

Pontificia Universidad Católica de Chile
Escuela de Ingeniería
Departamento de Ciencia de la Computación



Sistemas Urbanos Inteligentes

Aprendizaje Multitarea (Multitask Learning)

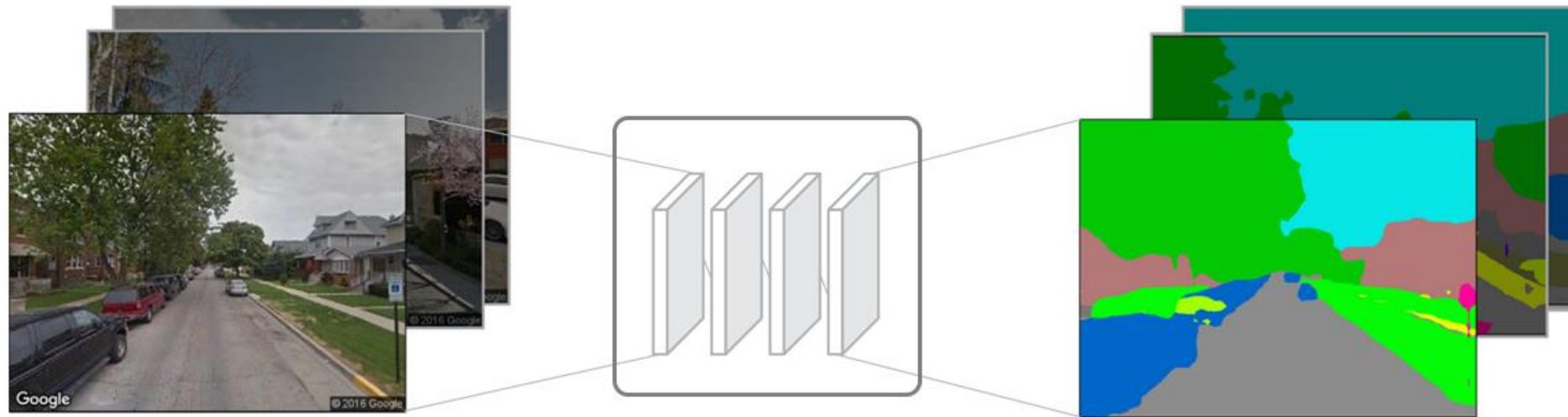
Hans Löbel

Dpto. Ingeniería de Transporte y Logística
Dpto. Ciencia de la Computación

¿Qué utilidad nos pueden entregar las CNN en contextos urbanos?



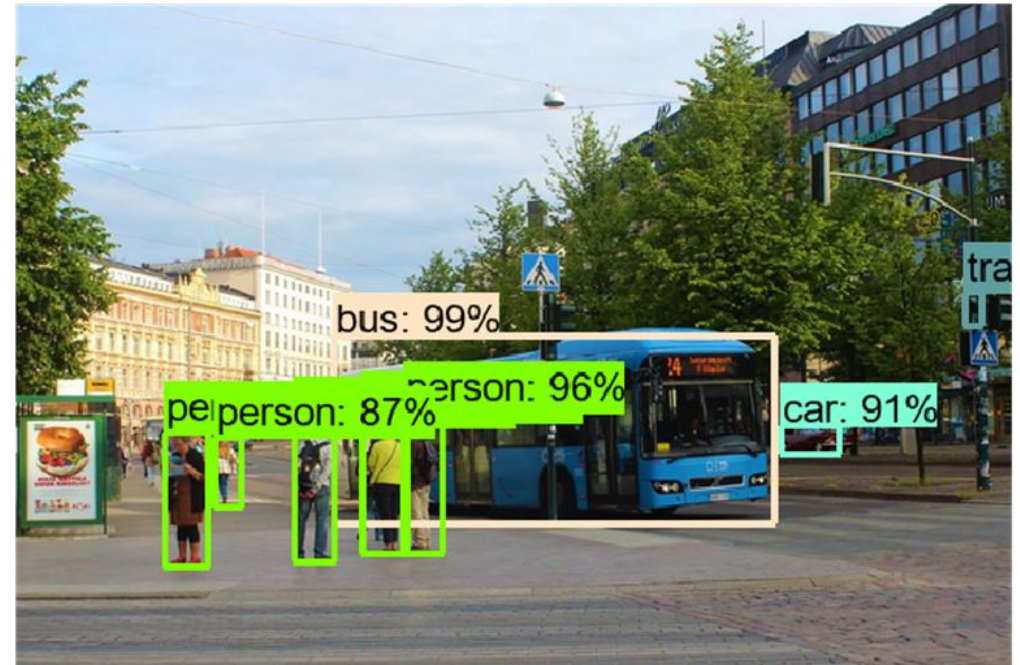
¿Cómo explicar la percepción?



Segmentación semántica de imágenes

¿Cómo mejorar la caracterización de las imágenes?

Segmentación + detección de objetos



Necesitamos subir en la escala de problemas de visión por computador, hasta la **detección de objetos**

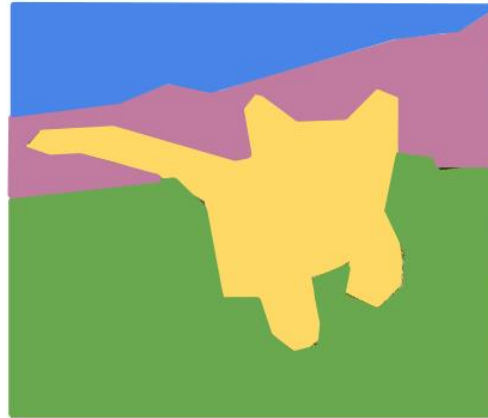
Classification



CAT

No spatial extent

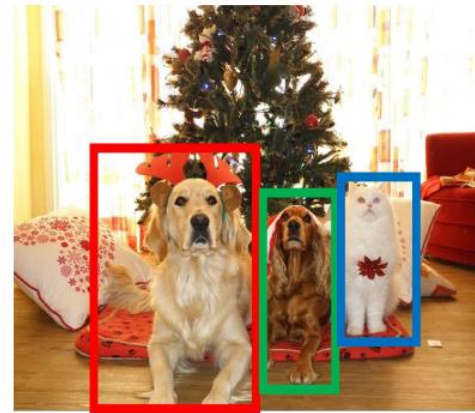
Semantic Segmentation



**GRASS, CAT,
TREE, SKY**

No objects, just pixels

Object Detection



DOG, DOG, CAT

Multiple Object

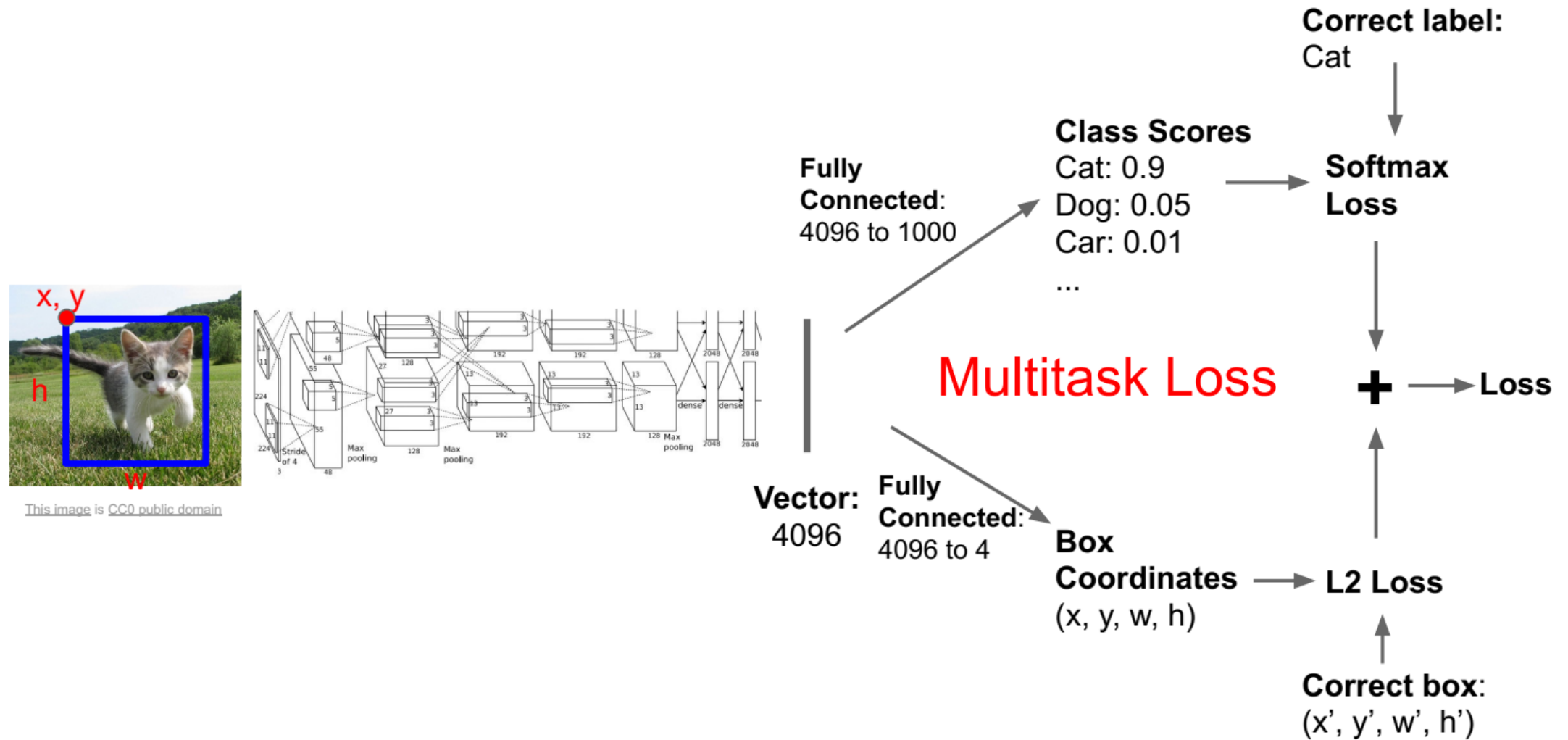
Instance Segmentation



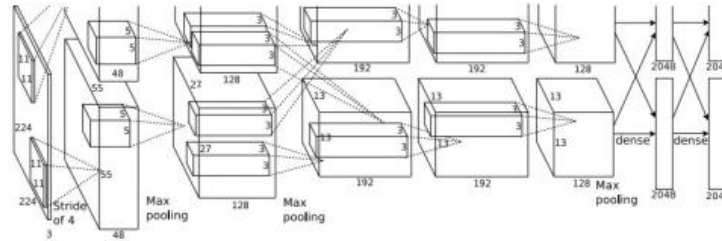
DOG, DOG, CAT

[This image is CC0 public domain](#)

