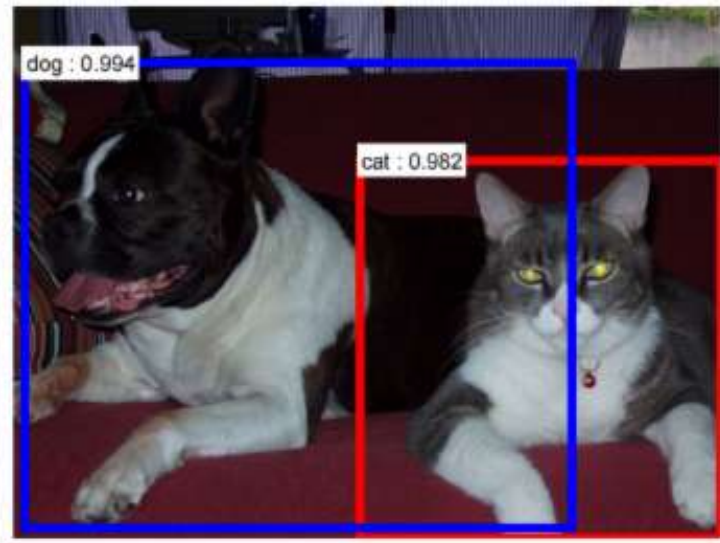
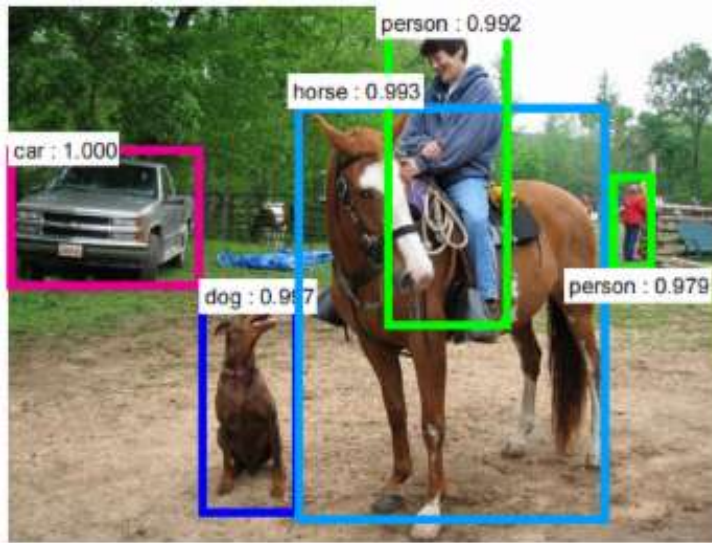
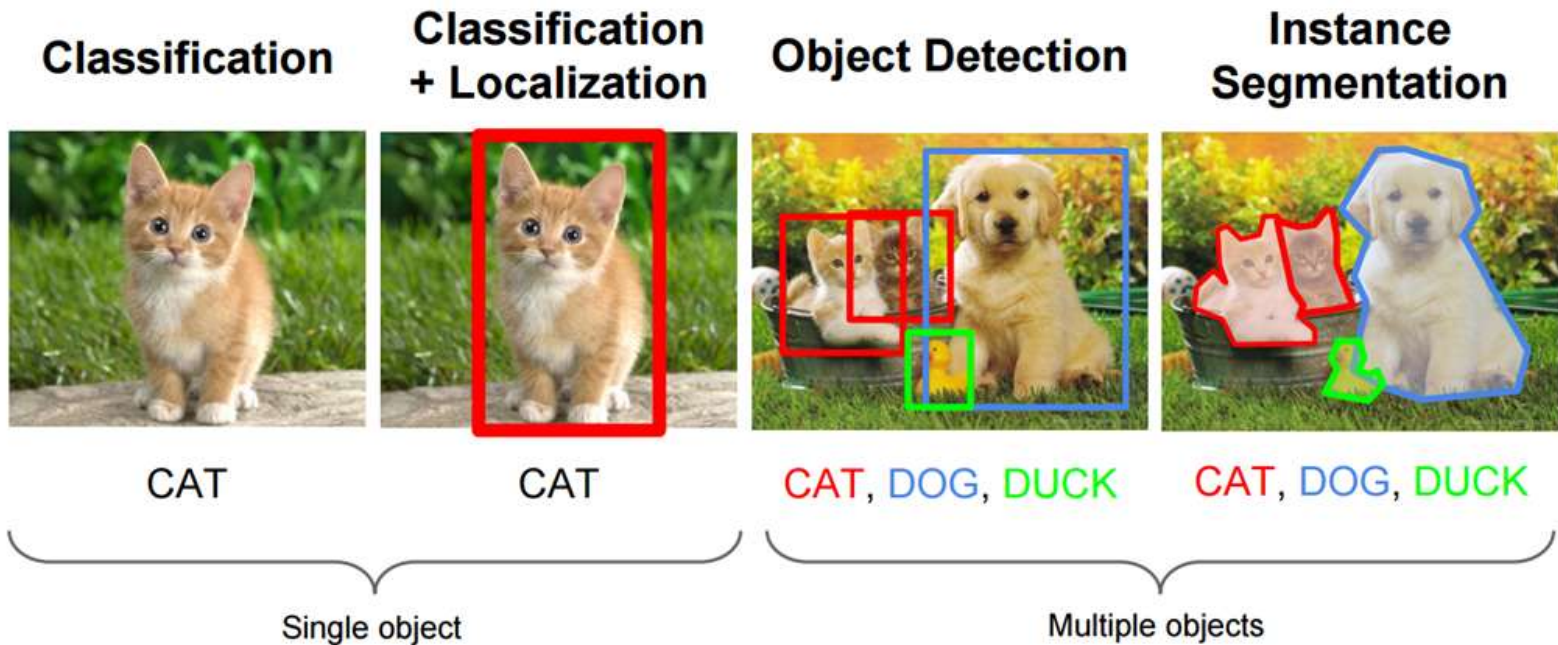


Object detection (deep methods)

Object detection

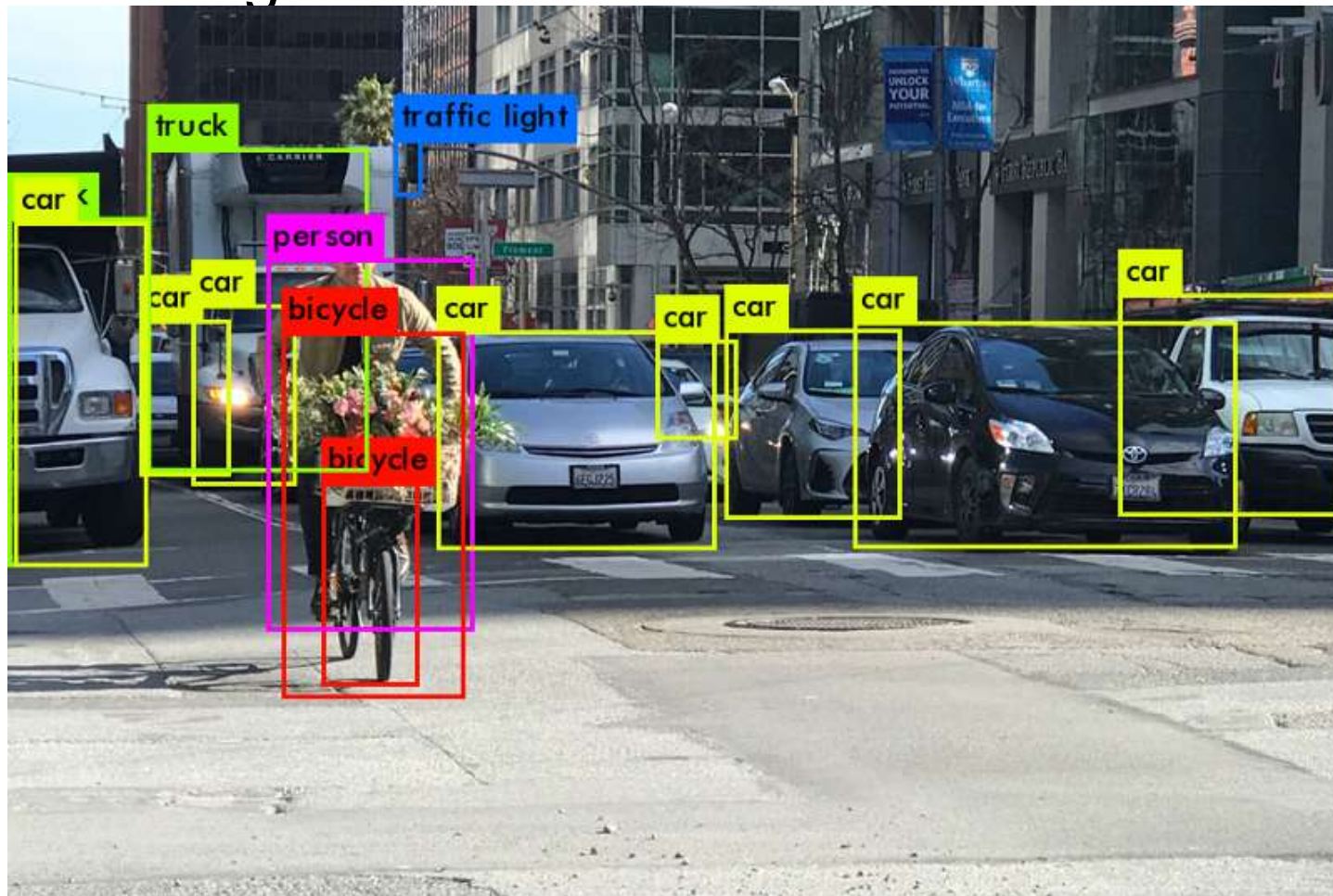


What is object detection



What are the challenges of object detection?

