

# Lecture 10: Object Detection

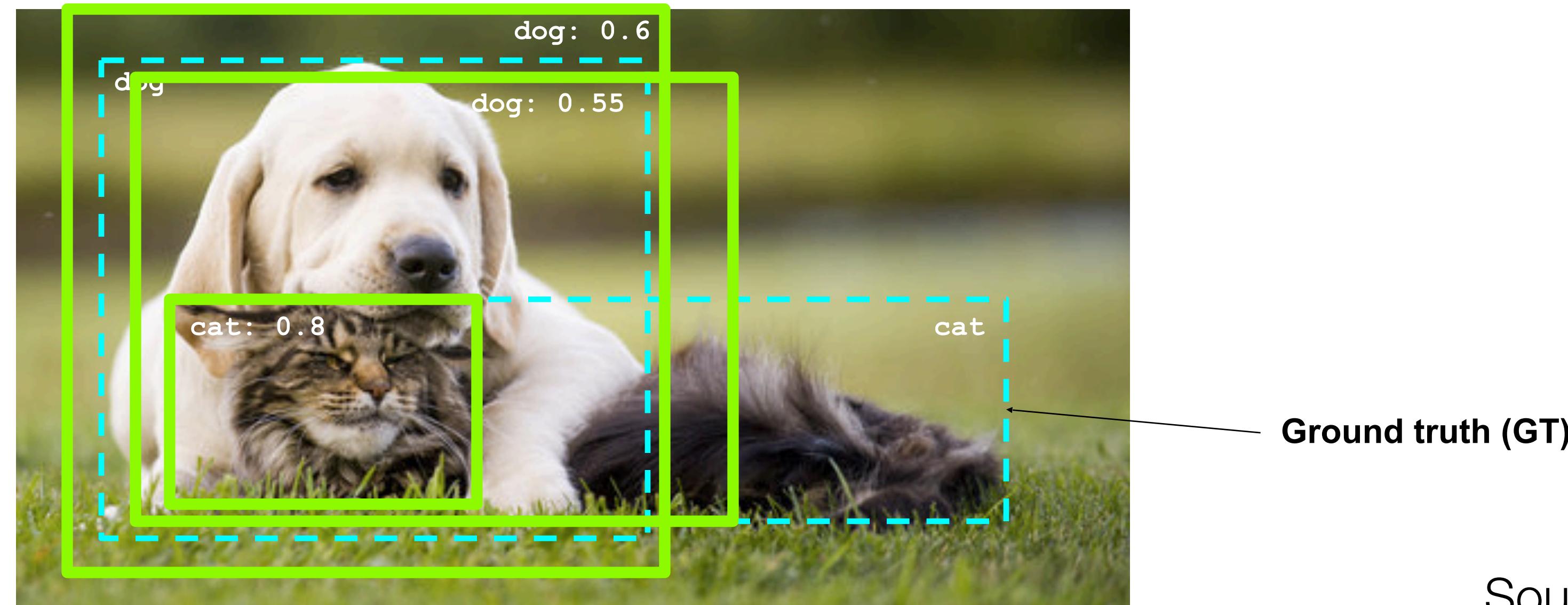
# Object detection with bounding boxes



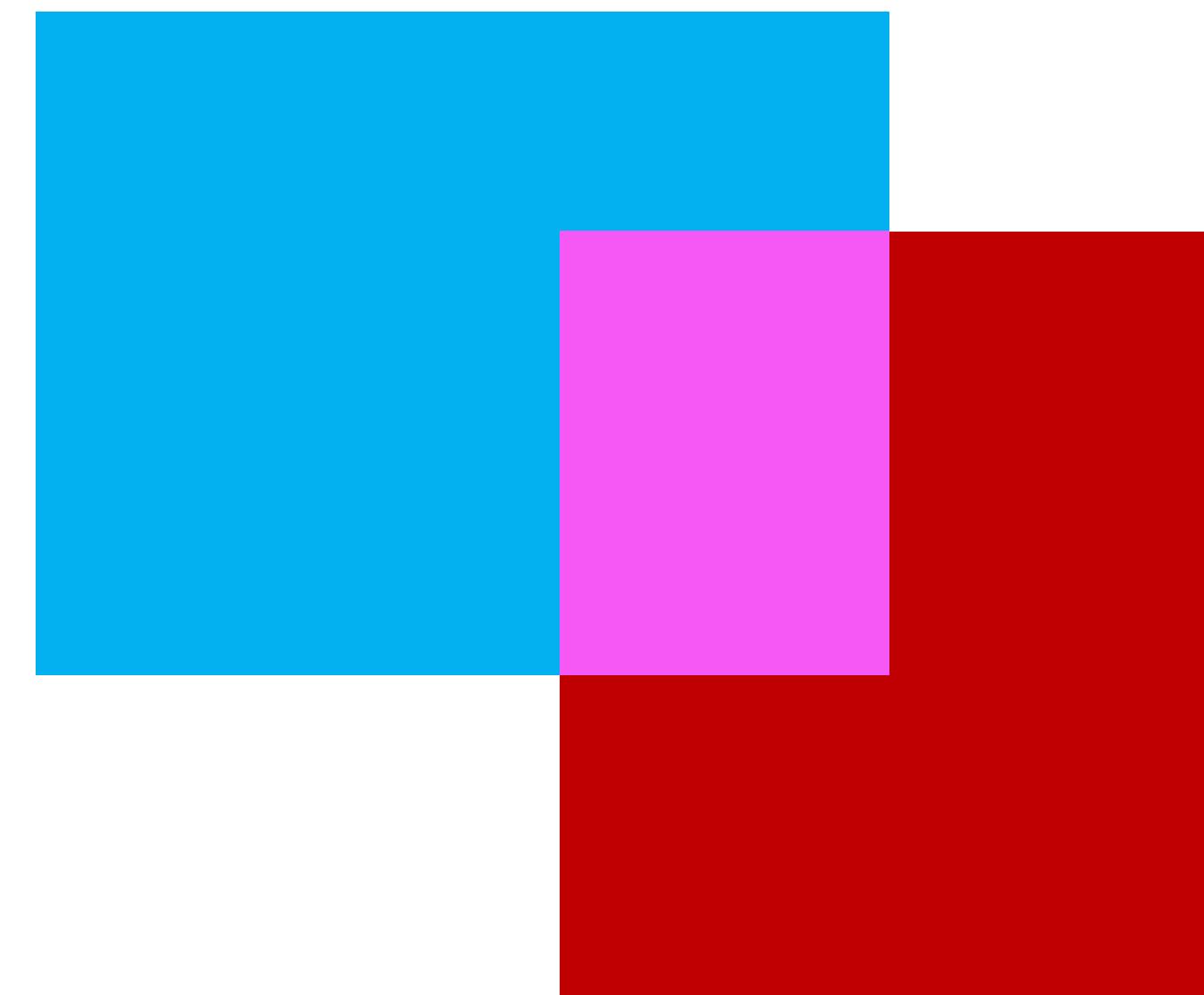
“Object detection”

# Evaluating an object detector

- At test time, predict bounding boxes, class labels, and confidence scores
- For each detection, determine whether it is a true or false positive
  - **Intersection over union (IoU):**  $\text{Area}(\text{GT} \cap \text{Det}) / \text{Area}(\text{GT} \cup \text{Det}) > 0.5$



# Evaluating an object detector

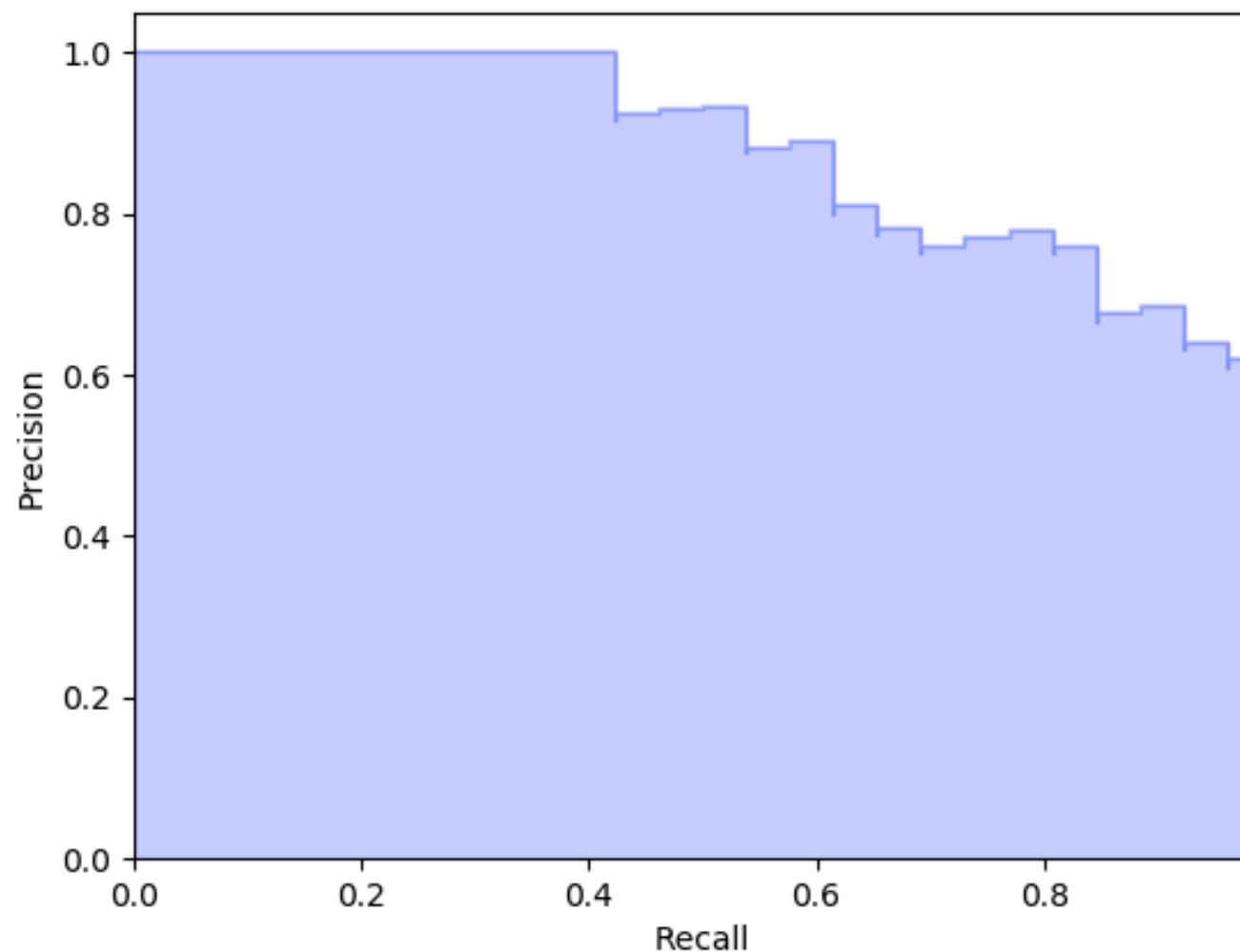


$$IoU(A, B) = \frac{|A \cap B|}{|A \cup B|}$$

Intersection over union (also known as Jaccard similarity)

# Evaluating an object detector

- For each class, plot Recall-Precision curve and compute Average Precision (area under the curve)
- Take mean of AP over classes to get mAP



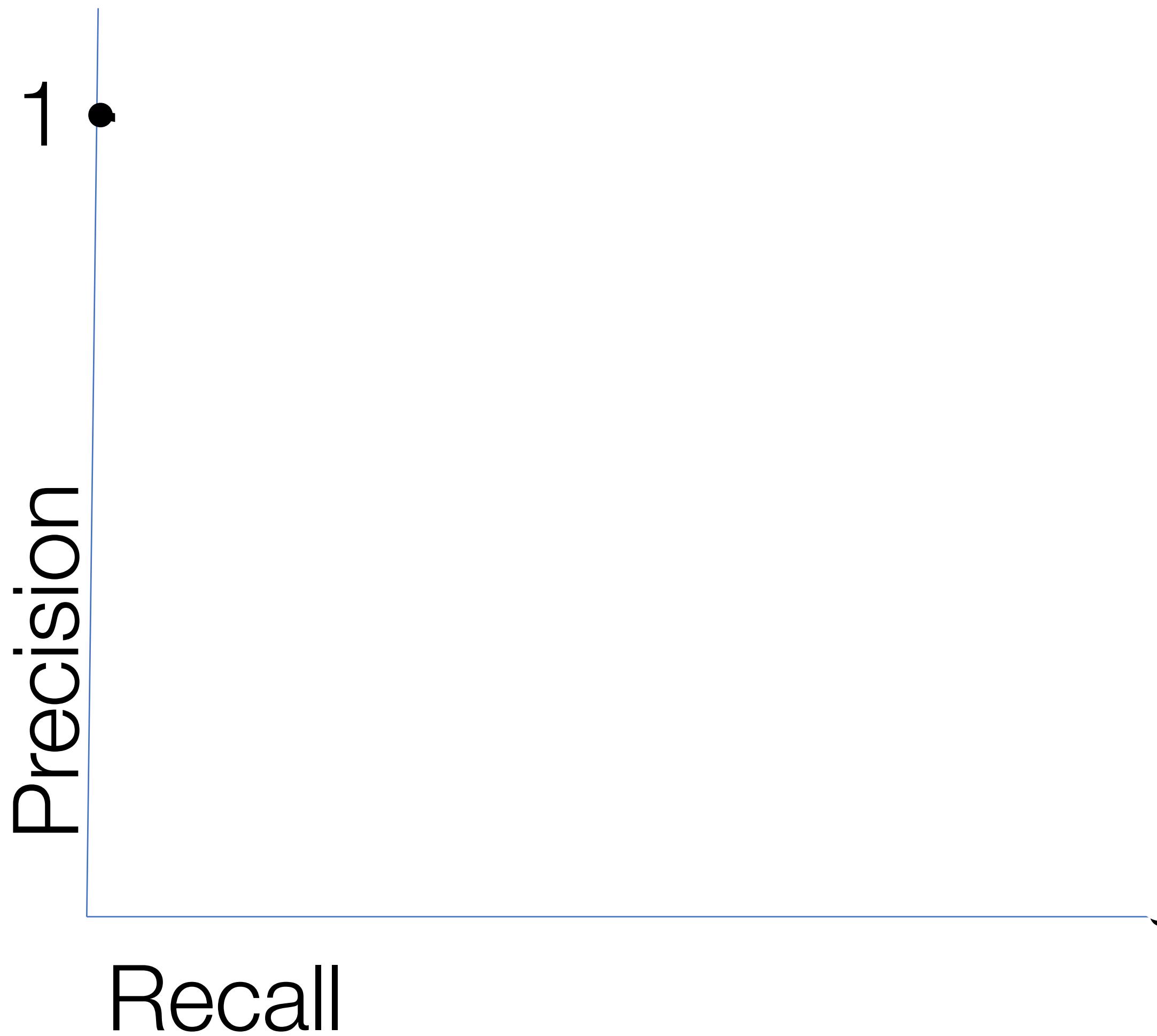
**Precision:**

true positive detections /  
total detections

**Recall:**

true positive detections /  
total positive test instances

# Average precision

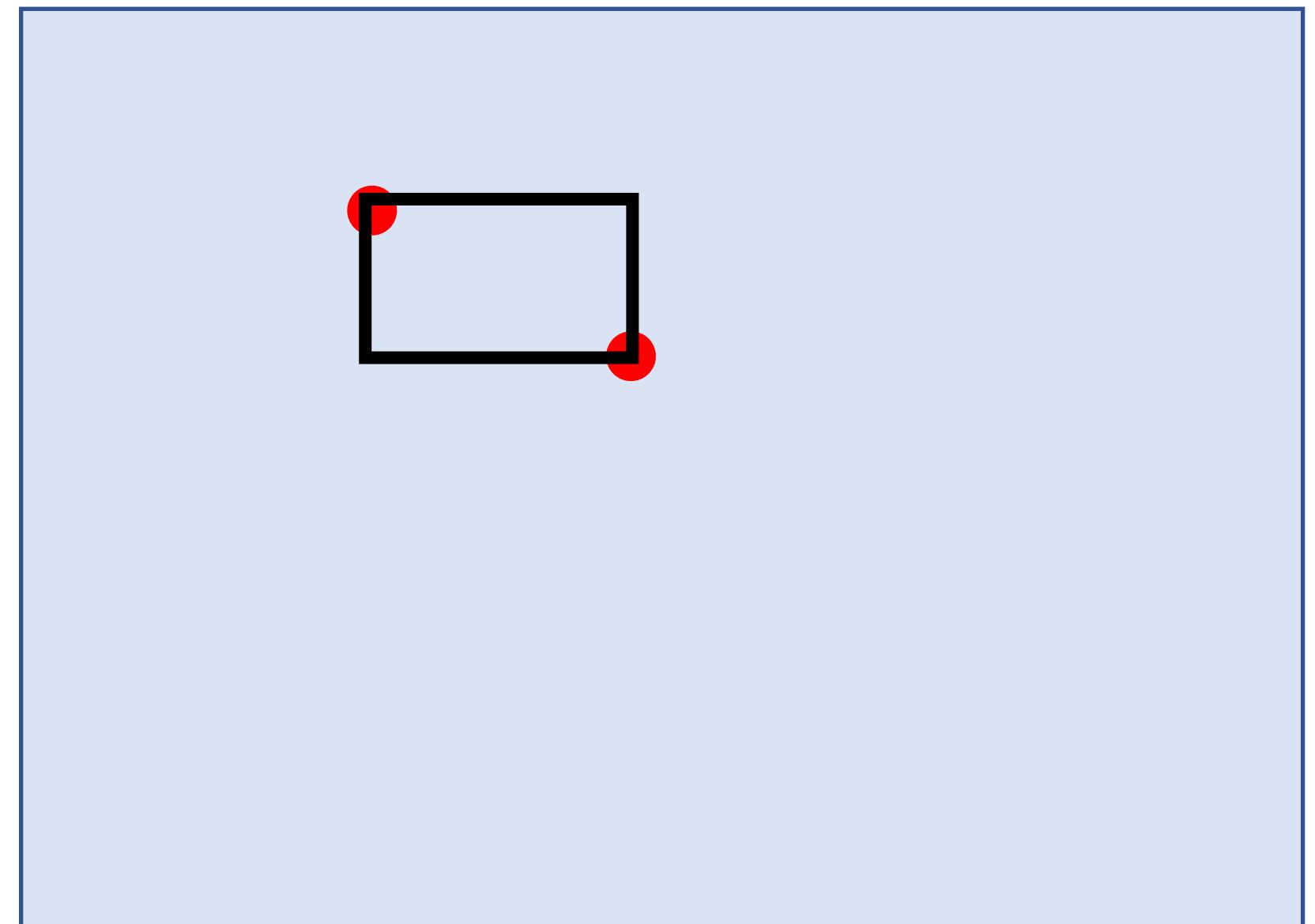


# Average precision



# Detection as classification

- Run through every possible box and classify
  - Well-localized object of class k or not?
- How many boxes?
  - Every pair of pixels = 1 box
  - $\binom{N}{2} = O(N^2)$
  - For  $300 \times 500$  image,  $N = 150K$
  - $2.25 \times 10^{10}$  boxes!
- Related challenge: almost all boxes are negative!

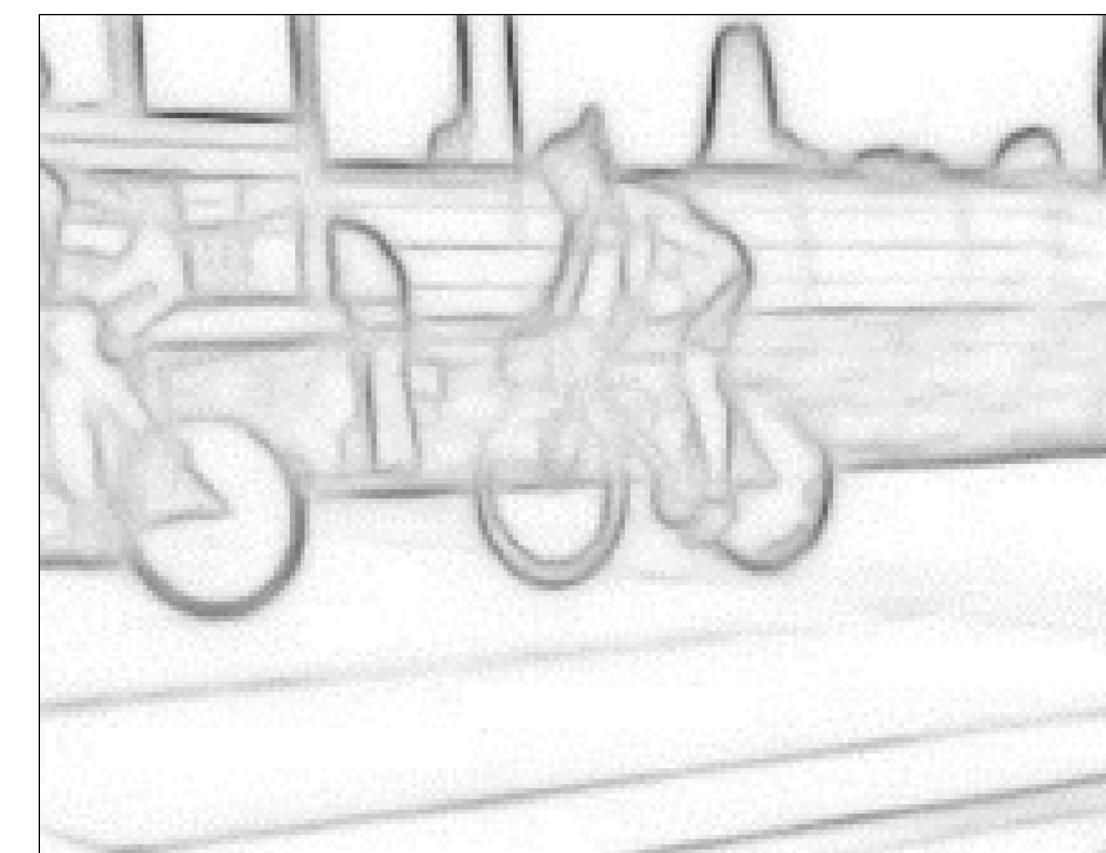


# Selective search

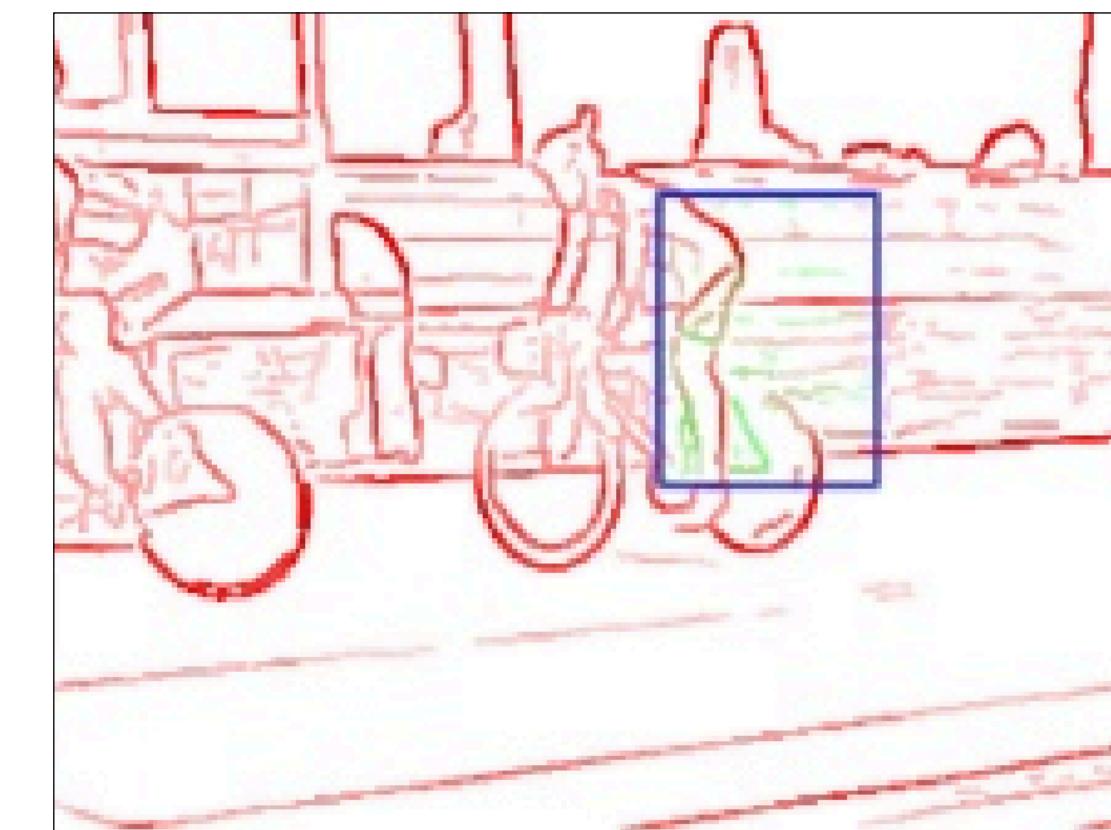
## Stage 1: generate candidate bounding boxes