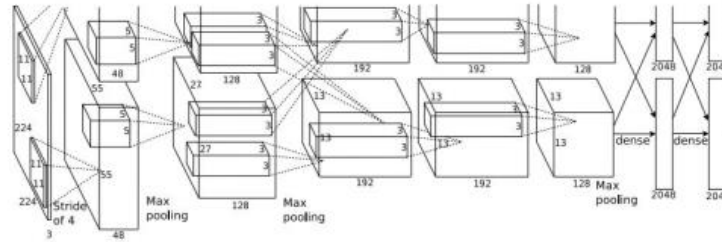
Para esto, necesitamos realizar dos predicciones de manera simultánea



Para esto, necesitamos realizar dos predicciones de manera simultánea



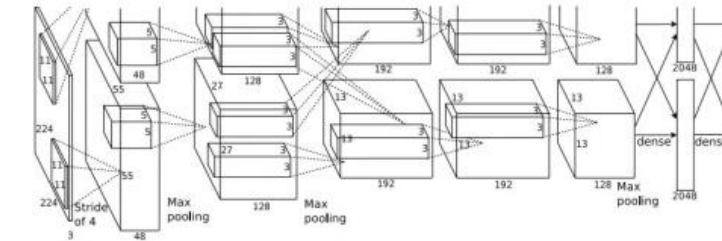
CAT: (x, y, w, h)



DOG: (x, y, w, h)

DOG: (x, y, w, h)

CAT: (x, y, w, h)



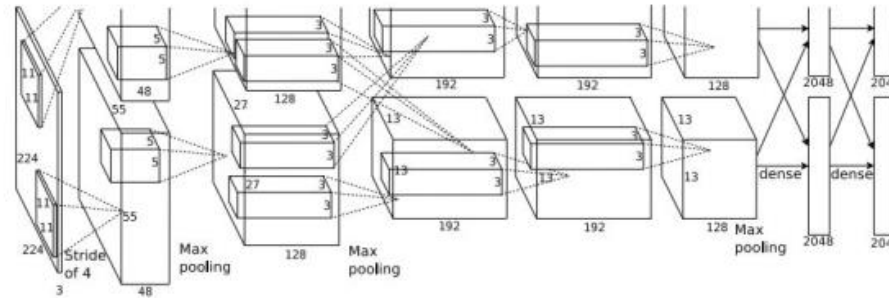
DUCK: (x, y, w, h)

DUCK: (x, y, w, h)

....

¿Qué problema nos genera este esquema de detección?

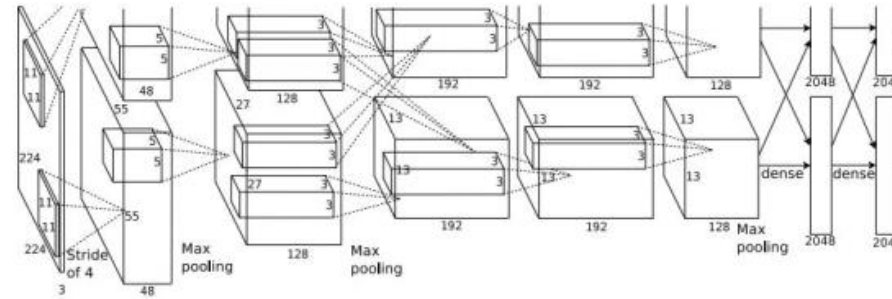
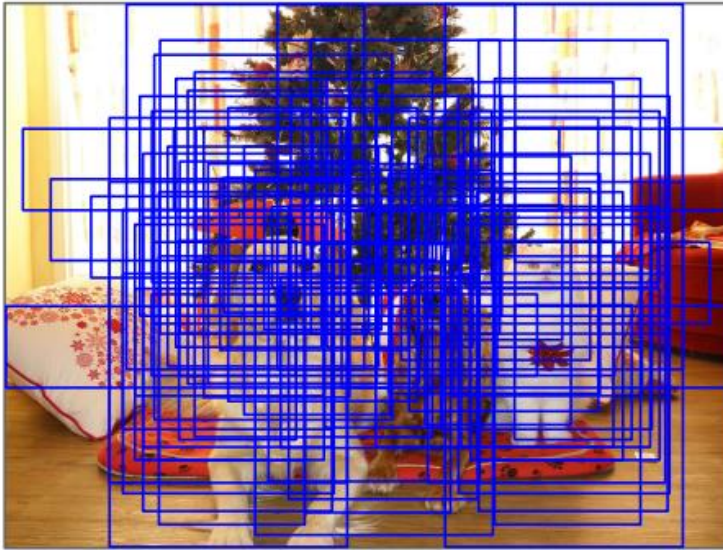
Podemos utilizar una ventana deslizante para analizar todas las posiciones



Dog? NO
Cat? NO
Background? YES

¿Qué problema nos genera este esquema de detección?

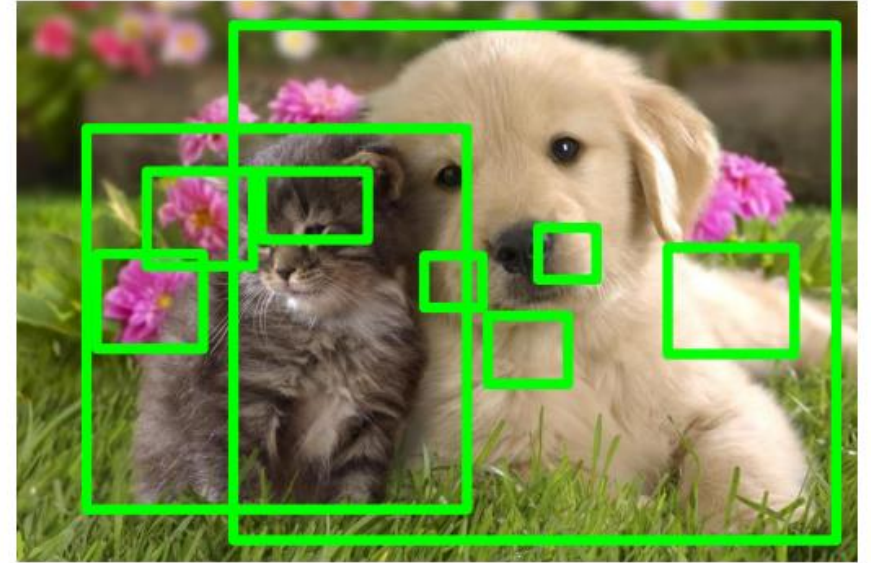
Podemos utilizar una ventana deslizante
para analizar todas las posiciones



Dog? NO
Cat? YES
Background? NO

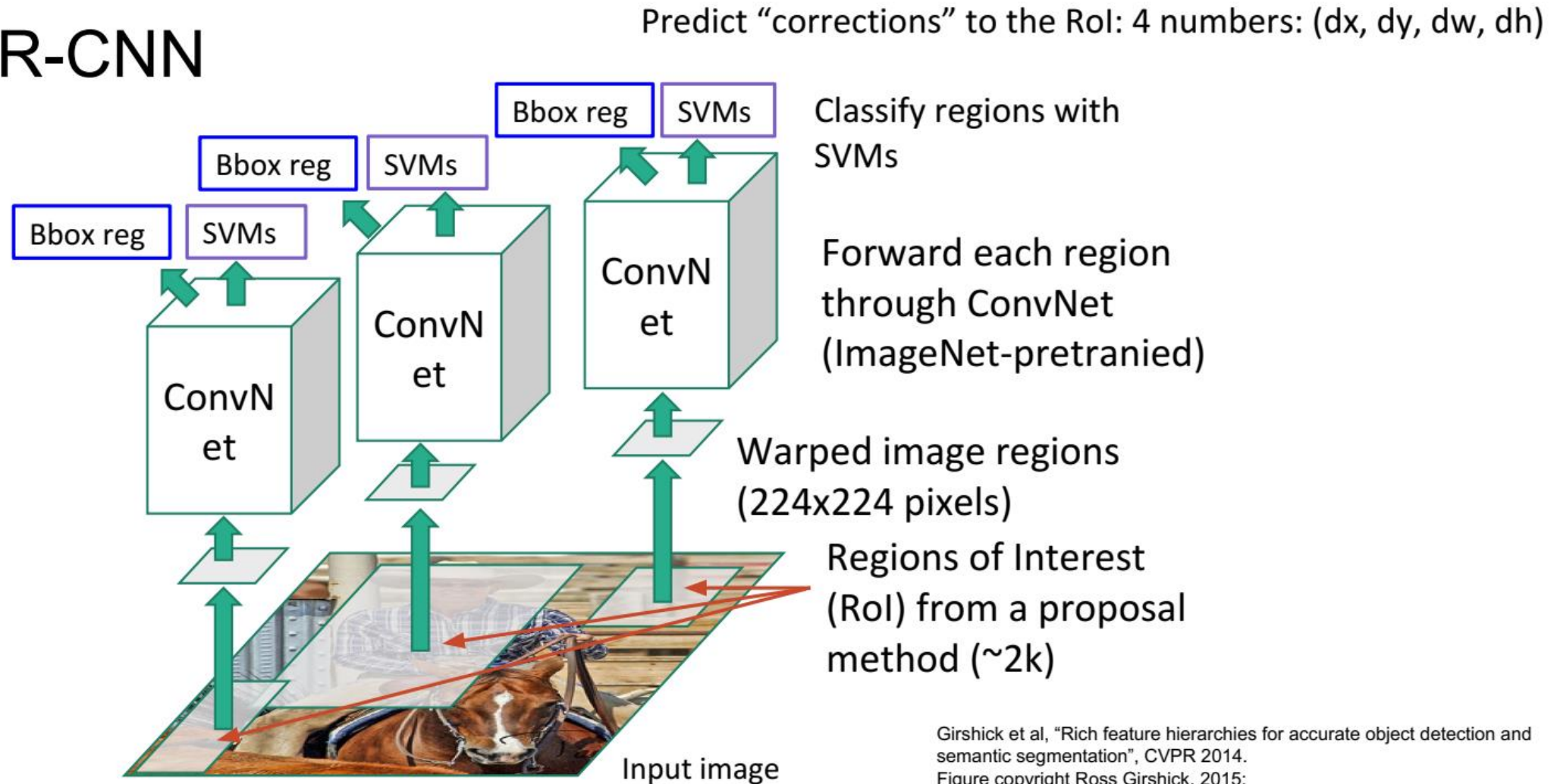
Necesitamos generar menos ventanas candidatas

