



# Object Detection

---

KOHI 2022

2022.09.17.

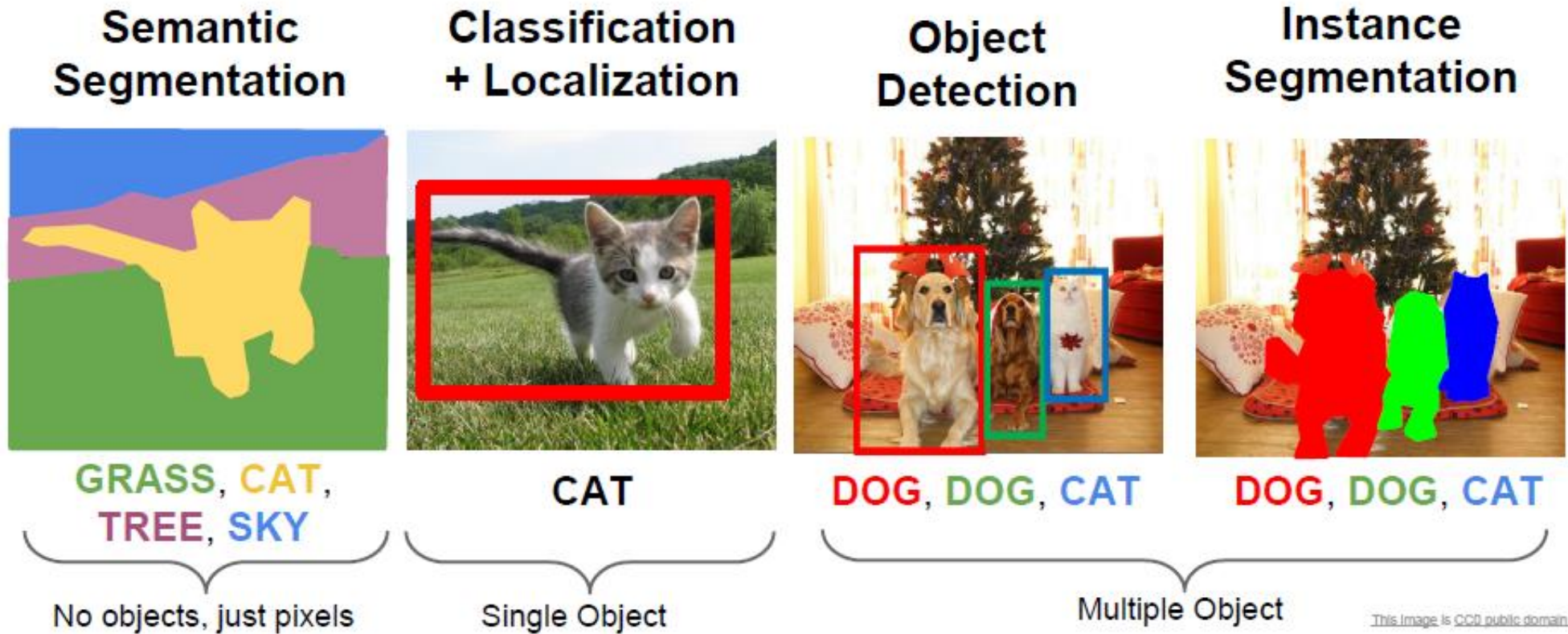
University of Ulsan, Asan Medical Center

Keewon Shin PhDc

# Contents

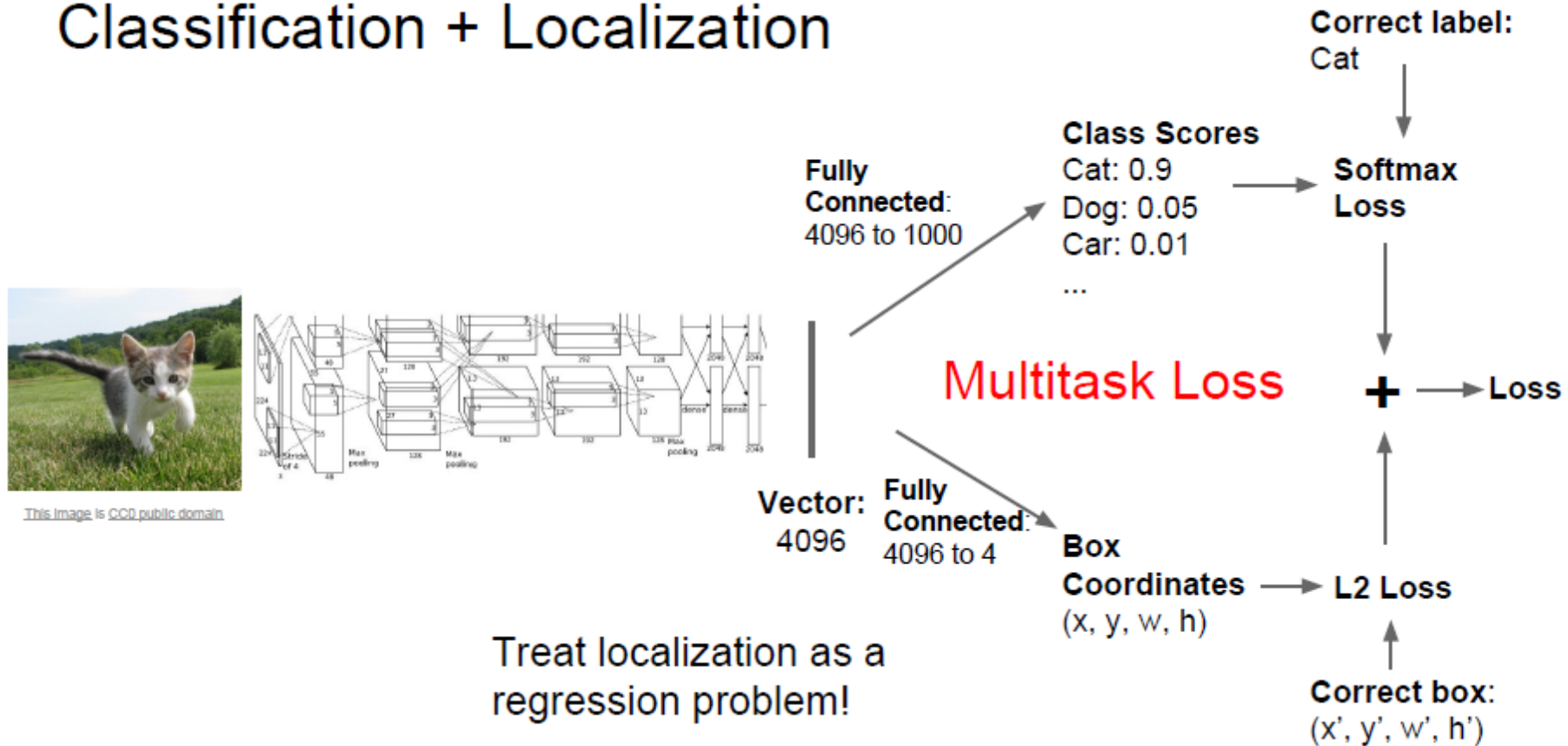
- **1. Computer vision tasks**
  - Object detection
  - Instance segmentation
  - Evaluation
  - Trend
  - Conclusion
- **2. Hands on : Mask-RCNN**  
(Instance segmentation and evaluation)

# Computer Vision Tasks



# Classification + Localization

## Classification + Localization



참고 : (x,y,w,h) or (x1,y1,x2,y2), Box를 정의하는 4개의 숫자면 됨

# Object Detection

## Object Detection: Impact of Deep Learning

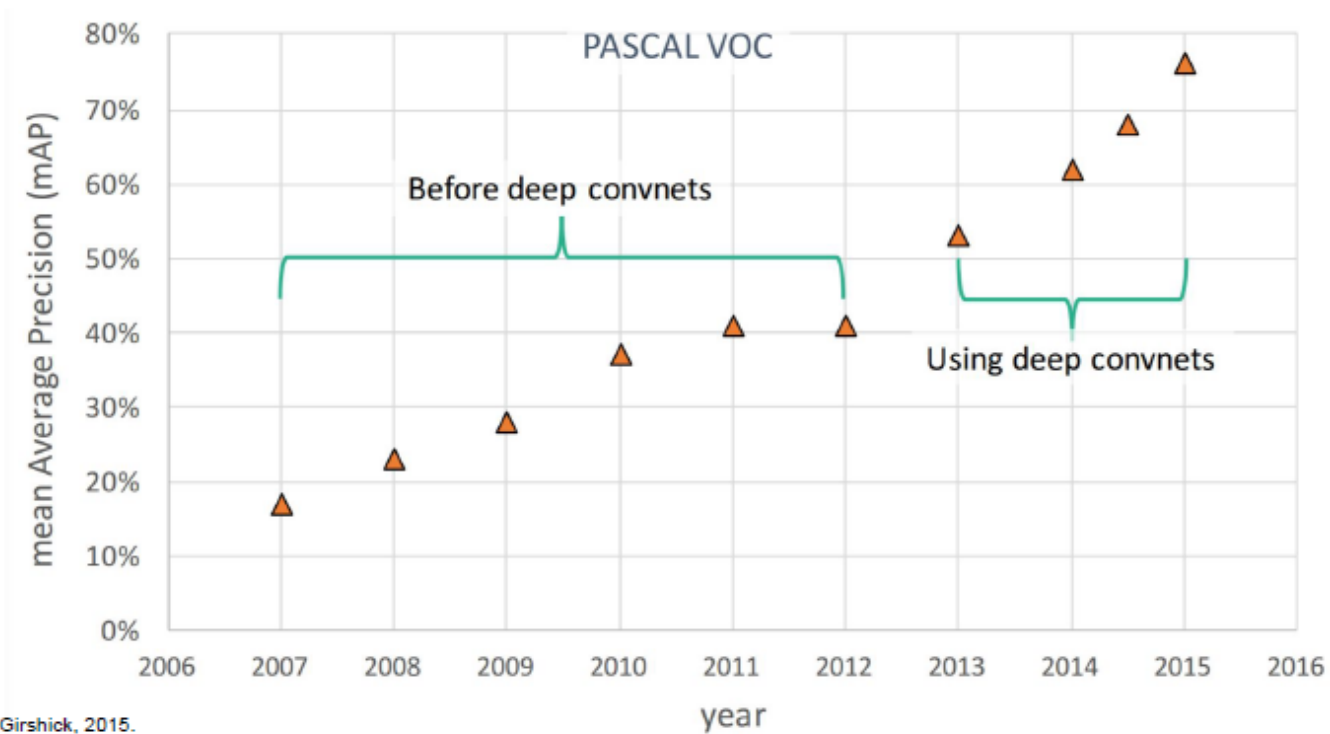
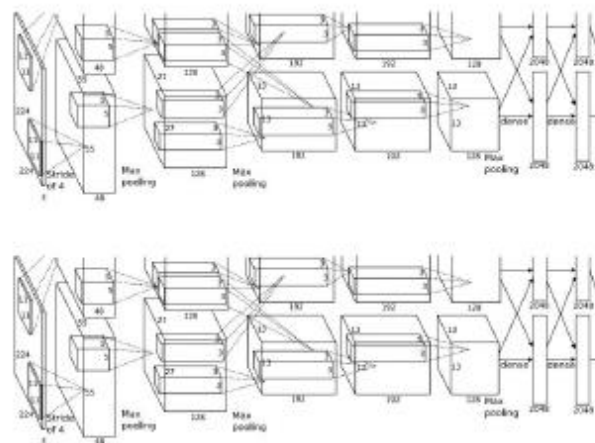


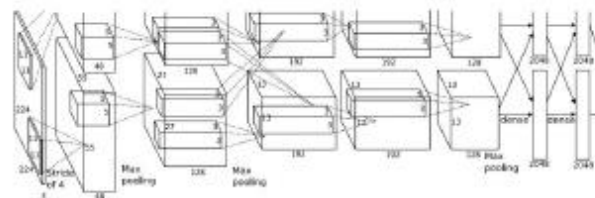
Figure copyright Ross Girshick, 2015.  
Reproduced with permission.