Input image



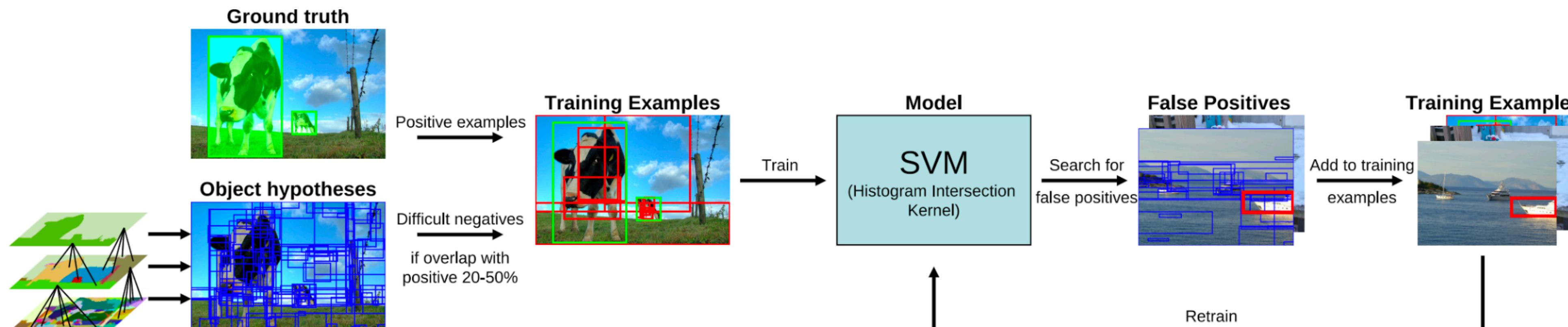
Edge detection



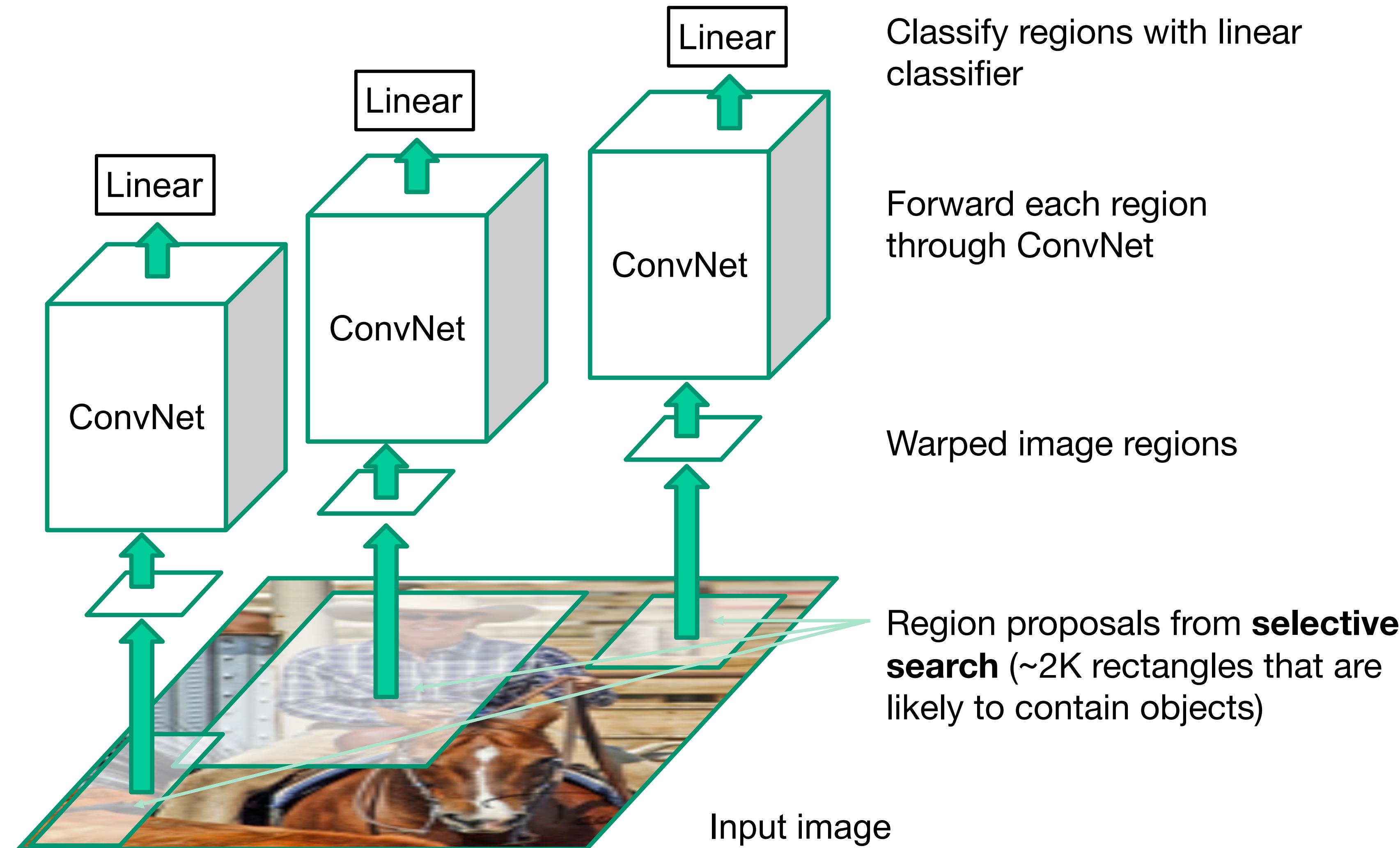
Bounding box proposal

## Stage 2: apply classifier only to each candidate bounding box

[Zitnick and Dollar, "Edge Boxes...", 2014]



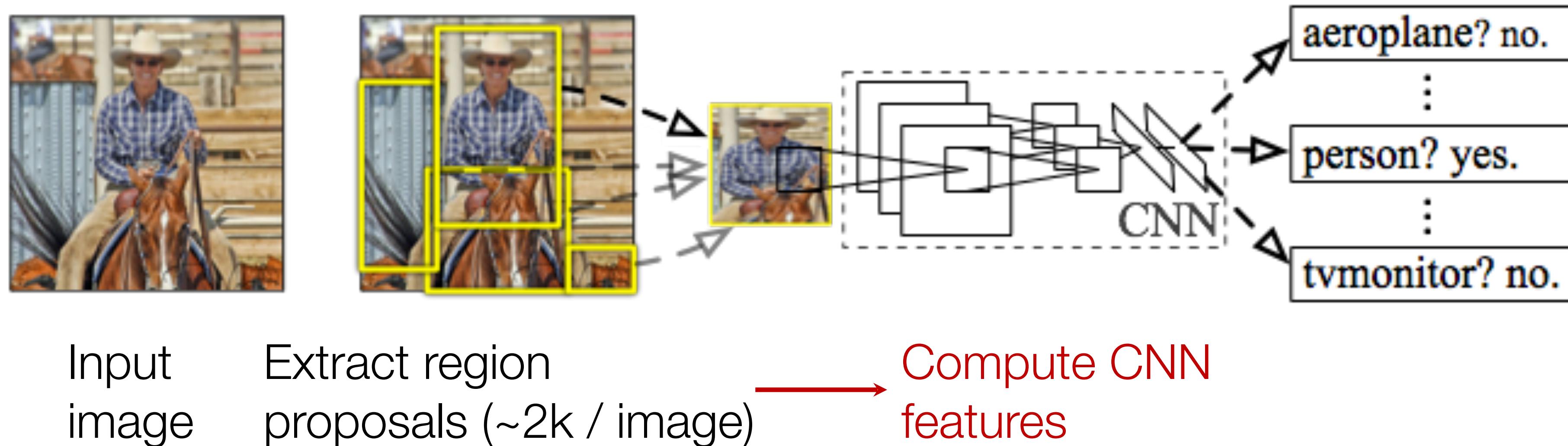
# R-CNN: Region proposals + CNN features



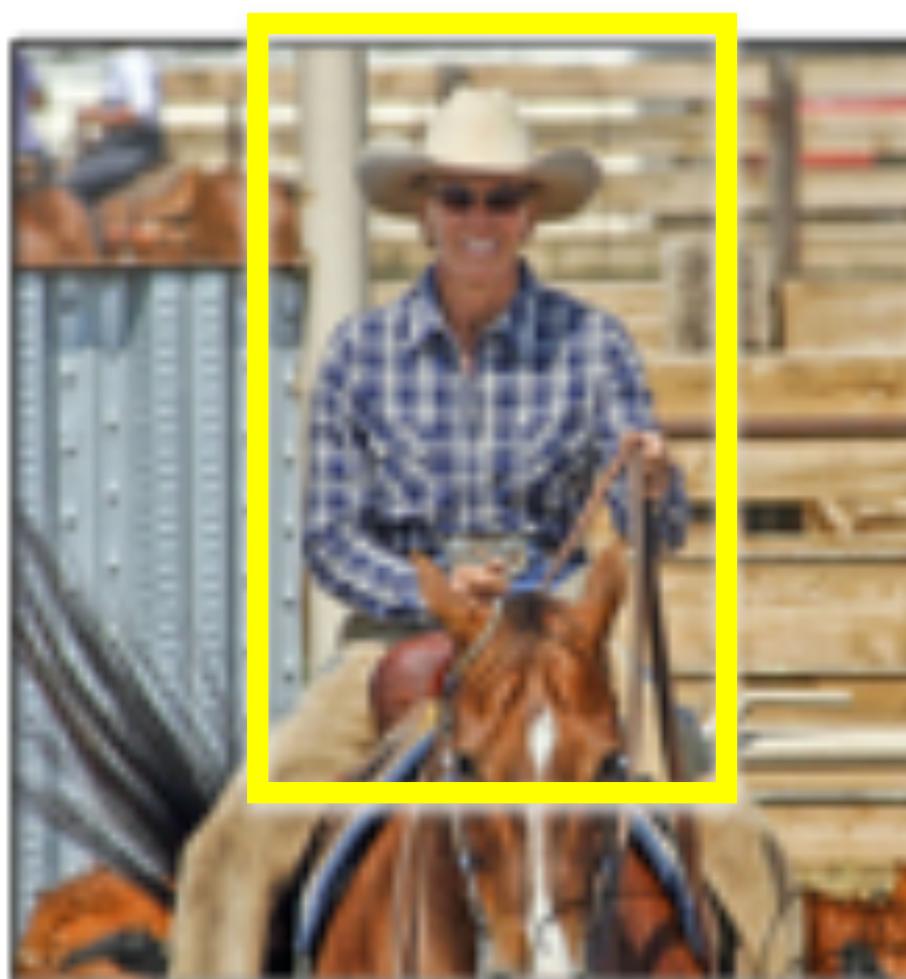
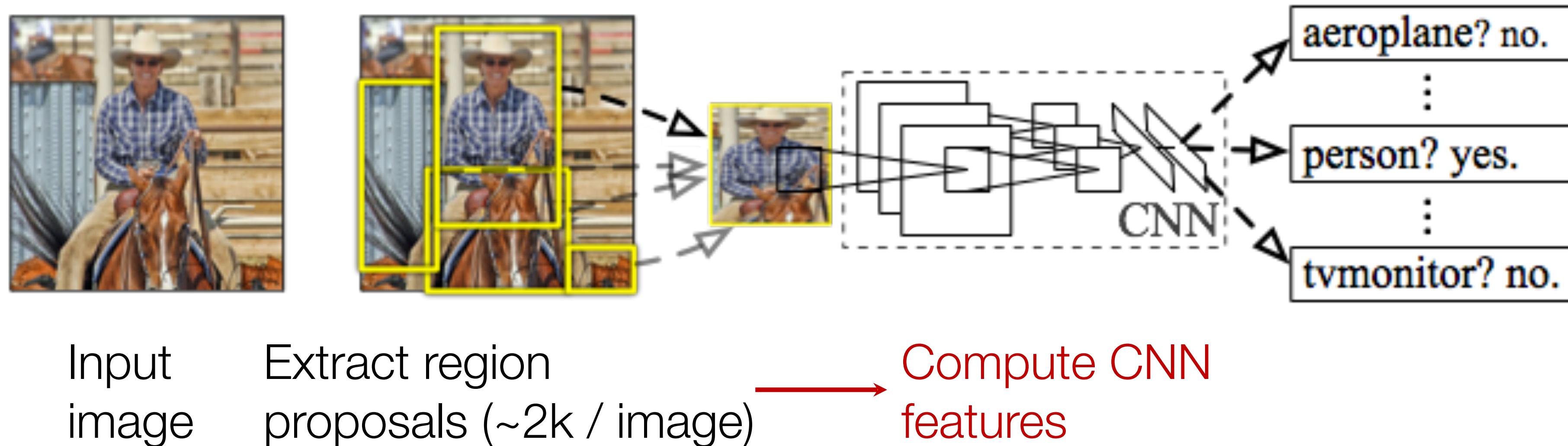
R. Girshick, J. Donahue, T. Darrell, and J. Malik, [Rich Feature Hierarchies for Accurate Object Detection and Semantic Segmentation](#), CVPR 2014.

Source: R. Girshick

# R-CNN at test time

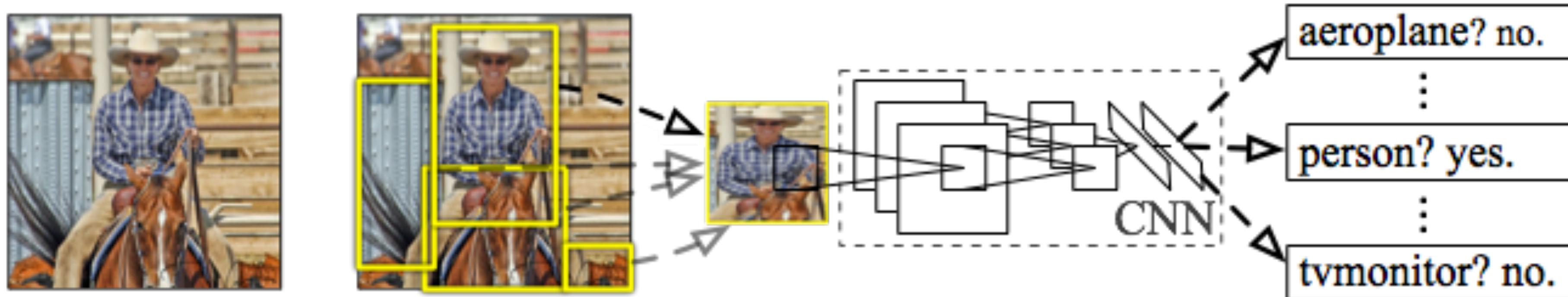


# R-CNN at test time



227 x 227

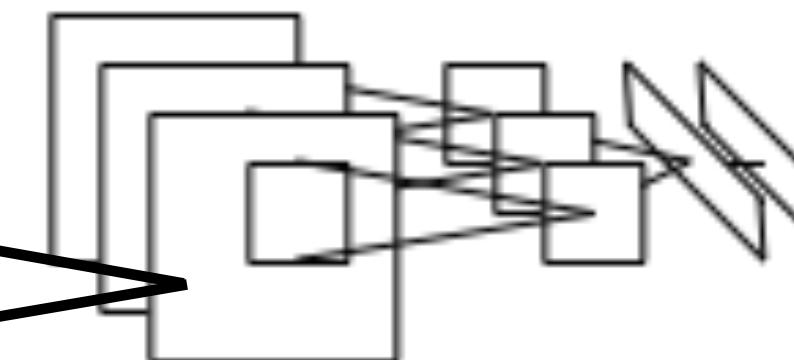
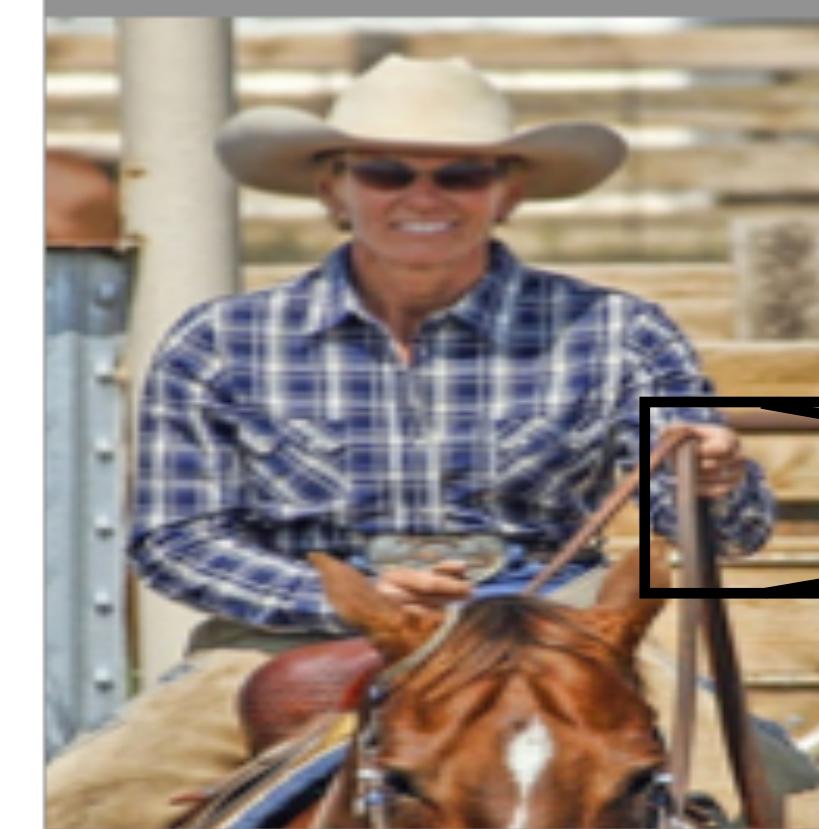
# R-CNN at test time



Input  
image

Extract region  
proposals (~2k / image)

Compute CNN  
features

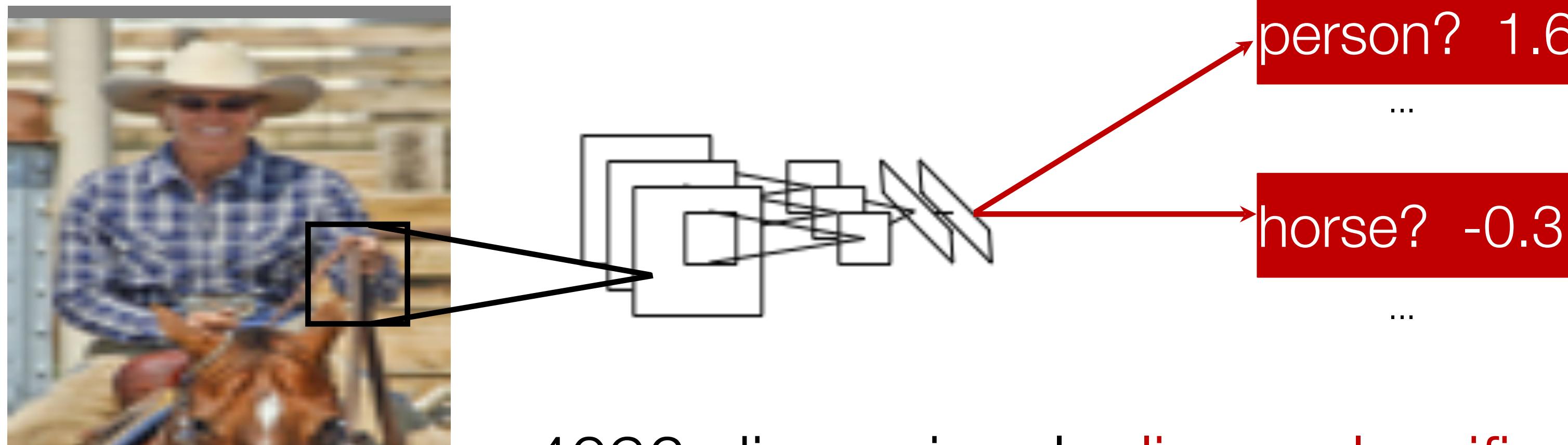
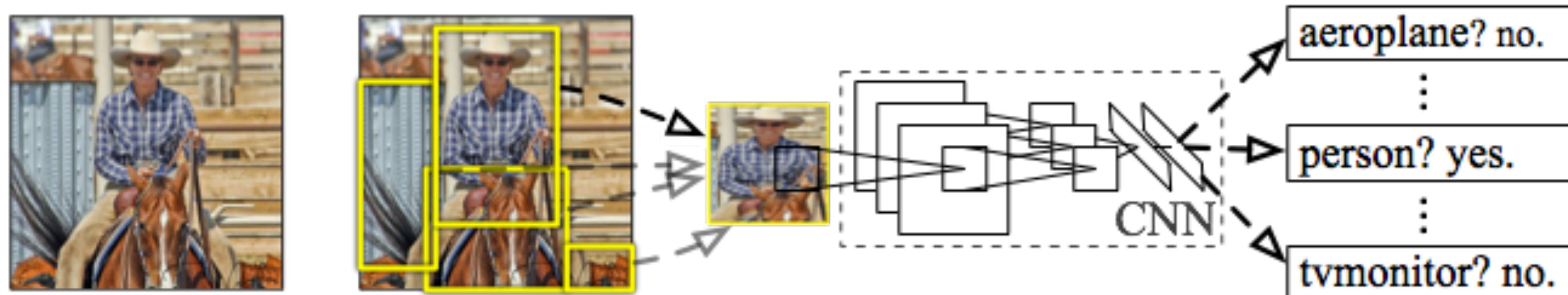