Region Proposals entregan rápidamente
ventanas donde podrían haber objetos



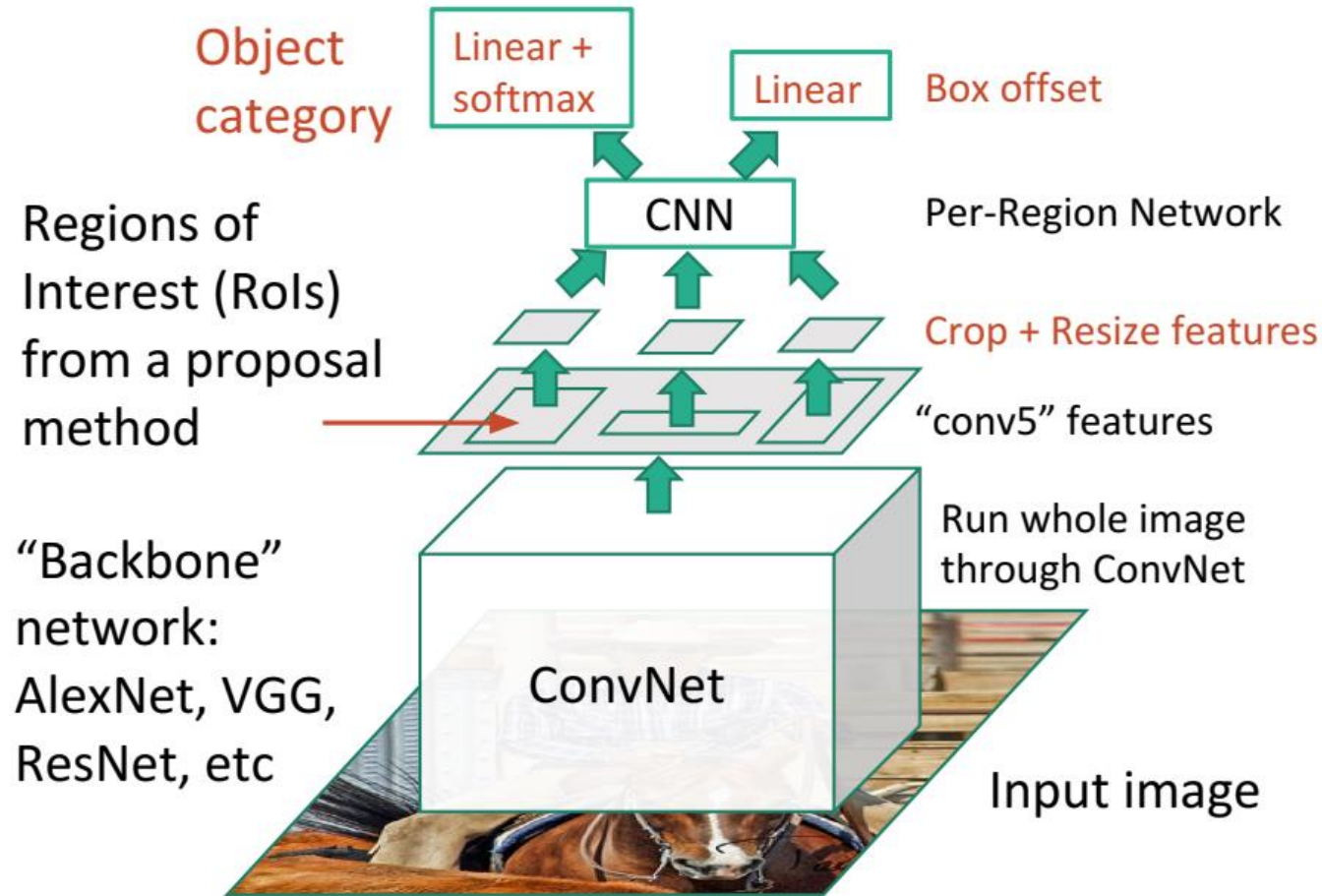
Alexe et al, "Measuring the objectness of image windows", TPAMI 2012
Uijlings et al, "Selective Search for Object Recognition", IJCV 2013
Cheng et al, "BING: Binarized normed gradients for objectness estimation at 300fps", CVPR 2014
Zitnick and Dollar, "Edge boxes: Locating object proposals from edges",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

R-CNN



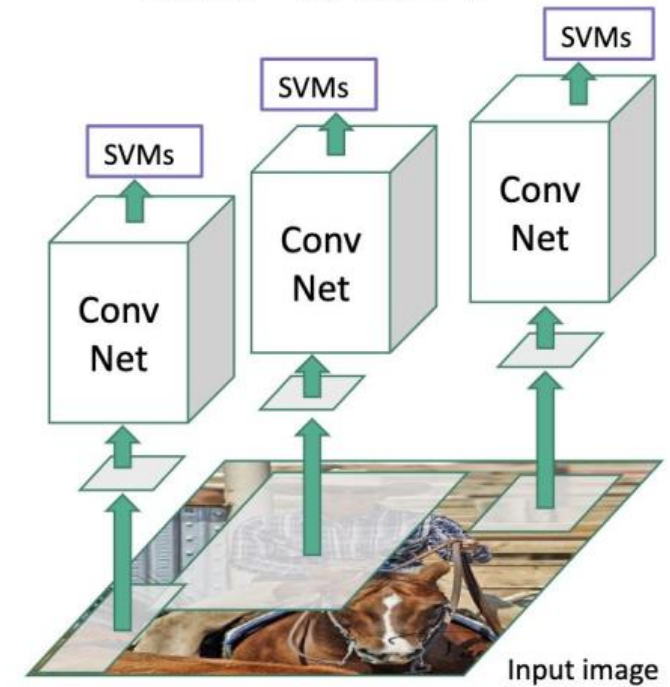
¿Cuál es el punto débil de este esquema?

Fast R-CNN

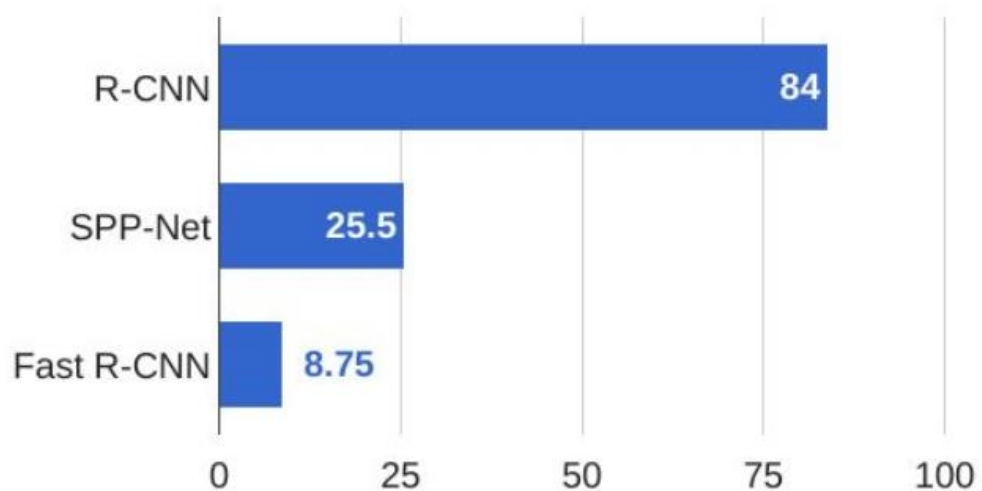


Girshick, "Fast R-CNN", ICCV 2015. Figure copyright Ross Girshick, 2015;