# Object Detection

## Object Detection as Regression?



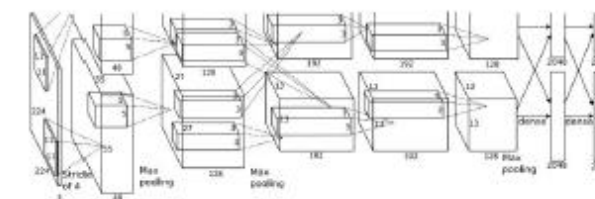
CAT: (x, y, w, h)



DOG: (x, y, w, h)

DOG: (x, y, w, h)

CAT: (x, y, w, h)



DUCK: (x, y, w, h)

DUCK: (x, y, w, h)

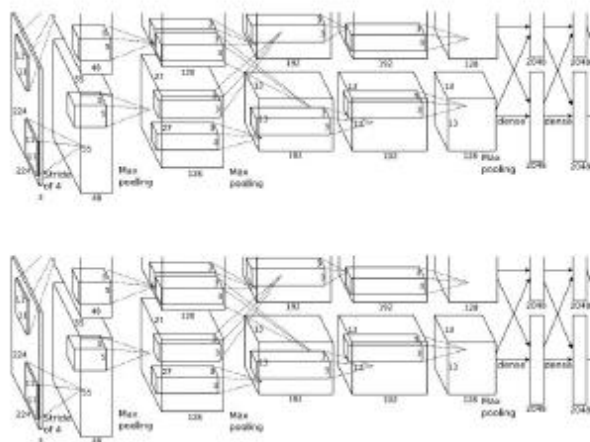
....



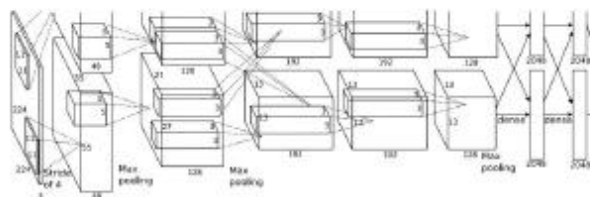
# Object Detection

## Object Detection as Regression?

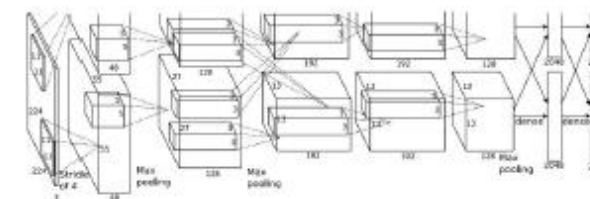
Each image needs a different number of outputs!



CAT: (x, y, w, h)      4 numbers



DOG: (x, y, w, h)  
DOG: (x, y, w, h) 16 numbers  
CAT: (x, y, w, h)

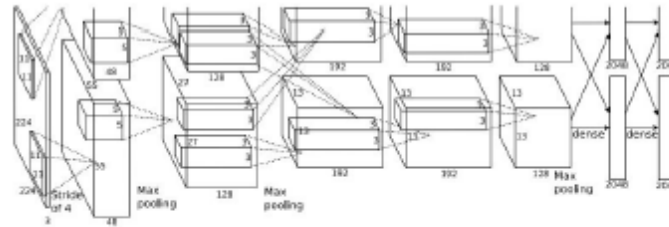


DUCK: (x, y, w, h) Many numbers!

# Object Detection

## Object Detection as Classification: Sliding Window

Apply a CNN to many different crops of the image, CNN classifies each crop as object or background



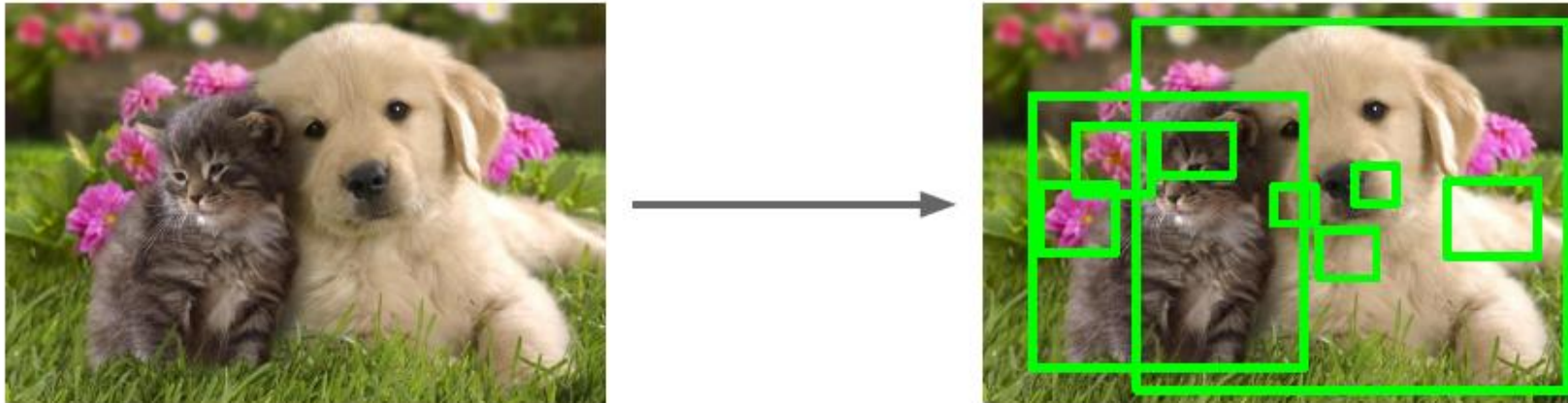
Dog? YES  
Cat? NO  
Background? NO



# Object Detection

## Region Proposals

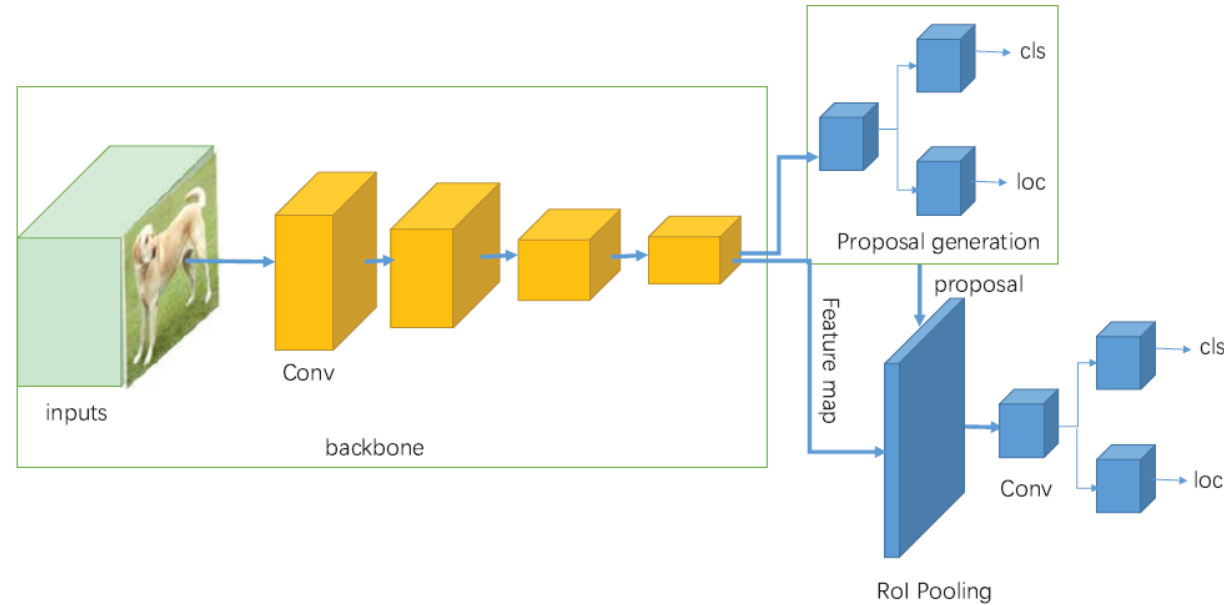
- Find “blobby” image regions that are likely to contain objects
- Relatively fast to run; e.g. Selective Search gives 1000 region proposals in a few seconds on CPU



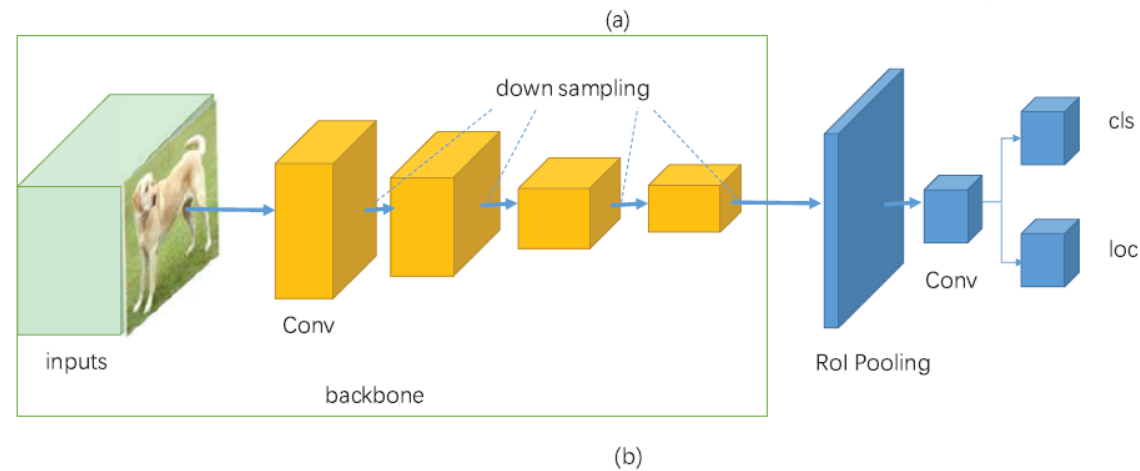
Alexe et al, "Measuring the objectness of image windows", TPAMI 2012  
Uijlings et al, "Selective Search for Object Recognition", IJCV 2013  
Cheng et al, "BING: Binarized normed gradients for objectness estimation at 300fps", CVPR 2014  
Zitnick and Dollar, "Edge boxes: Locating object proposals from edges", ECCV 2014