- Images may contain more than one class, multiple instances from the same class
- Bounding box localization & Evaluation

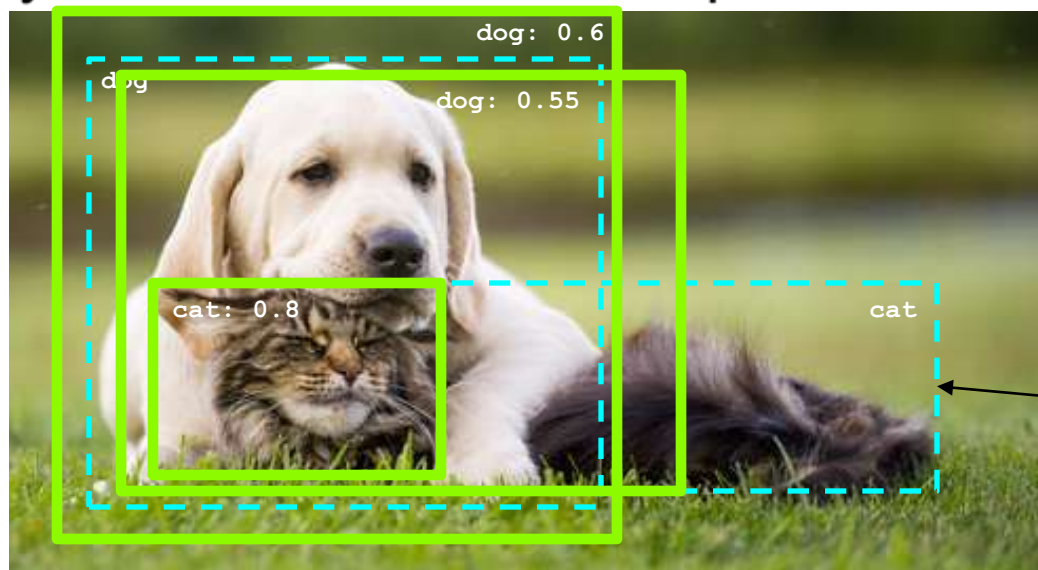


Outline

- Task definition and evaluation
- Generic object detection before deep learning
- Zoo of deep detection approaches
 - R-CNN
 - Fast R-CNN
 - Faster R-CNN
 - Yolo
 - SSD

Object detection evaluation

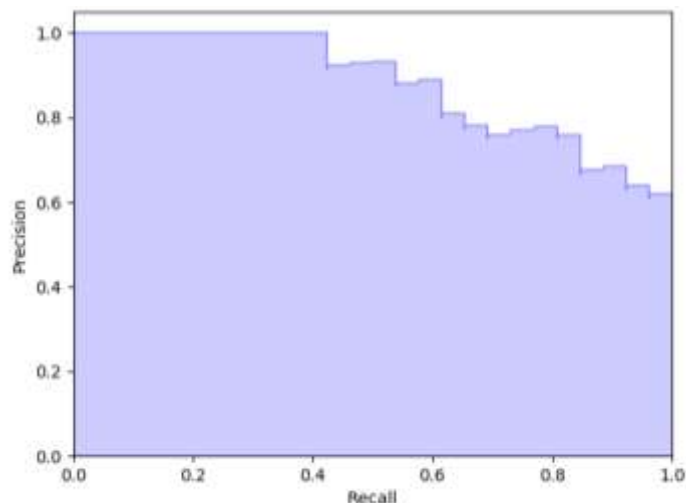
- At test time, predict bounding boxes, class labels, and confidence scores
- For each detection, determine whether it is a true or false positive
 - PASCAL criterion: $\text{Area}(\text{GT} \cap \text{Det}) / \text{Area}(\text{GT} \cup \text{Det}) > 0.5$
 - For multiple detections of the same ground truth box, only one considered a true positive



Ground truth (GT)

Object detection evaluation

- At test time, predict bounding boxes, class labels, and confidence scores
- For each detection, determine whether it is a true or false positive
- 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:

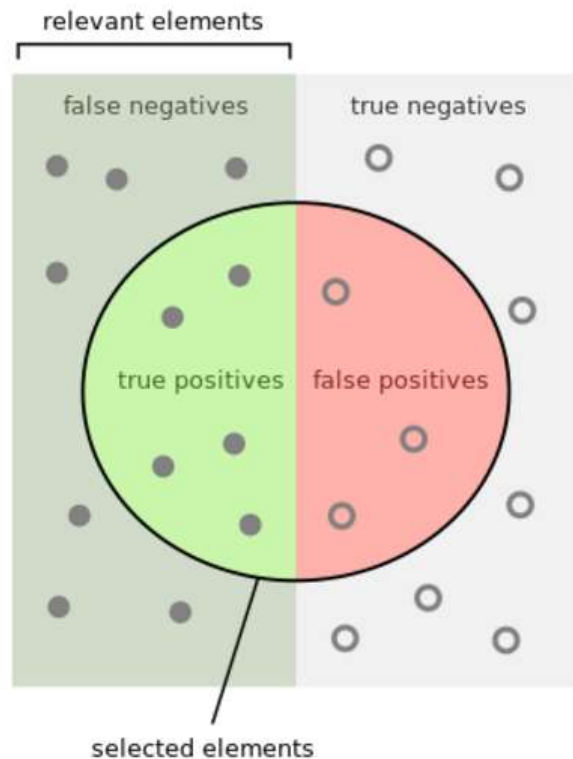
true positive detections /
total detections

Recall:

true positive detections /
total positive test instances

Terms

Recall, Precision, mAP, IoU

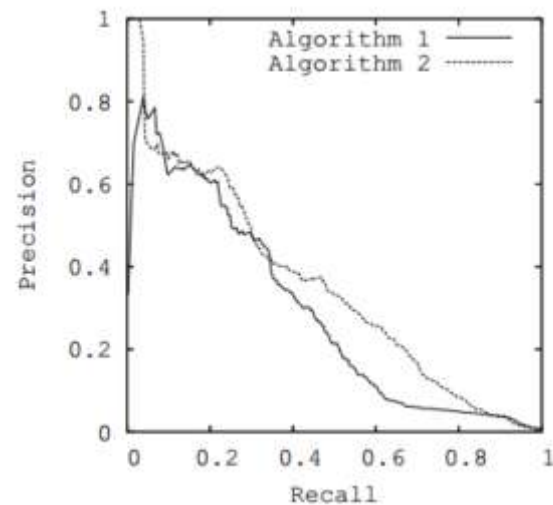


How many selected items are relevant?

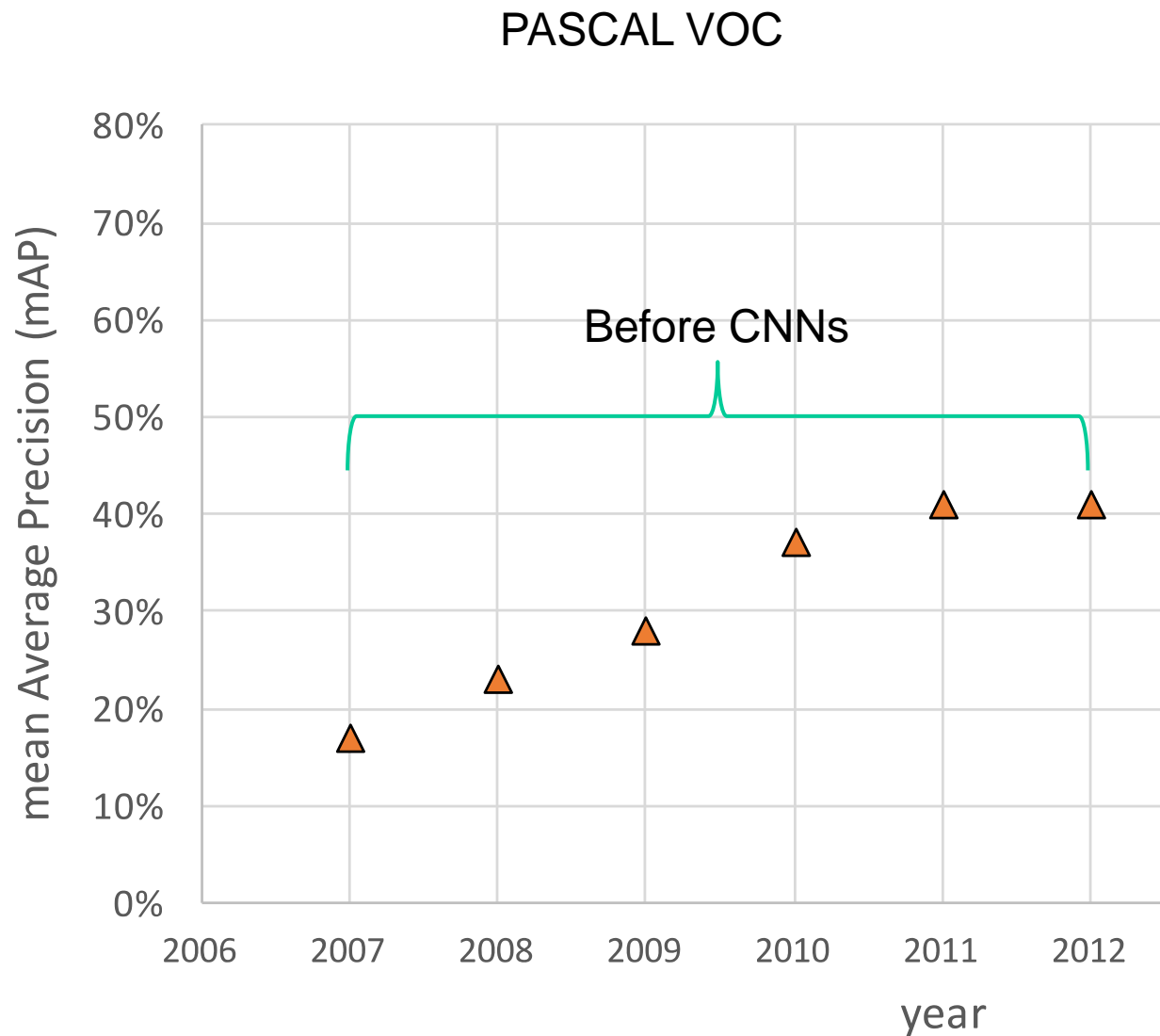
$$\text{Precision} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Positives}}$$

How many relevant items are selected?