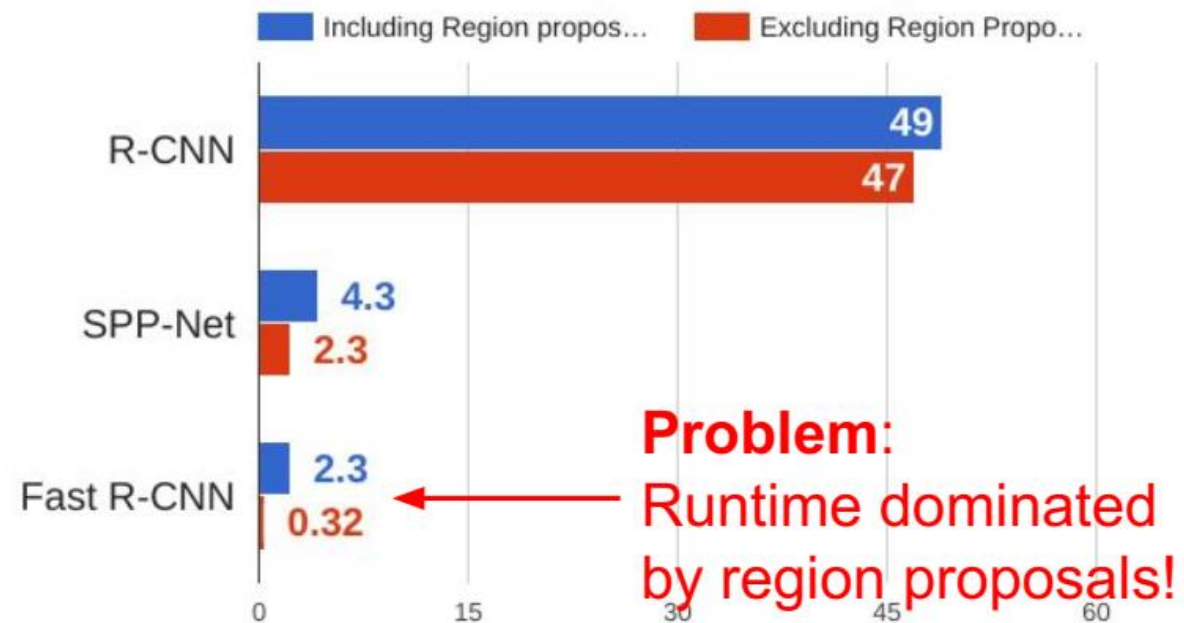
"Slow" R-CNN



Training time (Hours)



Test time (seconds)



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

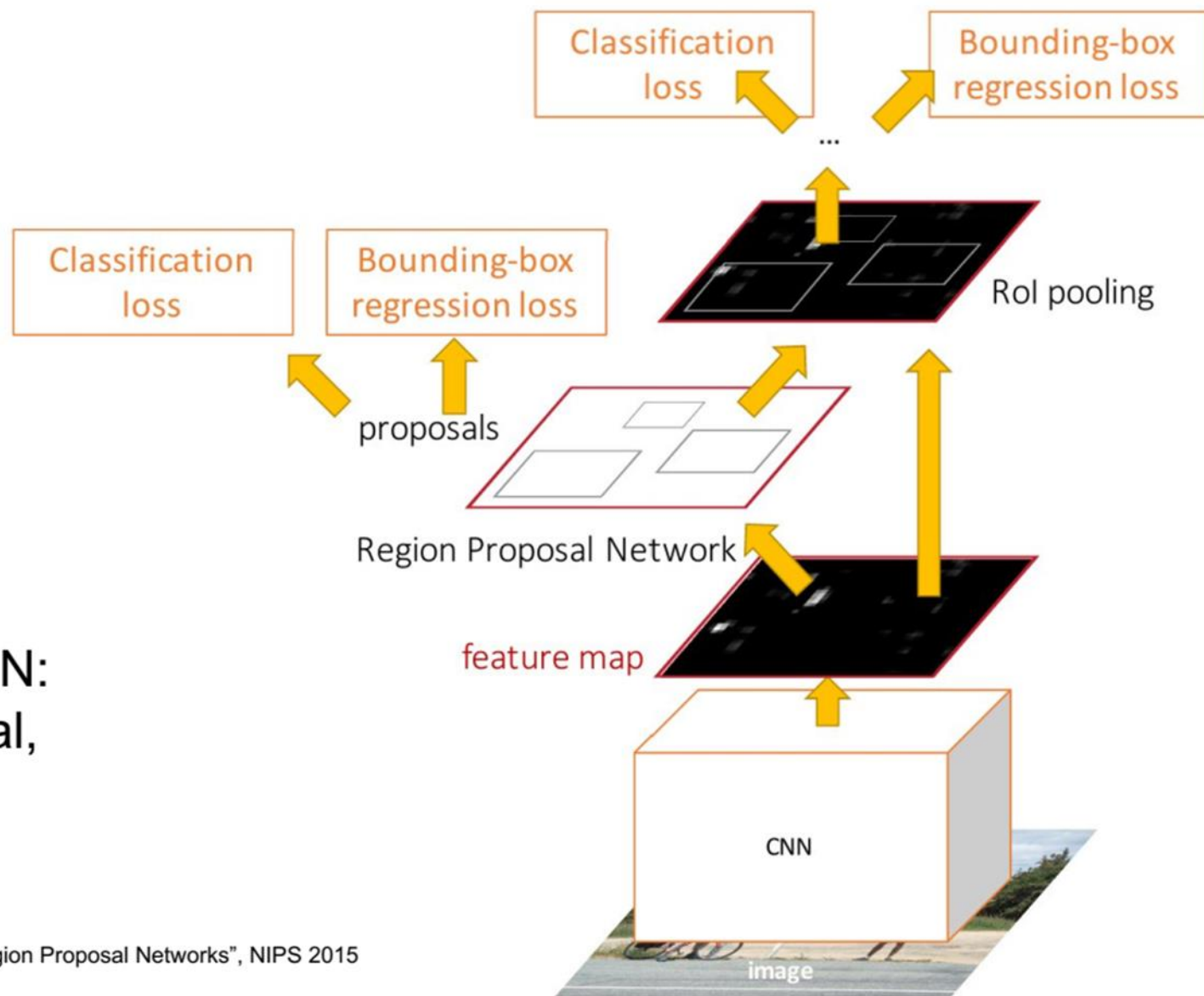
Girshick, "Fast R-CNN", ICCV 2015

Faster R-CNN:

Make CNN do proposals!

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

Otherwise same as Fast R-CNN:
Crop features for each proposal,
classify each one

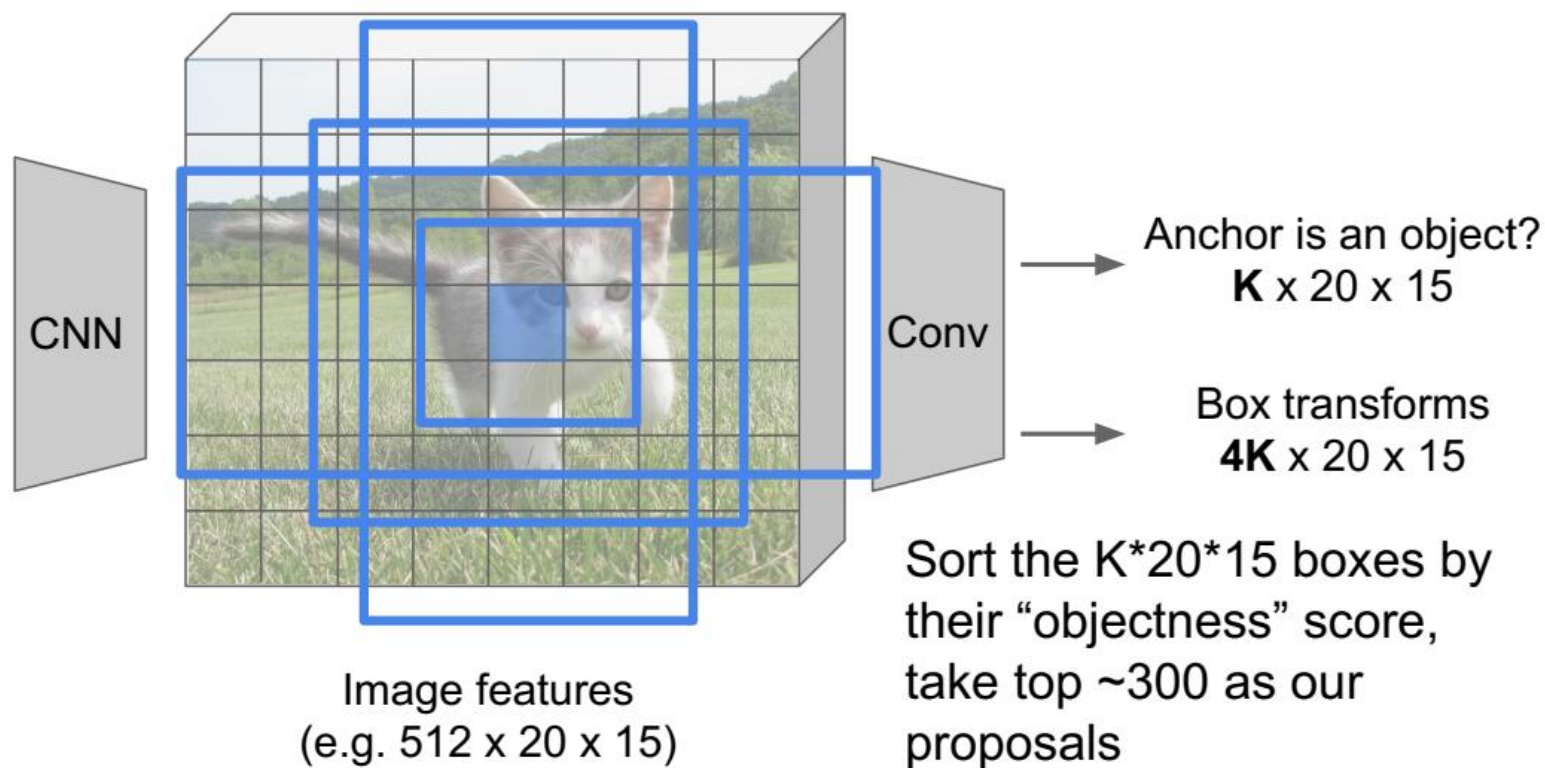


Region Proposal Network

In practice use K different anchor boxes of different size / scale at each point



Input Image
(e.g. 3 x 640 x 480)