# Object Detection

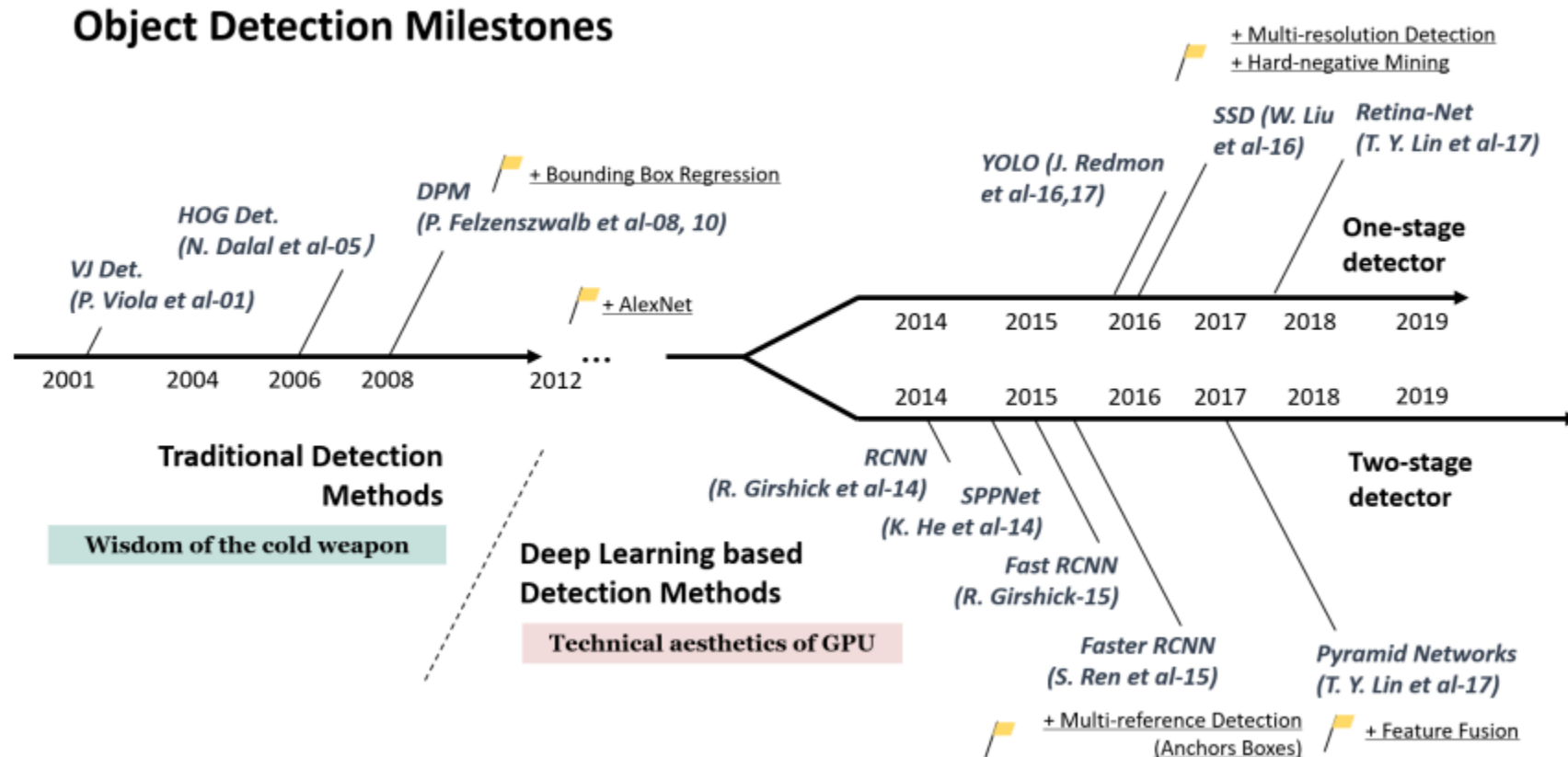
2 Stage approach



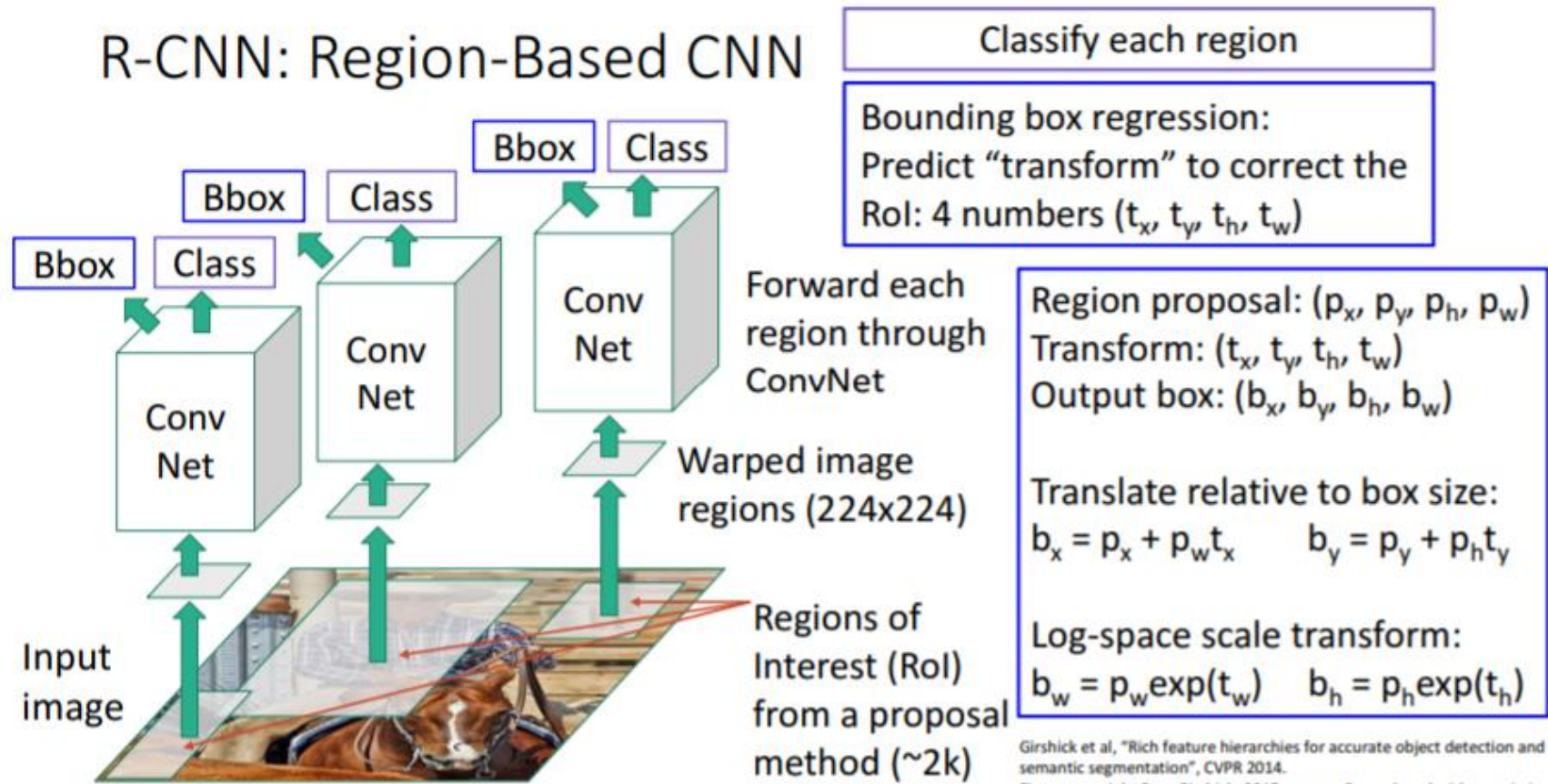
1 Stage approach



# Object Detection



# Object Detection

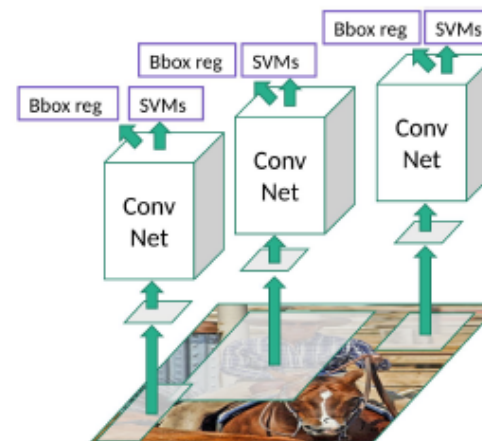


Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.  
Figure copyright Ross Girshick. 2015: [source](#). Reproduced with permission.