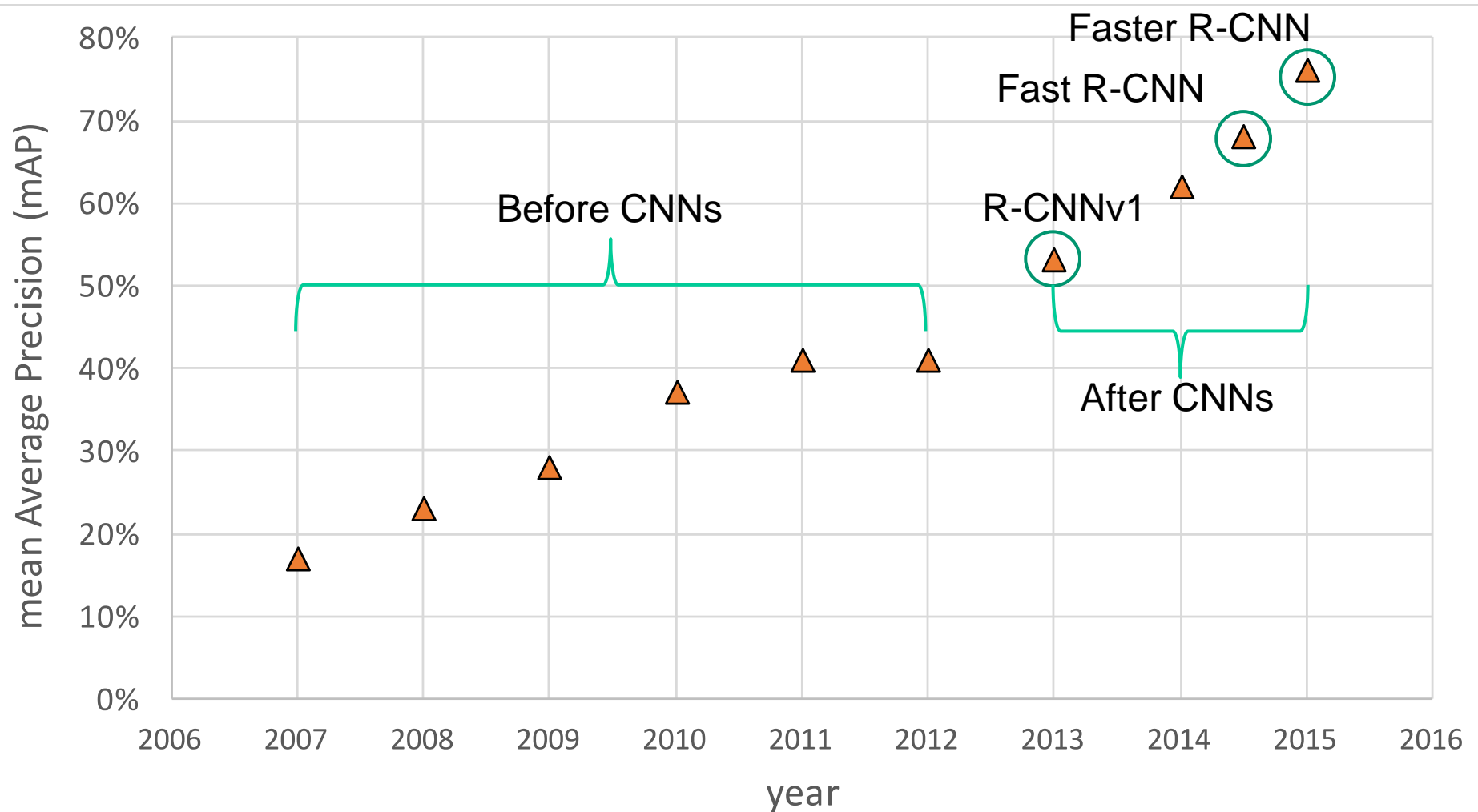
$$\text{Recall} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Negatives}}$$



Progress on PASCAL detection



Object detection progress



Conceptual approach: Sliding window detection



- Slide a window across the image and evaluate a detection model at each location
 - Thousands of windows to evaluate: efficiency and low false positive rates are essential
 - Difficult to extend to a large range of scales, aspect ratios

Histograms of oriented gradients (HOG)

- Partition image into blocks and compute histogram of gradient orientations in each block

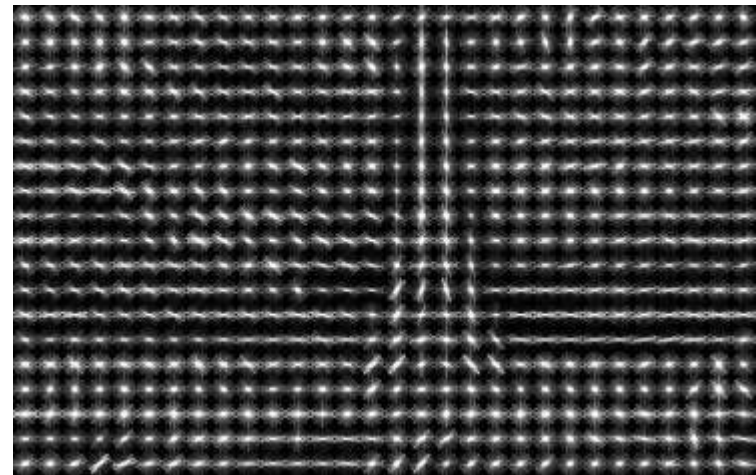
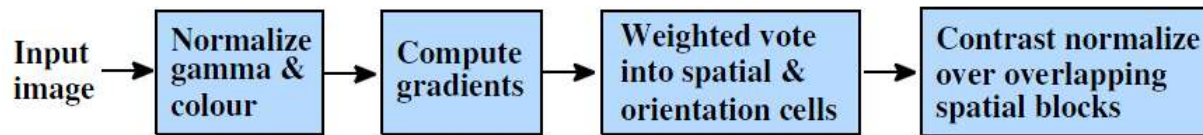


Image credit: N. Snavely

Pedestrian detection with HOG

- Train a pedestrian template using a linear support vector machine

positive training examples



negative training examples

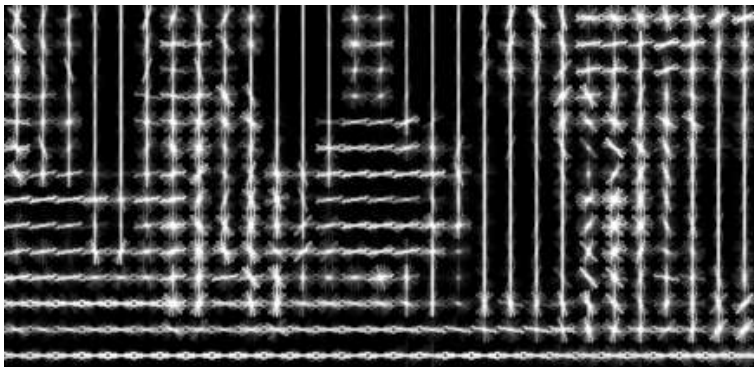


N. Dalal and B. Triggs, [Histograms of Oriented Gradients for Human Detection](#),
CVPR 2005

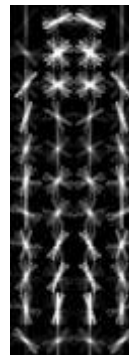
Pedestrian detection with HOG

- Train a pedestrian template using a linear support vector machine
- At test time, convolve feature map with template
- Find local maxima of response
- For multi-scale detection, repeat over multiple levels of a HOG *pyramid*