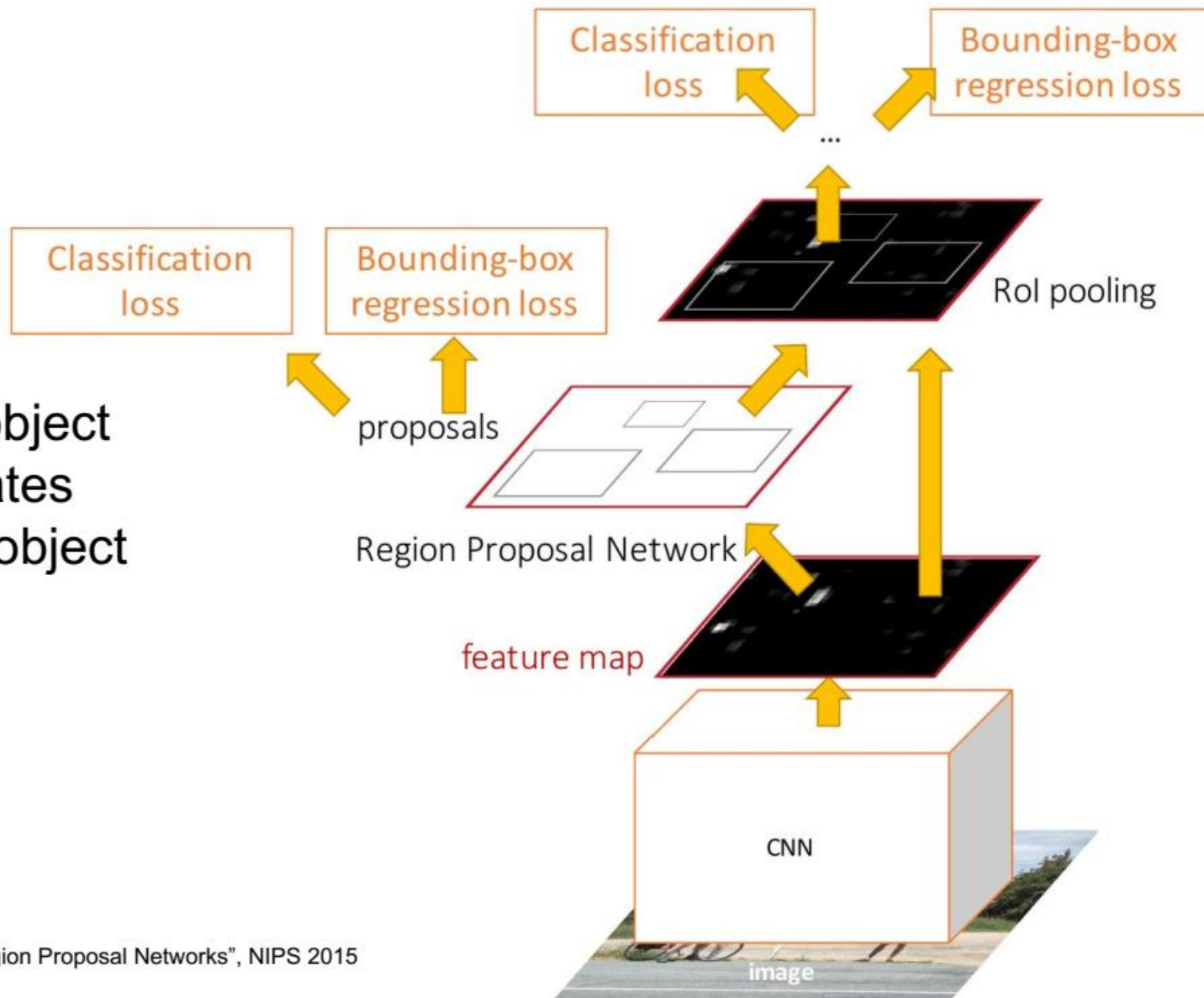


Faster R-CNN:

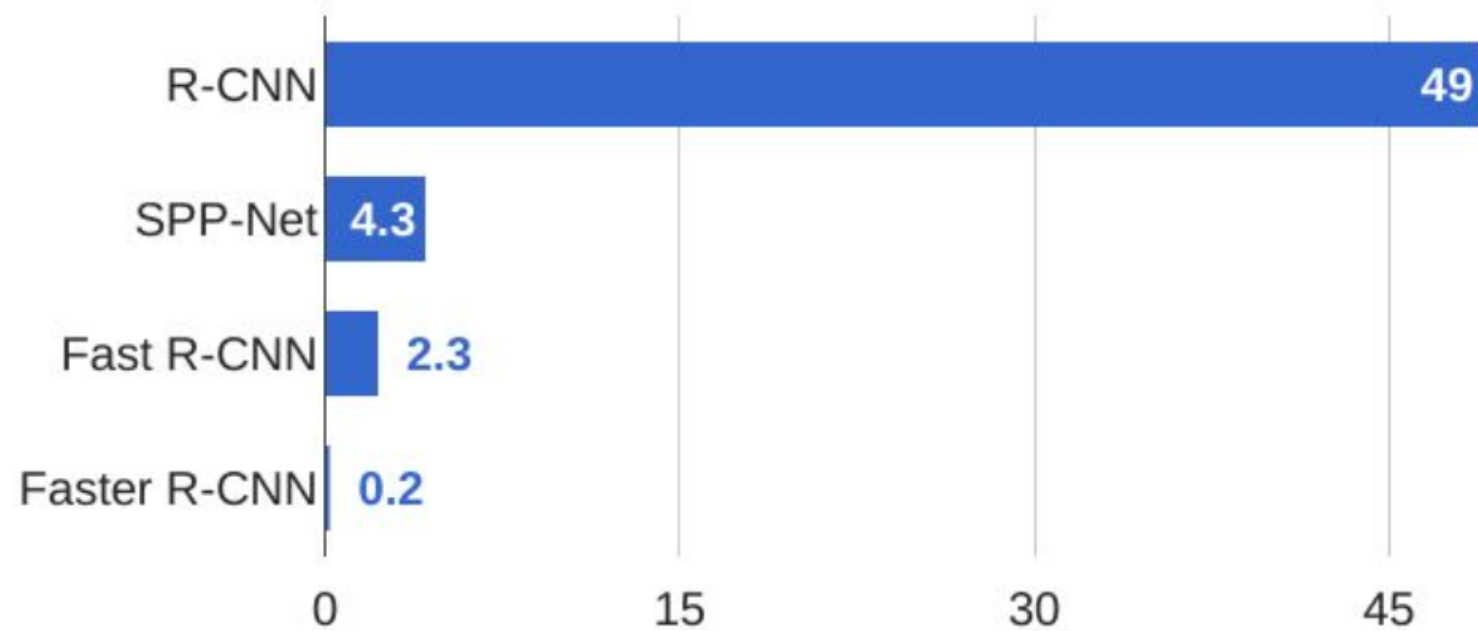
Make CNN do proposals!

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



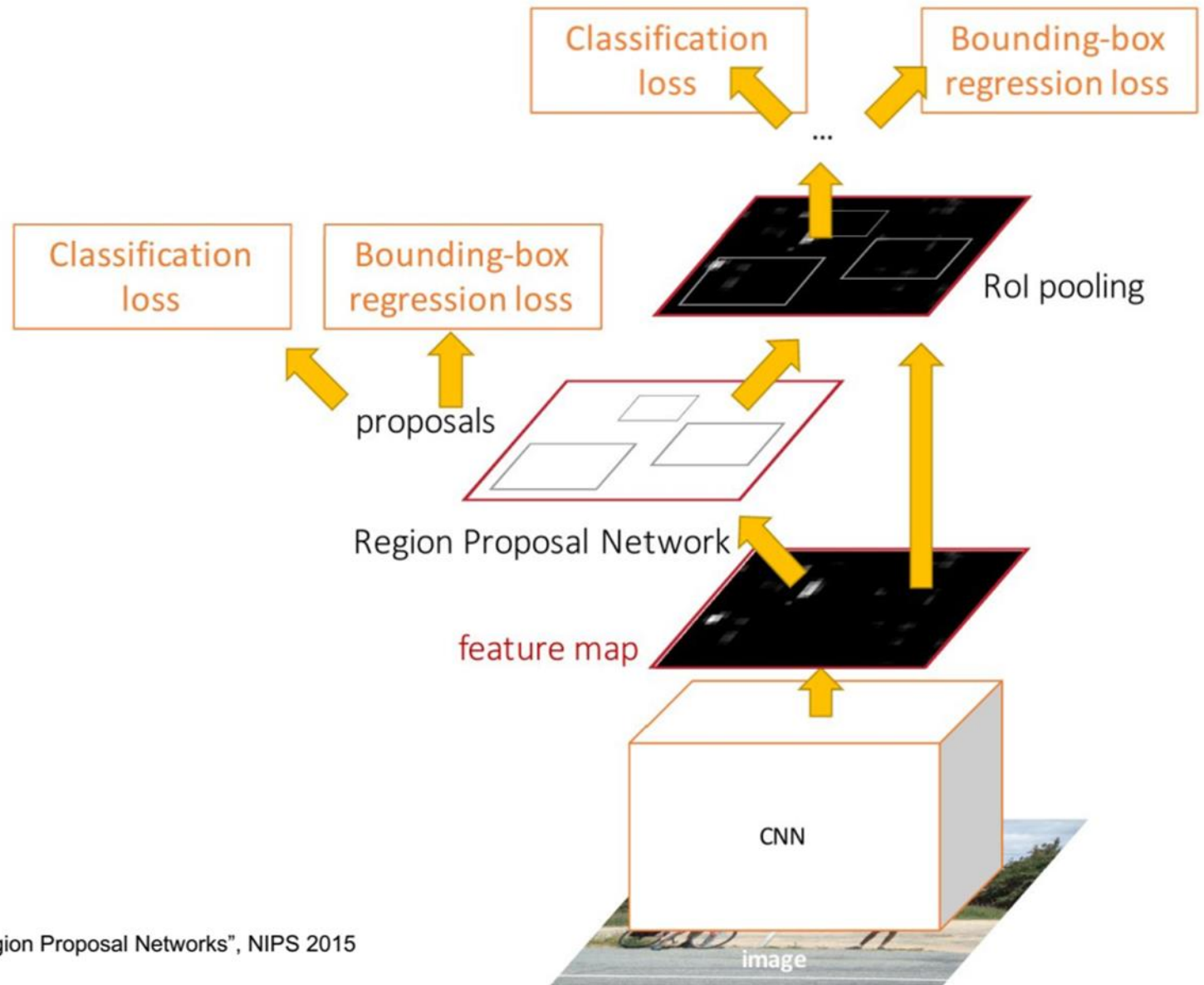
R-CNN Test-Time Speed



Faster R-CNN:

Make CNN do proposals!