# Object Detection

## R-CNN: Problems

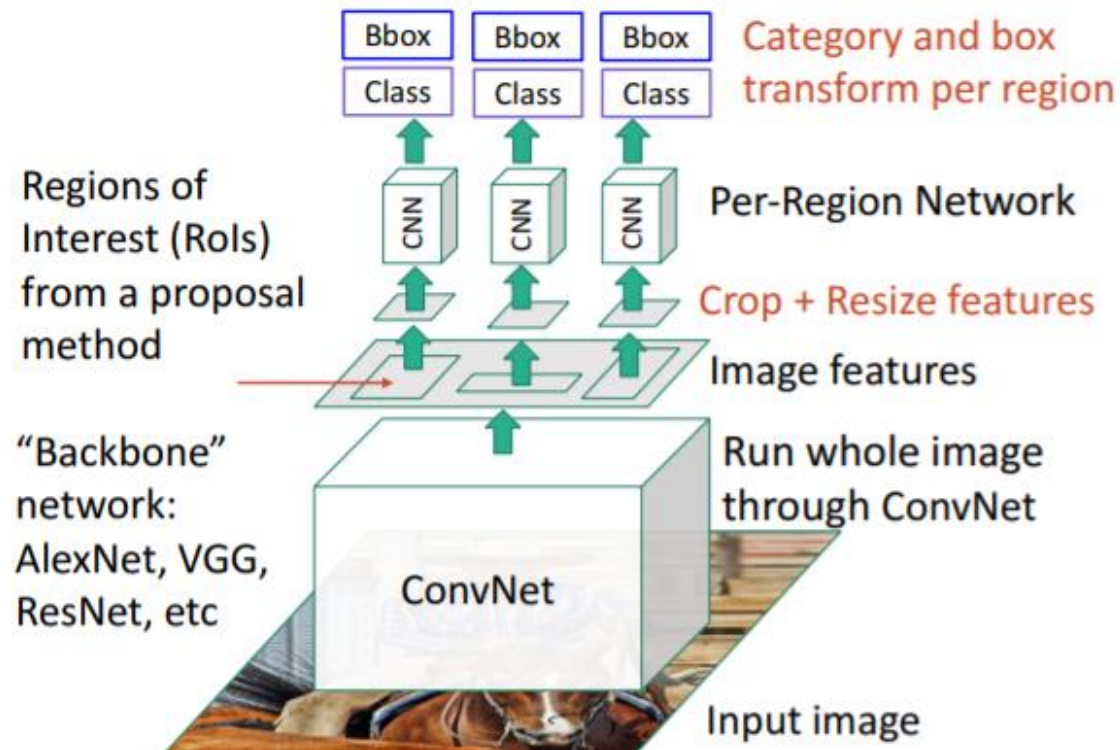
- Ad hoc training objectives
  - Fine-tune network with softmax classifier (log loss)
  - Train post-hoc linear SVMs (hinge loss)
  - Train post-hoc bounding-box regressions (least squares)
- Training is slow (84h), takes a lot of disk space
- Inference (detection) is slow
  - 47s / image with VGG16 [Simonyan & Zisserman. ICLR15]
  - Fixed by SPP-net [He et al. ECCV14]





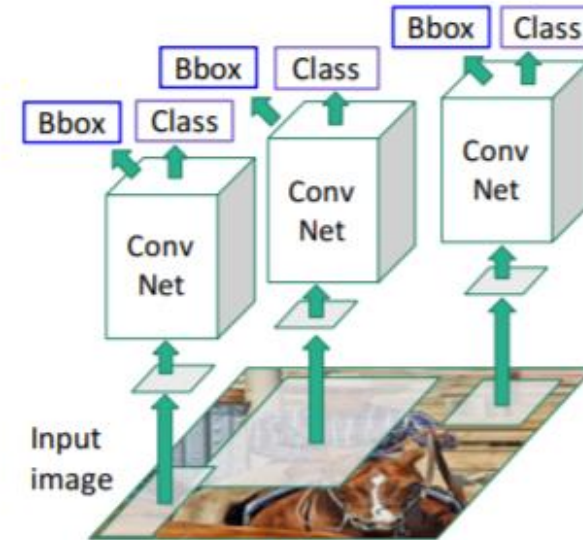
# Object Detection

## Faster R-CNN



## "Slow" R-CNN

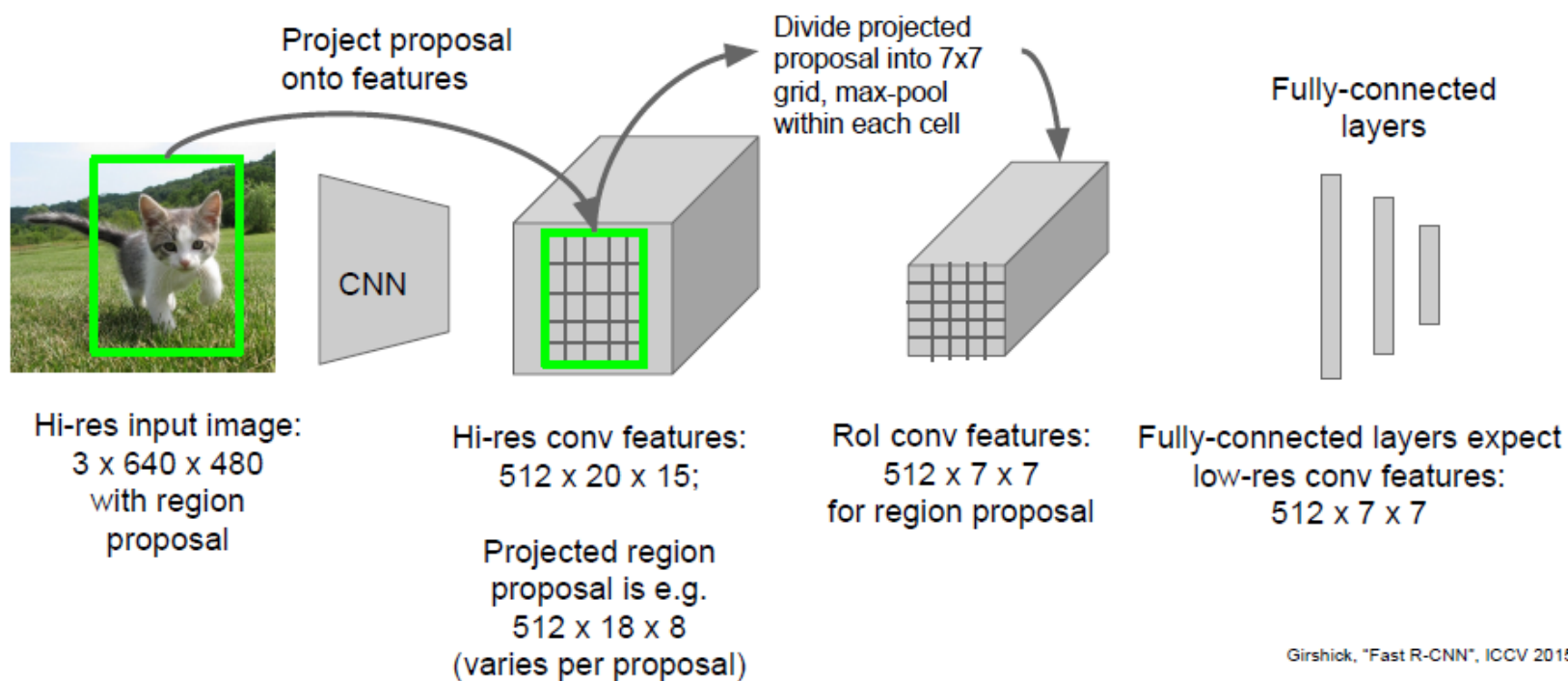
Process each region independently





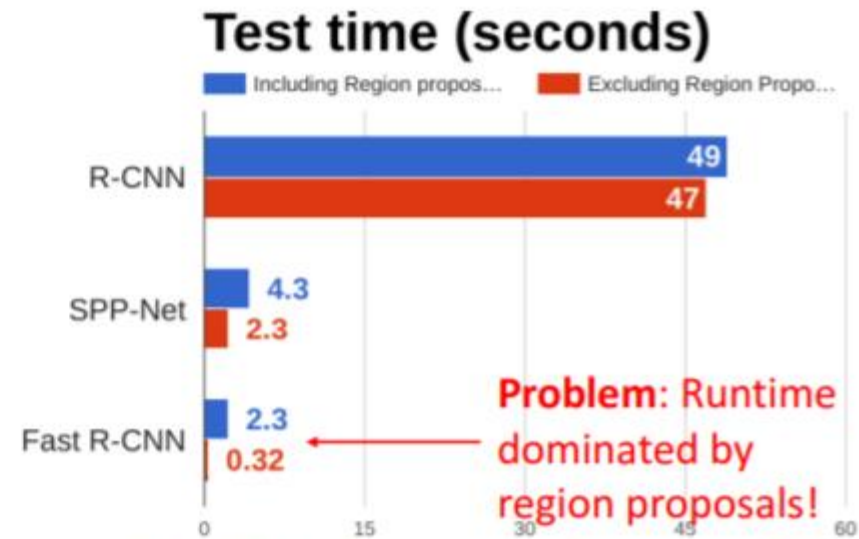
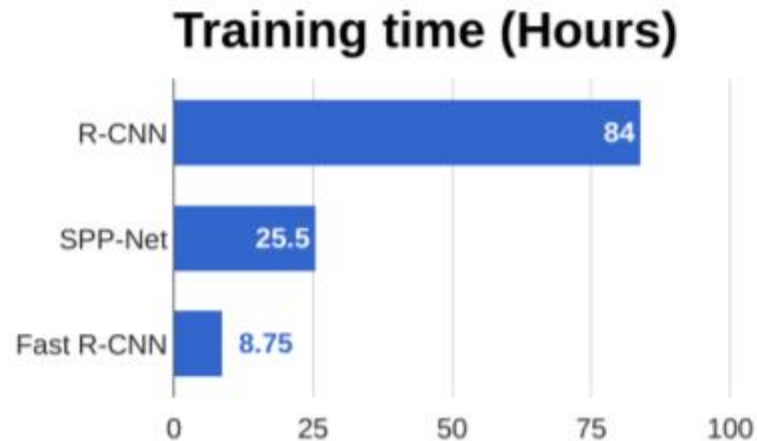
# Object Detection

## Faster R-CNN: RoI Pooling



# Object Detection

## R-CNN vs SPP vs Fast R-CNN



Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.  
He et al, "Spatial pyramid pooling in deep convolutional networks for visual recognition", ECCV 2014  
Girshick, "Fast R-CNN", ICCV 2015

**Recall:** Region proposals computed by heuristic "Selective Search" algorithm on CPU -- let's learn them with a CNN instead!

# Object Detection

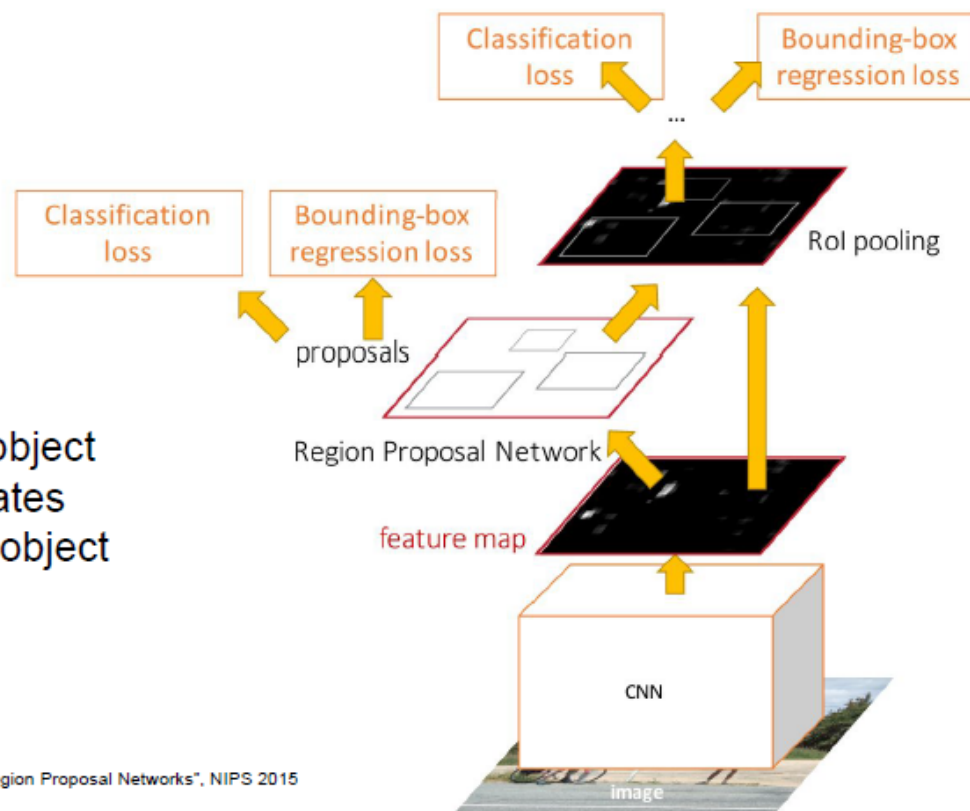
## Faster R-CNN:

Make CNN do proposals!