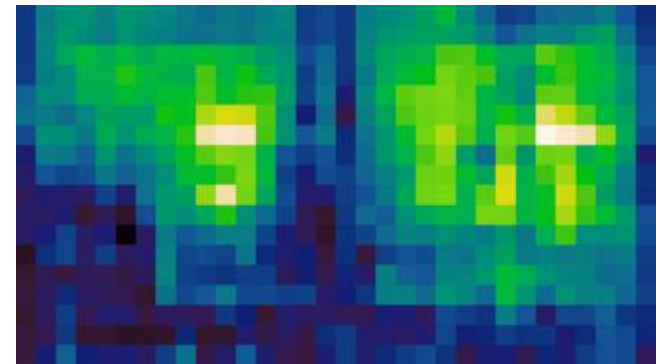
HOG feature map



Template



Detector response map



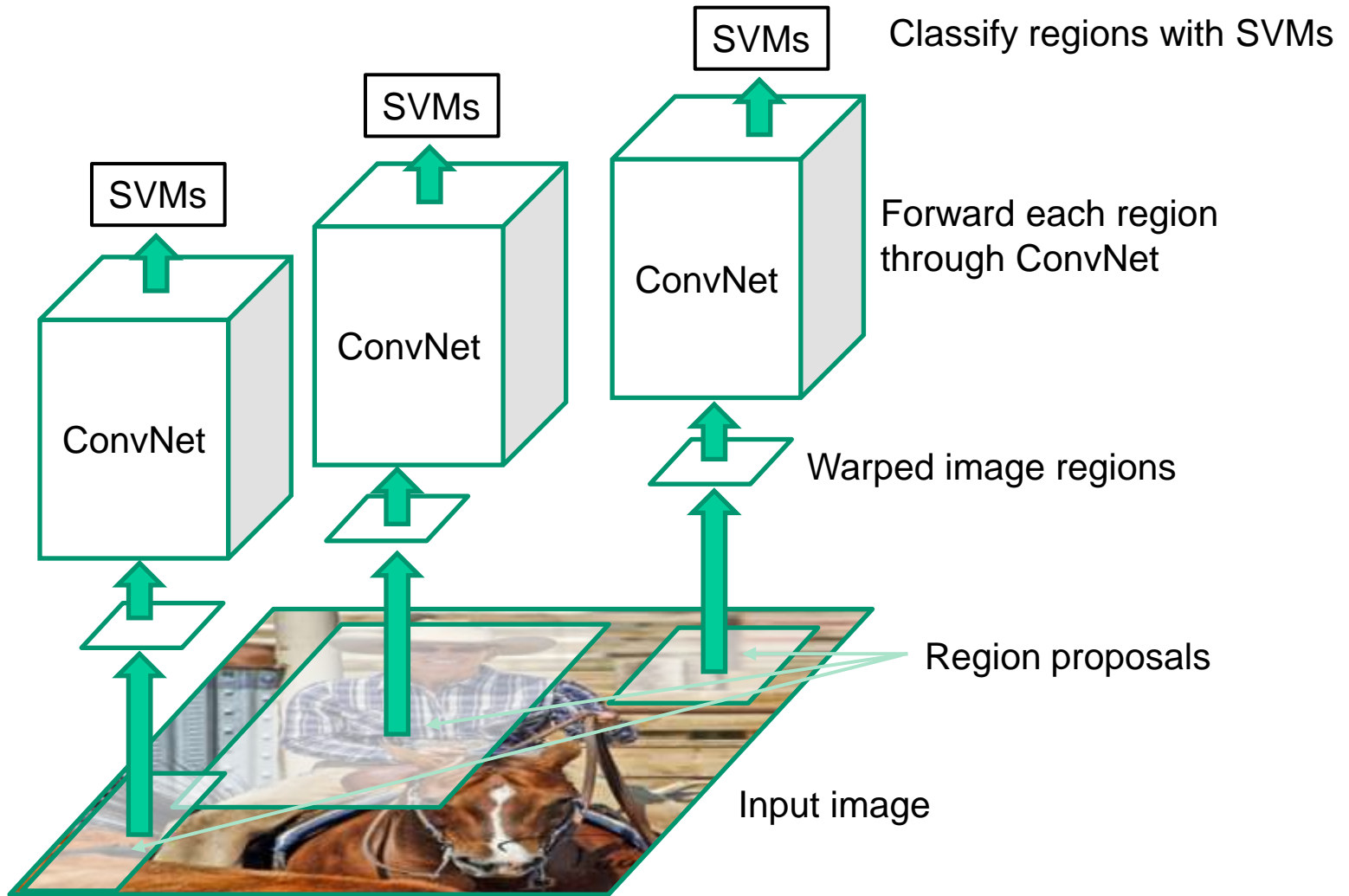
Example detections



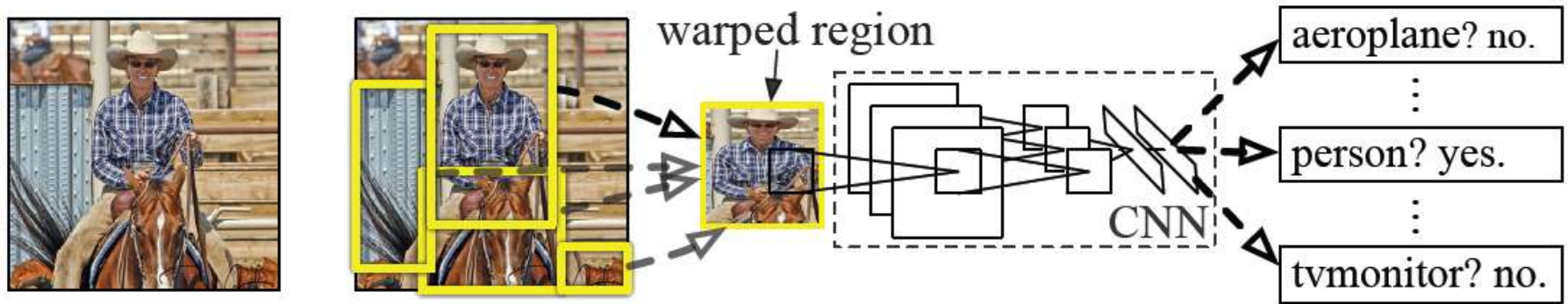
[Dalal and Triggs, CVPR 2005]

R-CNN: Region proposals + CNN features

Source: R. Girshick



R-CNN details



- **Regions:** ~2000 Selective Search proposals
- **Network:** AlexNet *pre-trained* on ImageNet (1000 classes), *fine-tuned* on PASCAL (21 classes)
- **Final detector:** warp proposal regions, extract fc7 network activations (4096 dimensions), classify with linear SVM
- **Bounding box regression** to refine box locations
- **Performance:** mAP of **53.7%** on PASCAL 2010 (vs. **35.1%** for Selective Search and **33.4%** for Deformable Part Models)

R-CNN pros and cons

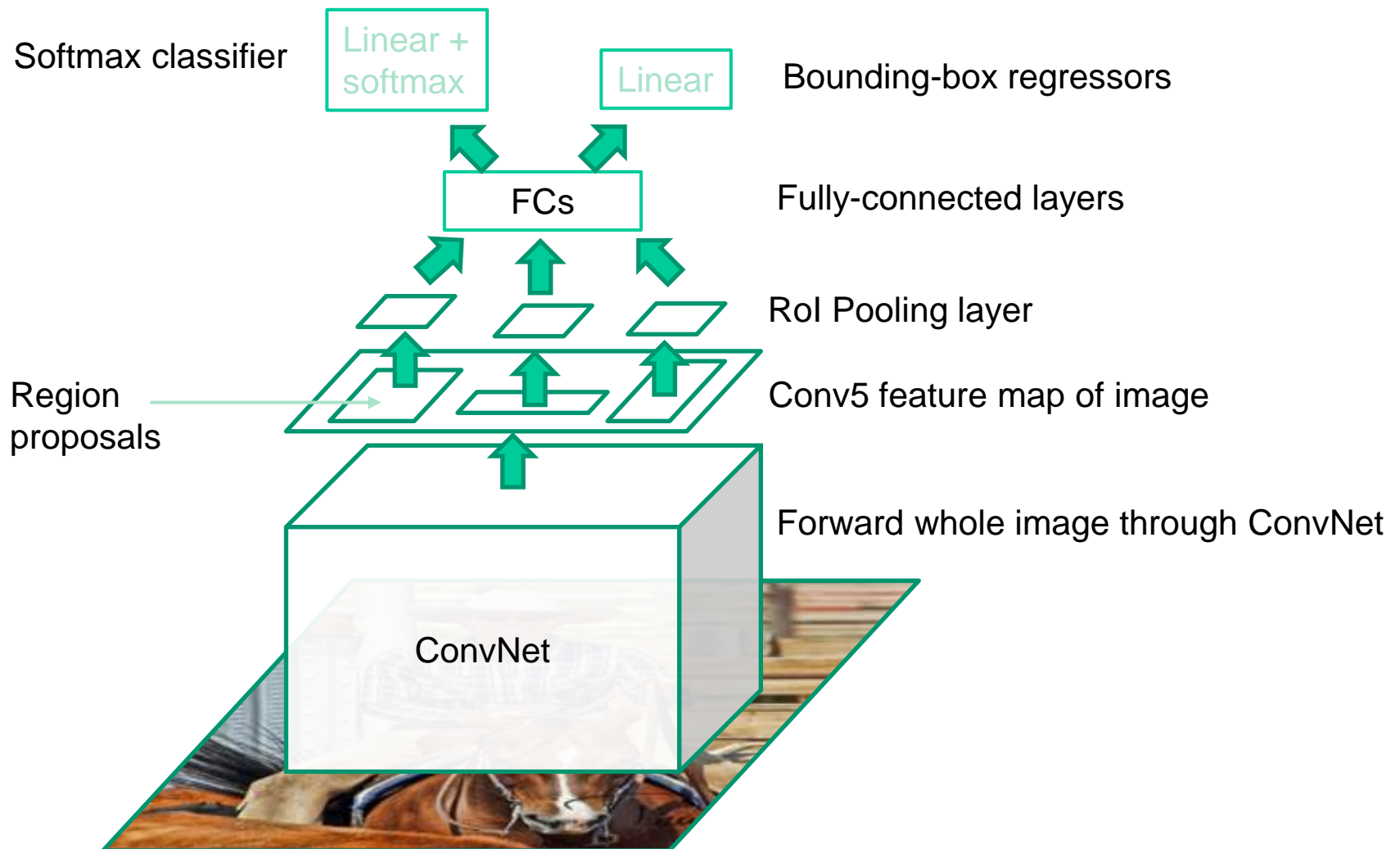
- **Pros**

- Accurate!
- Any deep architecture can immediately be “plugged in”

- **Cons**

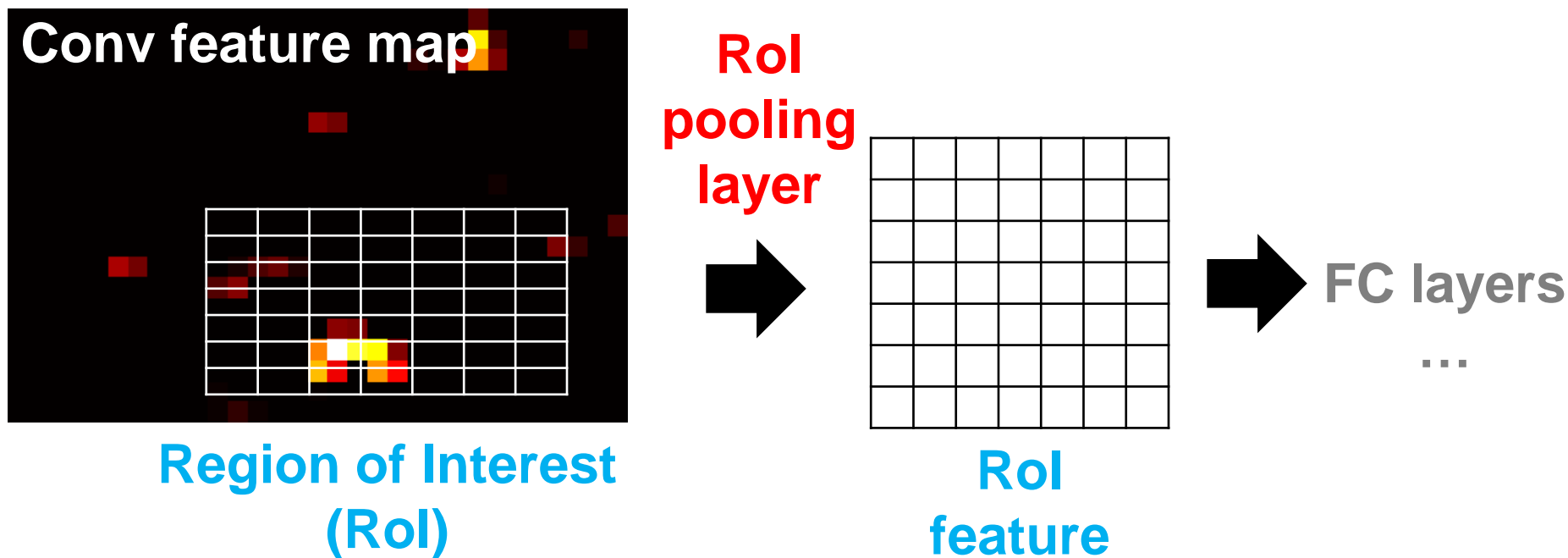
- Not a single end-to-end system
 - Fine-tune network with softmax classifier (log loss)
 - Train post-hoc linear SVMs (hinge loss)
 - Train post-hoc bounding-box regressions (least squares)
- Training is slow (84h), takes a lot of disk space
 - 2000 CNN passes per image
- Inference (detection) is slow (47s / image with VGG16)