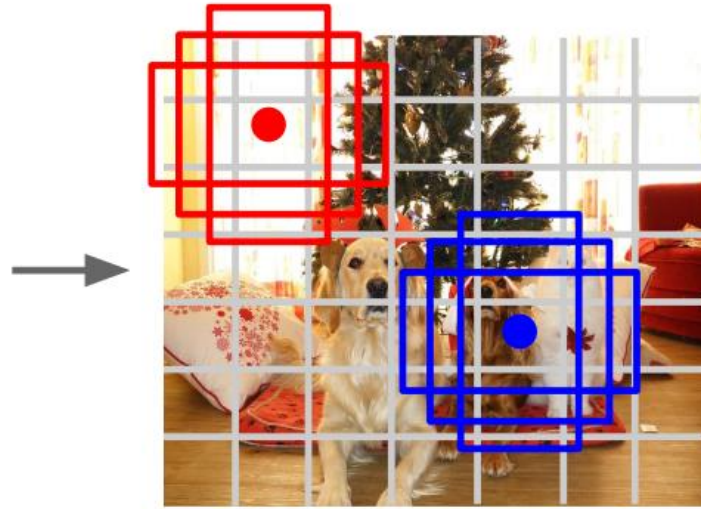
¿Necesitamos realmente la segunda subred?



Single-Stage Object Detectors: YOLO / SSD / RetinaNet



Input image
 $3 \times H \times W$



Divide image into grid
 7×7

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:
($dx, dy, dh, dw, \text{confidence}$)
- Predict scores for each of C classes (including background as a class)
- Looks a lot like RPN, but category-specific!

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
Lin et al, "Focal Loss for Dense Object Detection", ICCV 2017



<https://youtu.be/YDFf-TqJOFE>

Para lograr la segmentación a nivel de instancia, utilizamos las mismas ideas

Classification



CAT

No spatial extent

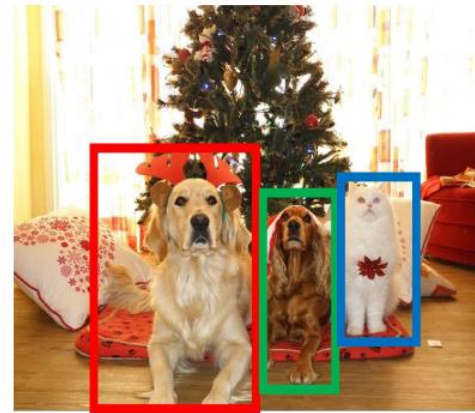
Semantic Segmentation



**GRASS, CAT,
TREE, SKY**

No objects, just pixels

Object Detection



DOG, DOG, CAT

Multiple Object

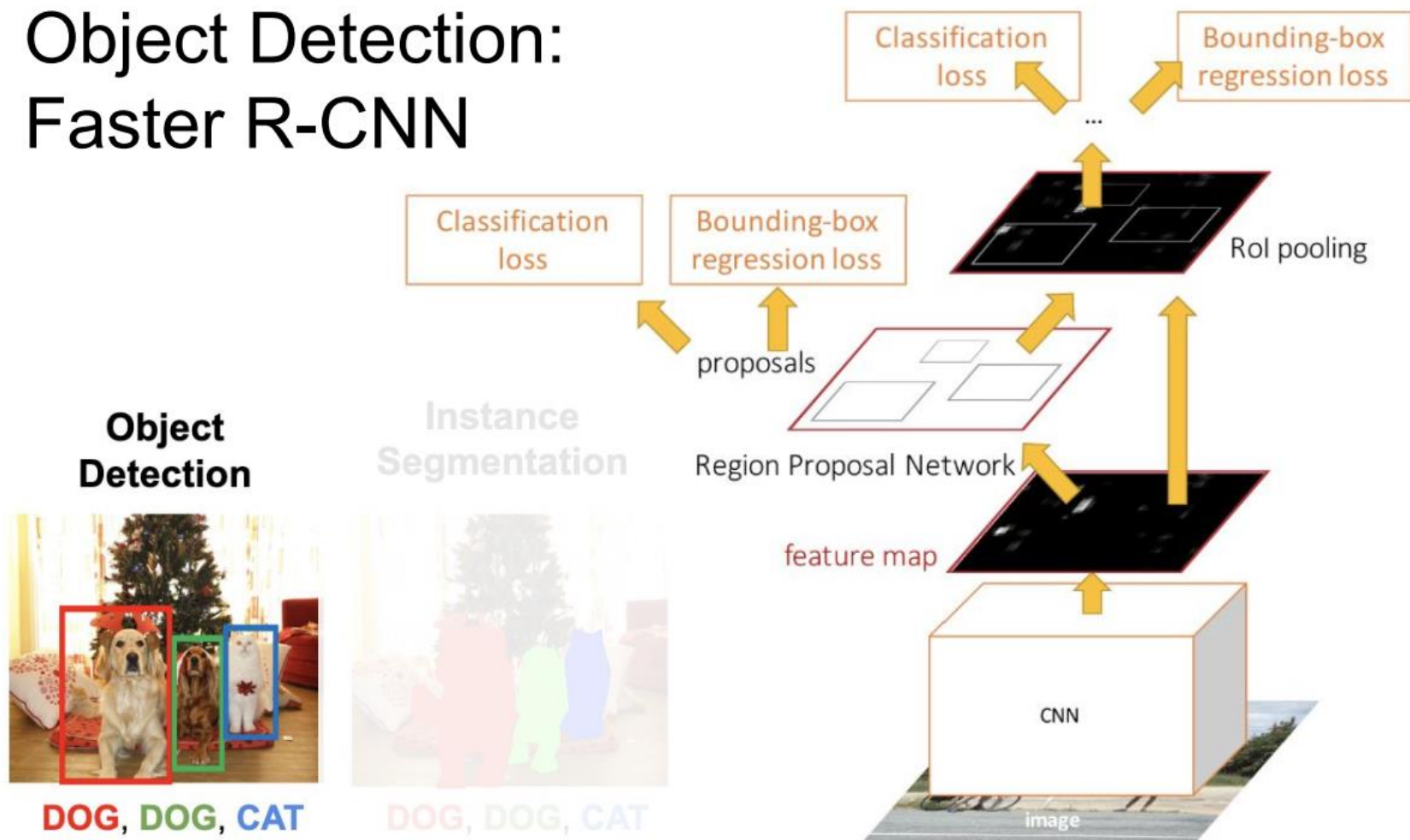
Instance Segmentation



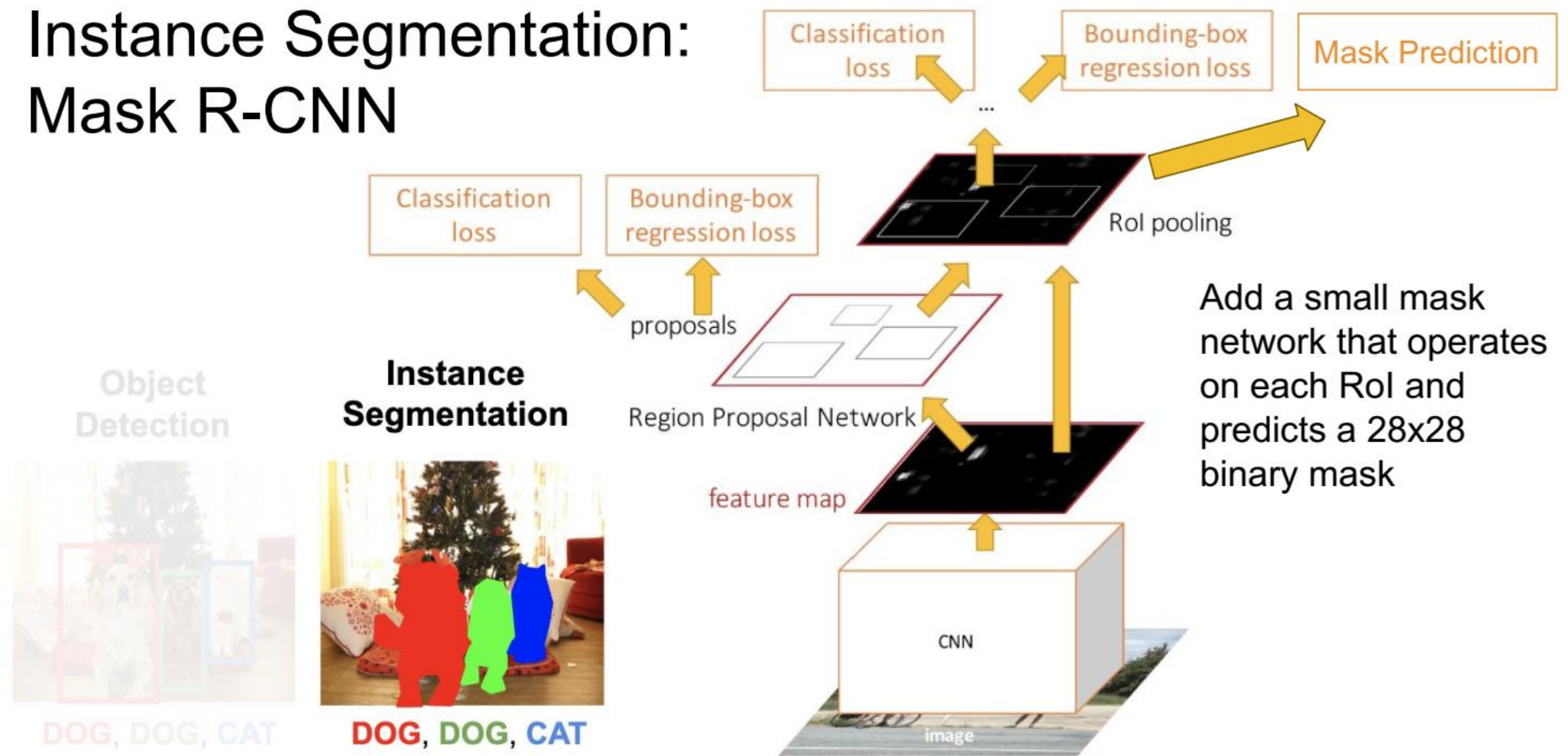
DOG, DOG, CAT

[This image is CC0 public domain](#)