Insert **Region Proposal Network (RPN)** to predict proposals from features

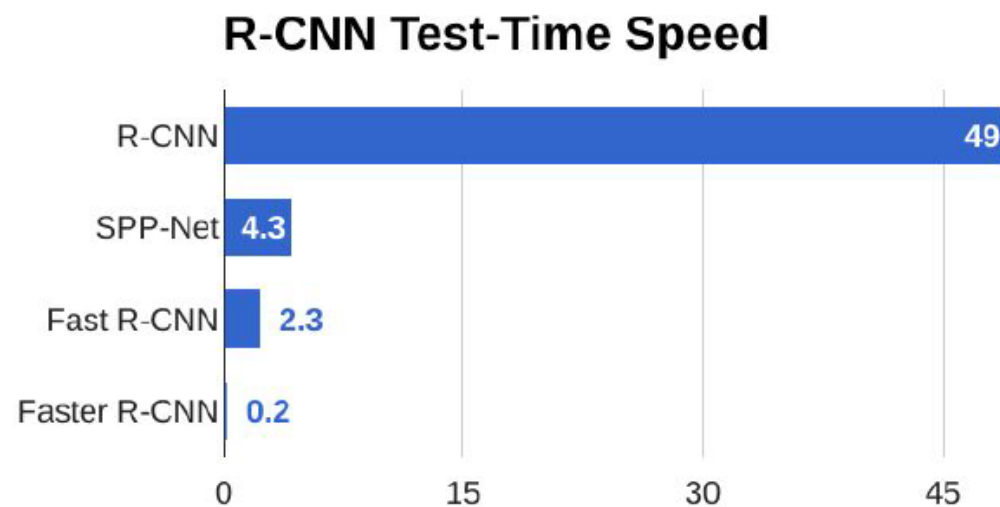
Jointly train with 4 losses:

1. RPN classify object / not object
2. RPN regress box coordinates
3. Final classification score (object classes)
4. Final box coordinates



# Object Detection

**Faster** R-CNN:  
Make CNN do proposals!



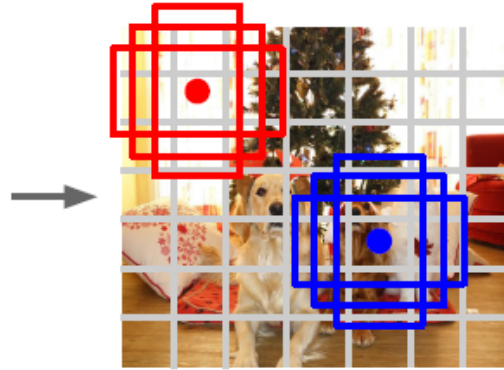
# Object Detection

## Detection without Proposals: YOLO / SSD

Go from input image to tensor of scores with one big convolutional network! →



Input image  
 $3 \times H \times W$



Divide image into grid  
 $7 \times 7$

Image a set of **base boxes**  
centered at each grid cell  
Here  $B = 3$

Within each grid cell:

- Regress from each of the  $B$  base boxes to a final box with 5 numbers:  
( $dx, dy, dh, dw, confidence$ )
- Predict scores for each of  $C$  classes (including background as a class)

Output:  
 $7 \times 7 \times (5 * B + C)$

Redmon et al, "You Only Look Once:  
Unified, Real-Time Object Detection", CVPR 2016  
Liu et al, "SSD: Single-Shot MultiBox Detector", ECCV 2016

# Object Detection : summary

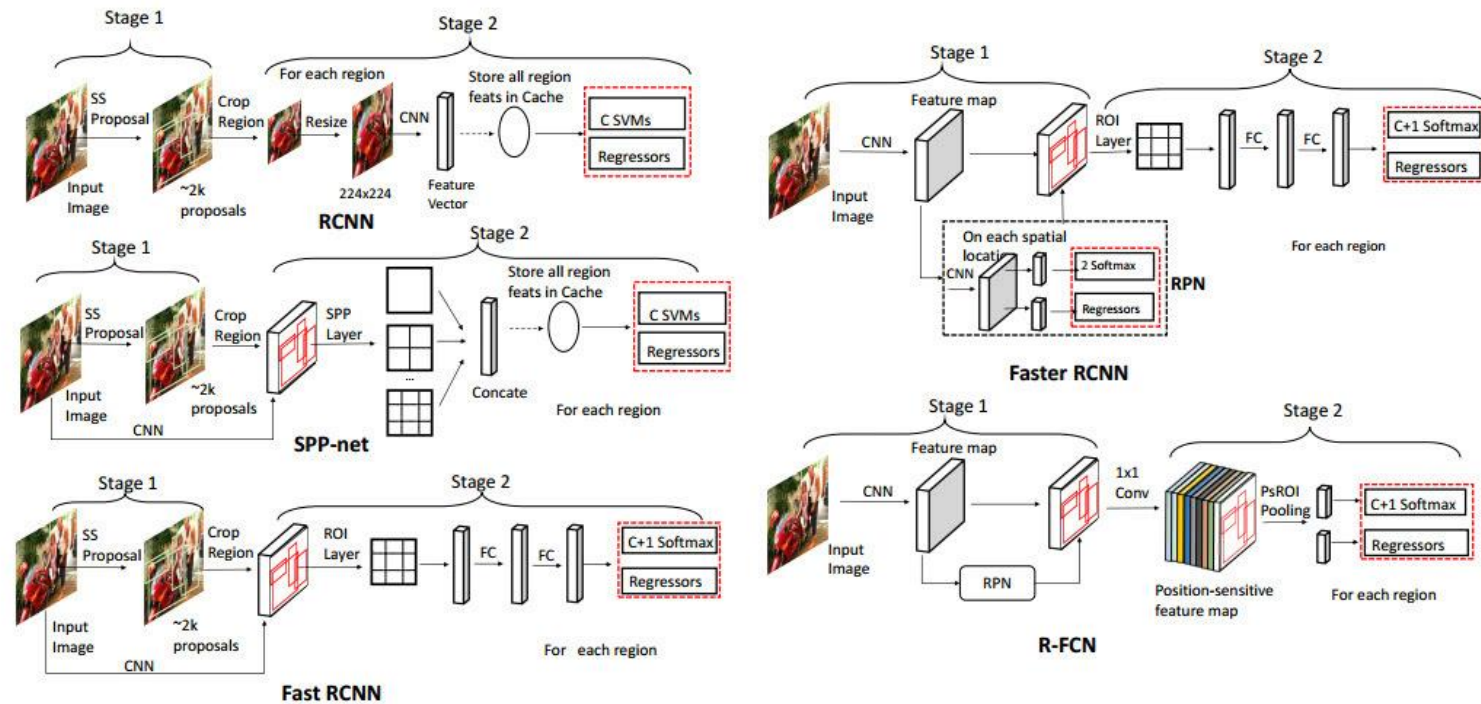


Figure 4: Overview of different two-stage detection frameworks for generic object detection. Red dotted rectangles denote the outputs that define the loss functions.