Fast R-CNN



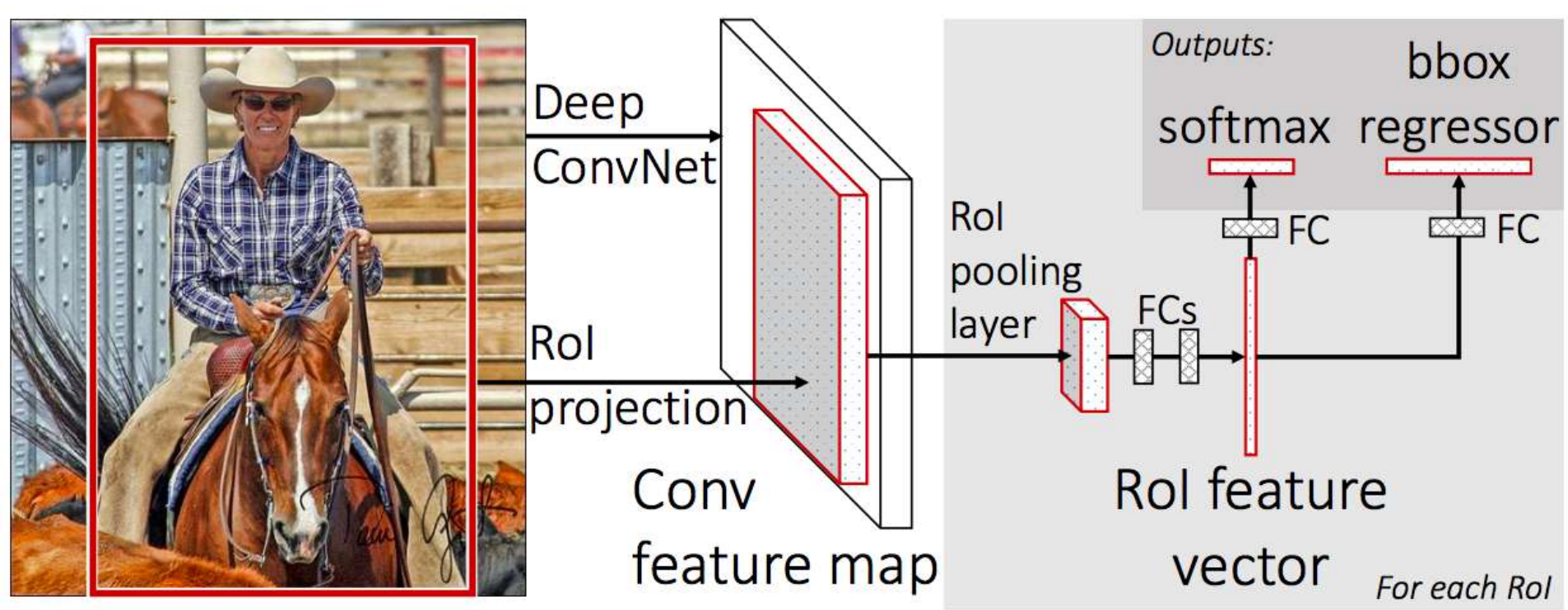
Rol pooling

- “Crop and resample” a fixed-size feature representing a region of interest out of the outputs of the last conv layer
 - Use nearest-neighbor interpolation of coordinates, max pooling

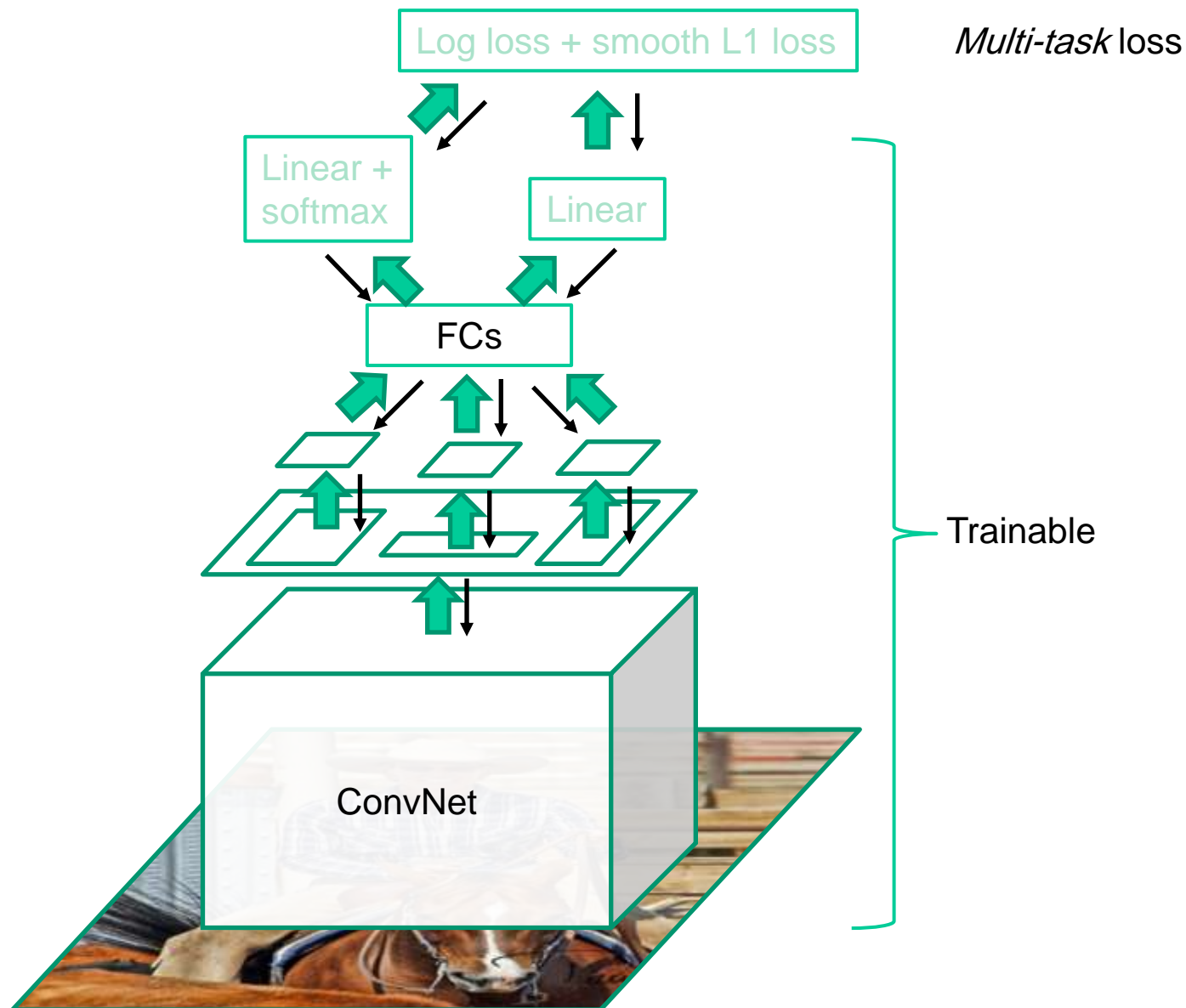


Prediction

- For each RoI, network predicts probabilities for $C+1$ classes (class 0 is background) and four bounding box offsets for C classes



Fast R-CNN training



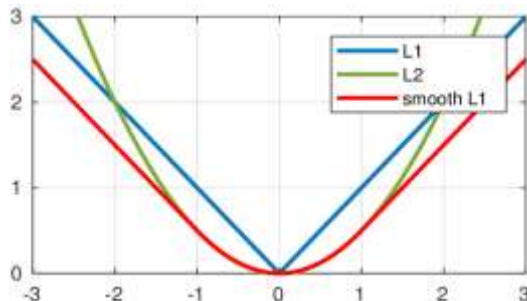
Multi-task loss

- Loss for ground truth class y , predicted class probabilities $P(y)$, ground truth box b , and predicted box \hat{b} :

$$L(y, P, b, \hat{b}) = \underbrace{-\log P(y)}_{\text{softmax loss}} + \lambda \mathbb{I}[y \geq 1] \underbrace{L_{\text{reg}}(b, \hat{b})}_{\text{regression loss}}$$

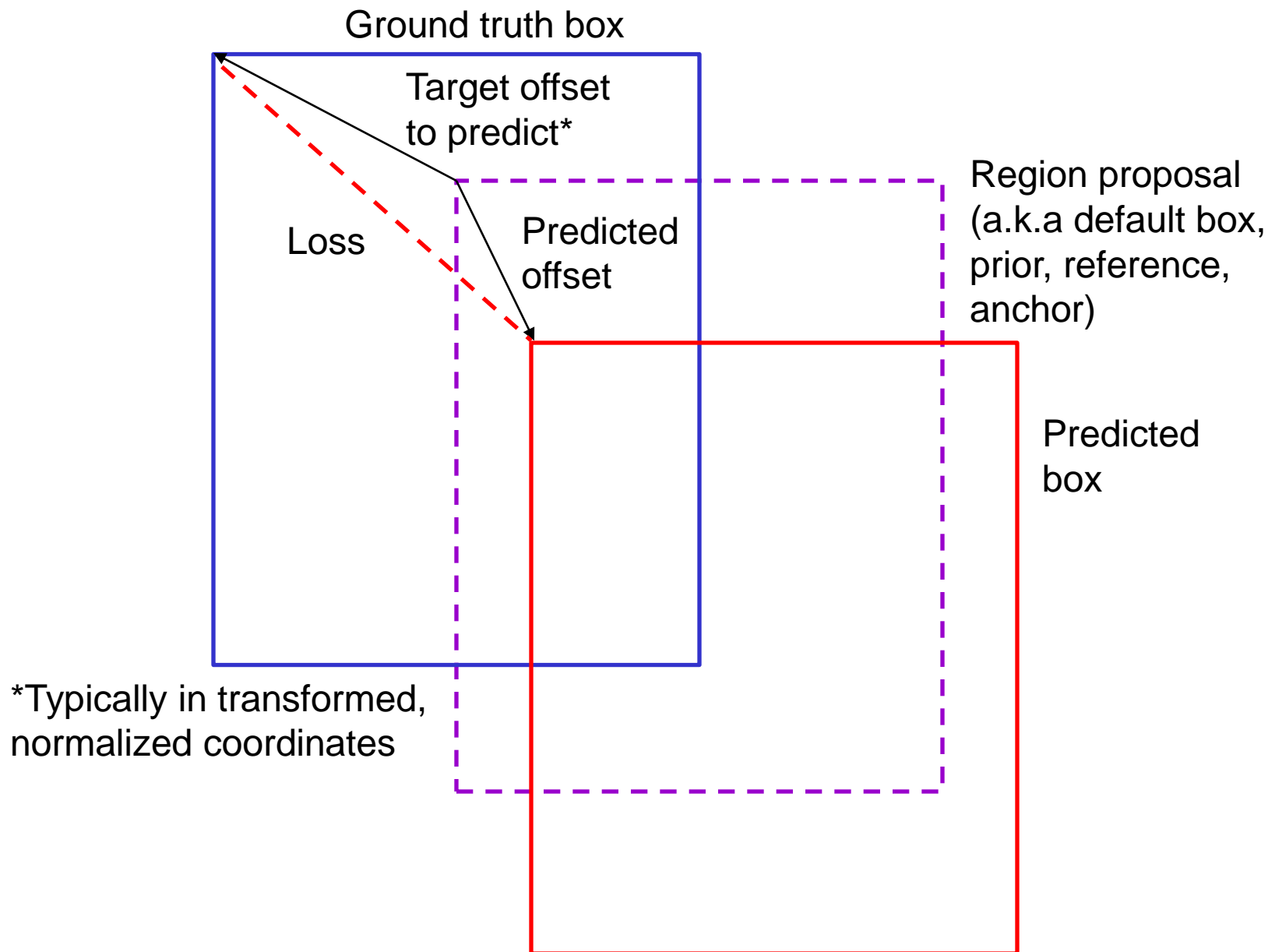
- Regression loss: *smooth L1 loss* on top of log space offsets relative to proposal

$$L_{\text{reg}}(b, \hat{b}) = \sum_{i=\{x,y,w,h\}} \text{smooth}_{L_1}(b_i - \hat{b}_i)$$



$$\text{smooth}_{L_1}(x) = \begin{cases} 0.5x^2 & \text{if } |x| < 1 \\ |x| - 0.5 & \text{otherwise} \end{cases}$$