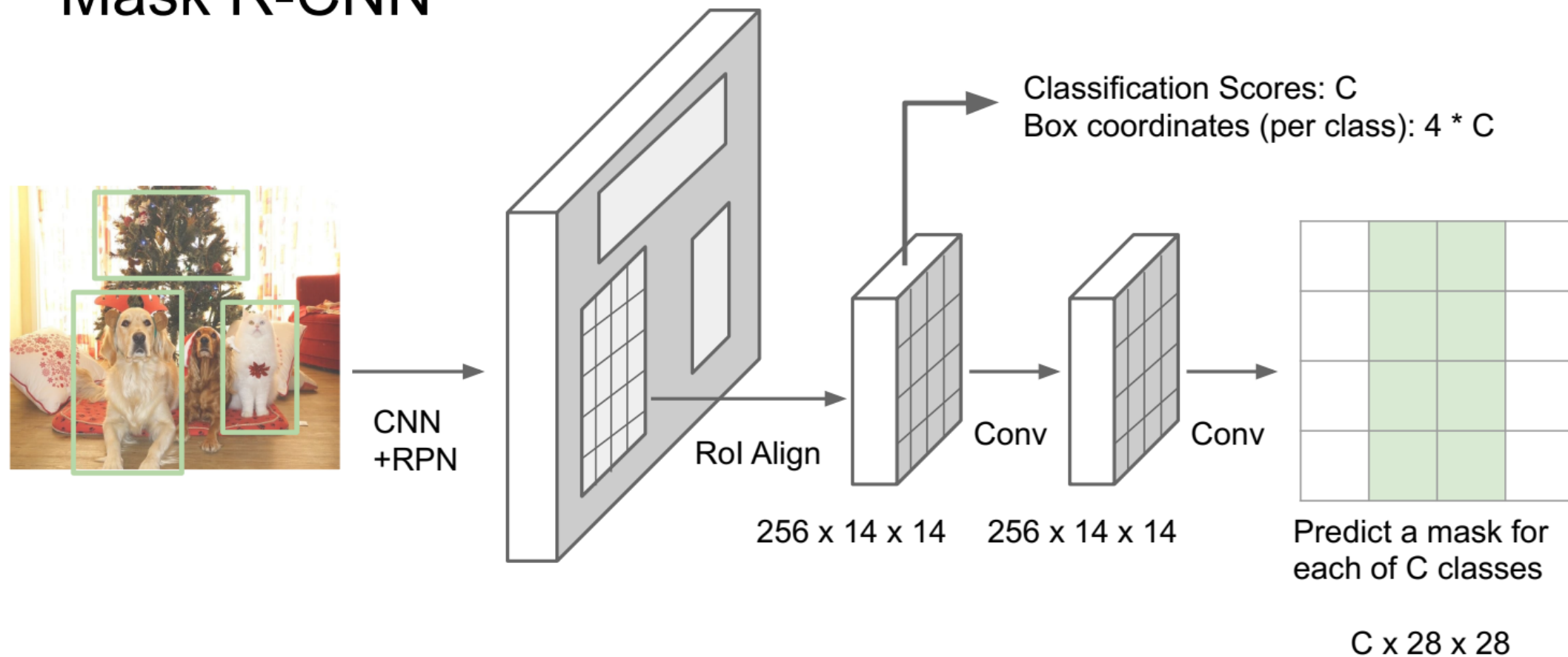
Object Detection: Faster R-CNN

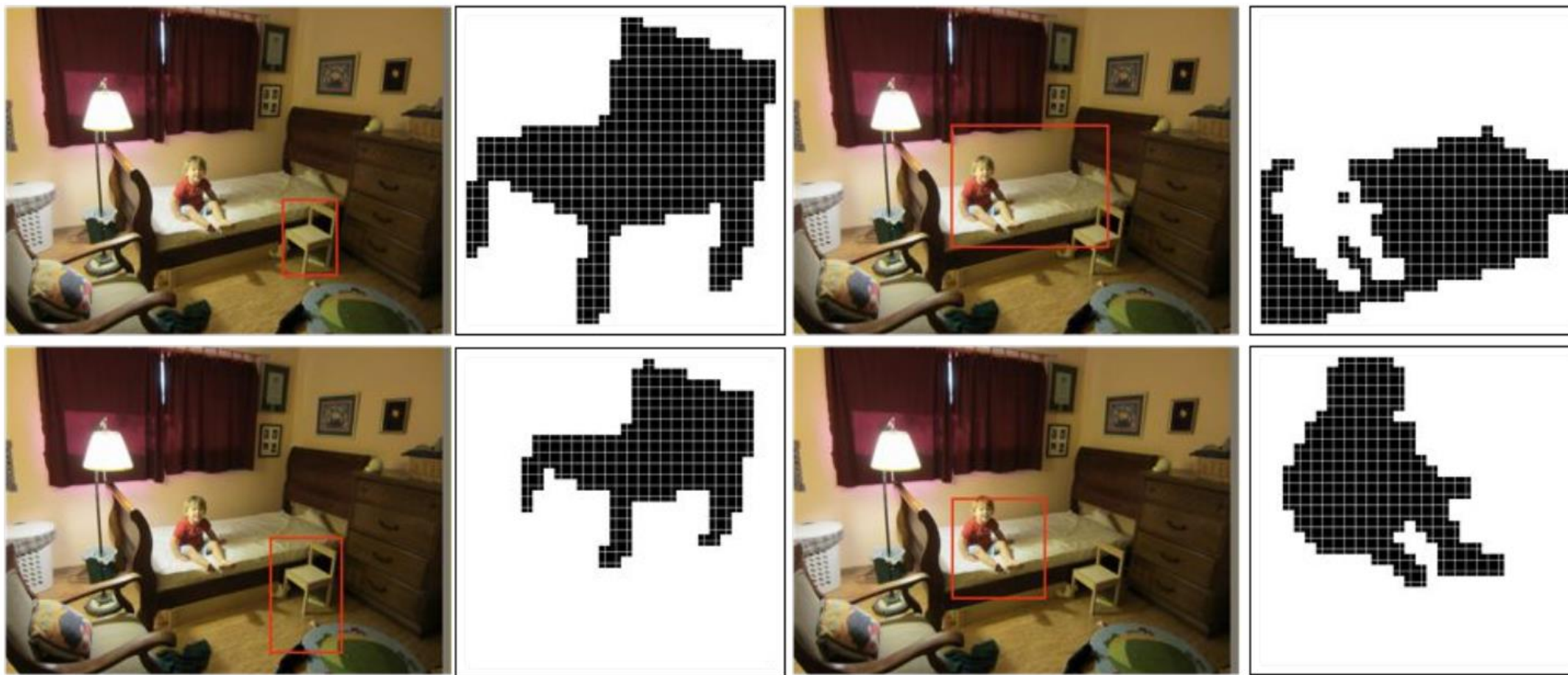


Instance Segmentation: Mask R-CNN

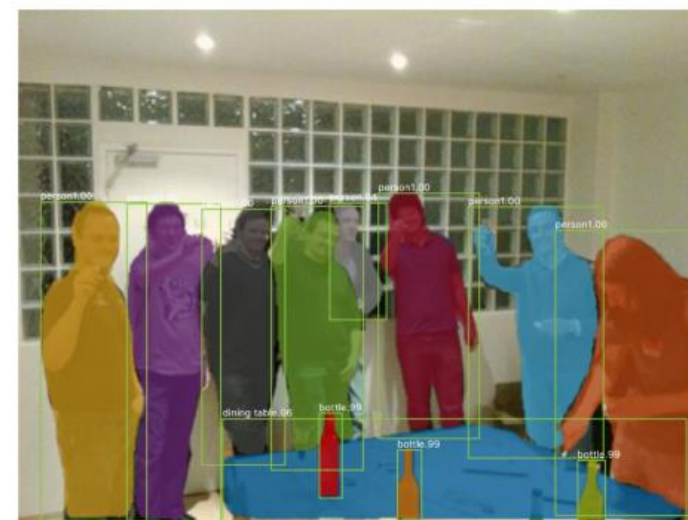
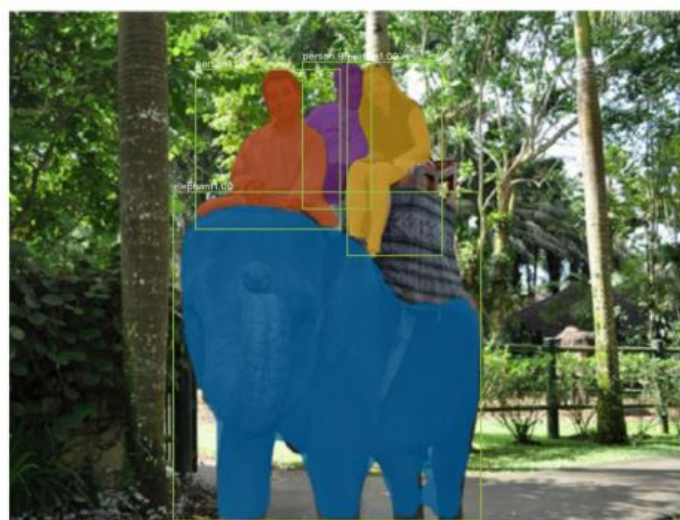