1. Crop

b. Scale (anisotropic)

c. Forward propagate  
Output: "fc7" features

# R-CNN at test time

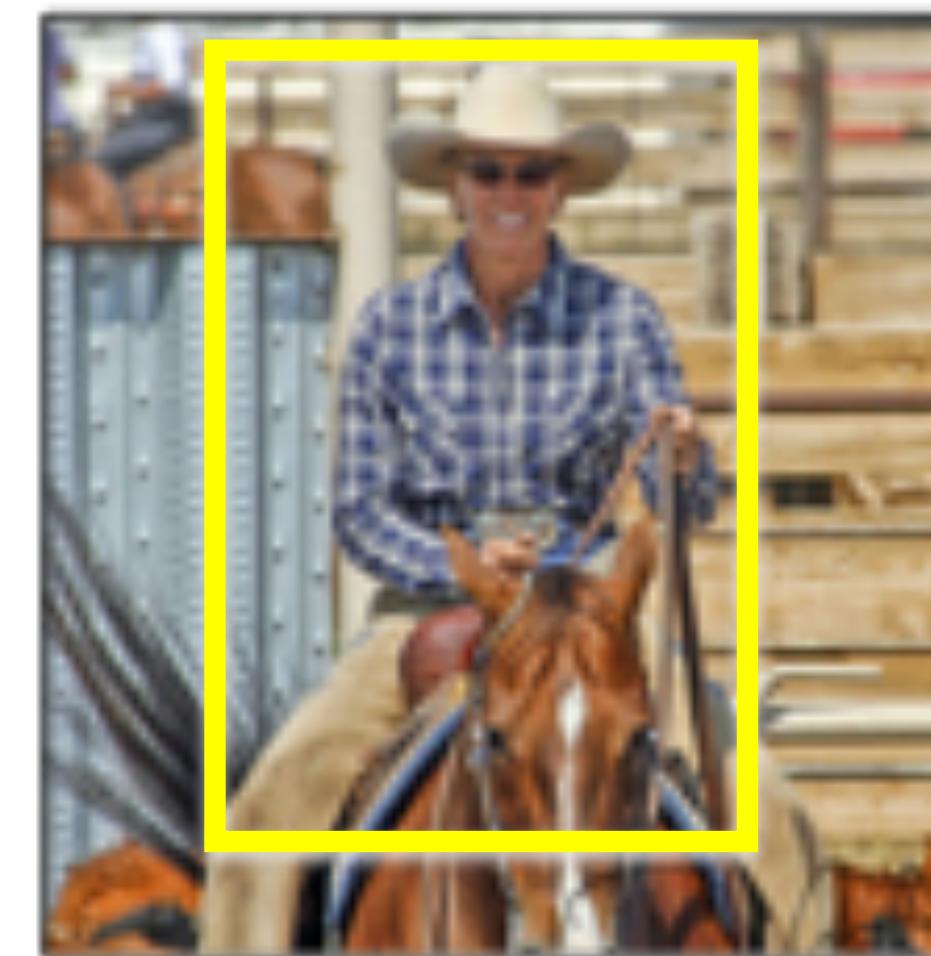


# R-CNN at test time: proposal refinement



Original  
proposal

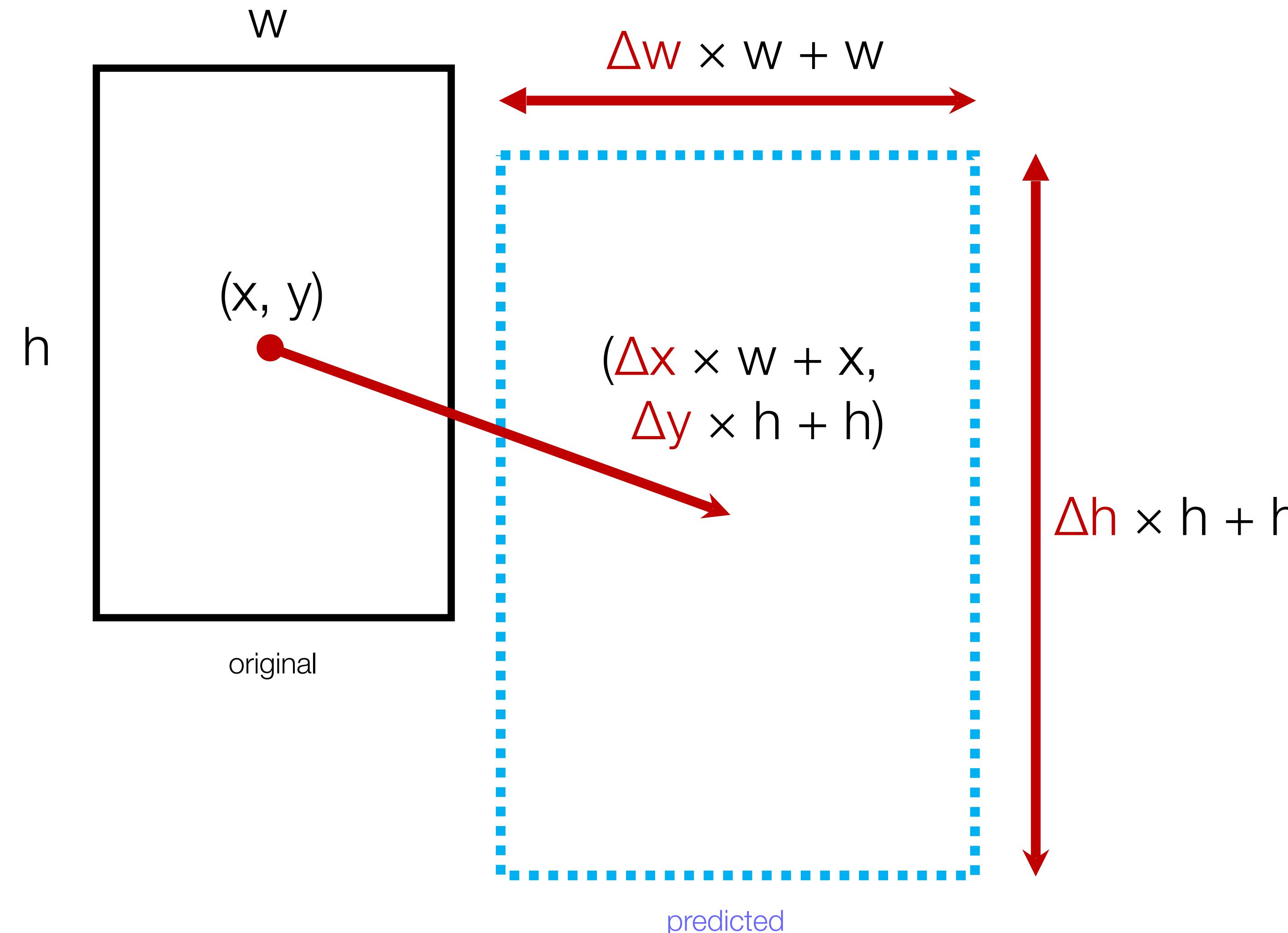
Linear regression  
on CNN features



Predicted  
object bounding box

Bounding-box regression

# Bounding-box regression

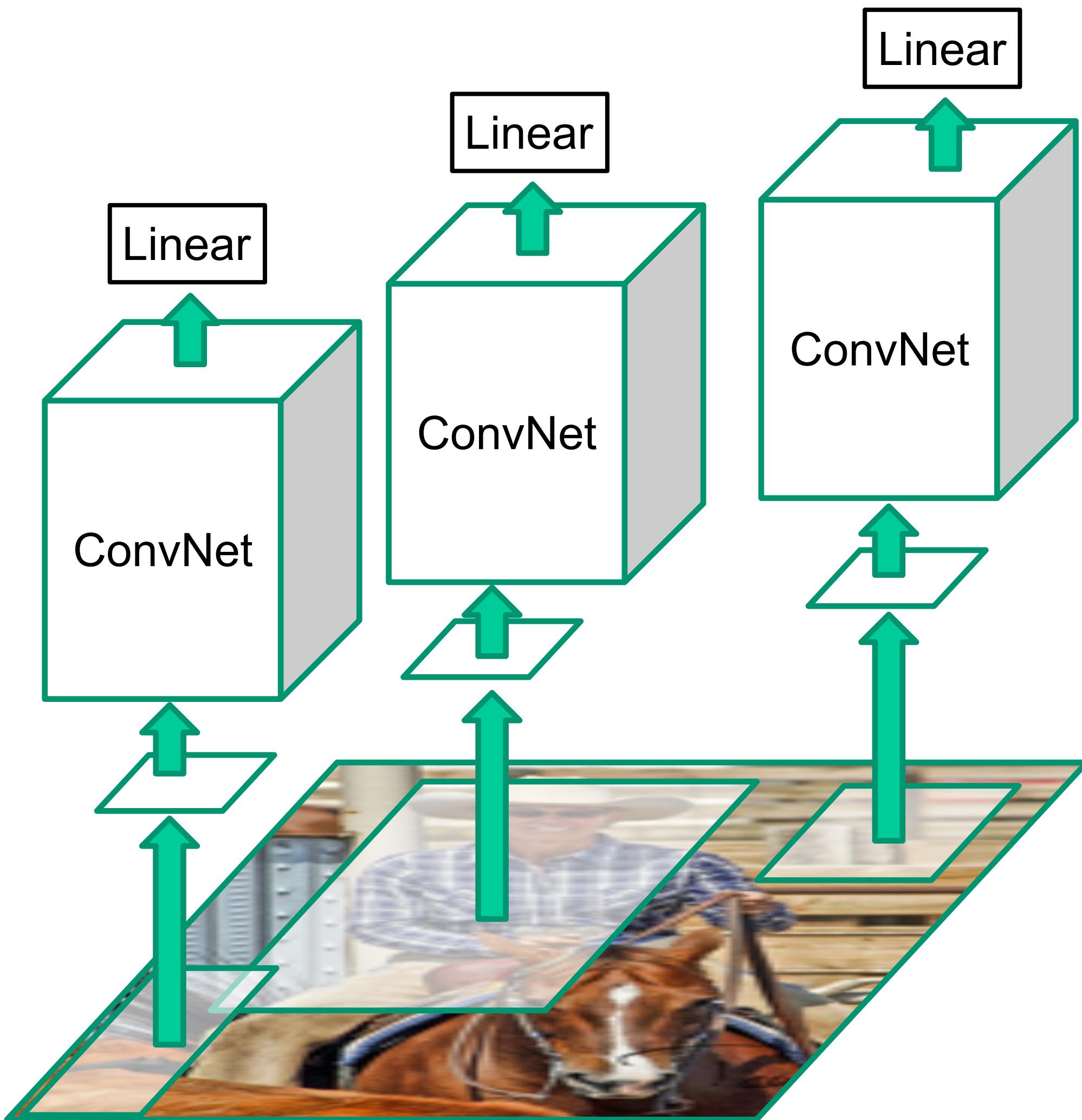


# Non-maximum suppression



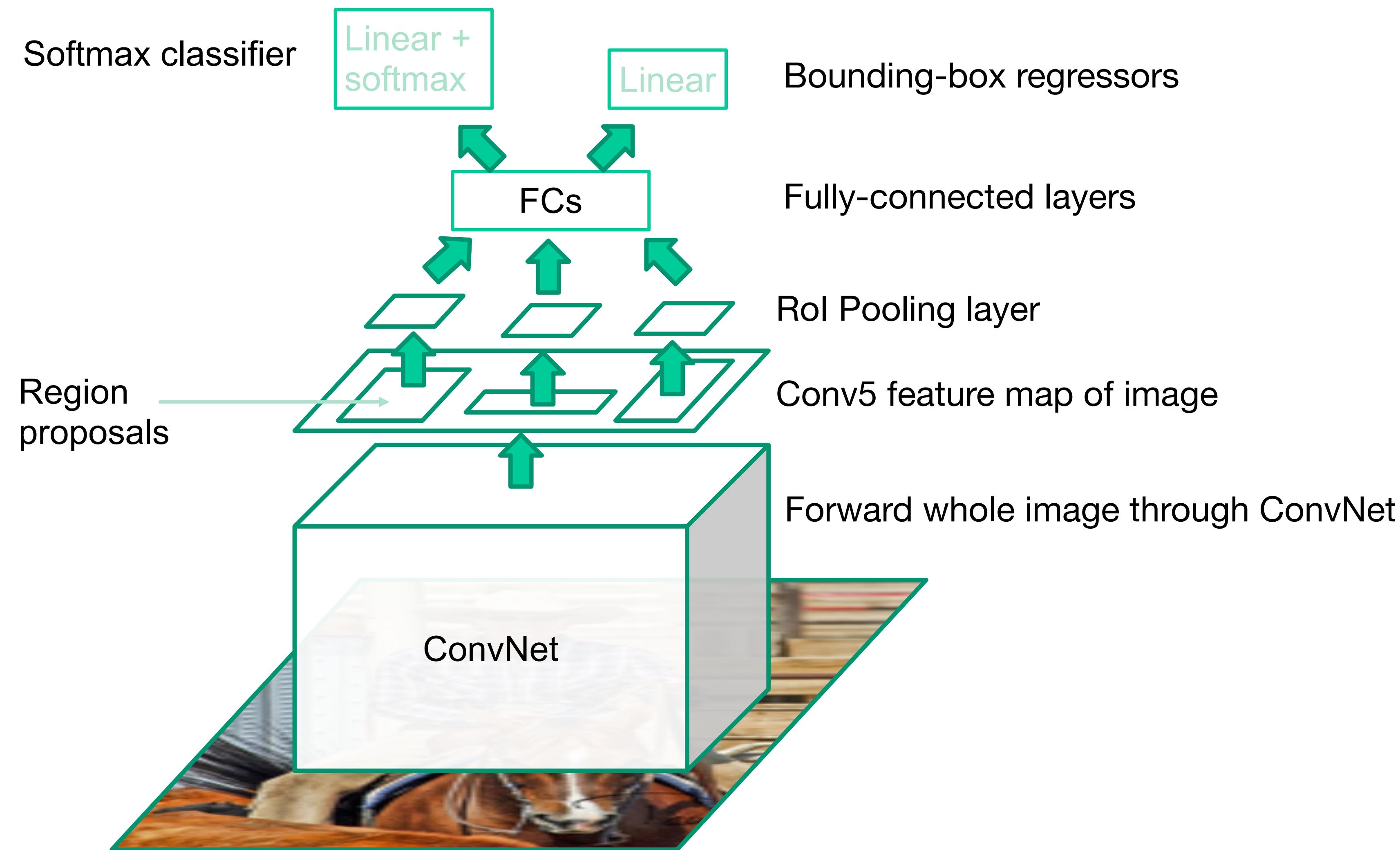
If two boxes overlap significantly (e.g.  $> 50\%$  IoU), drop the one with the lower score. Usually use greedy algorithm.

# Problems with R-CNN



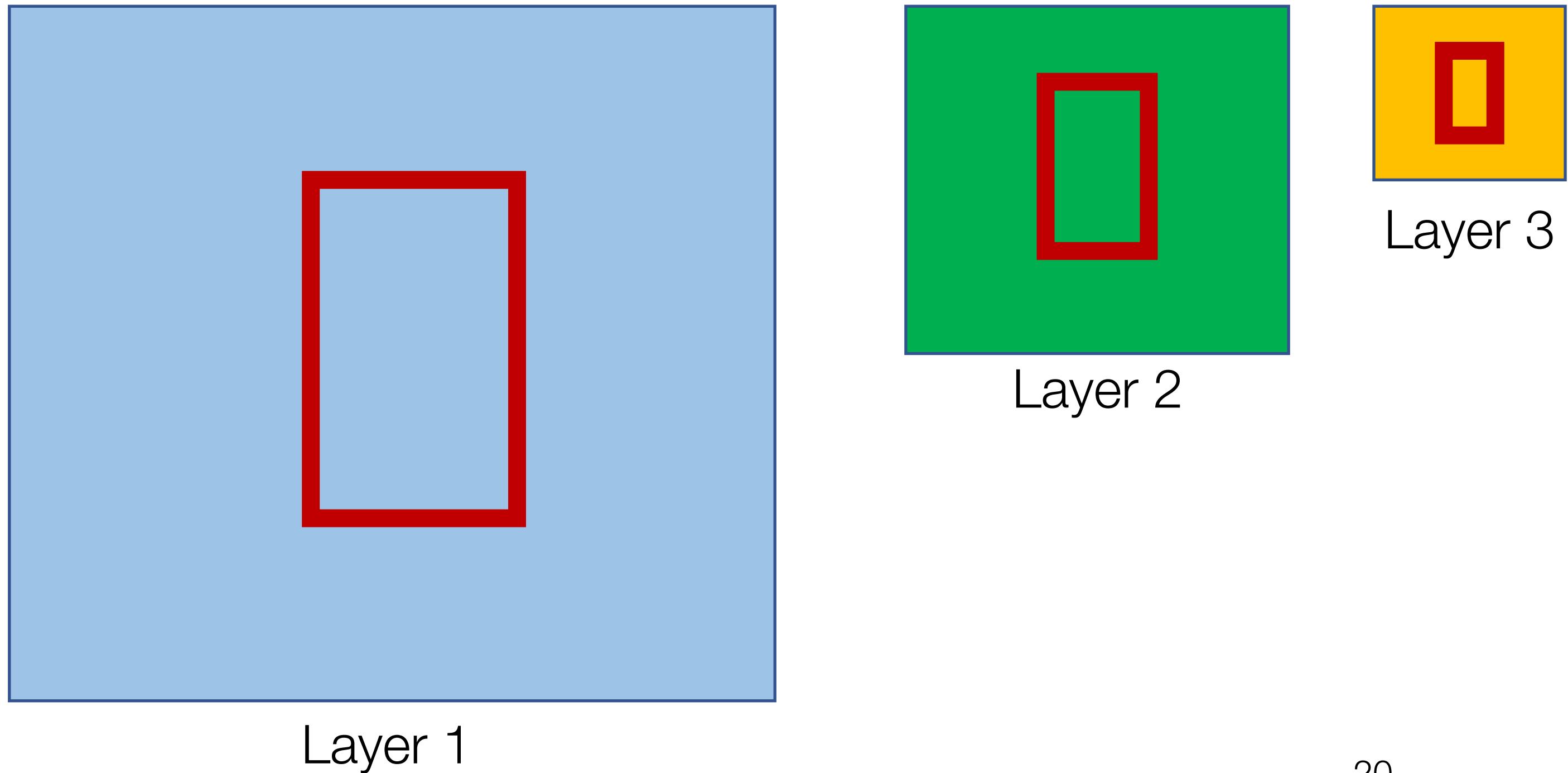
1. Slow! Have to run CNN per window
2. Hand-crafted mechanism for region proposal might be suboptimal.

# “Fast” R-CNN: reuse features between proposals



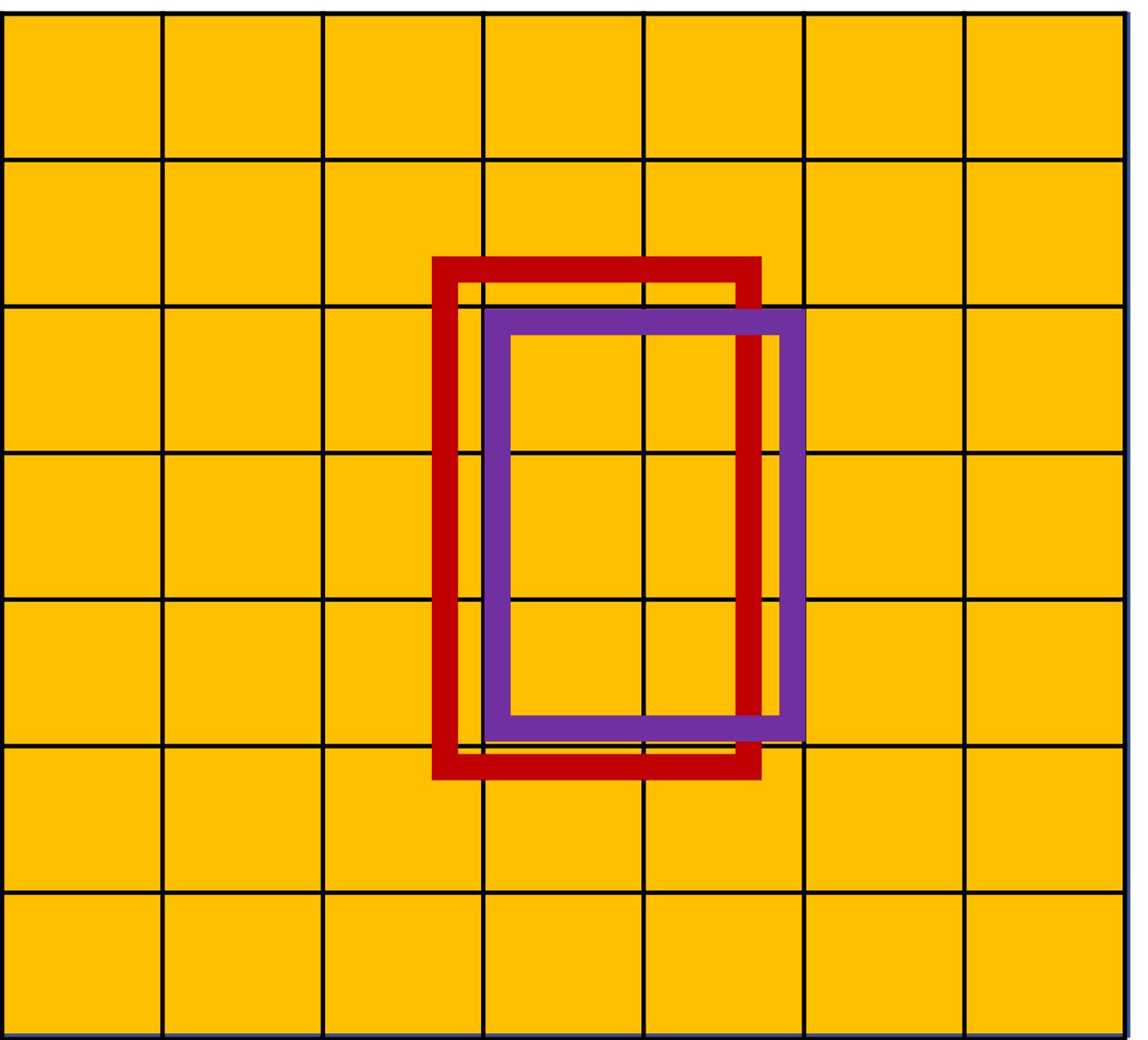
# ROI Pooling

- How do we crop from a feature map?
- Step 1: Resize boxes to account for subsampling