# Object Detection : summary

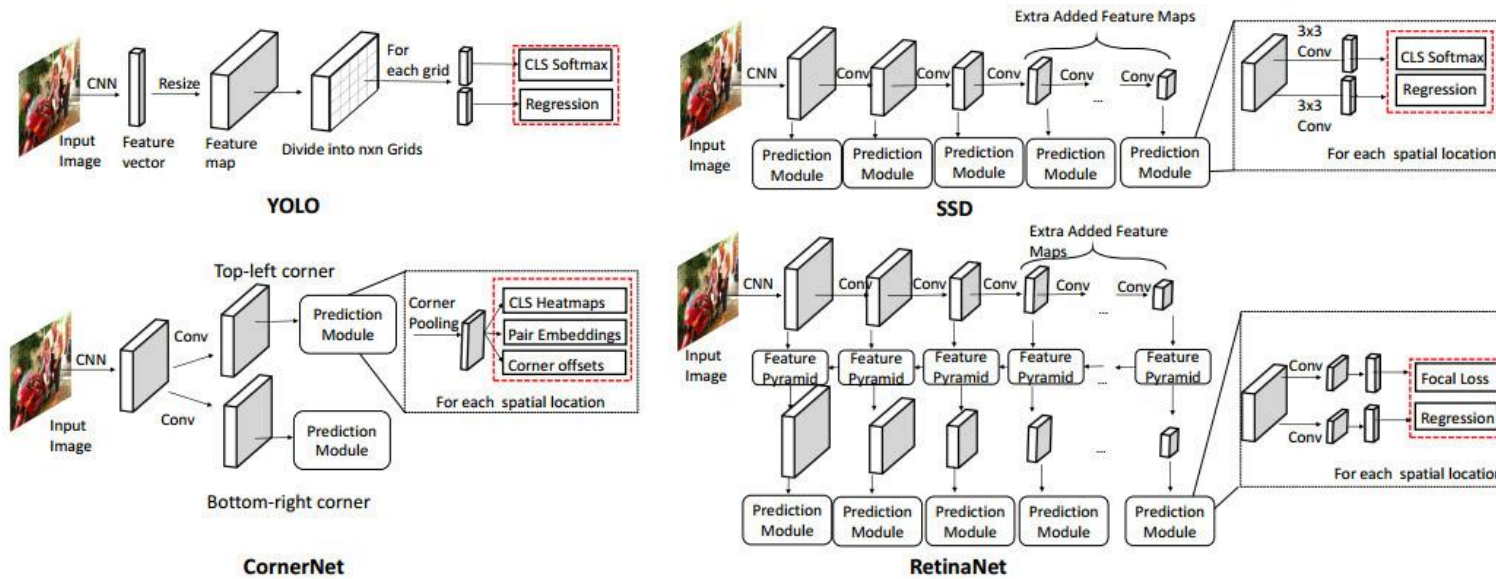
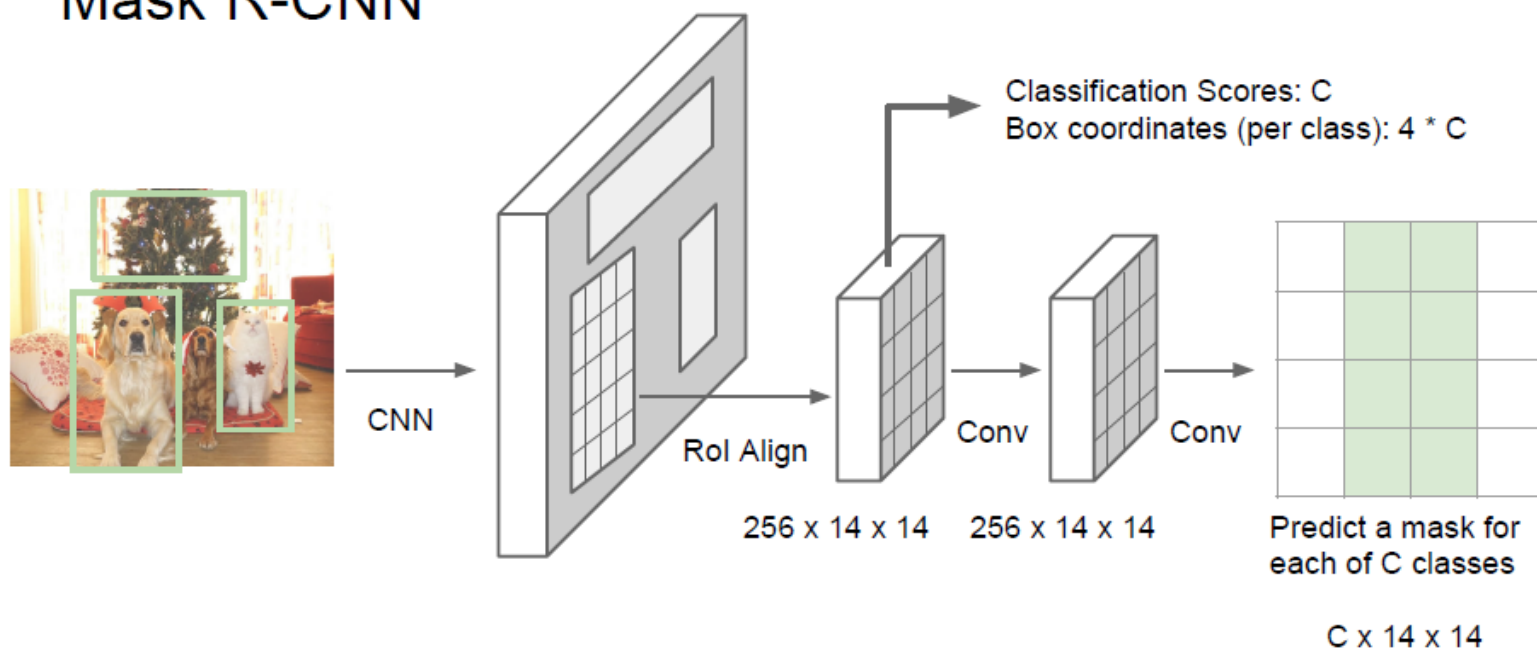


Figure 5: Overview of different one-stage detection frameworks for generic object detection. Red rectangles denotes the outputs that define the objective functions.

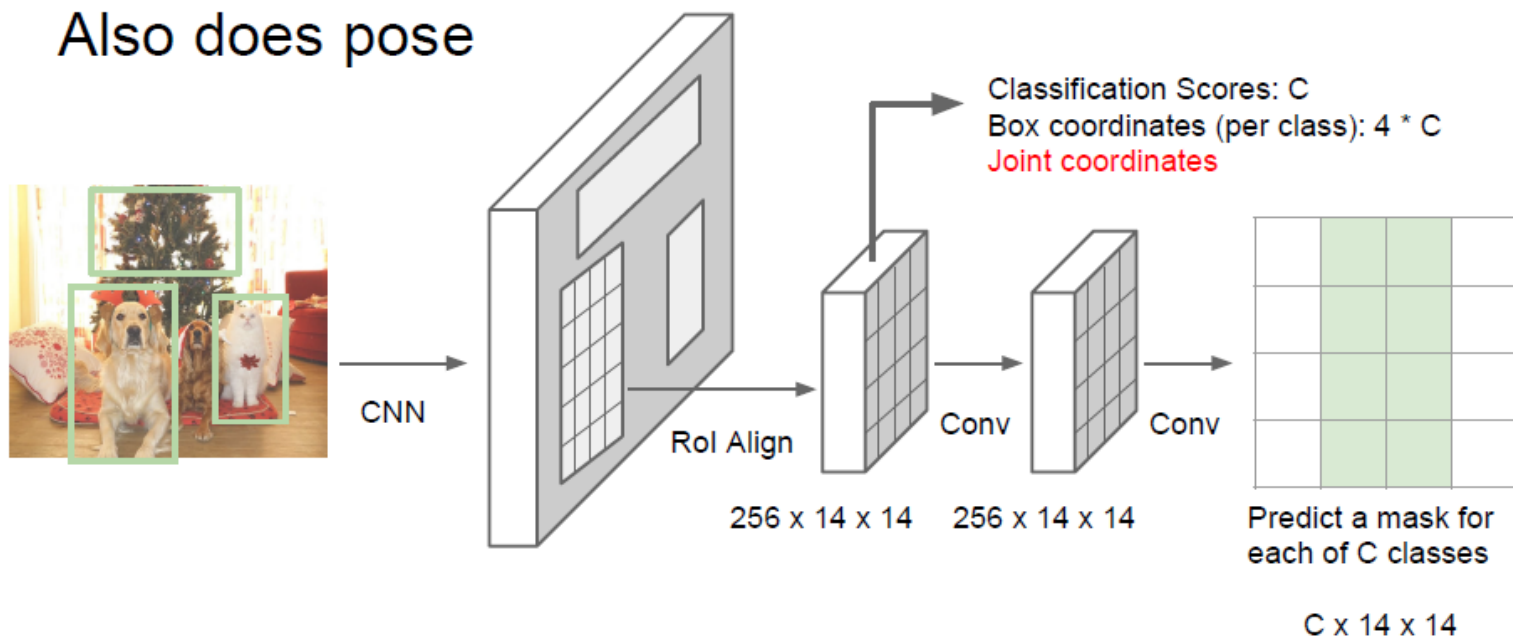
# Instance segmentation

## Mask R-CNN



# Instance segmentation

Mask R-CNN  
Also does pose



# Instance segmentation

Mask R-CNN: Very Good Results!

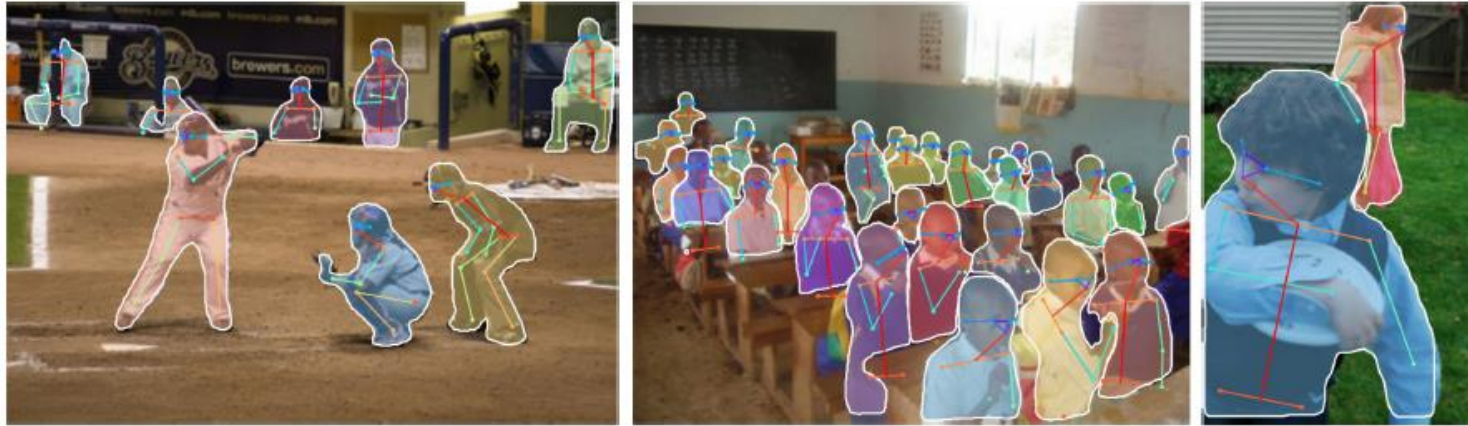


He et al, "Mask R-CNN", arXiv 2017  
Figures copyright Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick, 2017.  
Reproduced with permission.

# Instance segmentation

Mask R-CNN

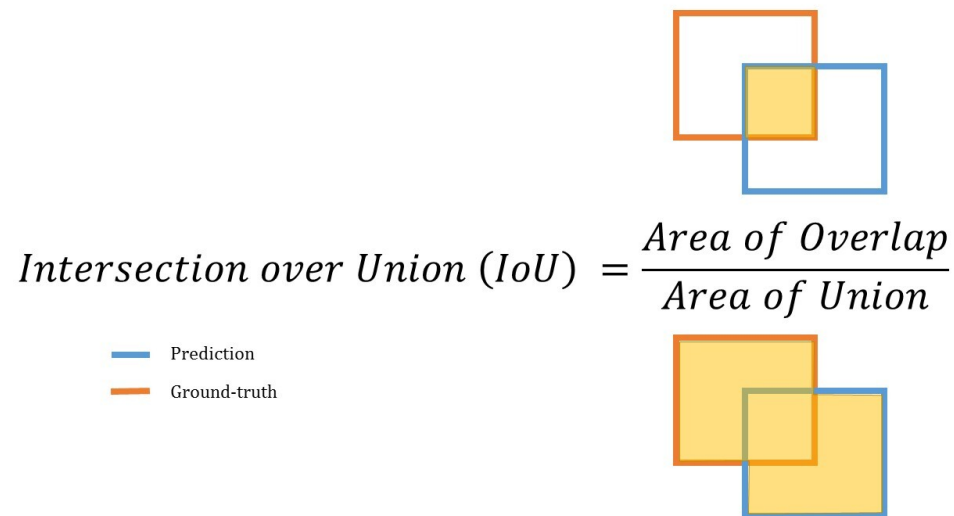
Also does pose



He et al, "Mask R-CNN", arXiv 2017  
Figures copyright Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick, 2017.  
Reproduced with permission.