Mask R-CNN

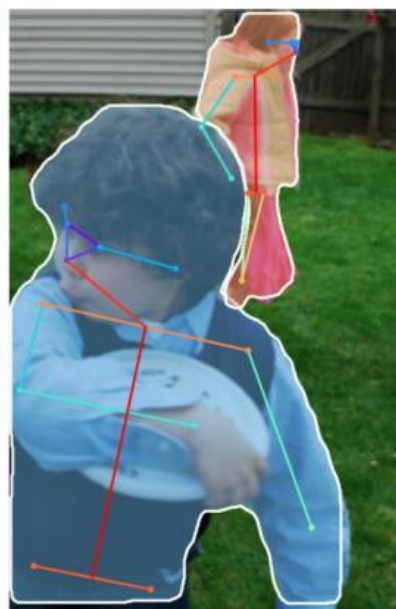
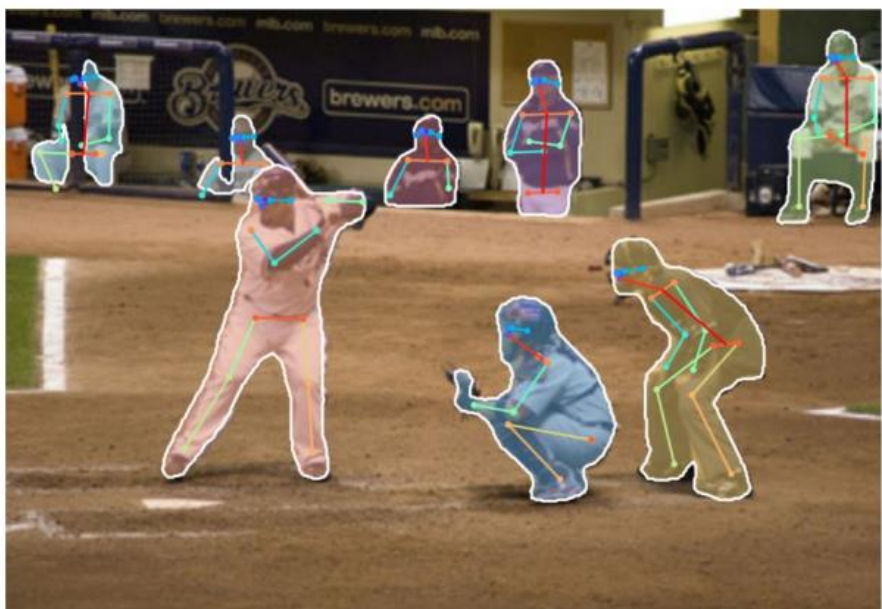




He et al, "Mask R-CNN", ICCV 2017



He et al, "Mask R-CNN", ICCV 2017

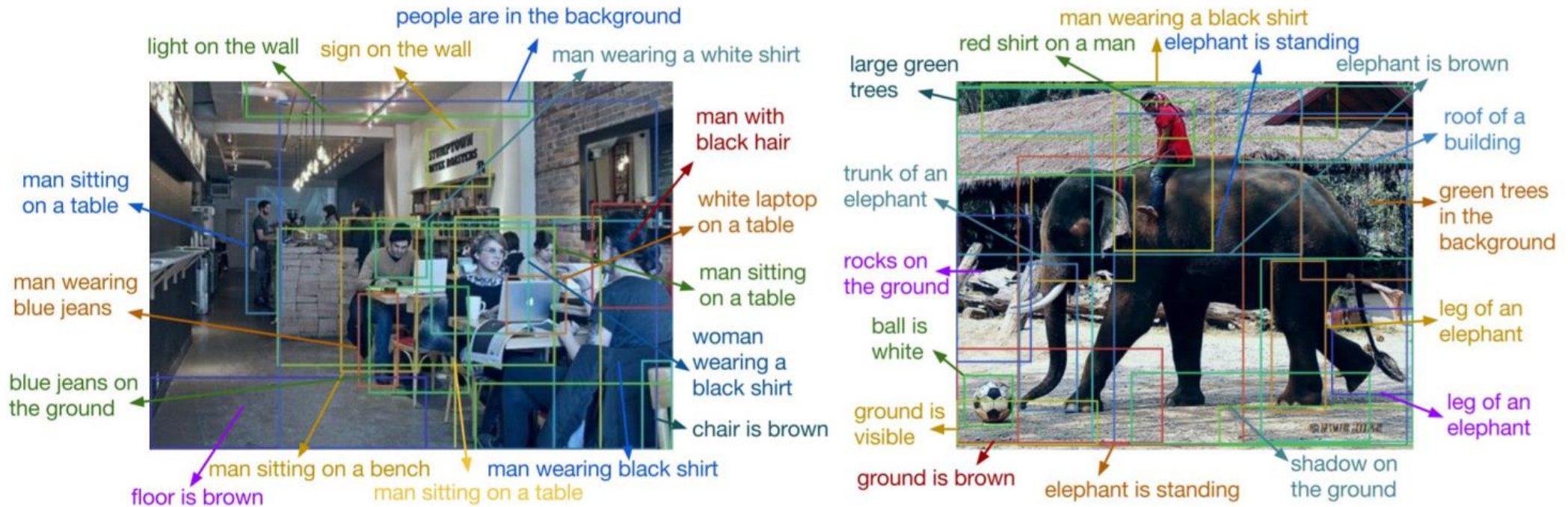


He et al, "Mask R-CNN", ICCV 2017

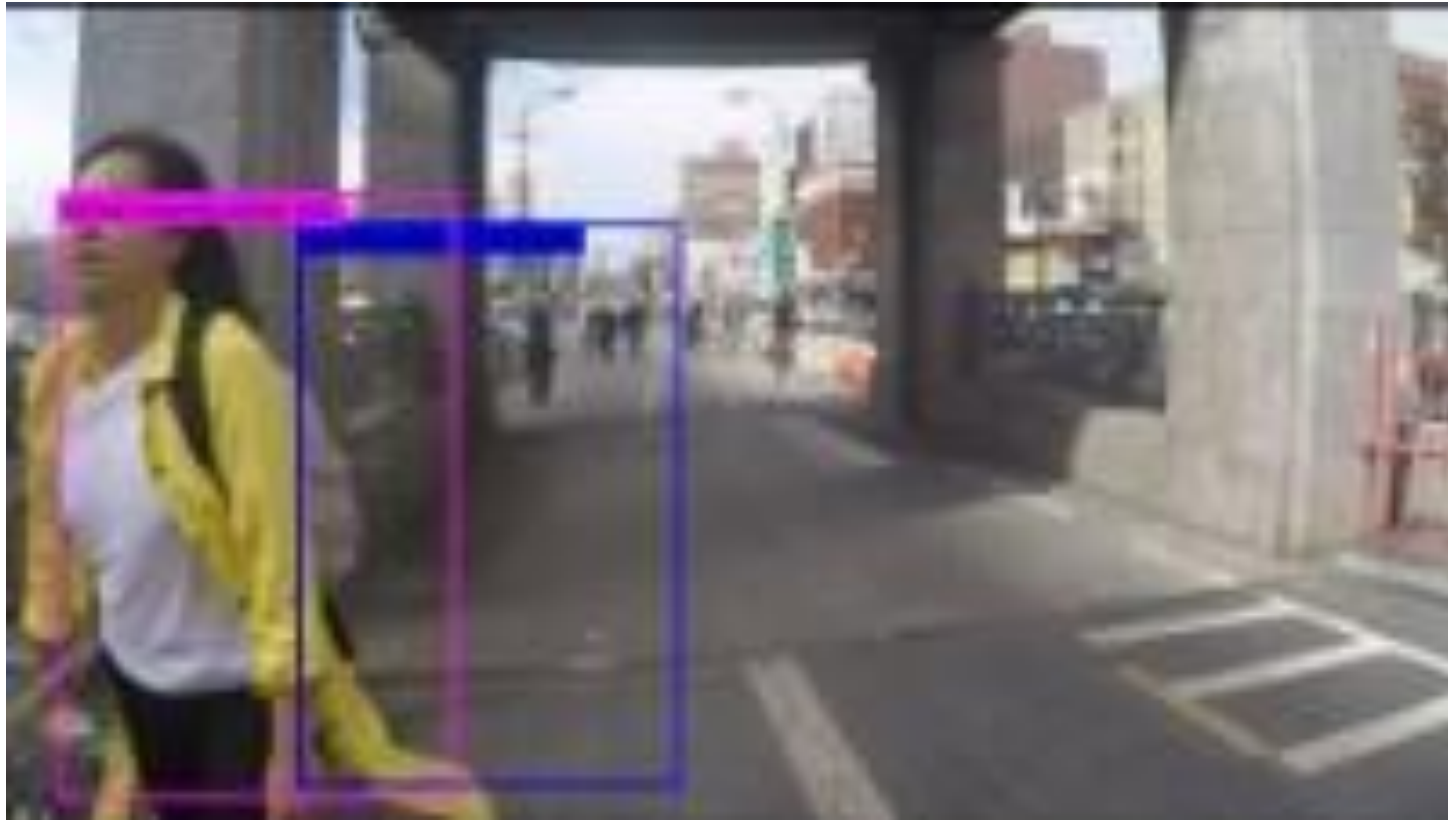


<https://youtu.be/OOT3UIXZztE>

Object Detection + Captioning = Dense Captioning



Johnson, Karpathy, and Fei-Fei, "DenseCap: Fully Convolutional Localization Networks for Dense Captioning", CVPR 2016
Figure copyright IEEE, 2016. Reproduced for educational purposes.



<https://youtu.be/25jyl67-poQ>

Pontificia Universidad Católica de Chile
Escuela de Ingeniería
Departamento de Ciencia de la Computación



Sistemas Urbanos Inteligentes

Aprendizaje Multitarea (Multitask Learning)

Hans Löbel

Dpto. Ingeniería de Transporte y Logística
Dpto. Ciencia de la Computación