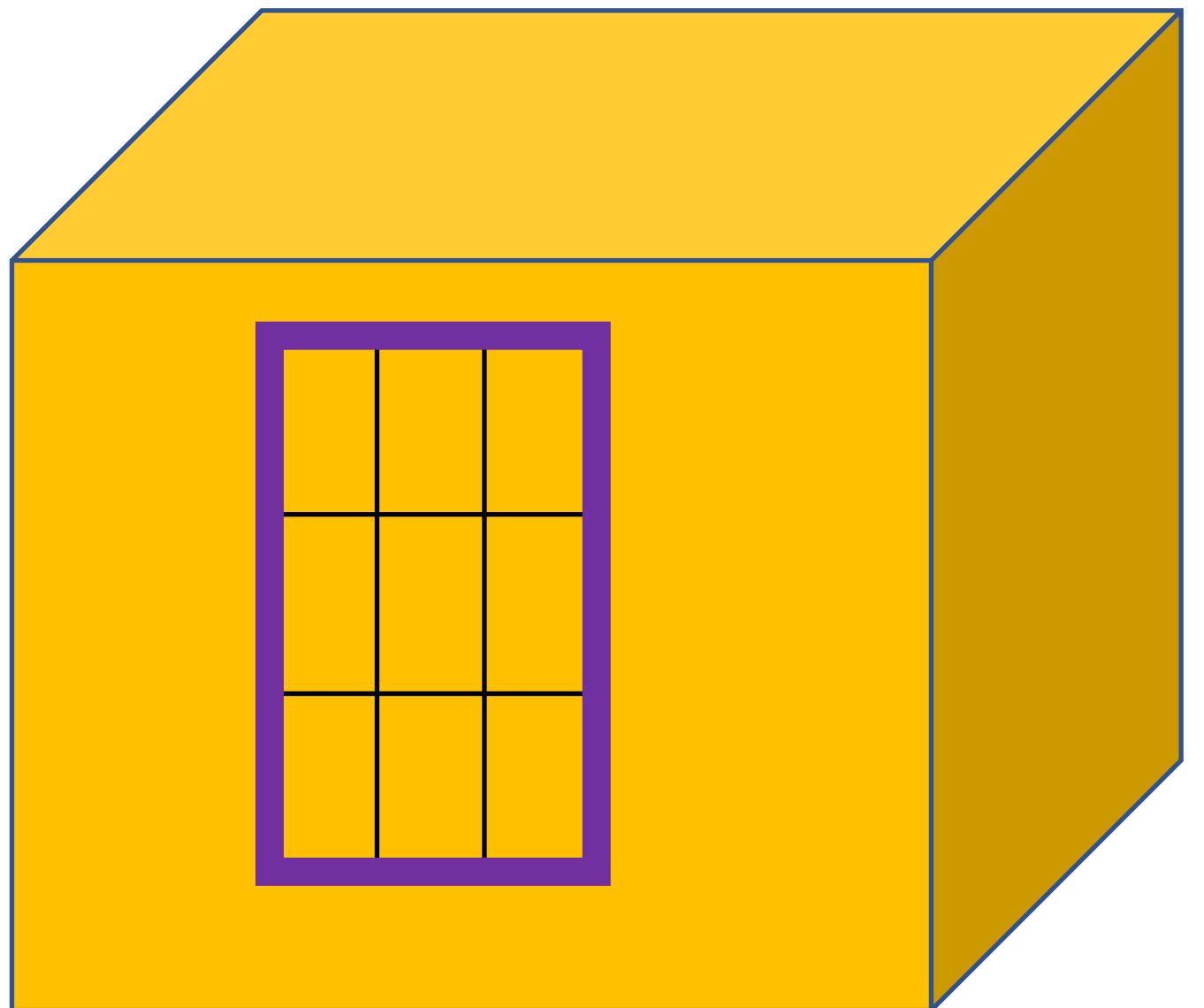
# ROI Pooling

- How do we crop from a feature map?
- Step 2: Snap to feature map grid



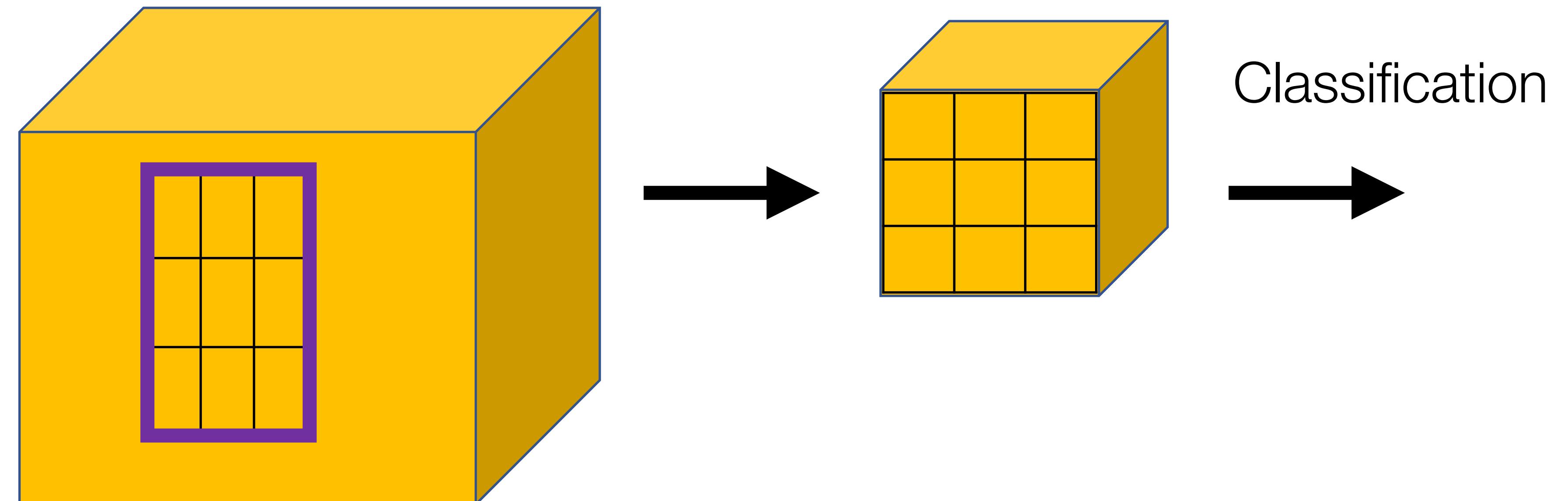
# ROI Pooling

- How do we crop from a feature map?
- Step 3: Overlay a new grid of fixed size



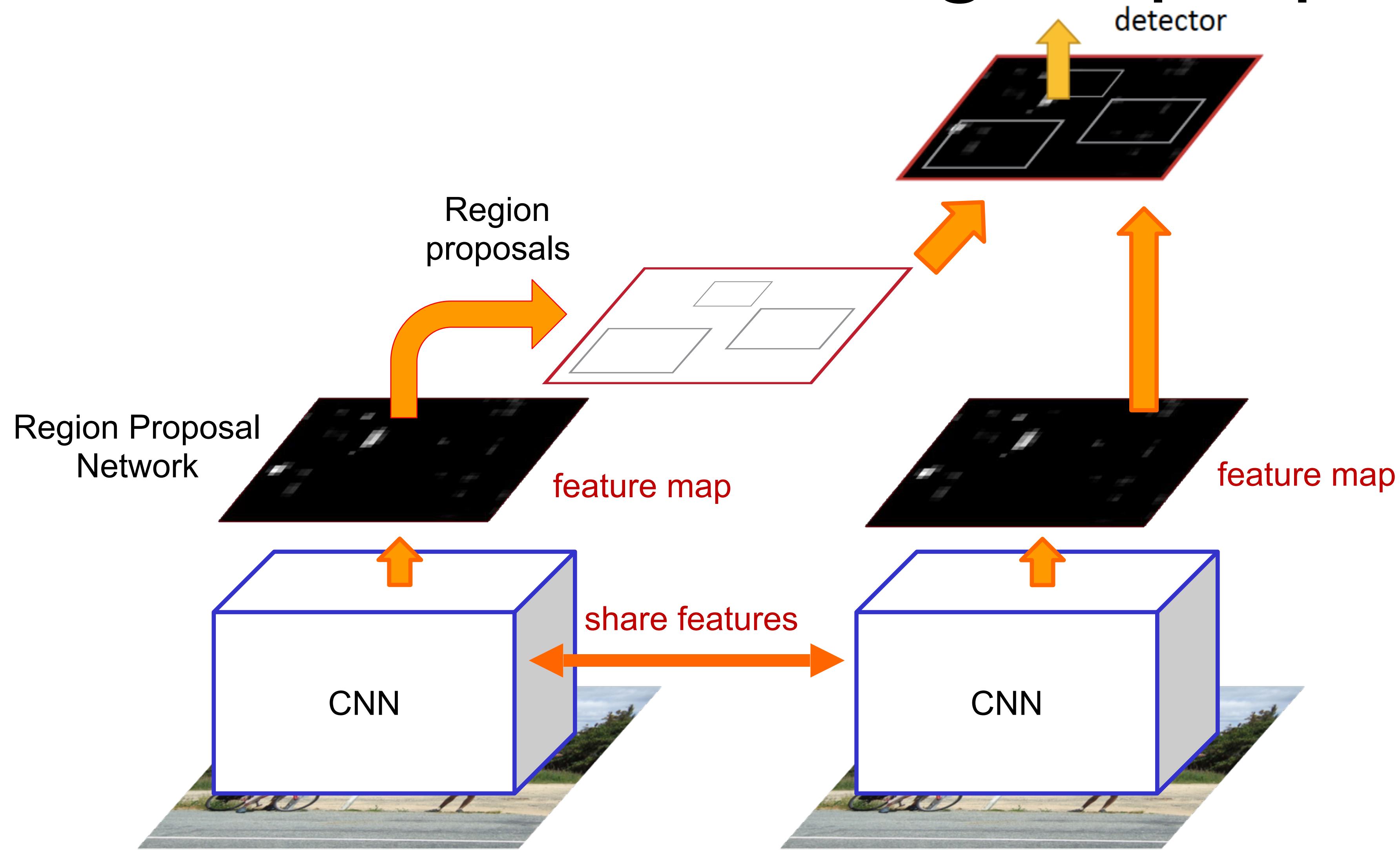
# ROI Pooling

- How do we crop from a feature map?
- Step 4: Take max in each cell



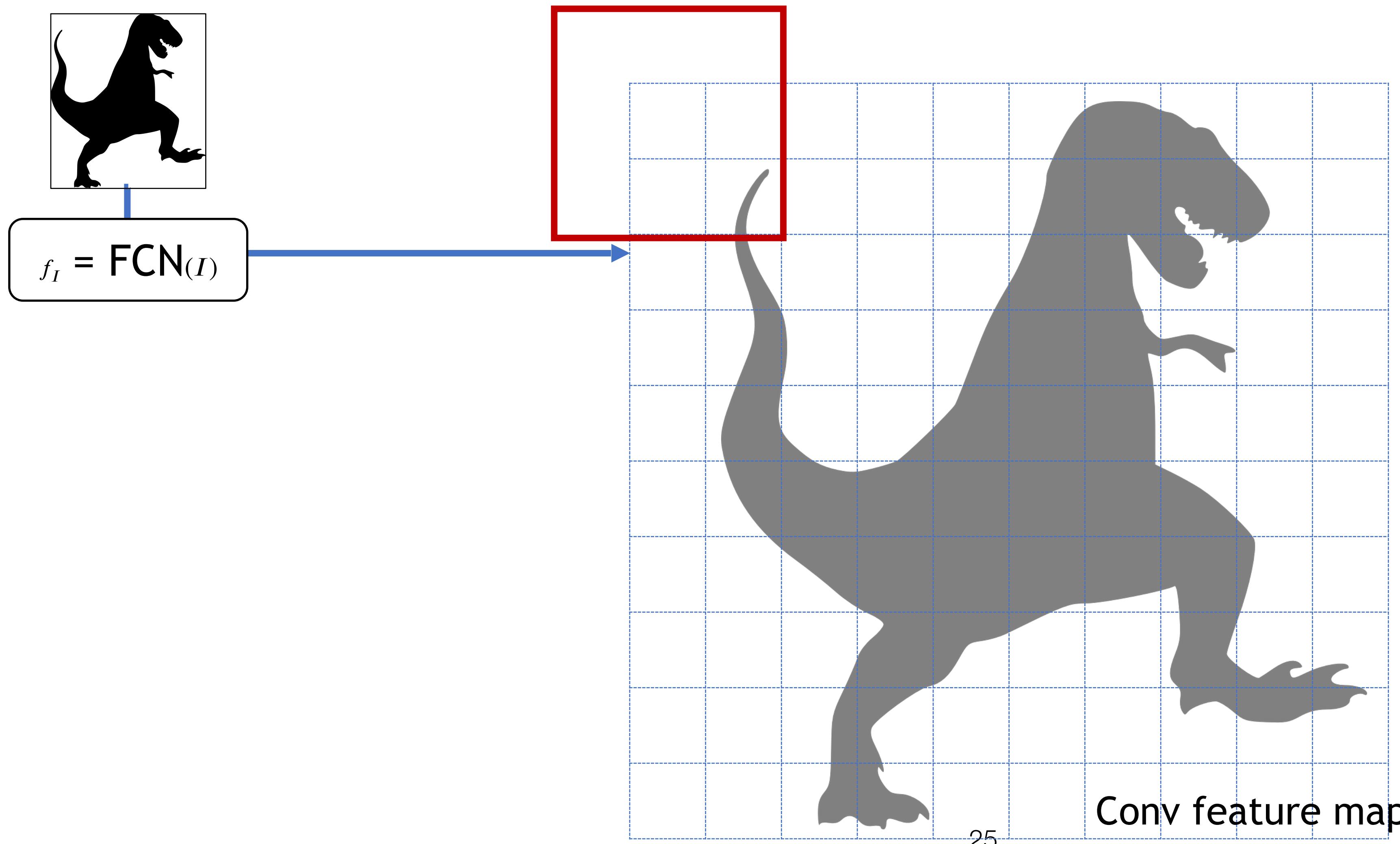
See more here: <https://deepsense.ai/region-of-interest-pooling-explained/>

