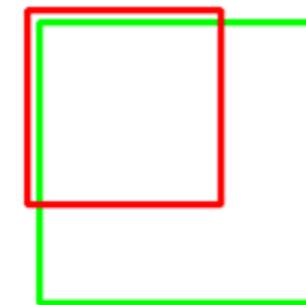


Ground truth  
Prediction

$$IoU = \frac{\text{area of overlap}}{\text{area of union}}$$

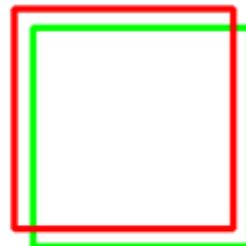


IoU: 0.4034



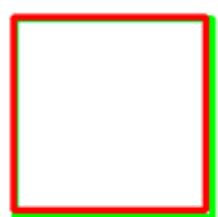
Poor

IoU: 0.7330



Good

IoU: 0.9264



Excellent