Bounding box regression



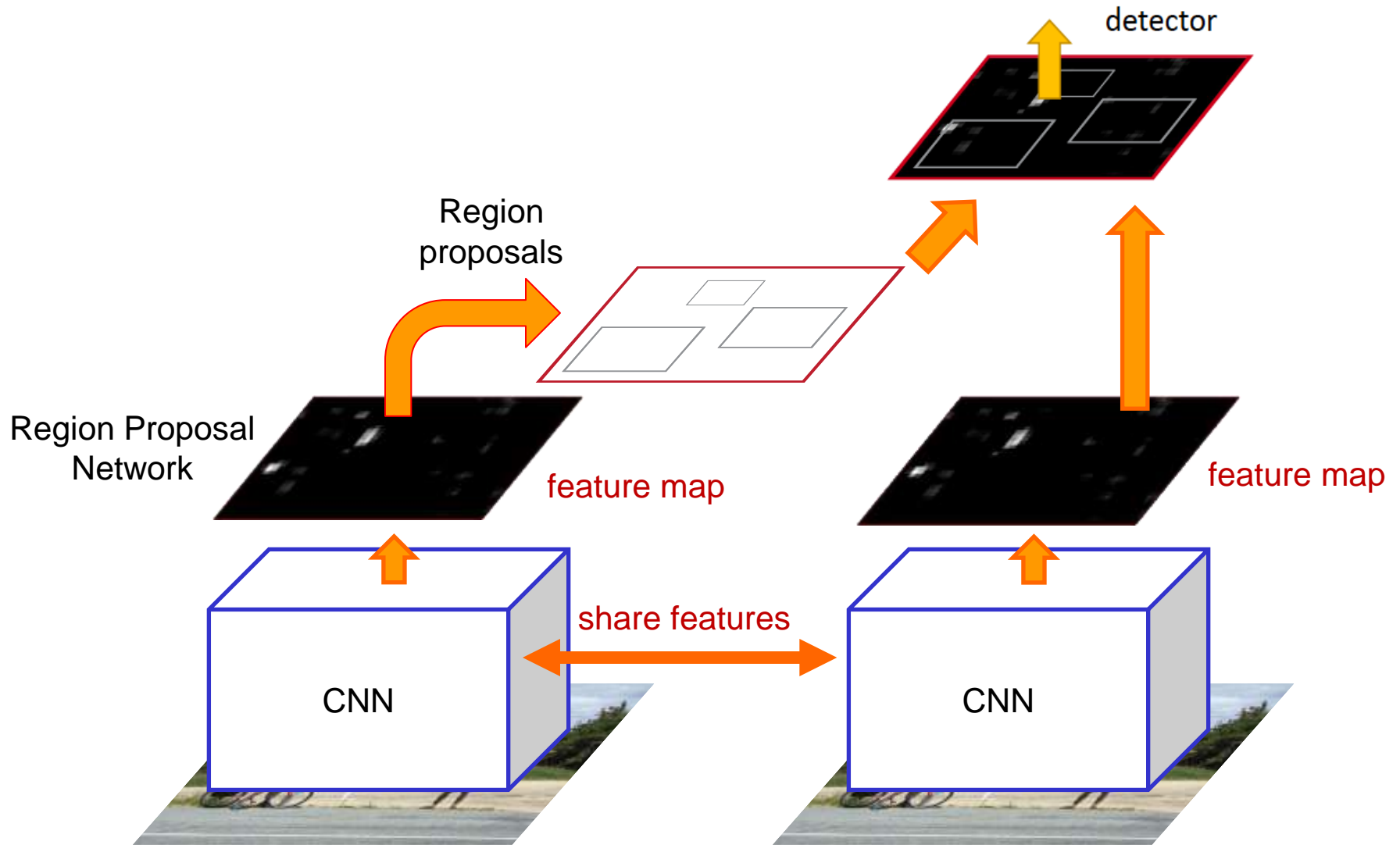
Fast R-CNN results

	Fast R-CNN	R-CNN
Train time (h)	9.5	84
- Speedup	8.8x	1x
Test time / image	0.32s	47.0s
Test speedup	146x	1x
mAP	66.9%	66.0%

(vs. 53.7% for AlexNet)

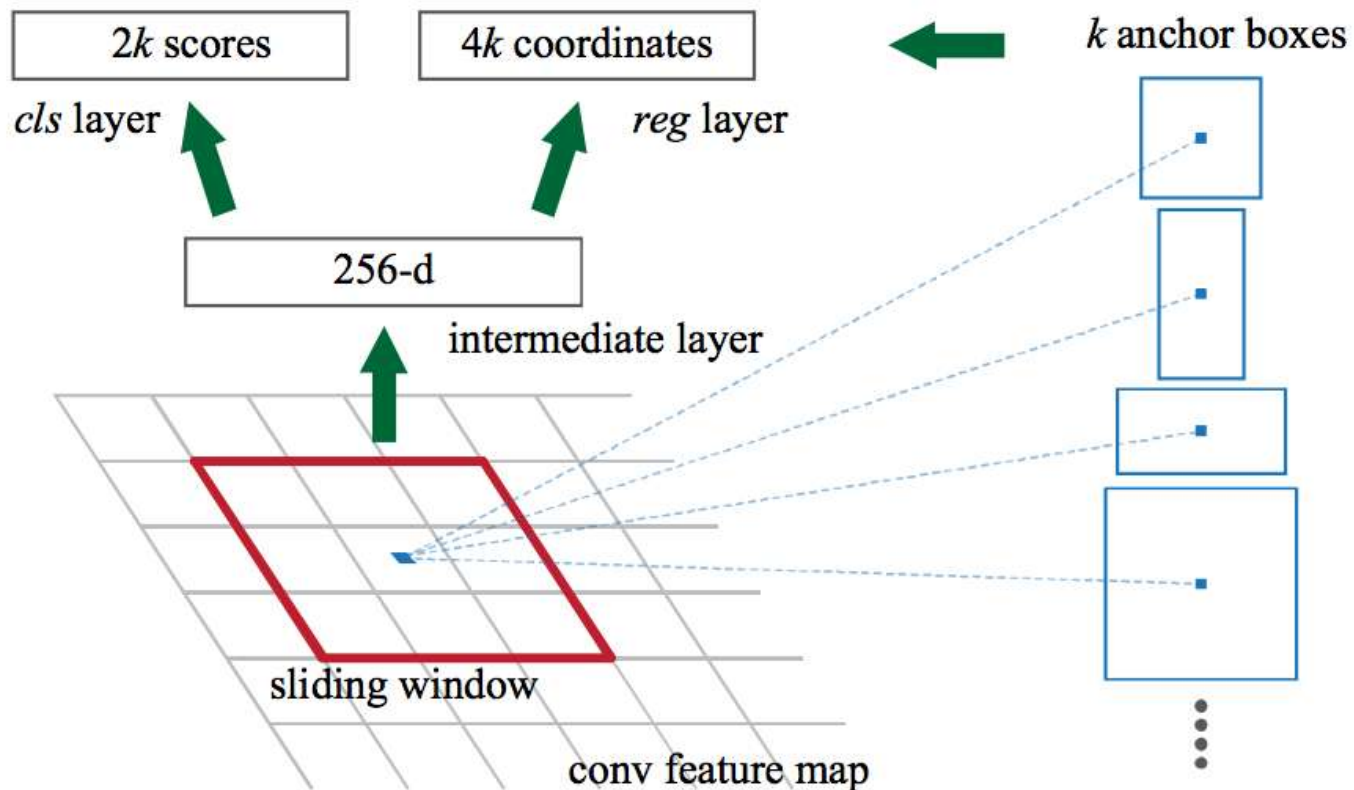
Timings exclude object proposal time, which is equal for all methods.
All methods use VGG16 from Simonyan and Zisserman.

Faster R-CNN

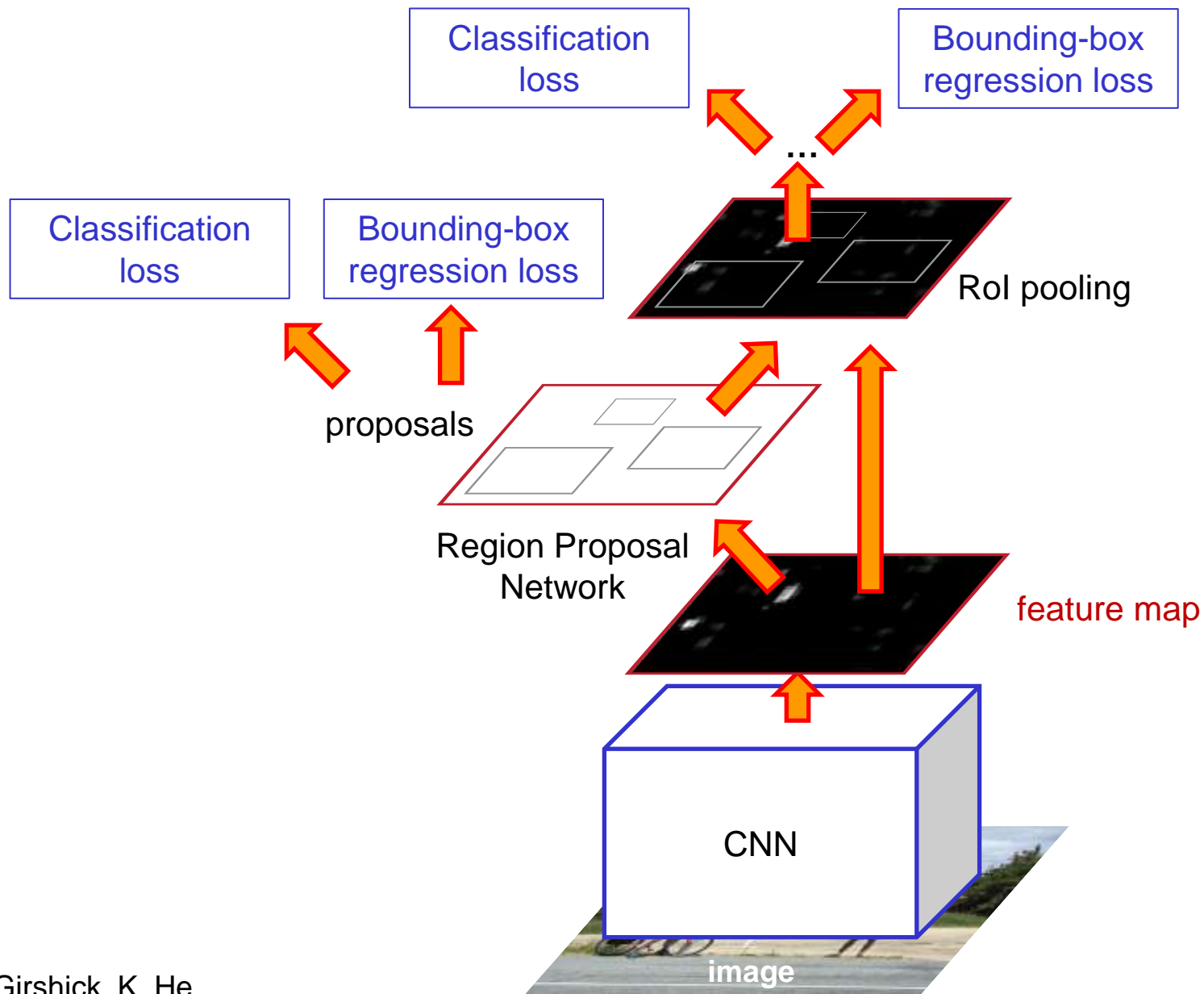


Region proposal network (RPN)