# “Faster” R-CNN: learn region proposals

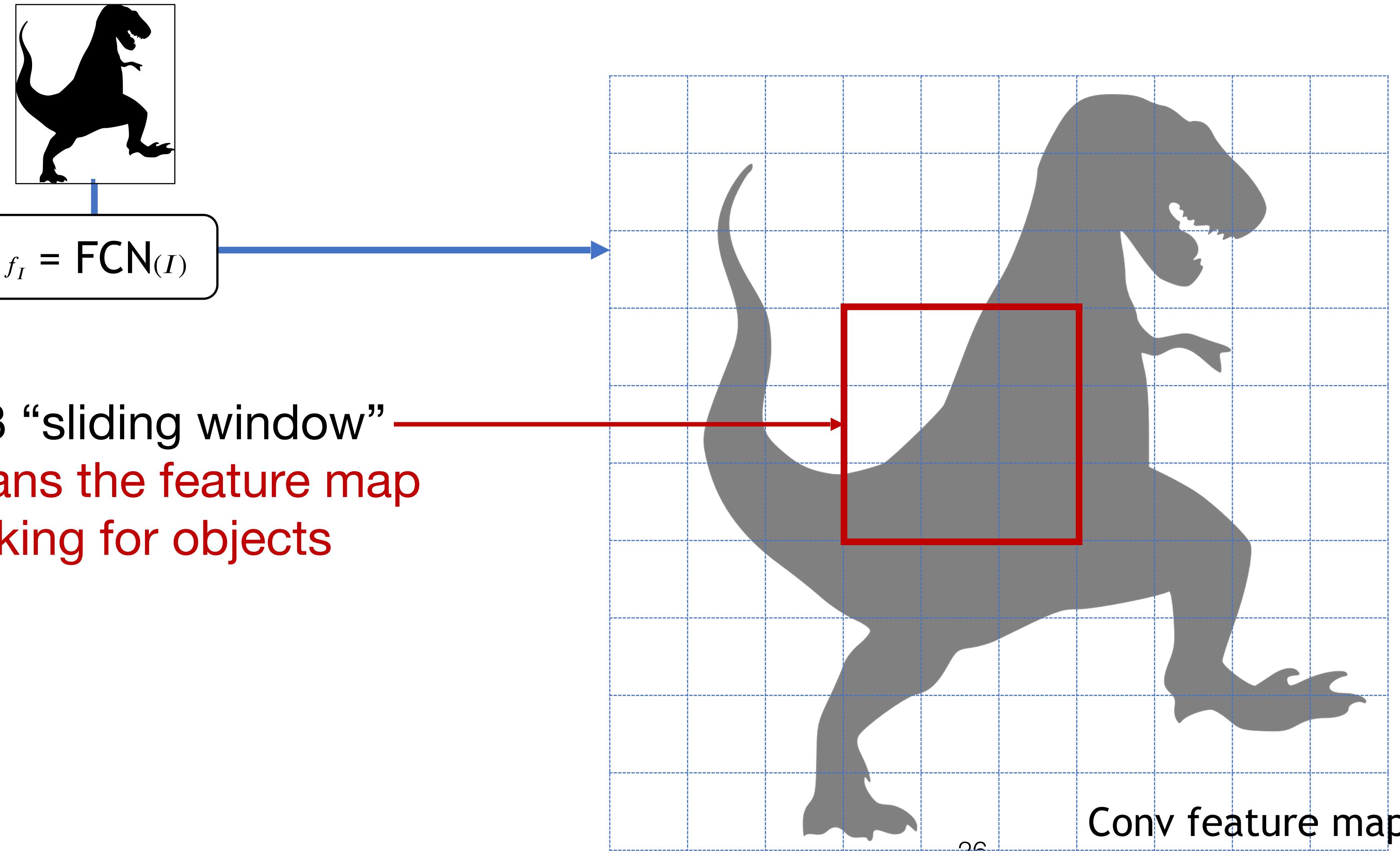


S. Ren, K. He, R. Girshick, and J. Sun, [Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks](#), NIPS 2015

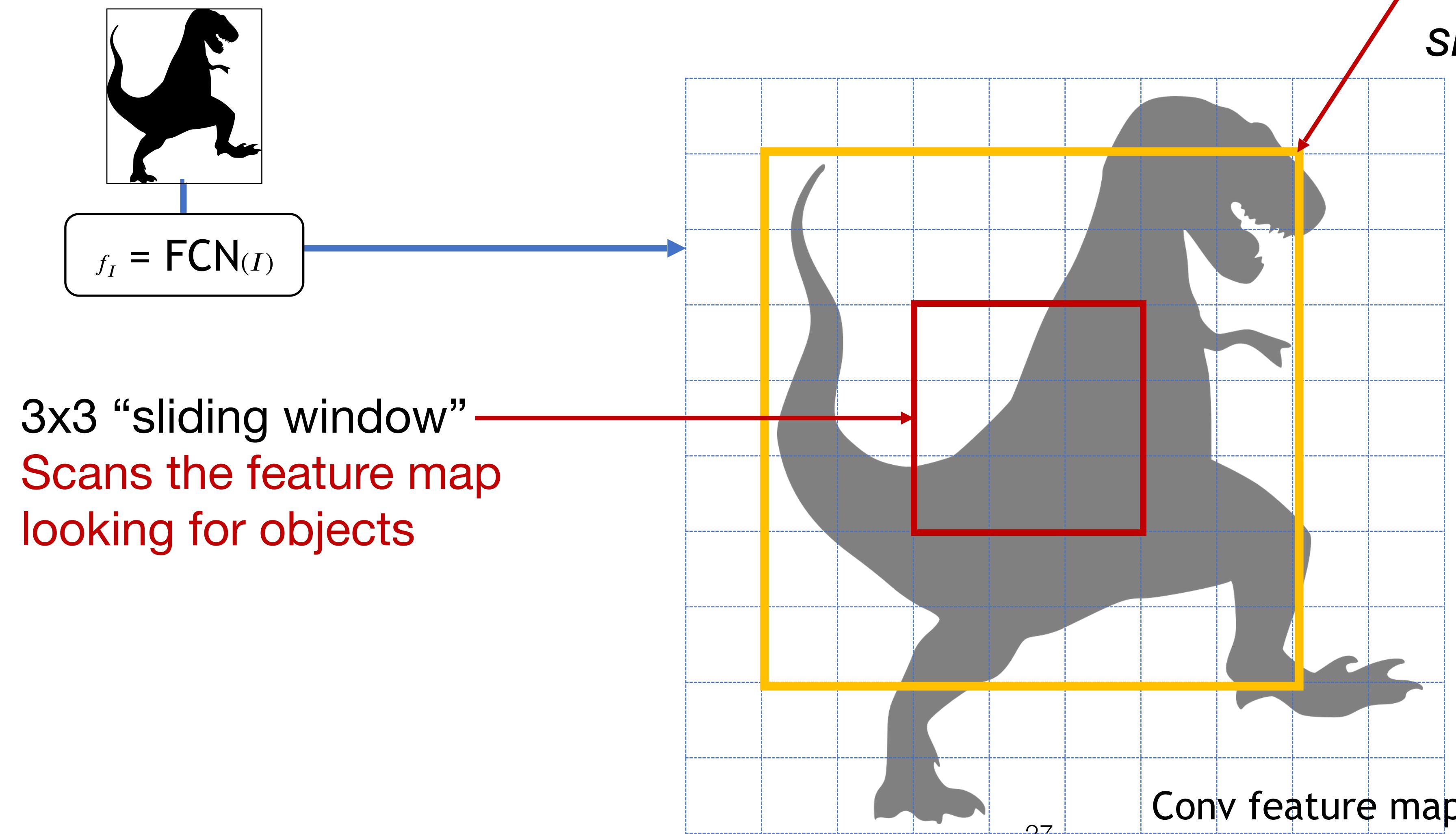
# RPN: Region Proposal Network



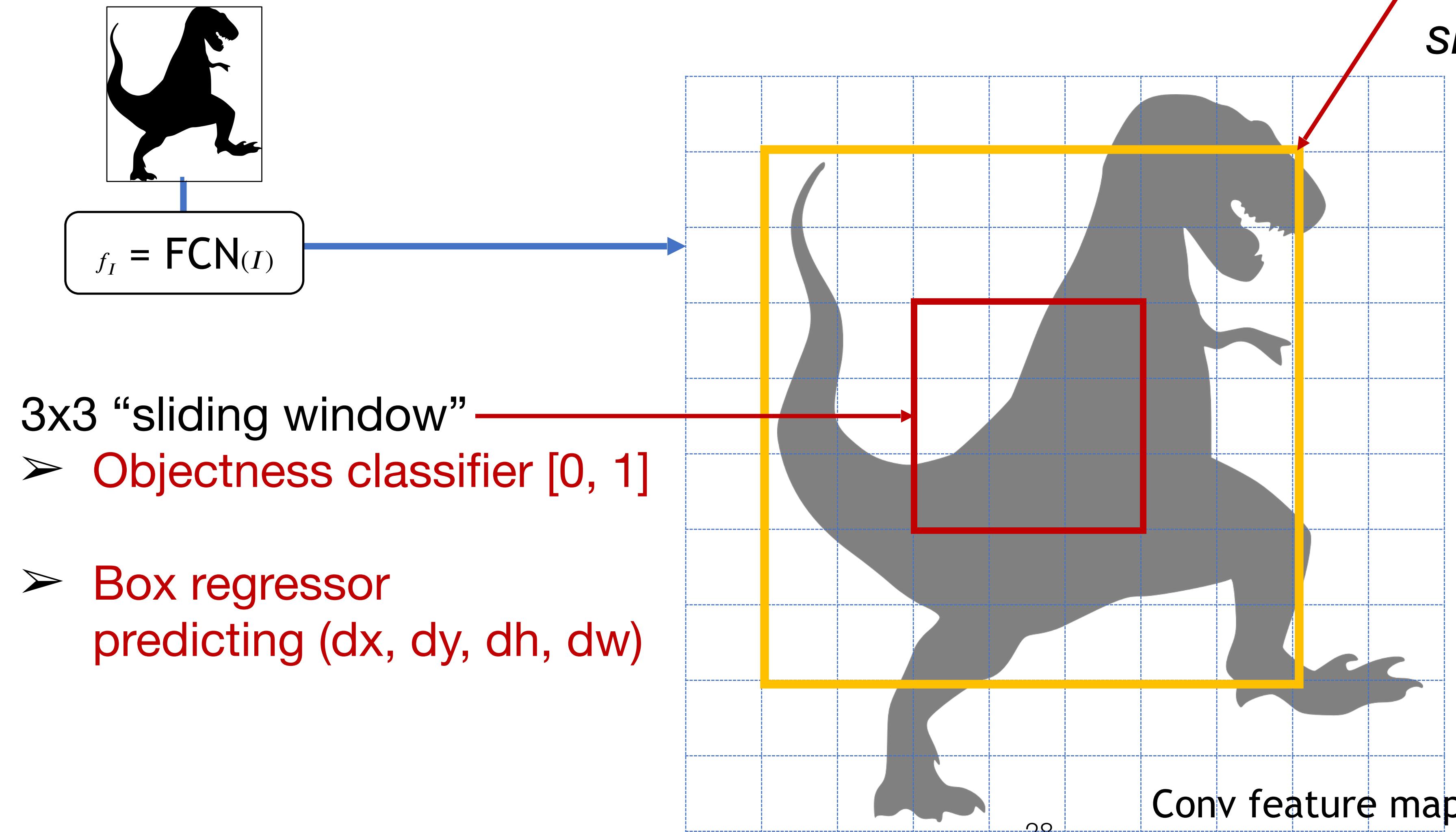
# RPN: Region Proposal Network



# RPN: Anchor Box



# RPN: Anchor Box

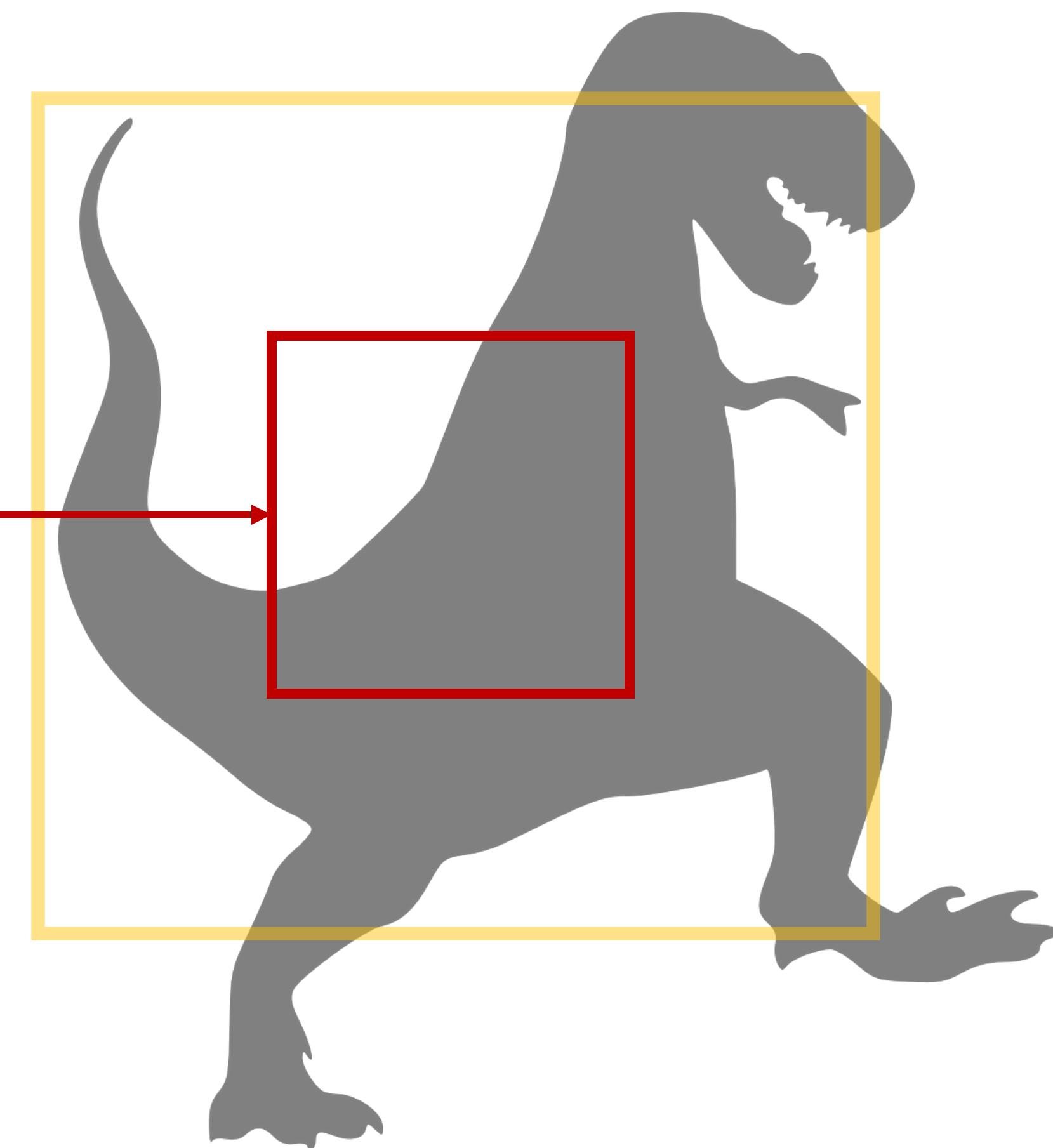


# RPN: Prediction (on object)

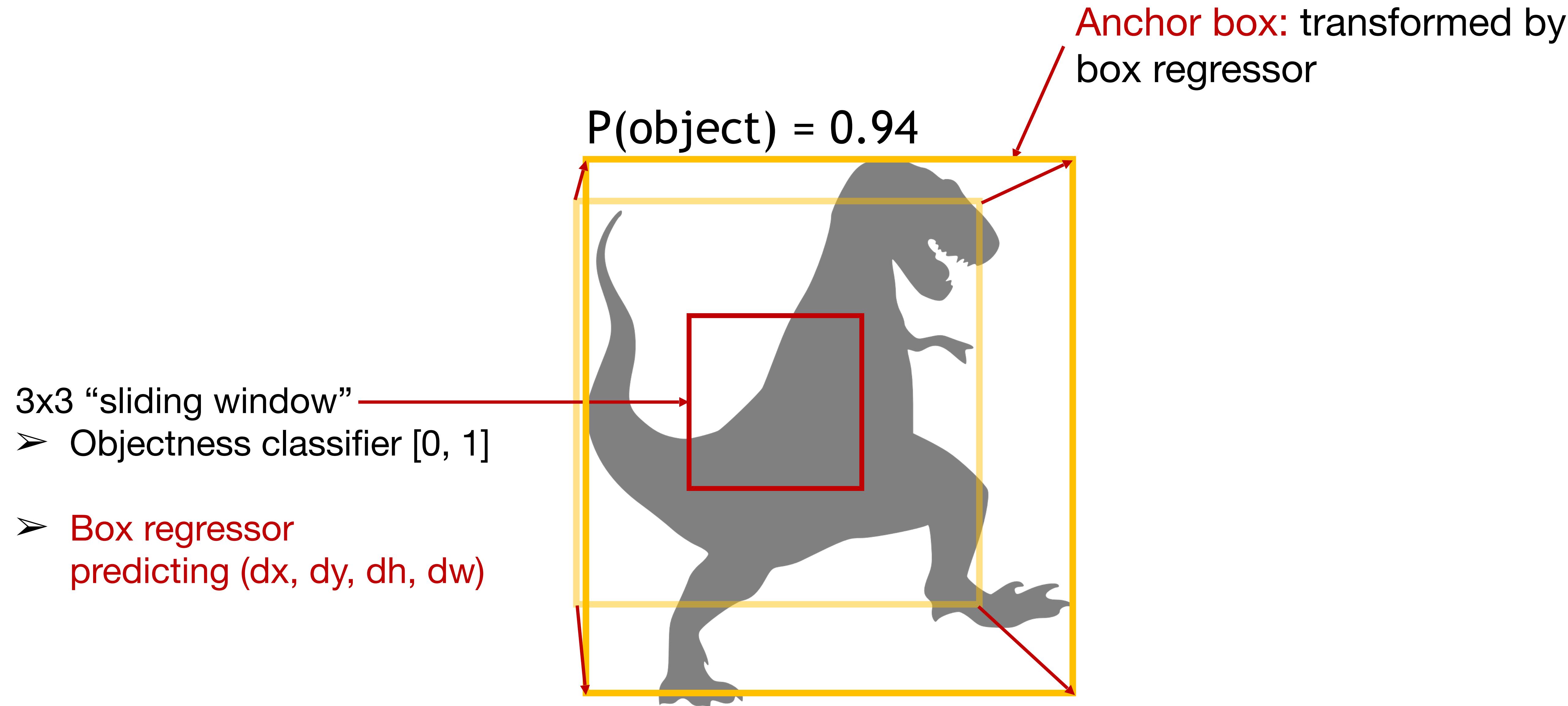
Objectness score

$$P(\text{object}) = 0.94$$

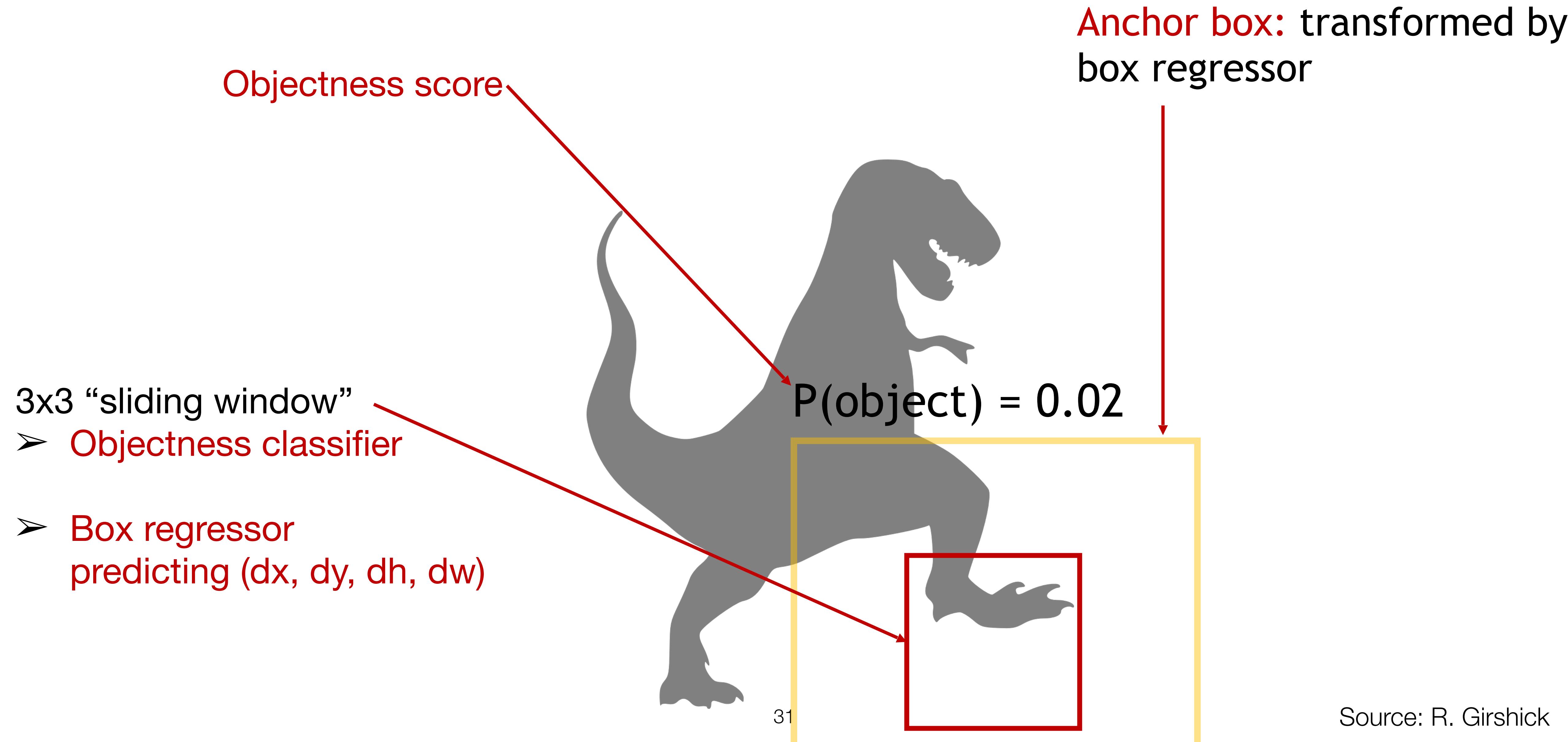
- 3x3 “sliding window”
- Objectness classifier [0, 1]
  - Box regressor  
predicting ( $dx$ ,  $dy$ ,  $dh$ ,  $dw$ )



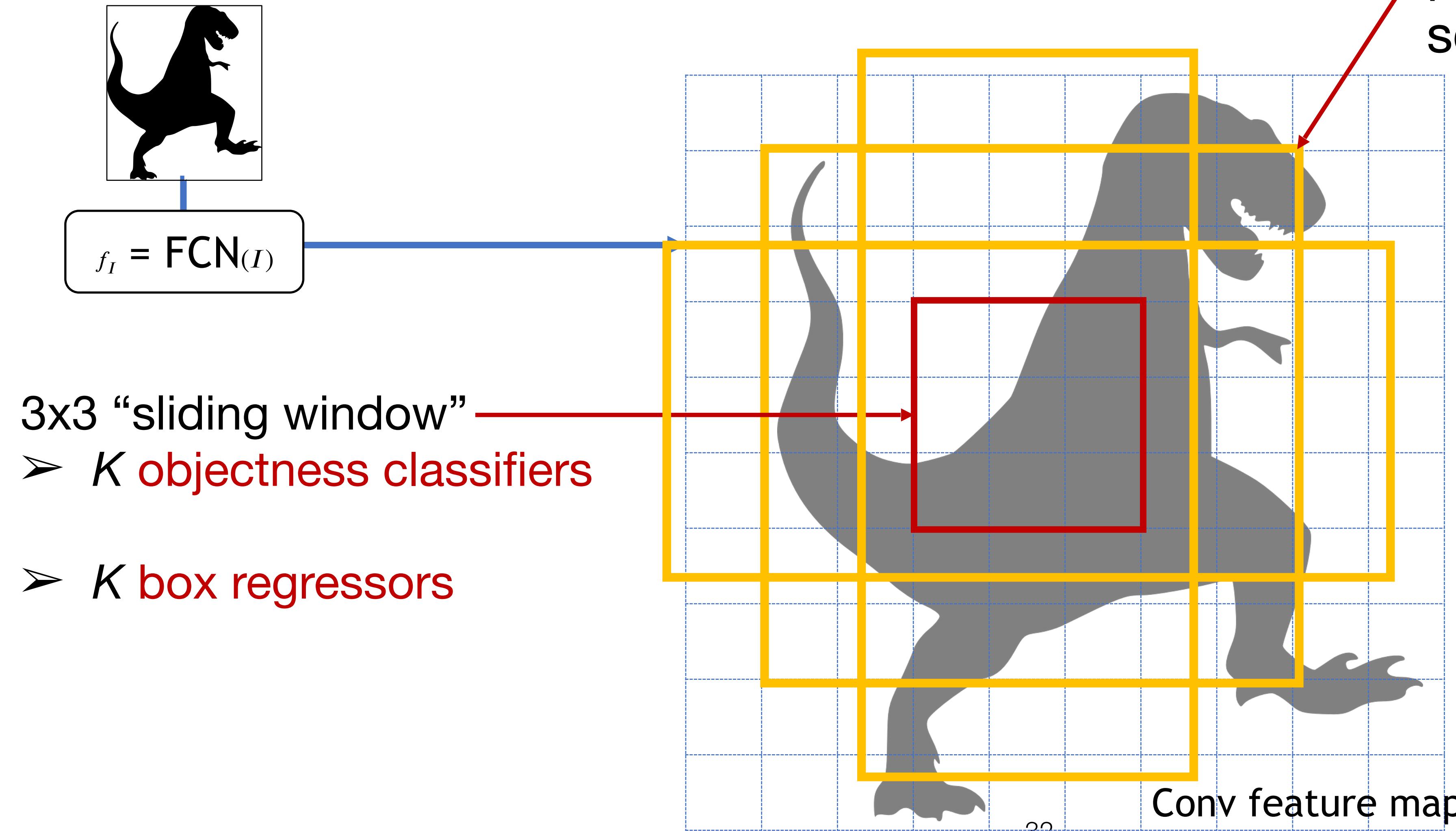
# RPN: Prediction (on object)



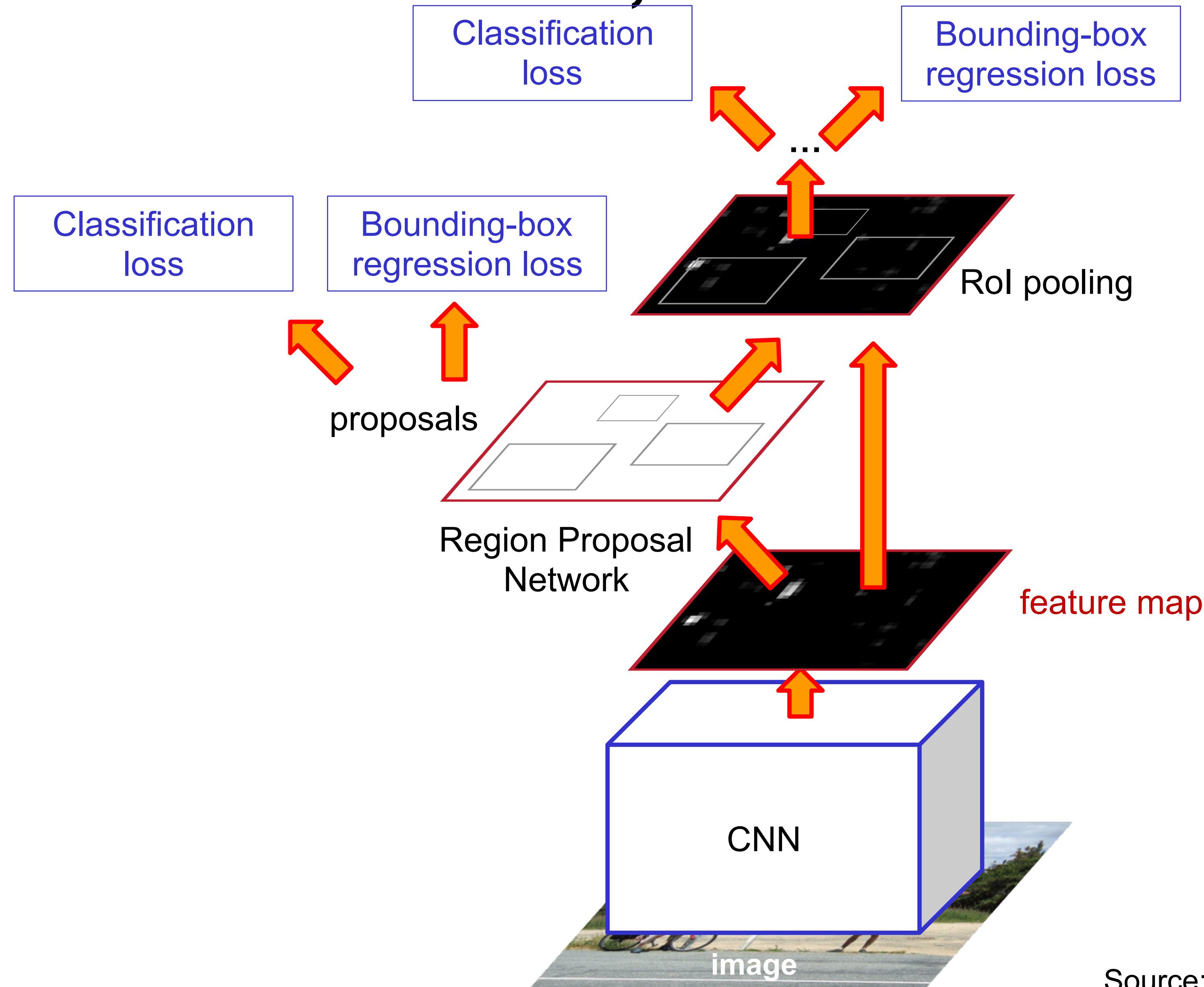
# RPN: Prediction (off object)



# RPN: Multiple Anchors



# One network, four losses



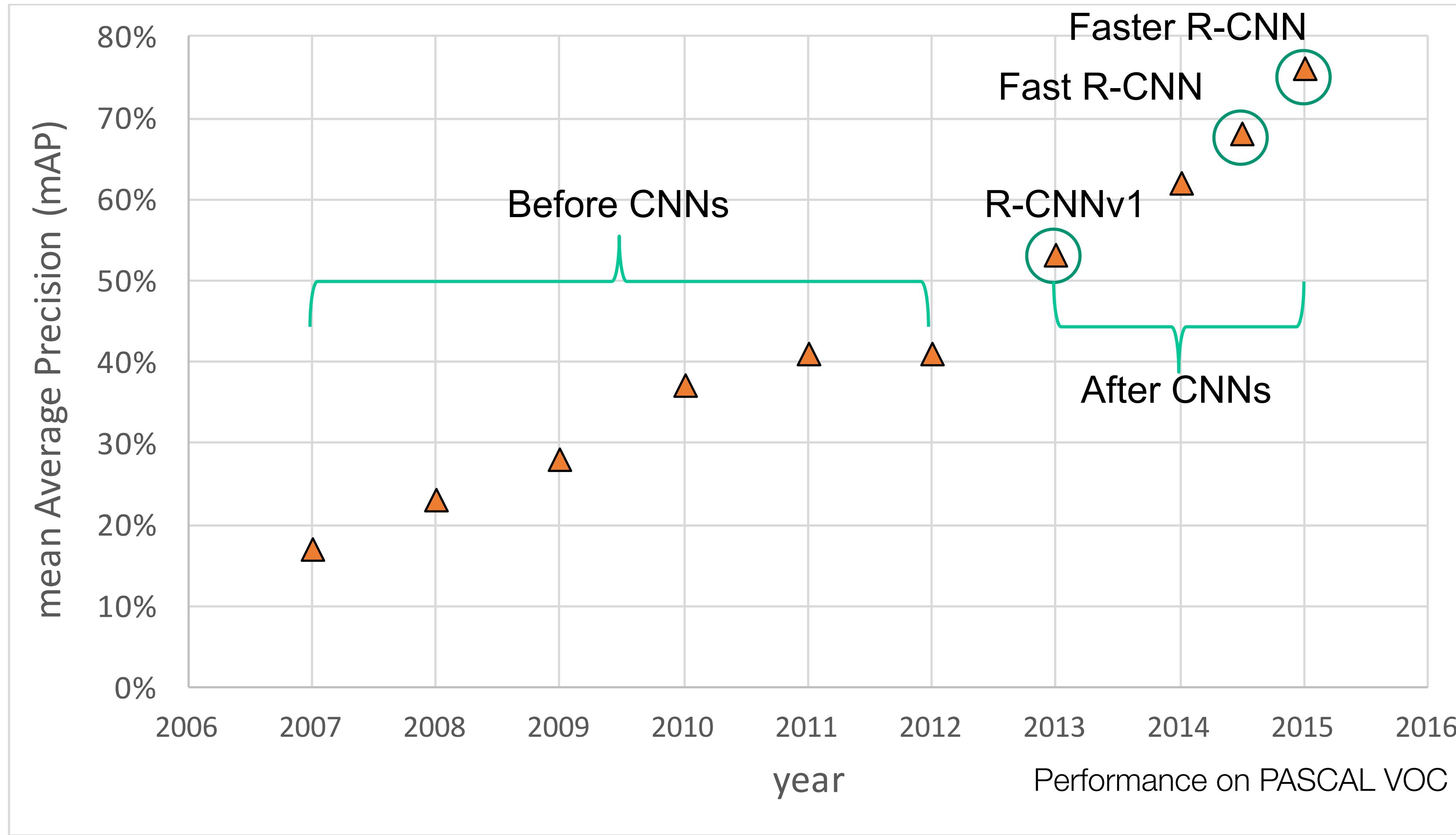
Source: R. Girshick, K. He, S. Lazebnik

# Faster R-CNN results

system	time	07 data	07+12 data
R-CNN	~50s	66.0	-
Fast R-CNN	~2s	66.9	70.0
Faster R-CNN	198ms	<b>69.9</b>	<b>73.2</b>

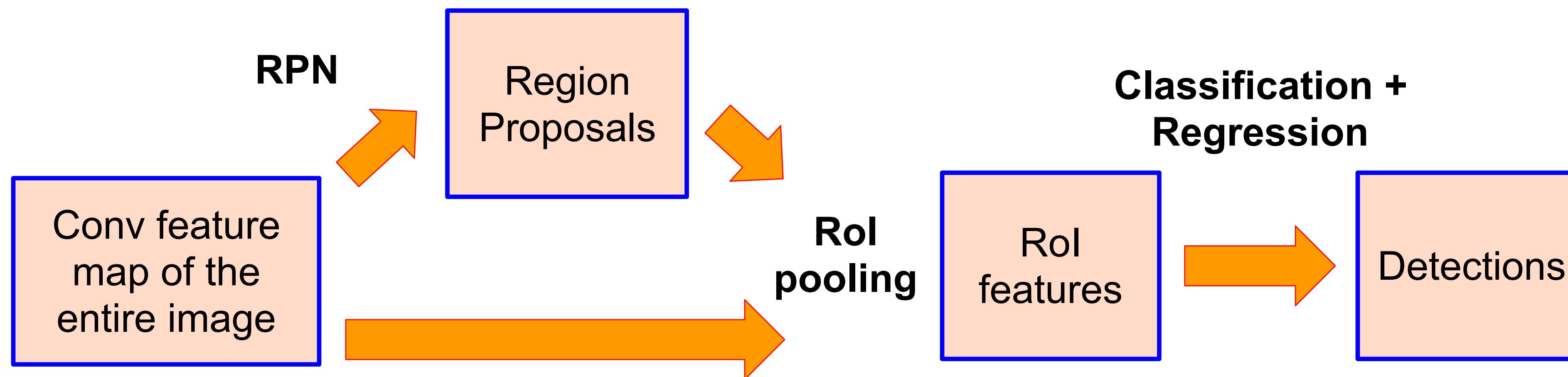
detection mAP on PASCAL VOC 2007, with VGG-16 pre-trained on ImageNet

# Object detection progress

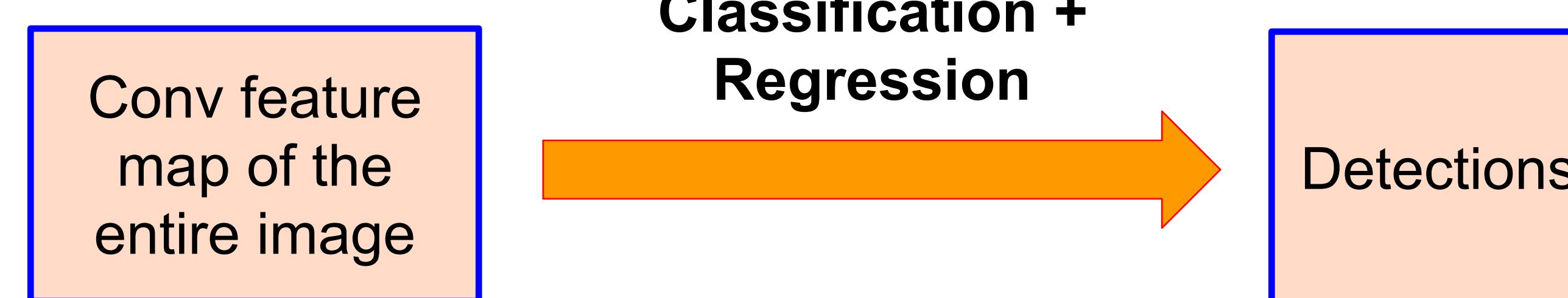


# Streamlined detection architectures

- The Faster R-CNN pipeline separates proposal generation and region classification:

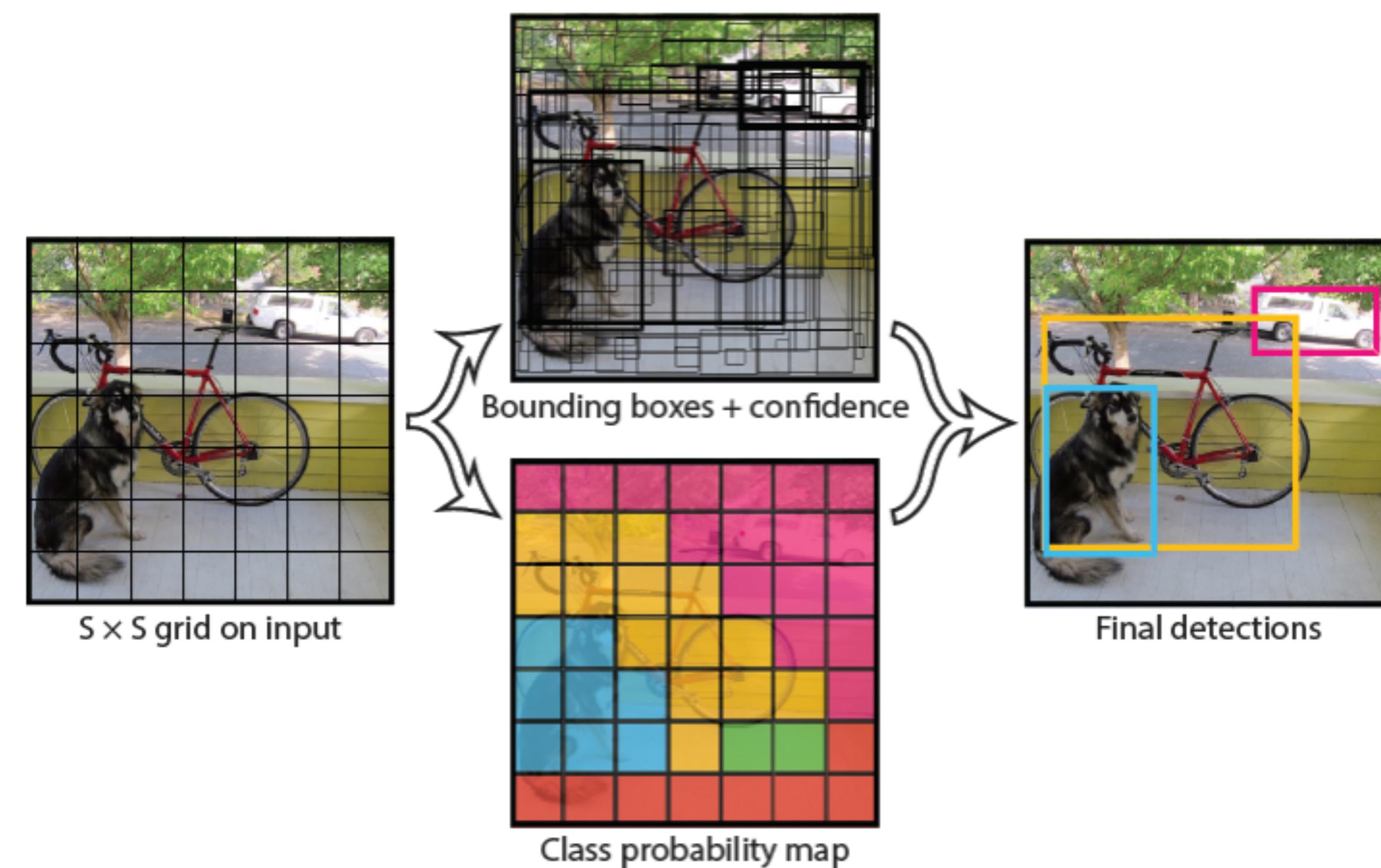


- Is it possible do detection in one shot?



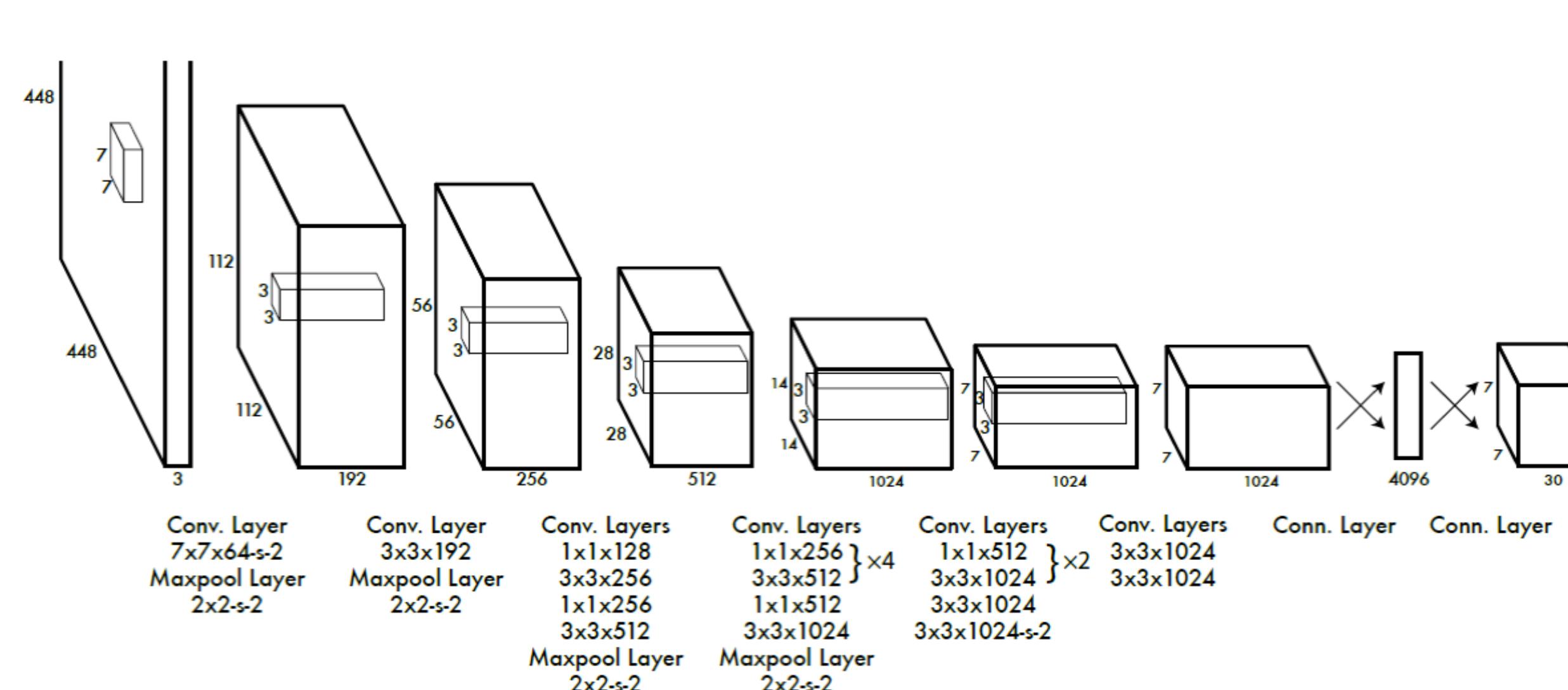
# Single-stage object detector

- Divide the image into a coarse grid and directly predict class label and a few candidate boxes for each grid cell



# YOLO detector

1. Take conv feature maps at 7x7 resolution
2. Predict, at each location, a score for each class and 2 bboxes w/ confidences
  - For PASCAL, output is 7x7x30 ( $30 = 20 + 2^*(4+1)$ )
  - 7x speedup over Faster R-CNN (45-155 FPS vs. 7-18 FPS)  
but less accurate (e.g. 65% vs. 72 mAP%)



# Challenges in object detection

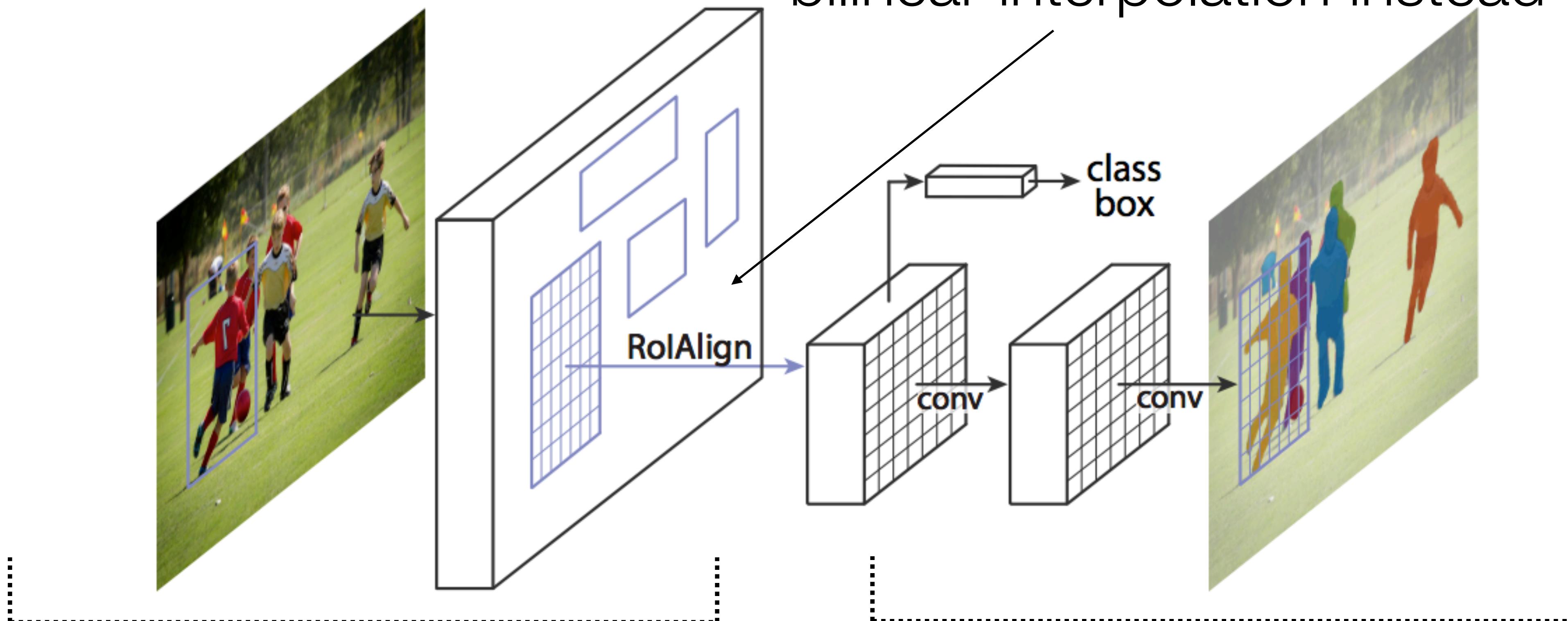
# Beyond bounding boxes: instance segmentation



Predict segmentation mask for each object  
From COCO [Lin et al., 2014]

# Instance segmentation

ROI pooling with tiny change:  
bilinear interpolation instead of max



Faster R-CNN

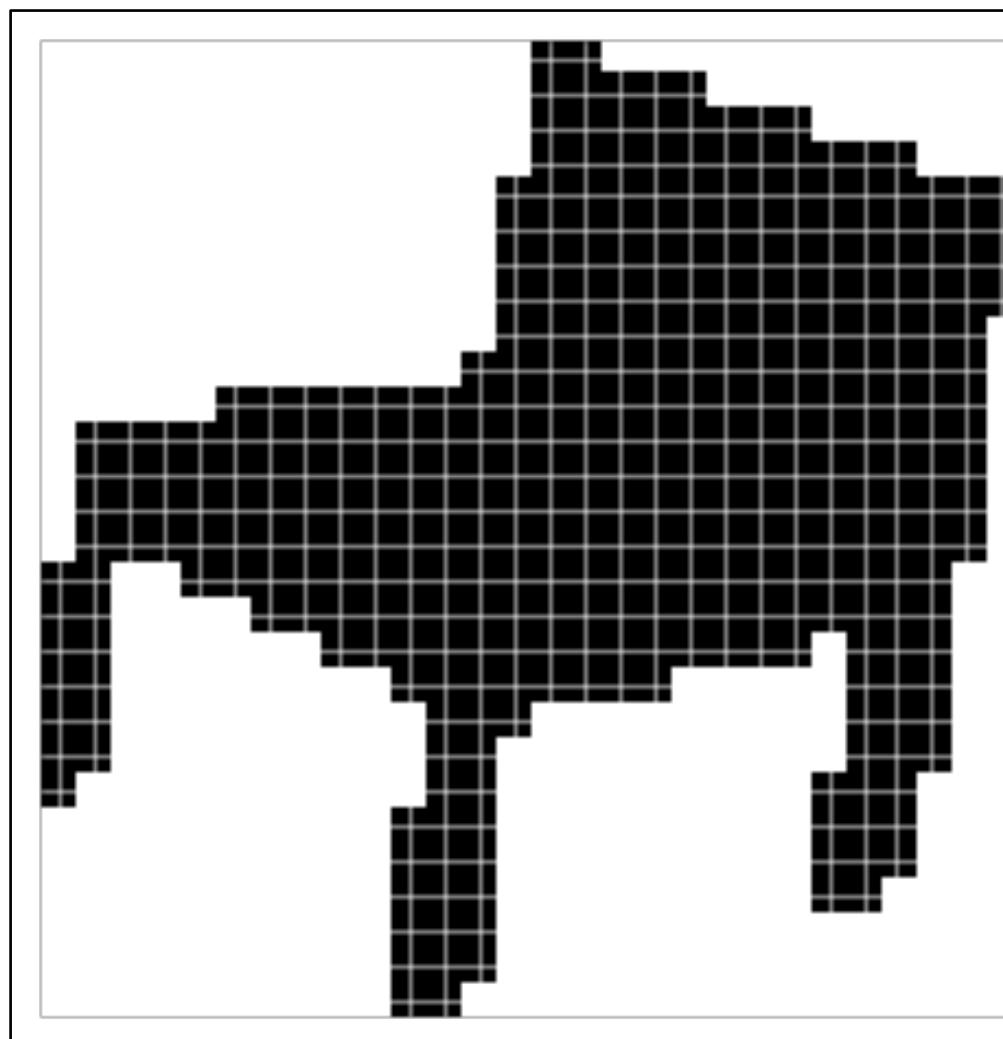
Extra “head” on network  
predicts binary mask