# Evaluation

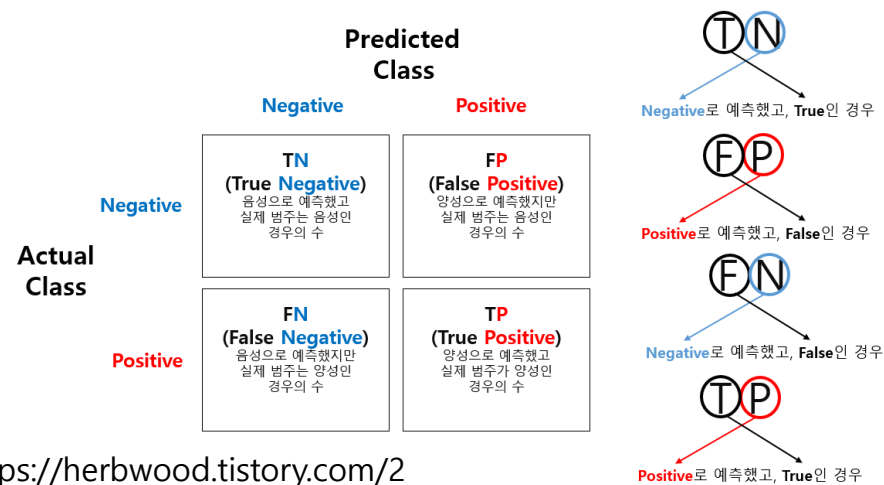
- IoU(Intersection over union)



- Precision and Recall

$$\text{Precision} = \frac{TP}{TP + FP} \quad \text{Recall} = \frac{TP}{TP + FN} \quad \text{F1score} = \frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$$

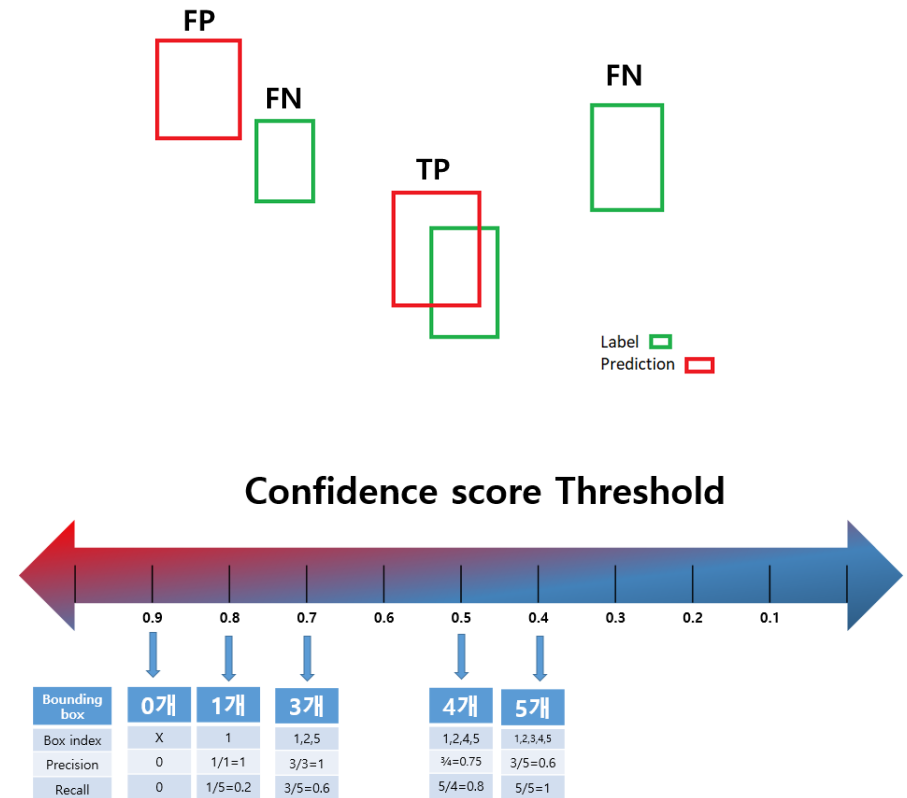
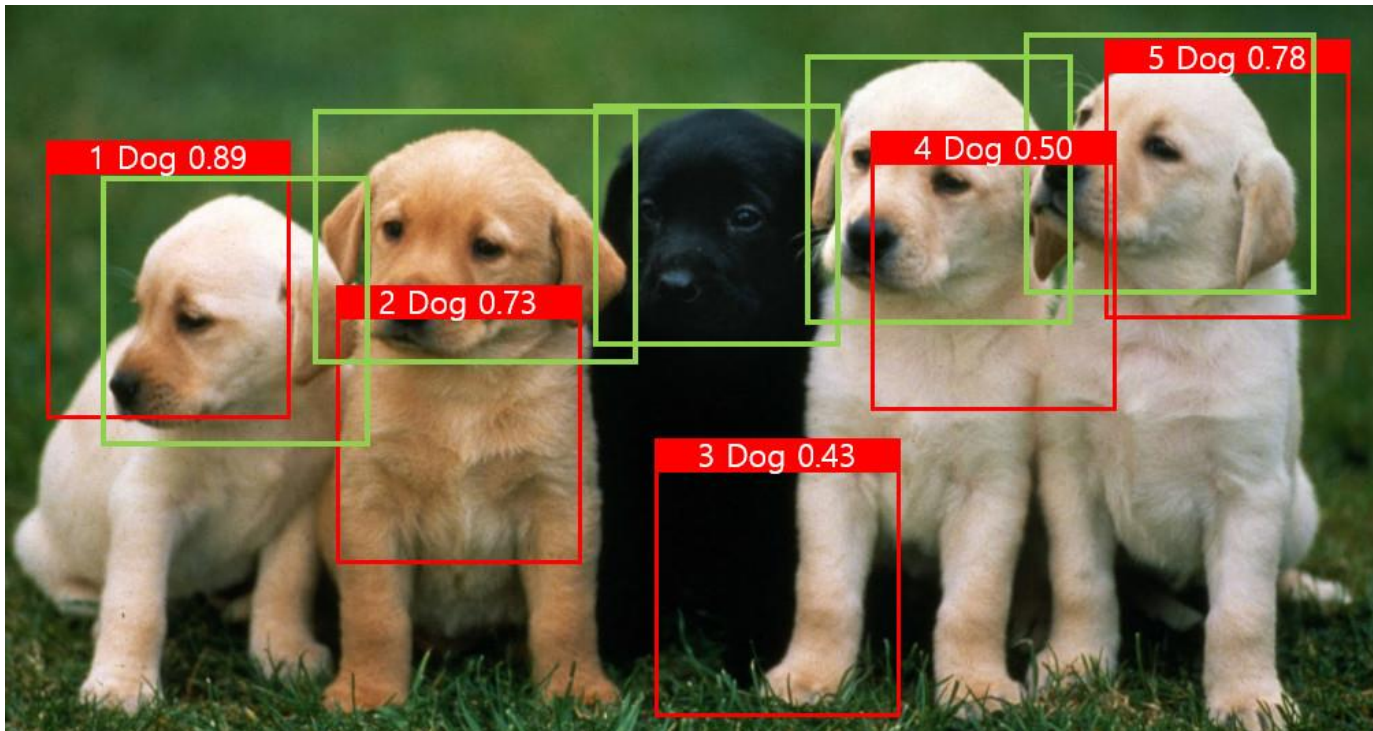
where,  
 TP be the number of true positives,  
 FP be the number of false positives,  
 FN be the number of false negatives





# Evaluation

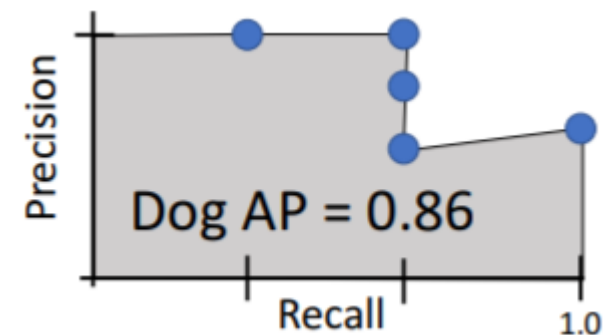
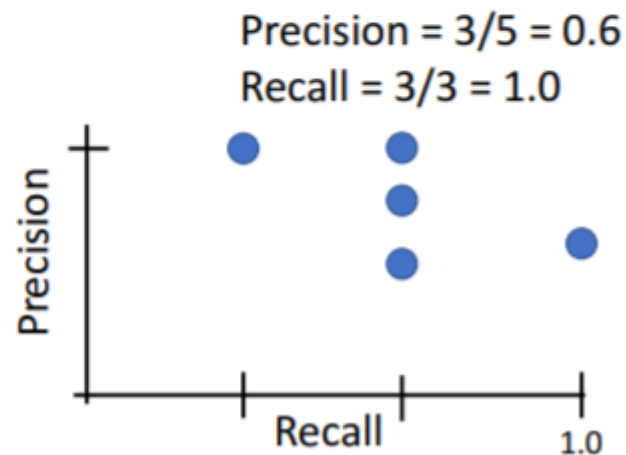
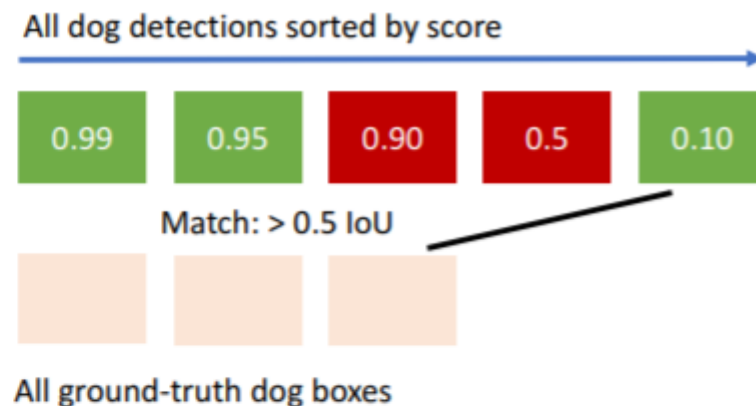
- Confidence score, distance



# Evaluation

- Evaluating Object Detectors:  
Mean Average Precision (mAP)