- Slide a small window (3x3) over the conv5 layer
 - Predict object/no object
 - Regress bounding box coordinates with reference to *anchors* (3 scales x 3 aspect ratios)



One network, four losses



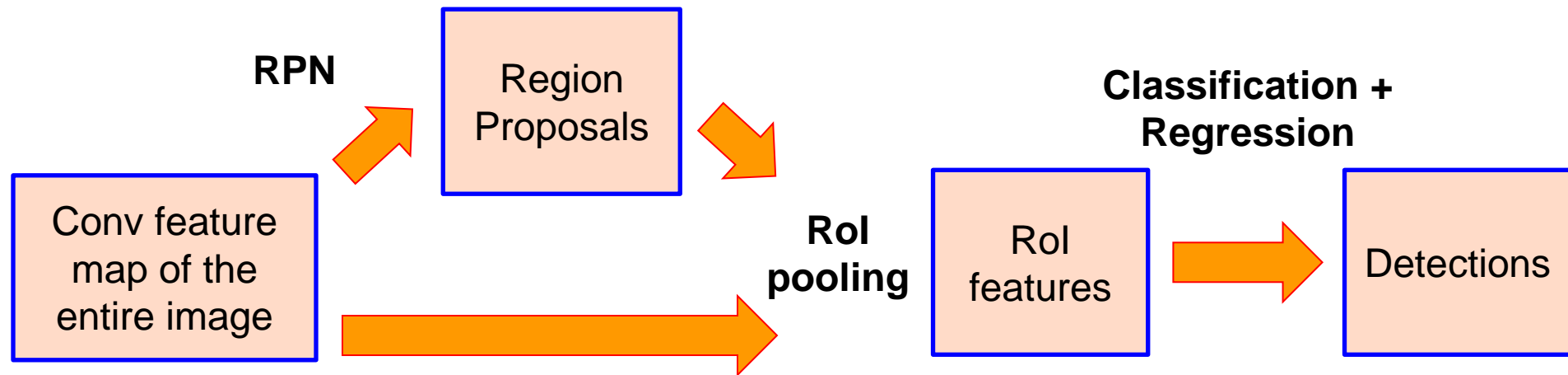
Faster R-CNN results

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

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

Streamlined detection architectures

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



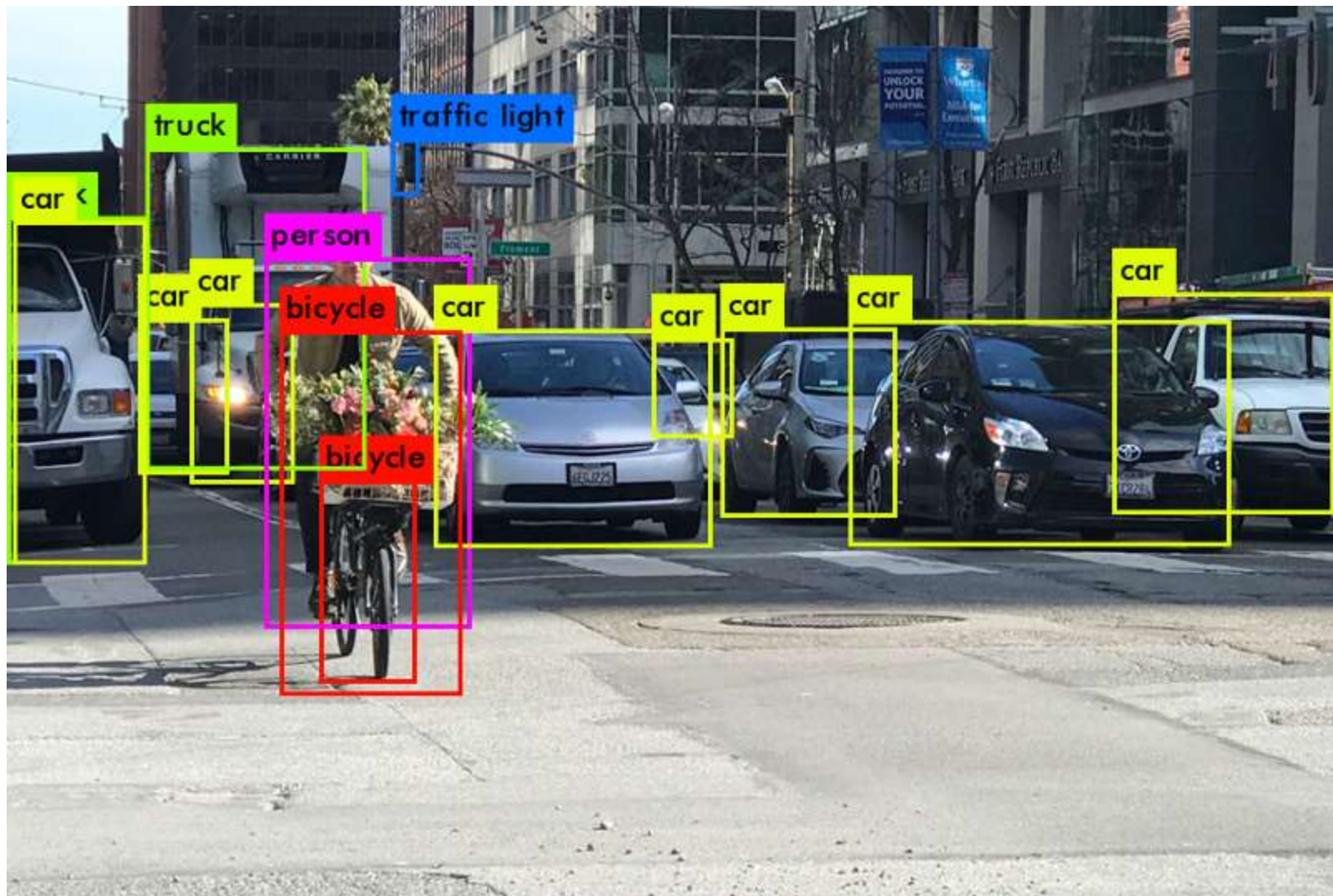
- Is it possible do detection in one shot?



You Only Look Once



Objects detected by YOLO



Grid cell



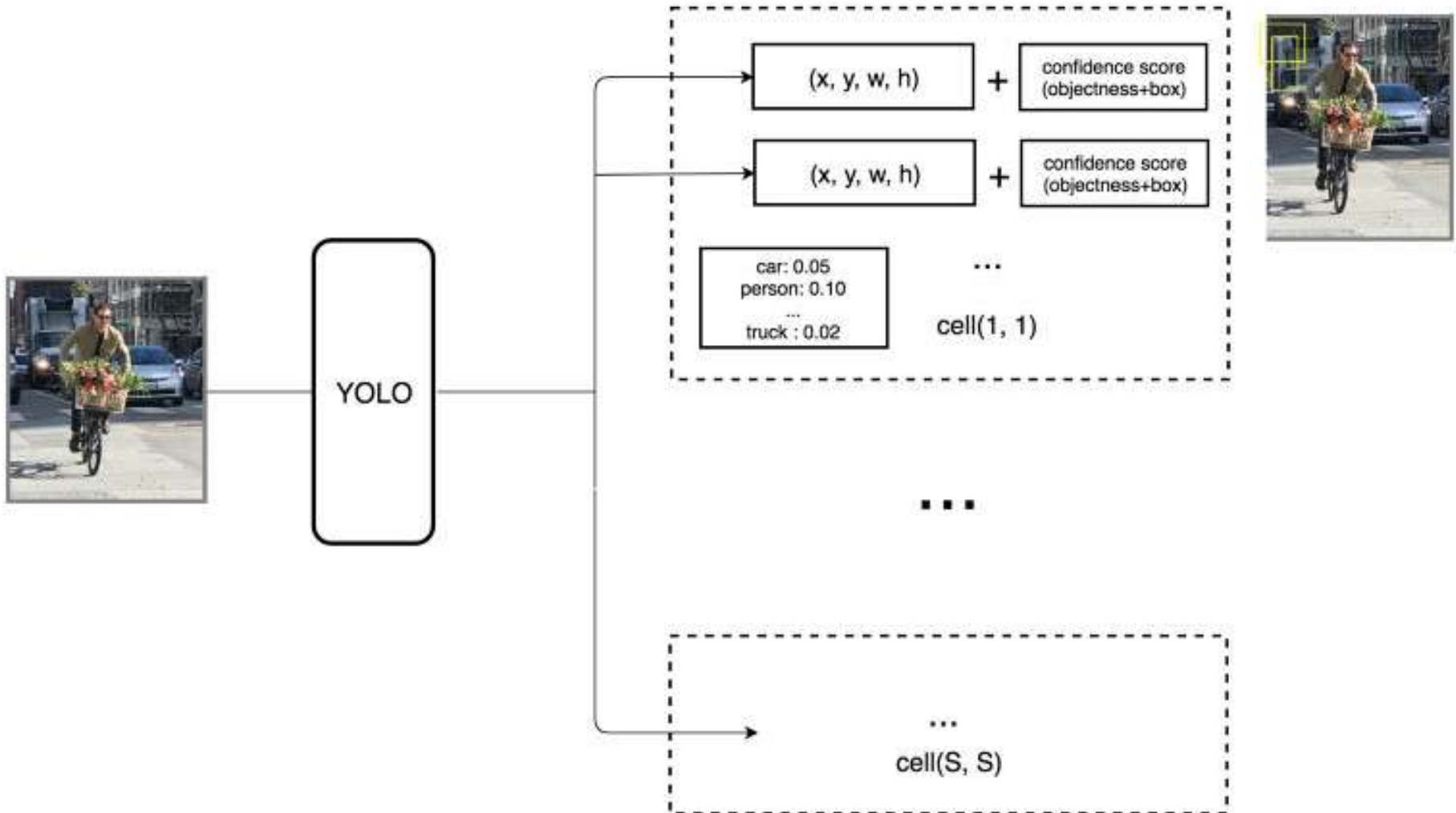
YOLO divides the input image into an $S \times S$ grid. Each grid cell predicts only one object.