# Example Mask Training Targets

Image with training proposal

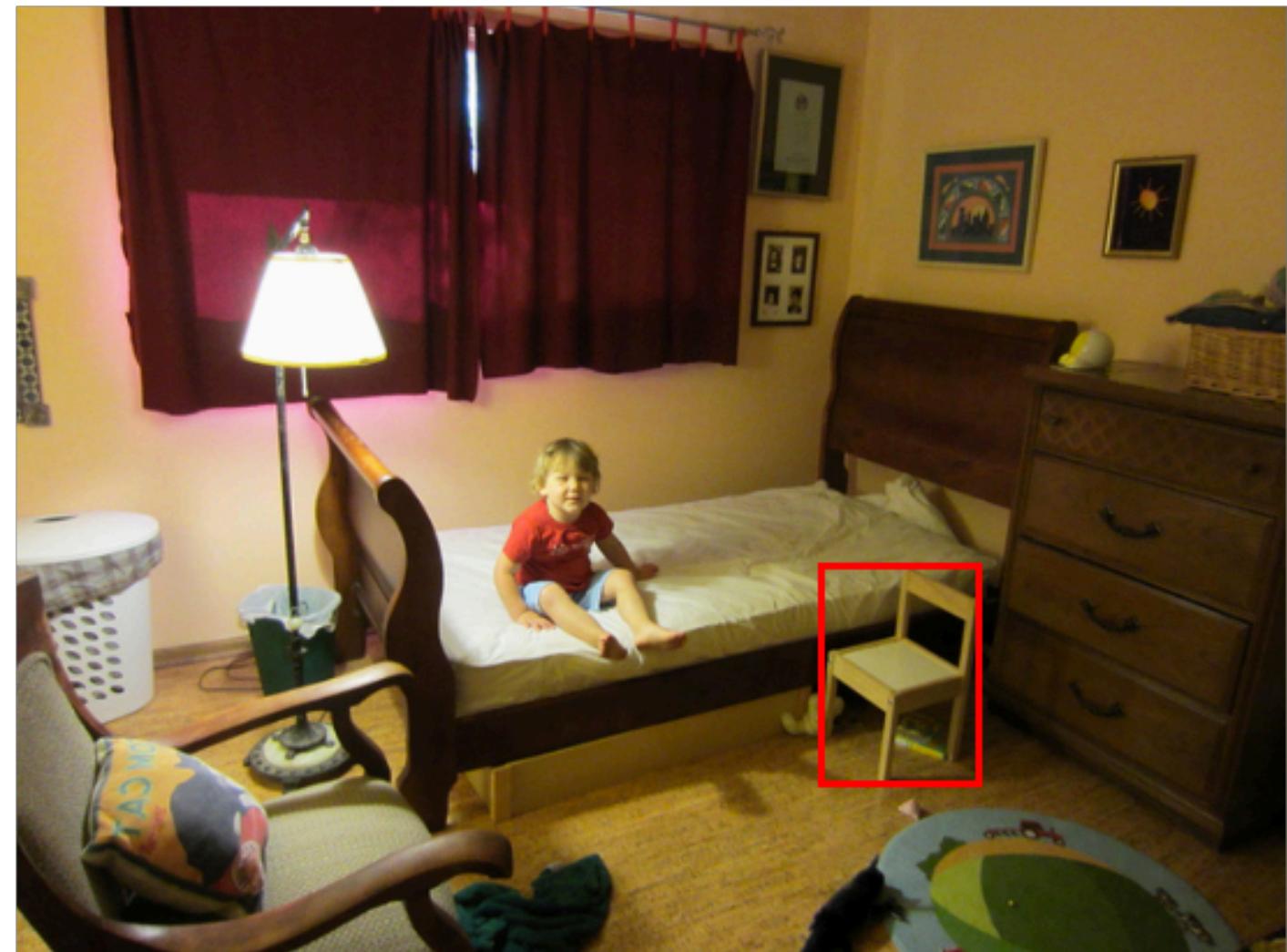


28x28 mask target

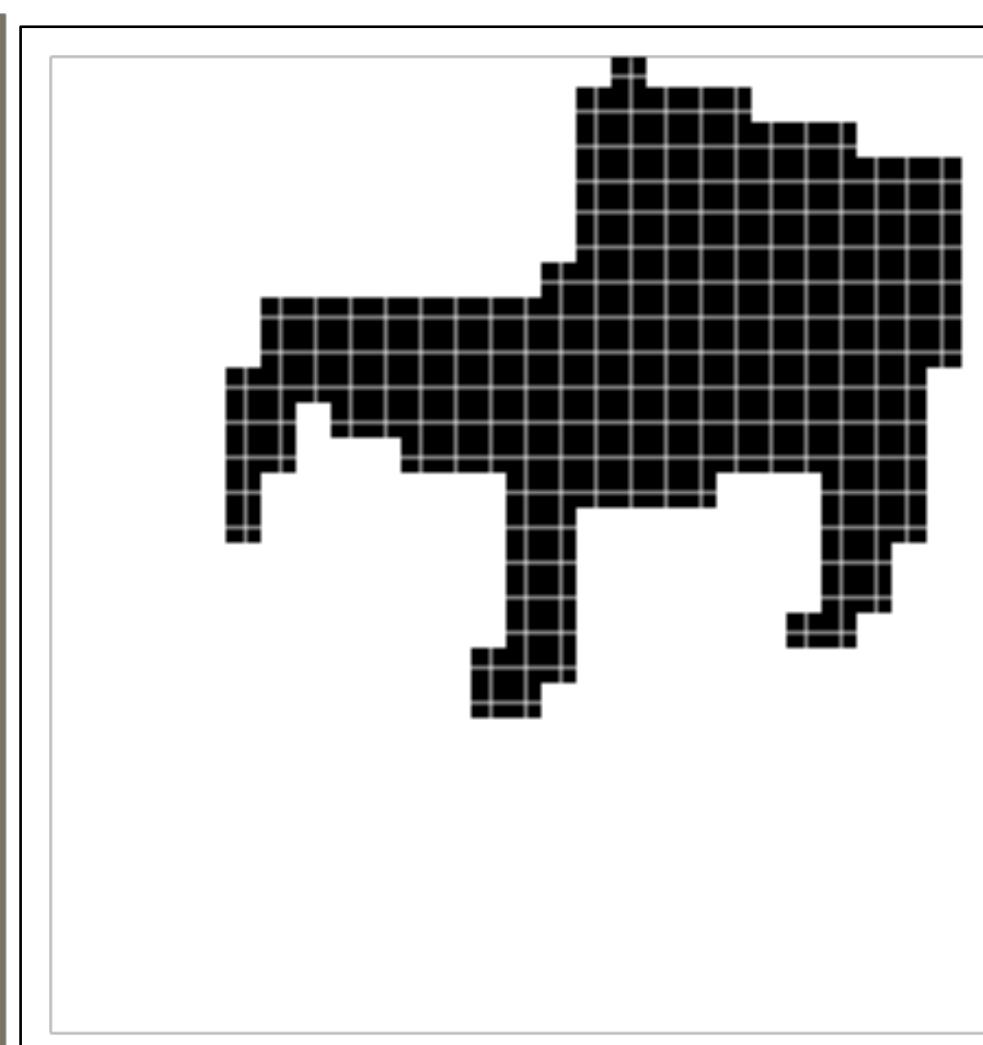
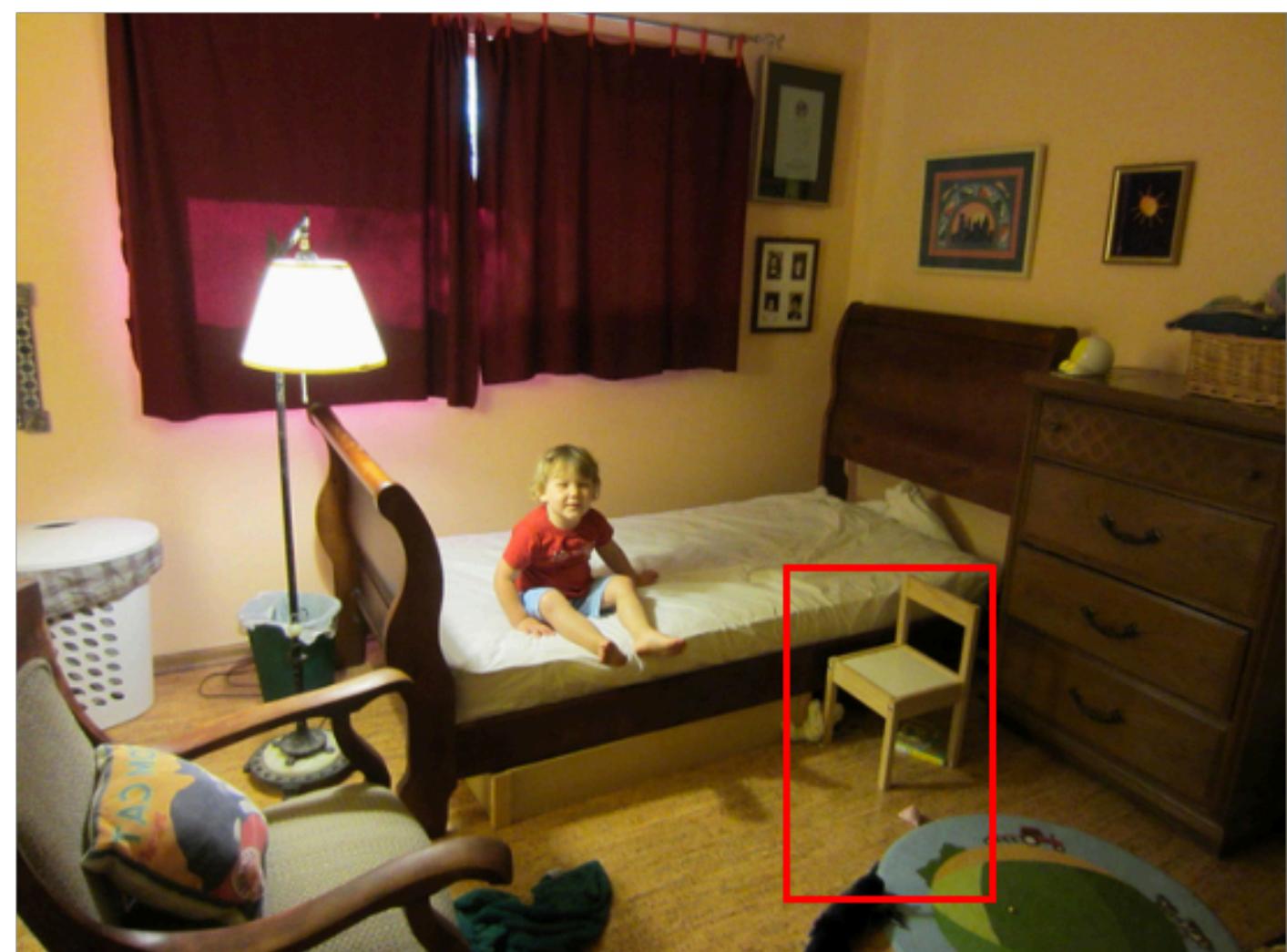
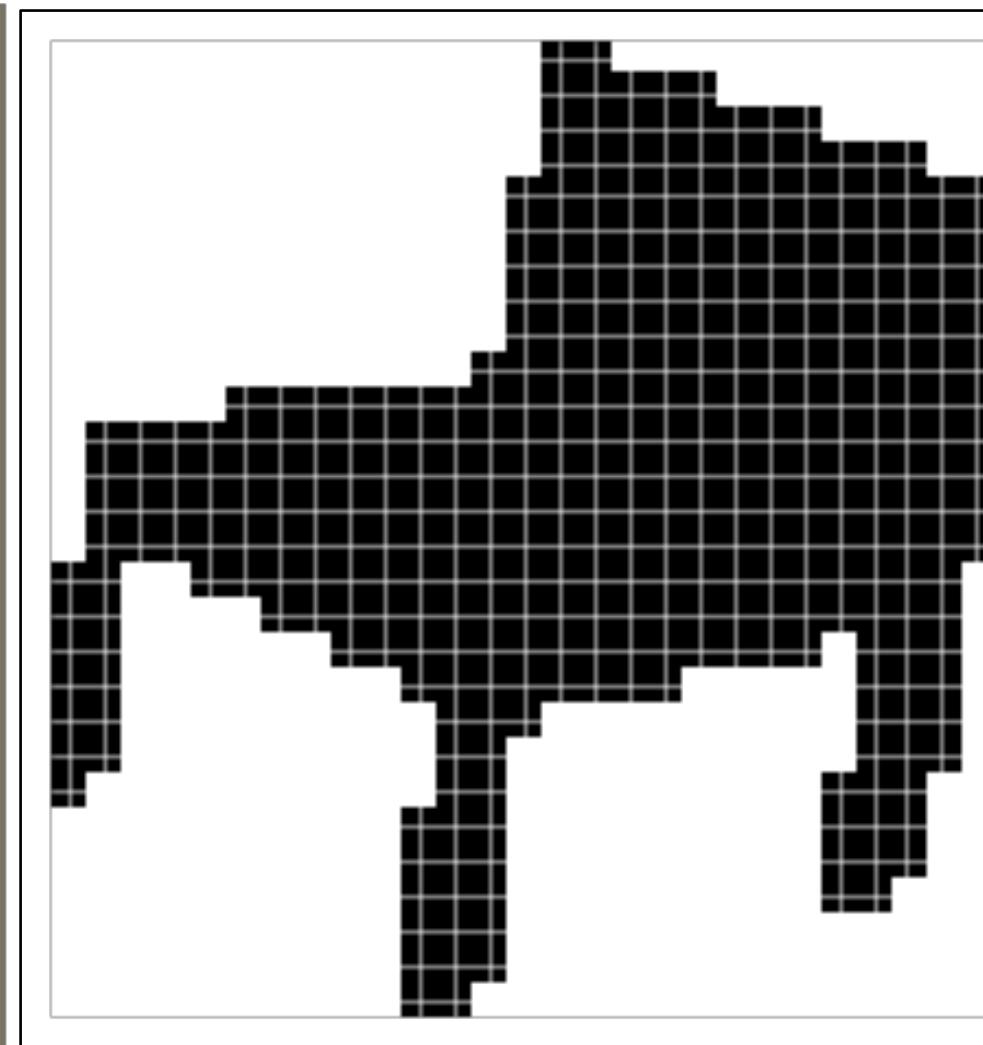


# Example Mask Training Targets

Image with training proposal



28x28 mask target



# Example Mask Training Targets

Image with training proposal



28x28 mask target

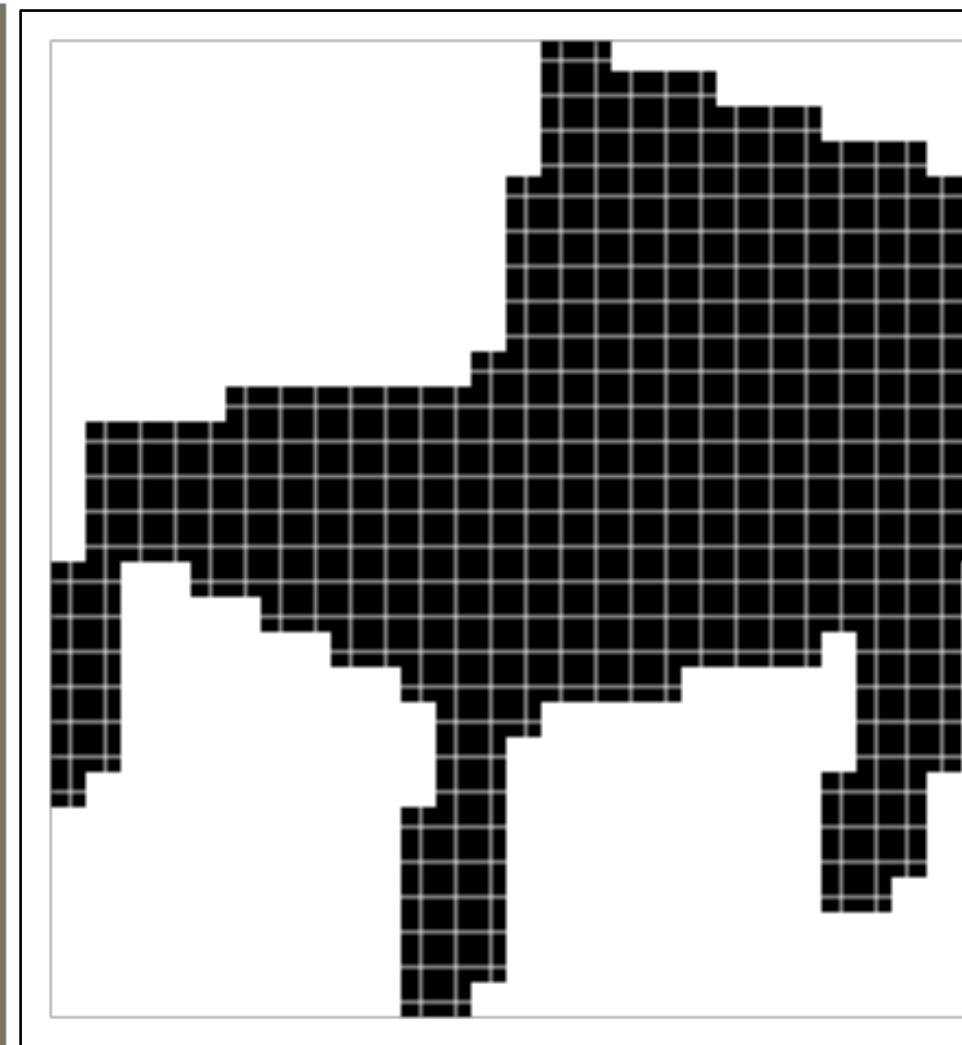
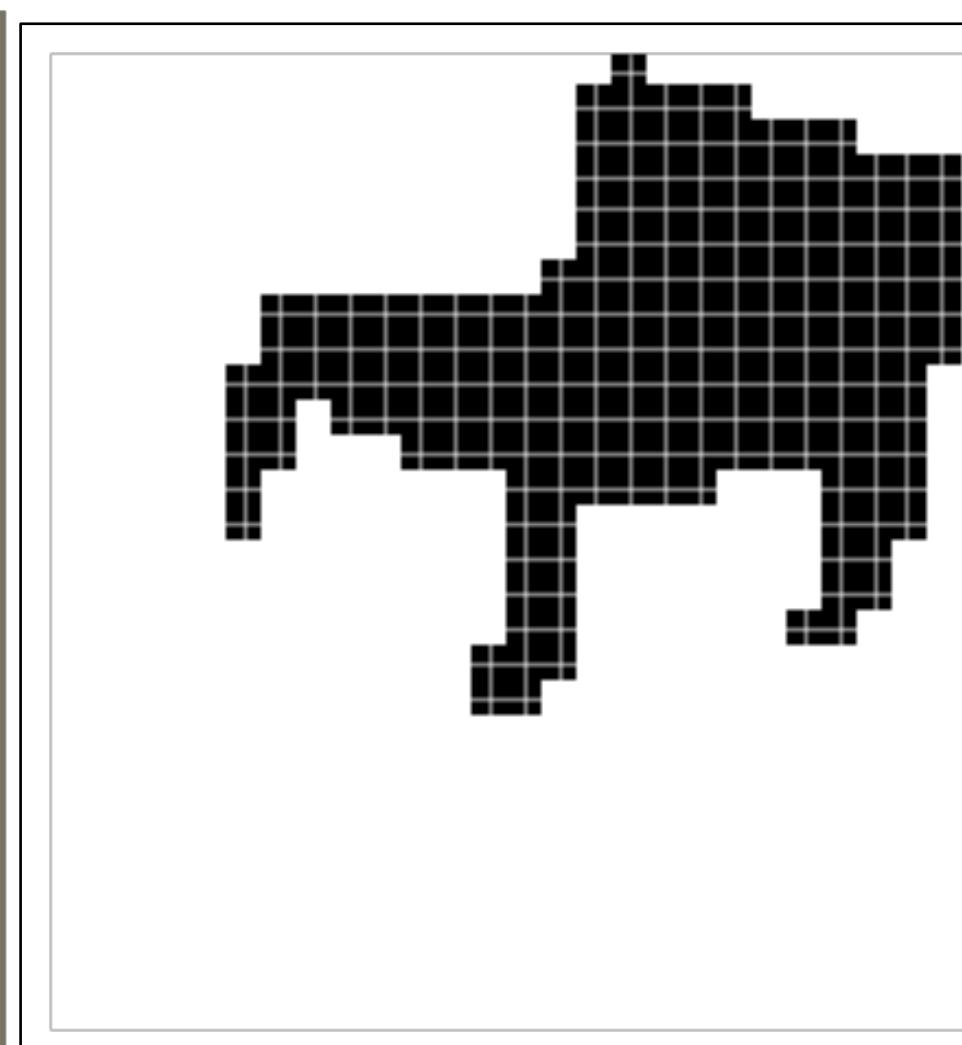
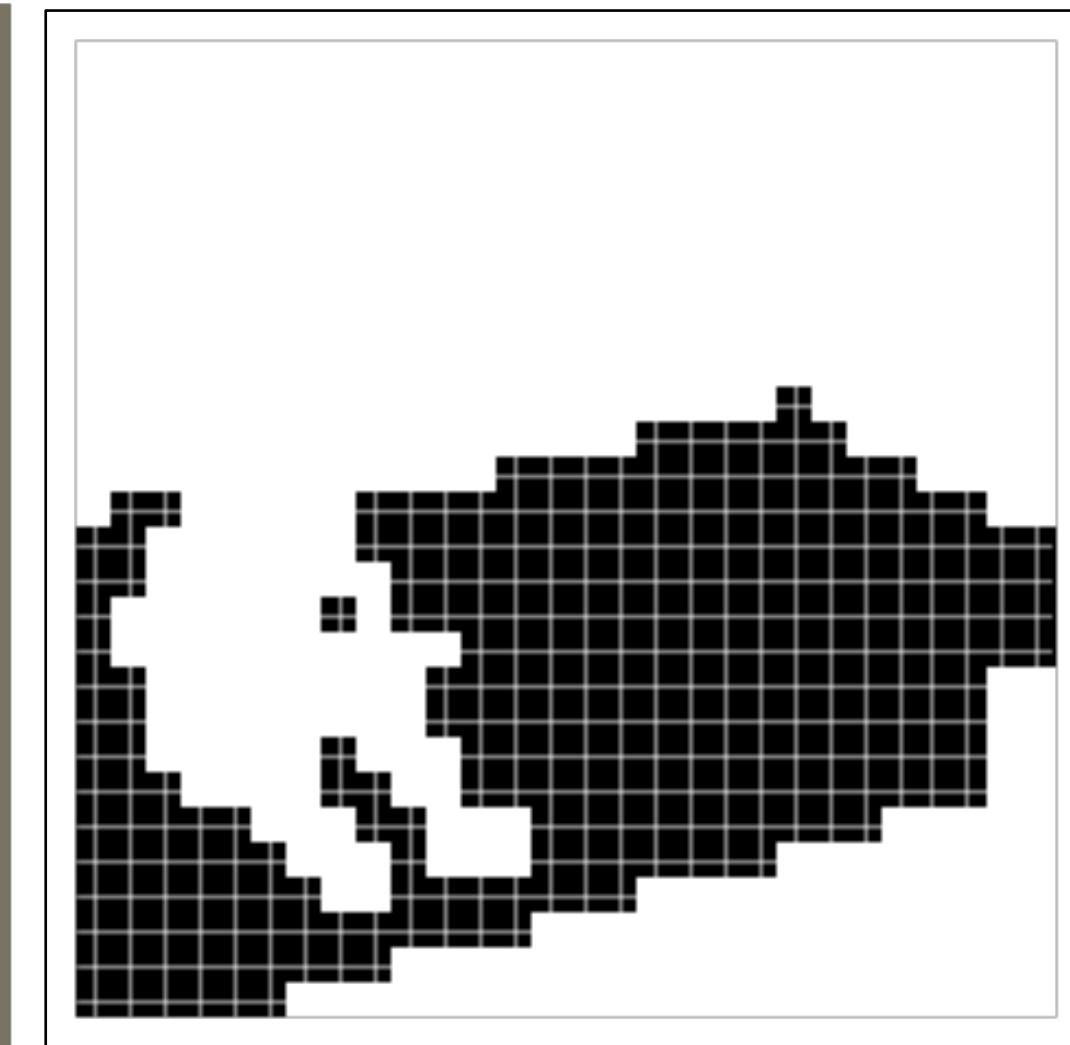


Image with training proposal



28x28 mask target



# Example Mask Training Targets

Image with training proposal



28x28 mask target

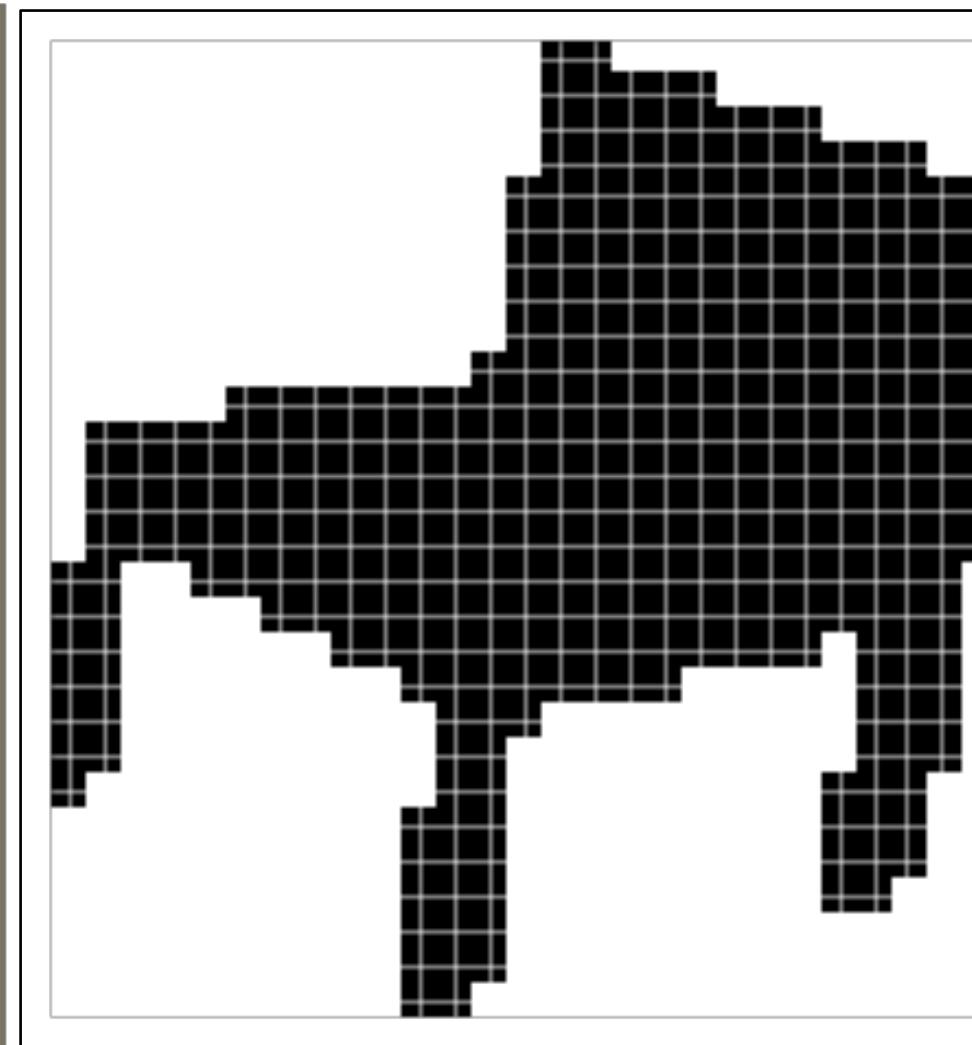
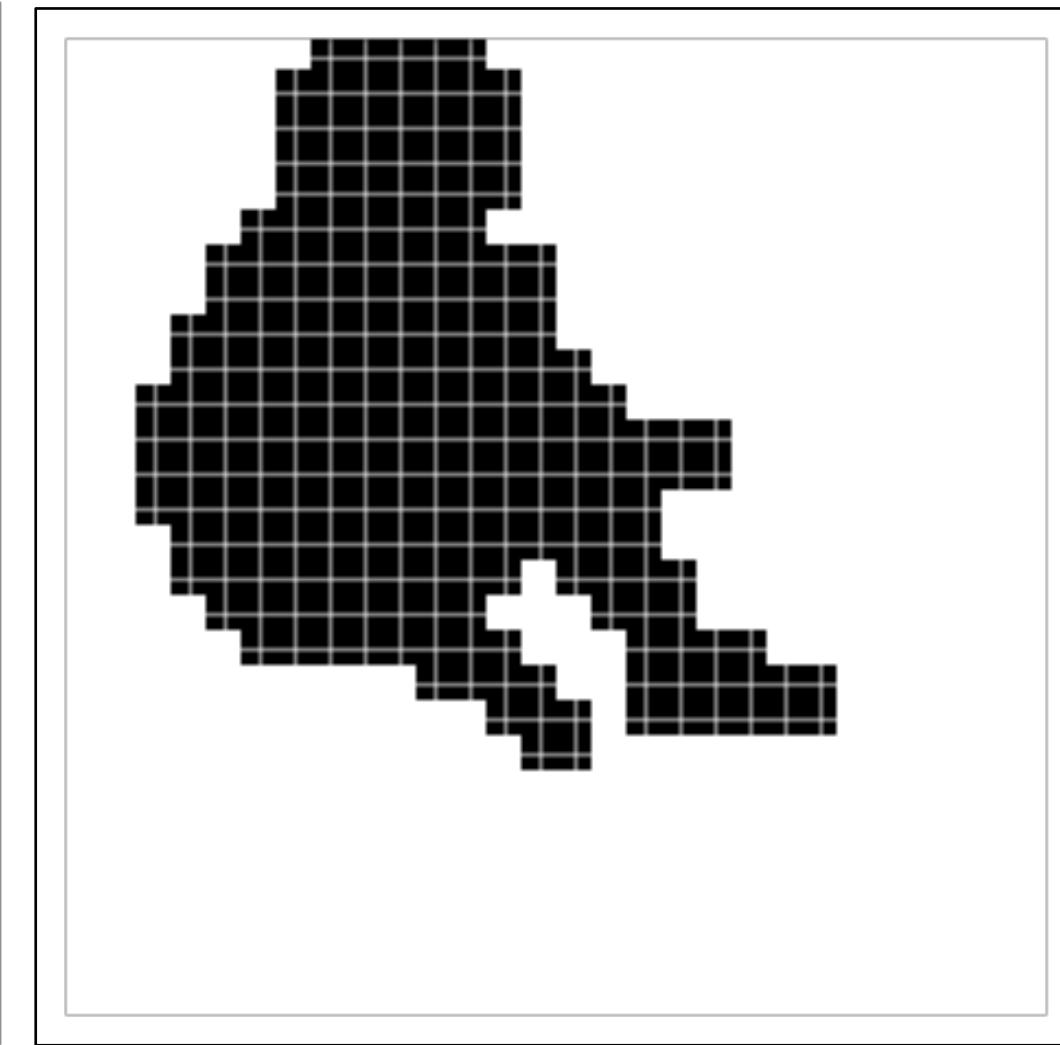
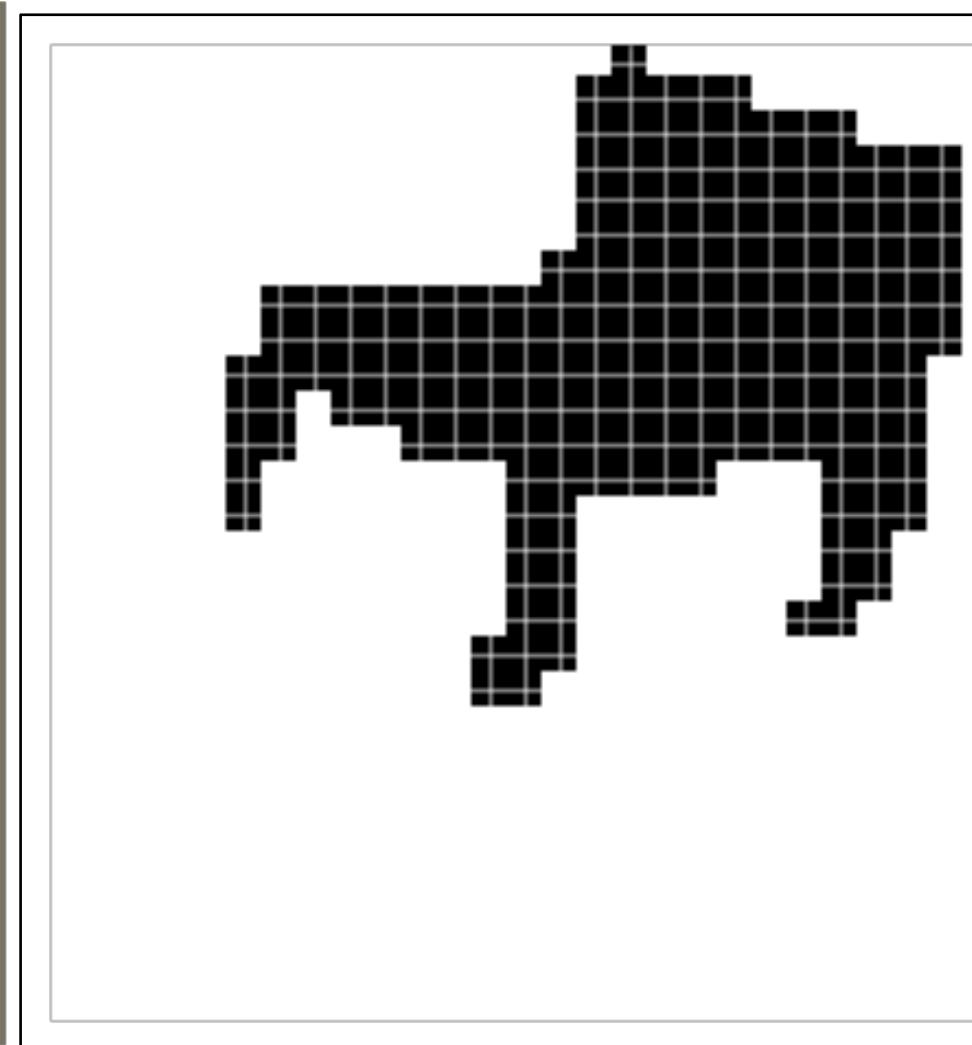
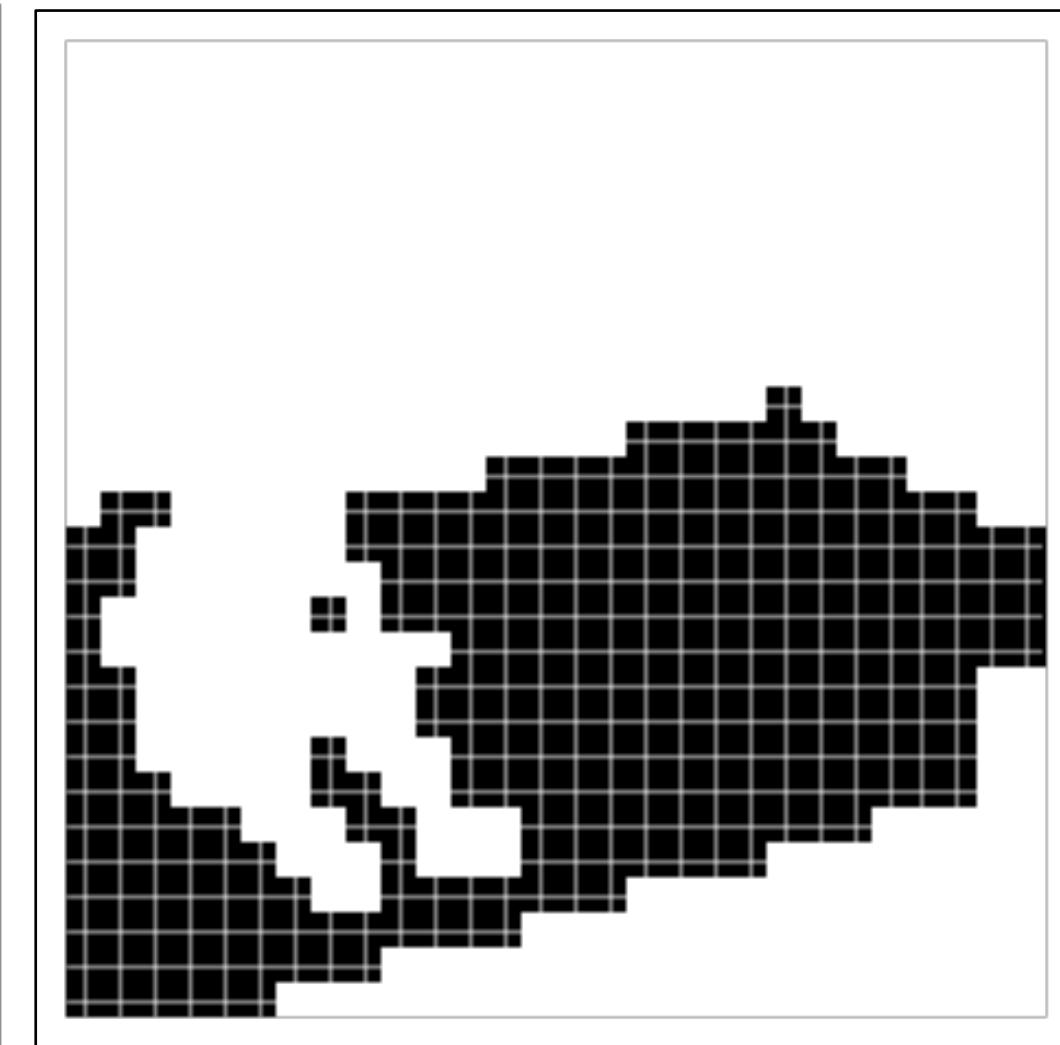