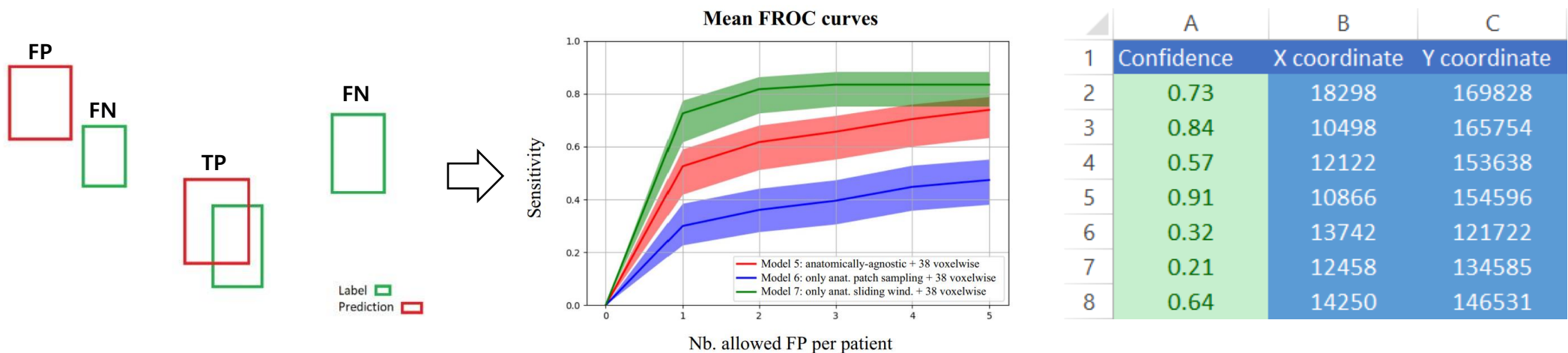
1. Run object detector on all test images (with NMS)
2. For each category, compute Average Precision (AP) = area under Precision vs Recall Curve
  1. For each detection (highest score to lowest score)
    1. If it matches some GT box with IoU > 0.5, mark it as positive and eliminate the GT
    2. Otherwise mark it as negative
    3. Plot a point on PR Curve



# Evaluation

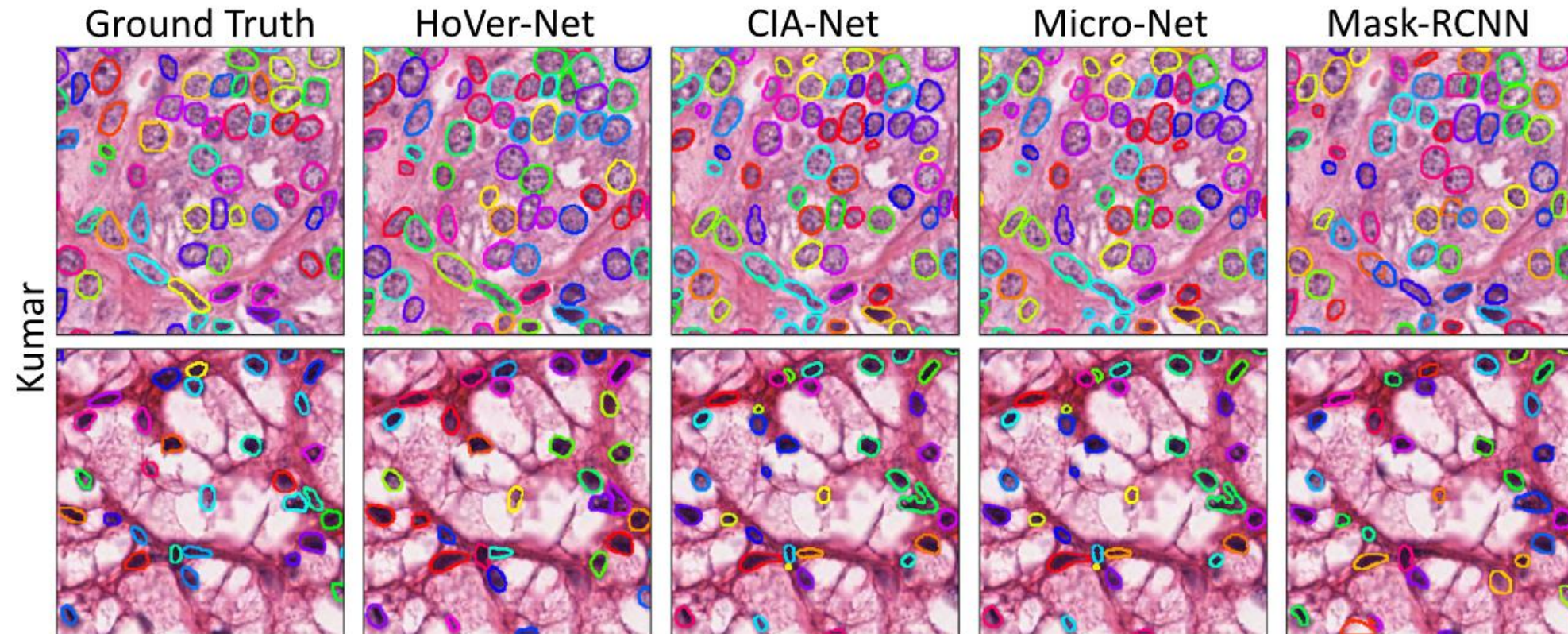
- Free-response receiver operating characteristic (FROC)

1. Slide-based Evaluation: The merits of the algorithms will be assessed for discriminating between slides containing metastasis and normal slides. Receiver operating characteristic (ROC) analysis at the slide level will be performed and the measure used for comparing the algorithms will be the area under the ROC curve (AUC).
2. **Lesion-based Evaluation:** For the lesion-based evaluation, free-response receiver operating characteristic (FROC) curve will be used. The FROC curve is defined as the plot of sensitivity versus the average number of false-positives per image.





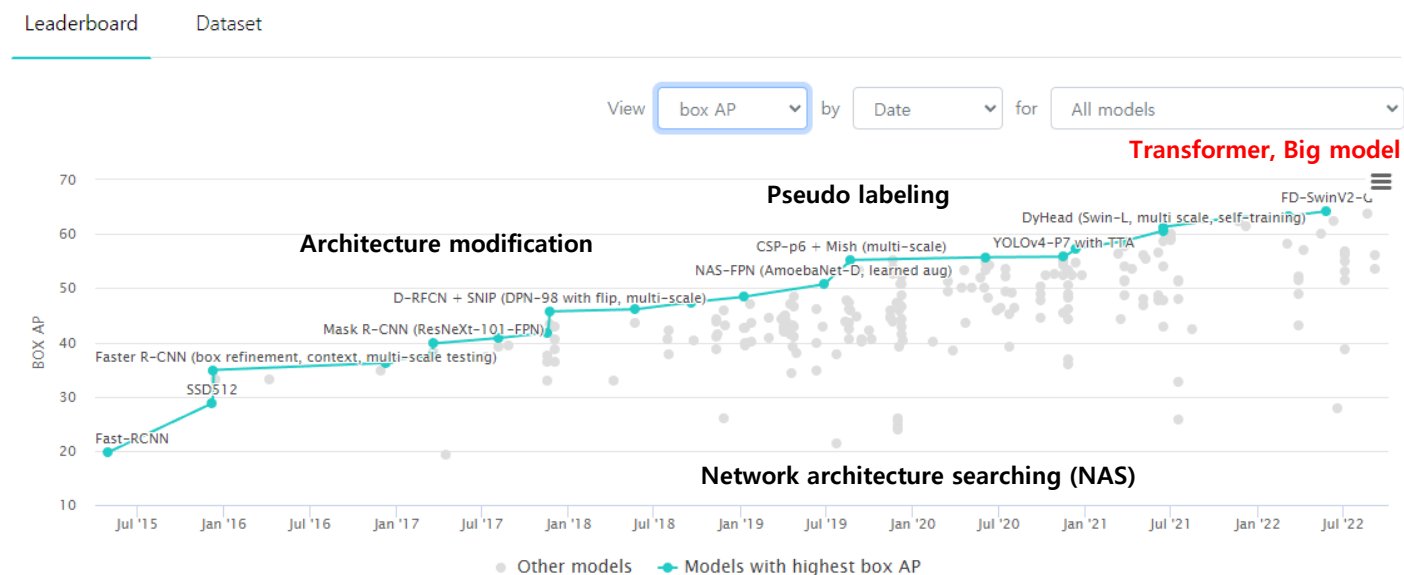
# Evaluation



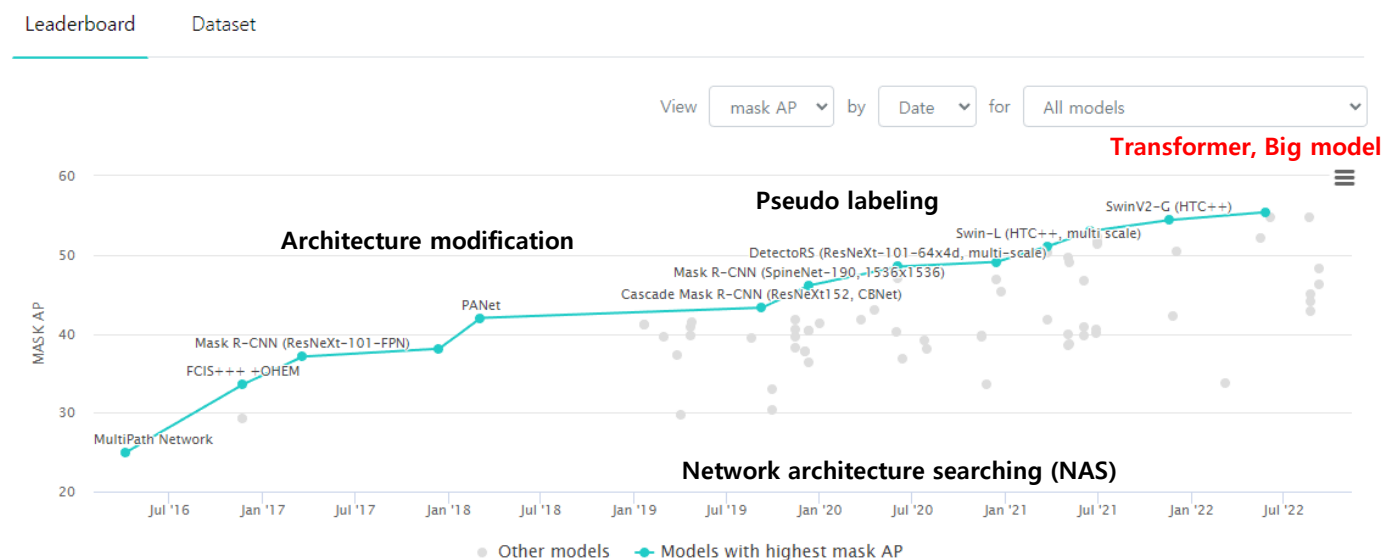
**Hover-Net: Simultaneous segmentation and classification of nuclei in multi-tissue histology images,**  
**Medical image analysis, 2019**

# Trend

## Object Detection on COCO test-dev



## Instance Segmentation on COCO test-dev



# Conclusion

- Object detection 모델 성능은 특정 트렌드를 가지며, 지금도 계속 향상되고 있다.
- 모델 개발자 입장 : 모델 한계점을 파악하고, 개선함으로써 성능, 속도를 향상
- 모델 사용자 입장 : 우리 데이터셋 분석 및 현재 테스트에 맞는 모델 선정이 중요!
- Segmentation vs object detection vs instance segmentation



# Hands on

- 1. Google colab (<https://colab.research.google.com/>)
- 2. 노트 열기 – Github – 아래 깃허브 주소 입력  
[https://github.com/kevinkwshin/Handson\\_detection](https://github.com/kevinkwshin/Handson_detection)