Image with training proposal



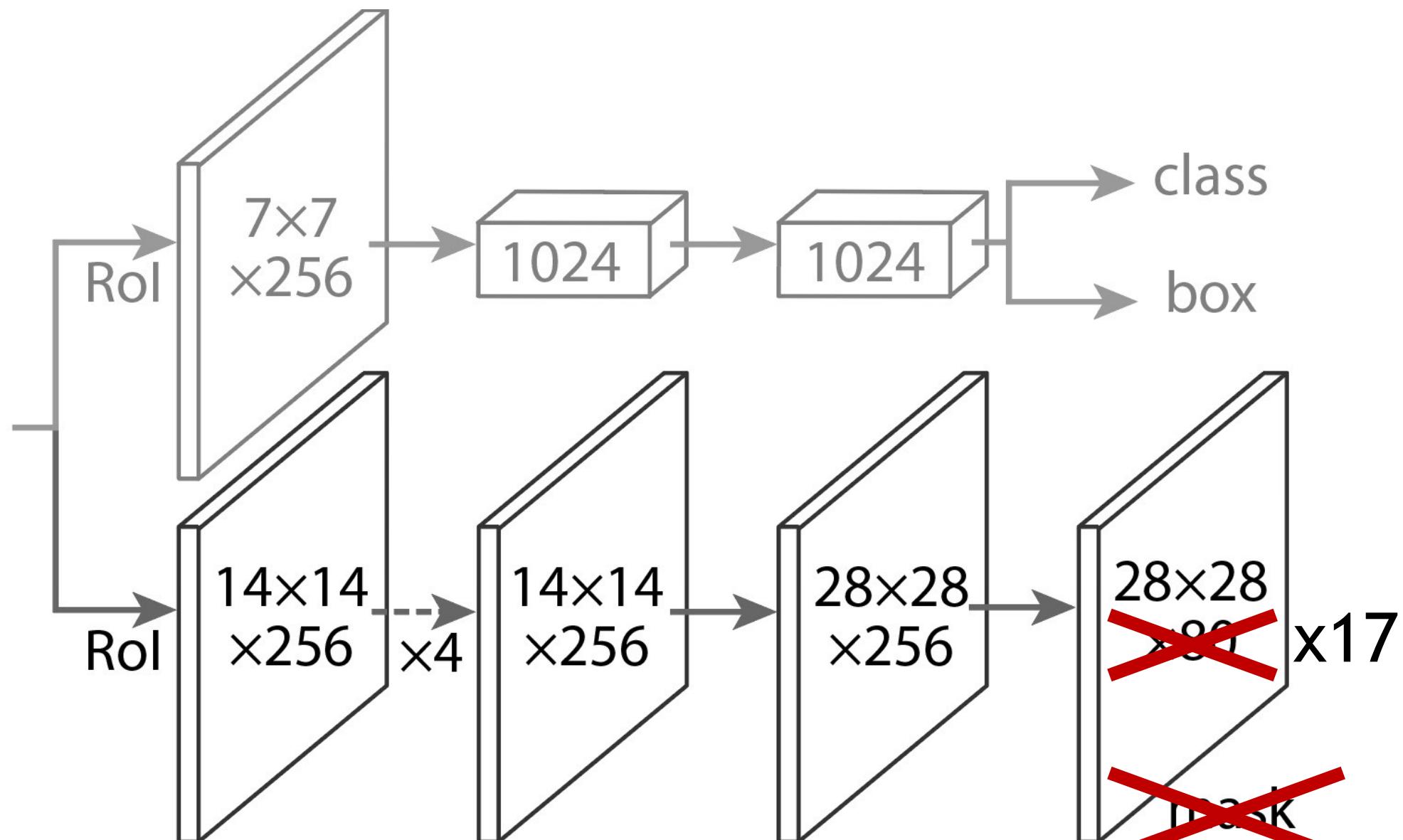
28x28 mask target





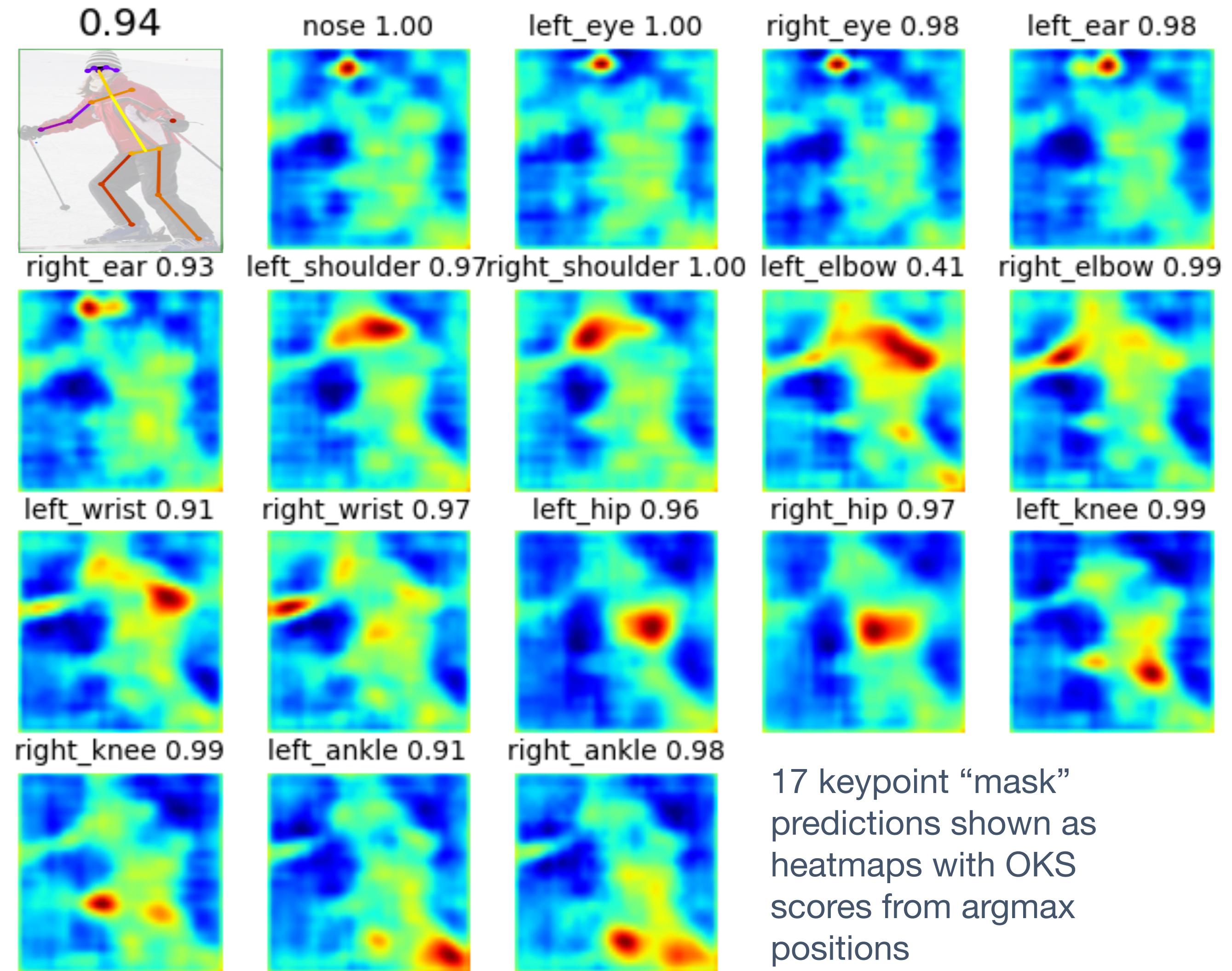


# Human Pose



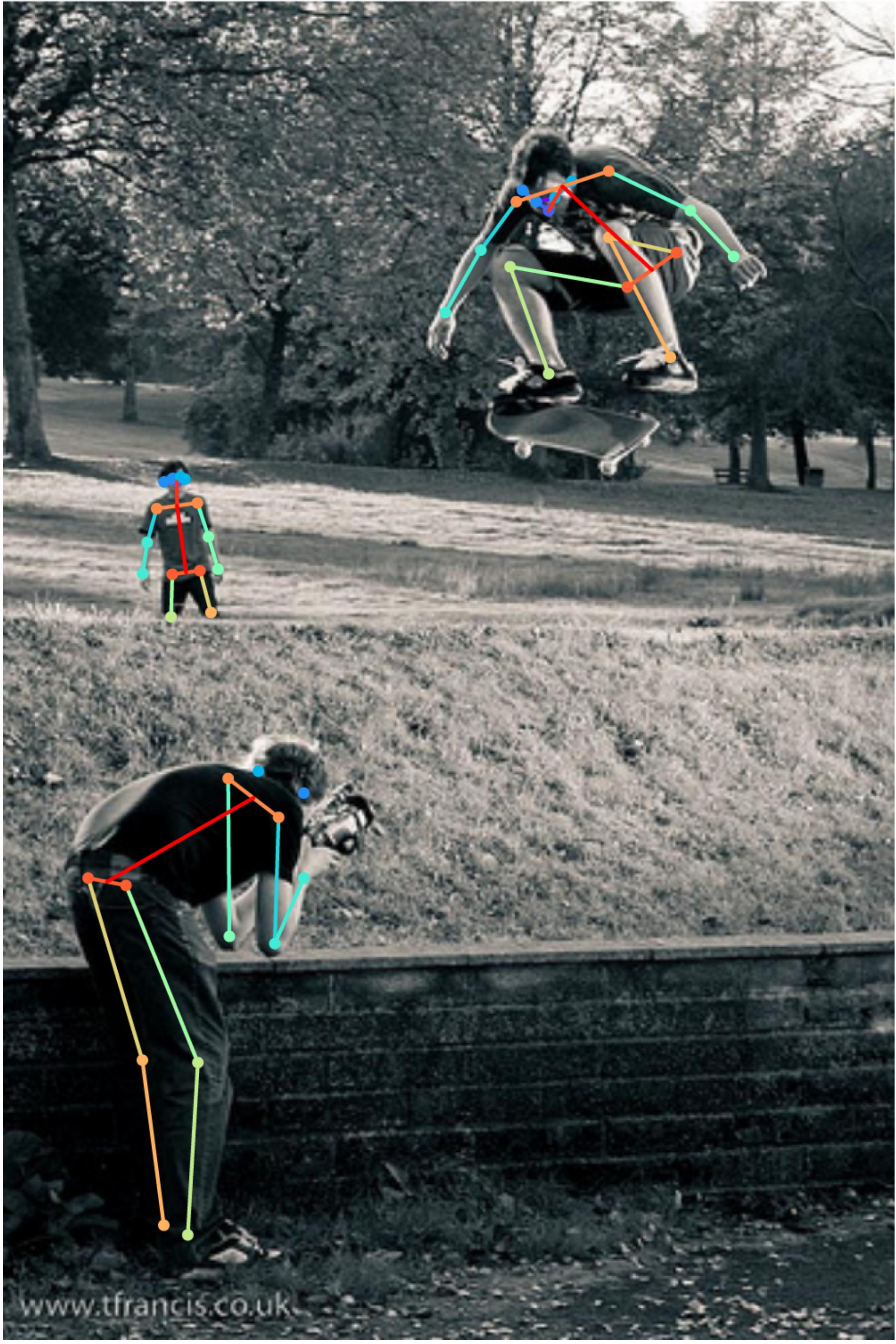
(Not shown: Head architecture is slightly different for keypoints)

- Add keypoint head (28x28x17)
- Predict one “mask” for each keypoint
- Softmax over spatial locations (encodes one keypoint per mask “prior”)



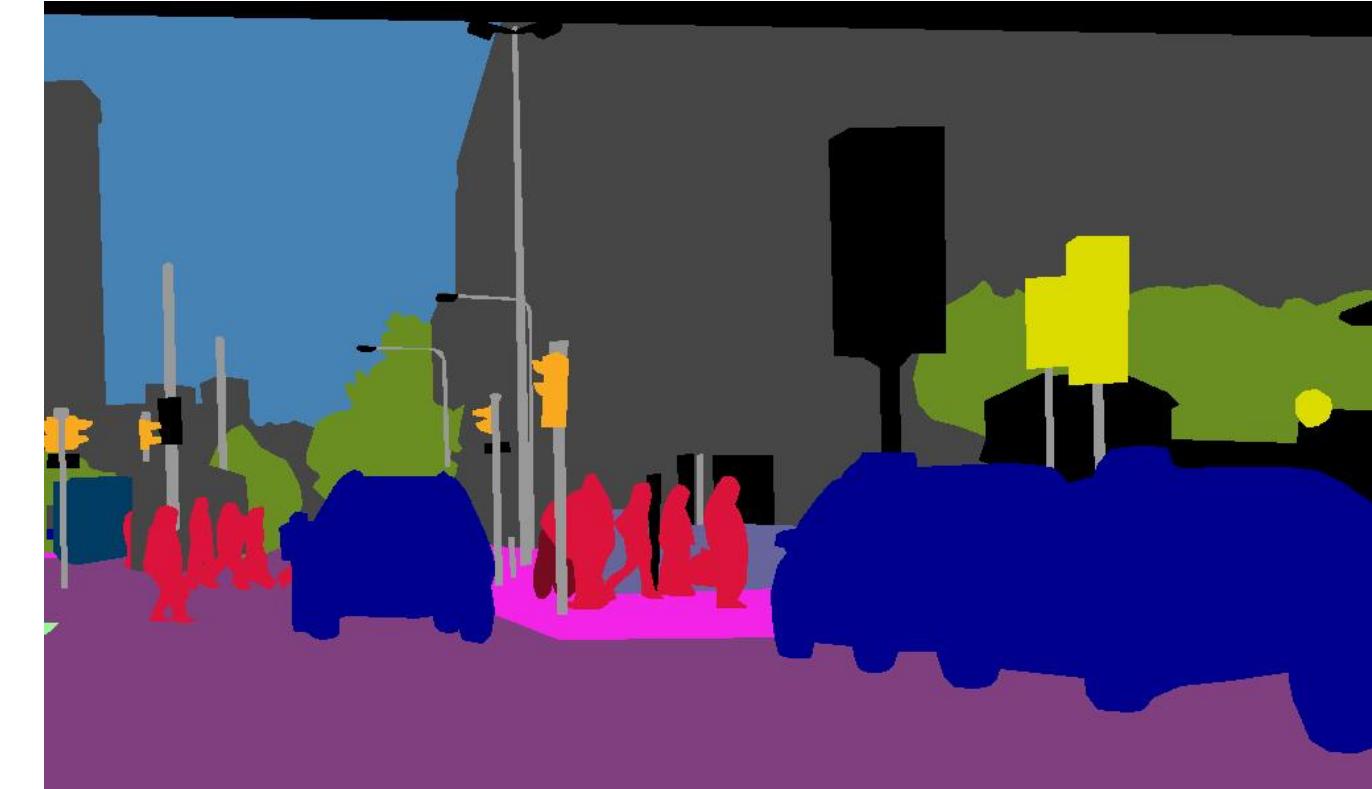
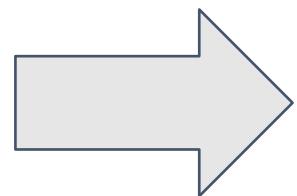
17 keypoint “mask”  
predictions shown as  
heatmaps with OKS  
scores from argmax  
positions

Source: R. Girshick

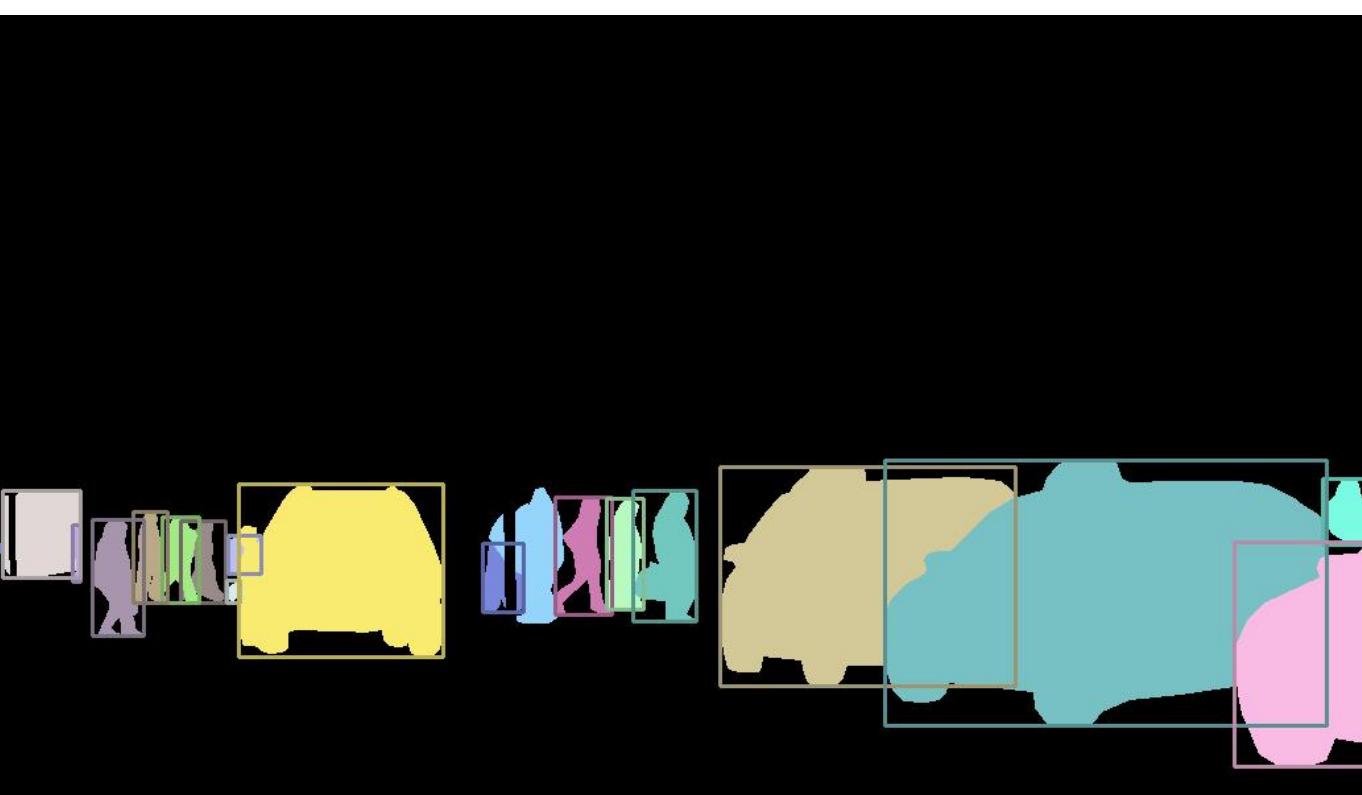




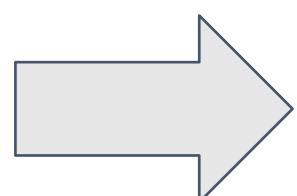
# Panoptic Segmentation



Semantic segmentation “stuff”



Instance detection, “things”



panoptic segmentation [Krilov et al. 2018]  
predict label + instance id per pixel

UPSNNet-101-M: Cityscapes