The yellow grid cell tries to predict the “person” object whose center (the blue dot) falls inside the grid cell.

For each grid cell,

- it predicts **B** boundary boxes and each box has one box **confidence score**,
- it detects one object only regardless of the number of boxes **B**,
- it predicts **C conditional class probabilities** (one per class for the likeliness of the object class).

YOYO make **S x S** predictions with **B** boundary boxes.



Boundary box

- Each boundary box contains 5 elements:
 (x, y, w, h) and a **box confidence score**.
- The confidence score reflects how likely the box contains an object (**objectness**) and how accurate is the boundary box.
- We normalize the bounding box width w and height h by the image width and height. x and y are offsets to the corresponding cell.

$$\begin{aligned}\text{box confidence score} &\equiv P_r(\text{object}) \cdot \text{IoU} \\ \text{conditional class probability} &\equiv P_r(\text{class}_i | \text{object}) \\ \text{class confidence score} &\equiv P_r(\text{class}_i) \cdot \text{IoU} \\ &= \text{box confidence score} \times \text{conditional class probability}\end{aligned}$$

where

$P_r(\text{object})$ is the probability the box contains an object.

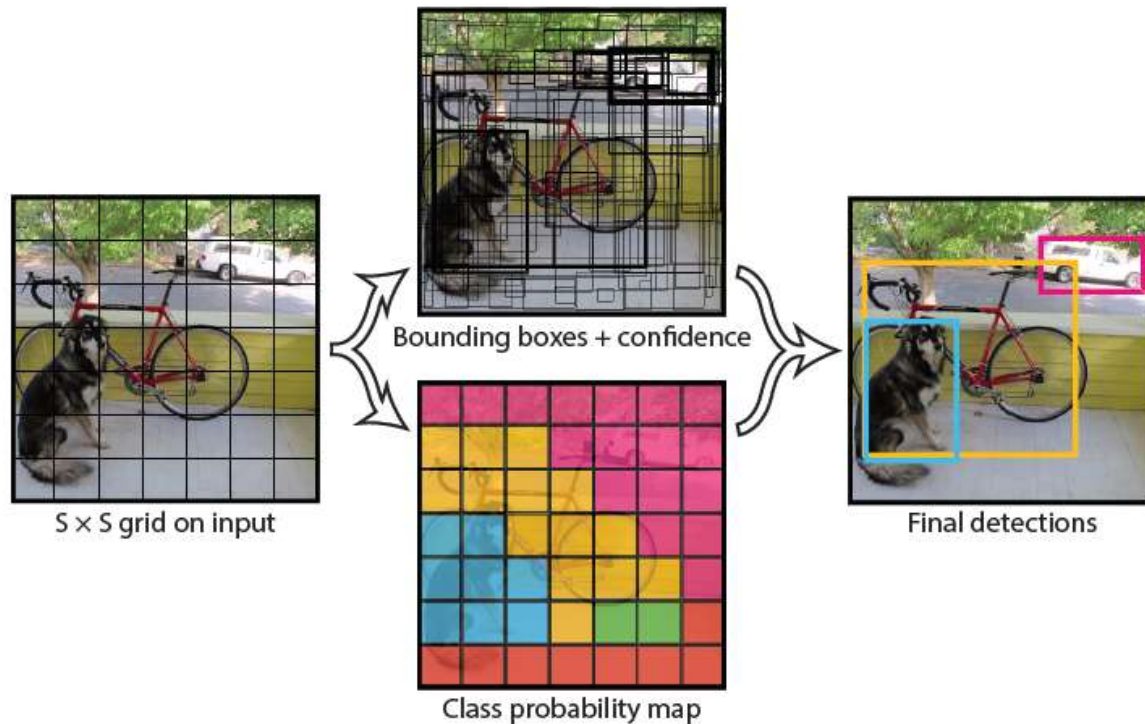
IoU is the IoU (intersection over union) between the predicted box and the ground truth.

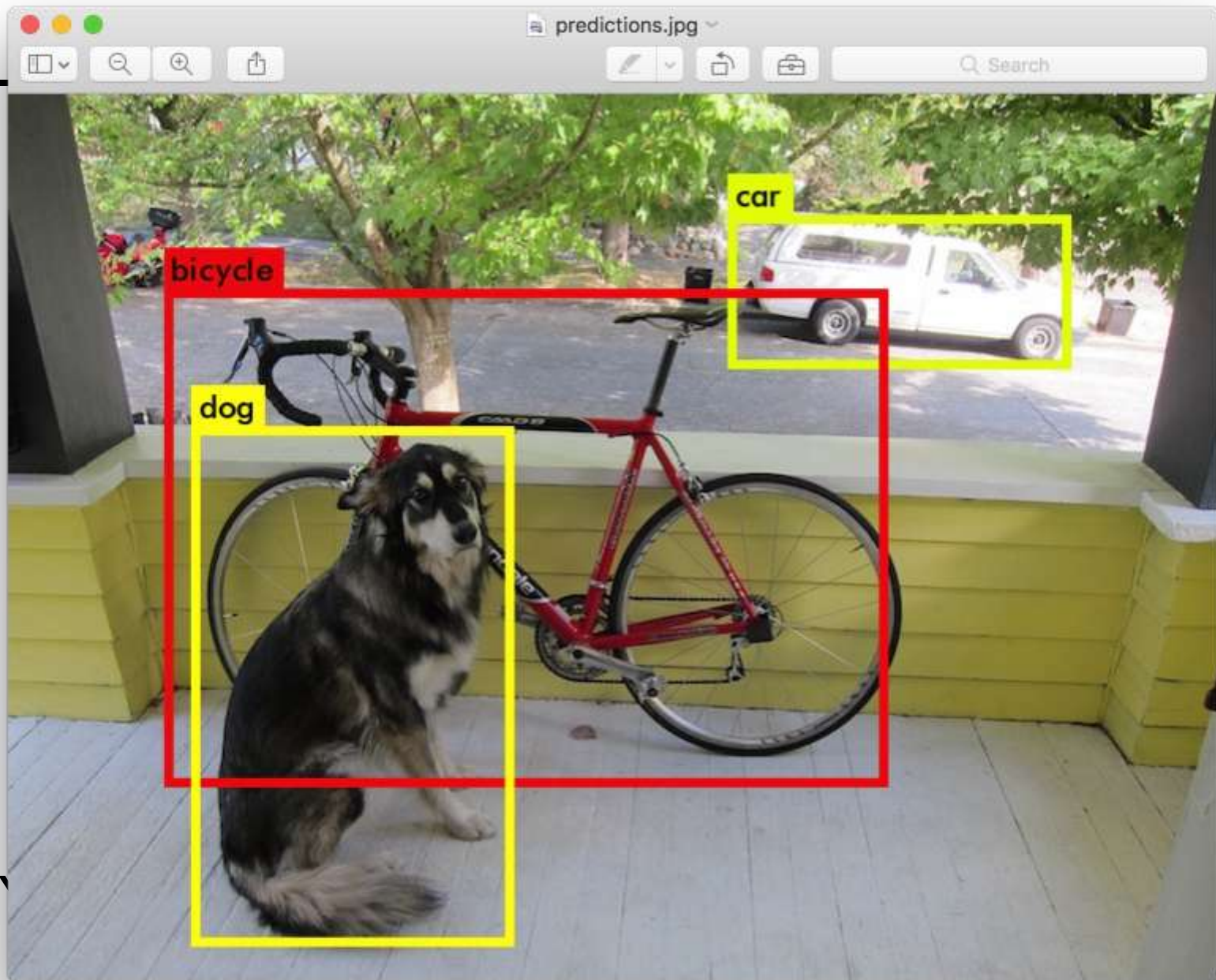
$P_r(\text{class}_i | \text{object})$ is the probability the object belongs to class_i given an object is presence.

$P_r(\text{class}_i)$ is the probability the object belongs to class_i

YOLO

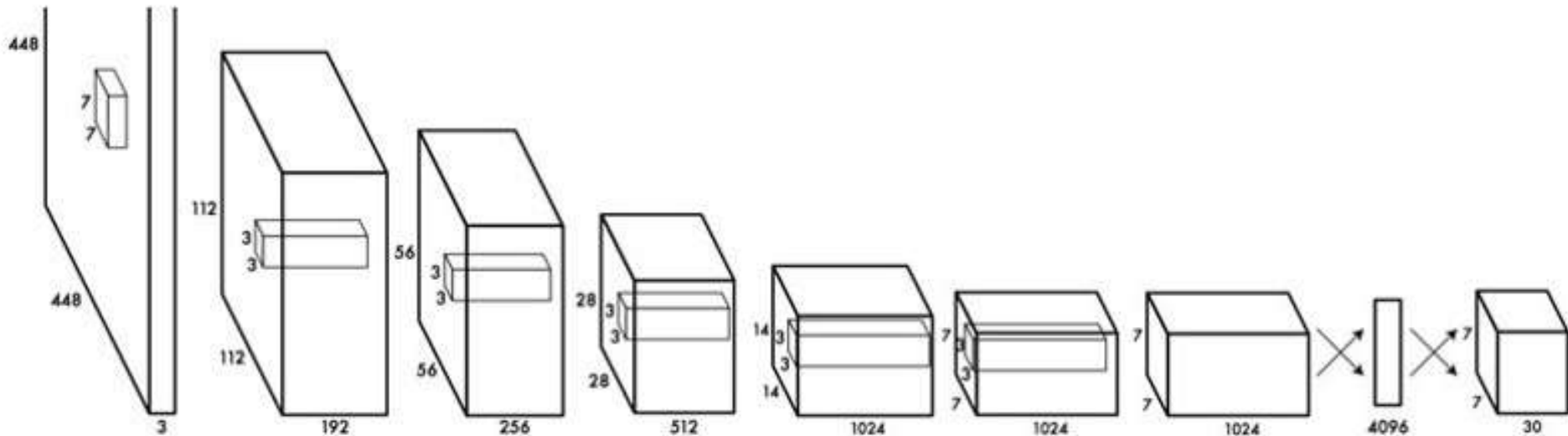
- Divide the image into a coarse grid and directly predict class label and a few candidate boxes for each grid cell





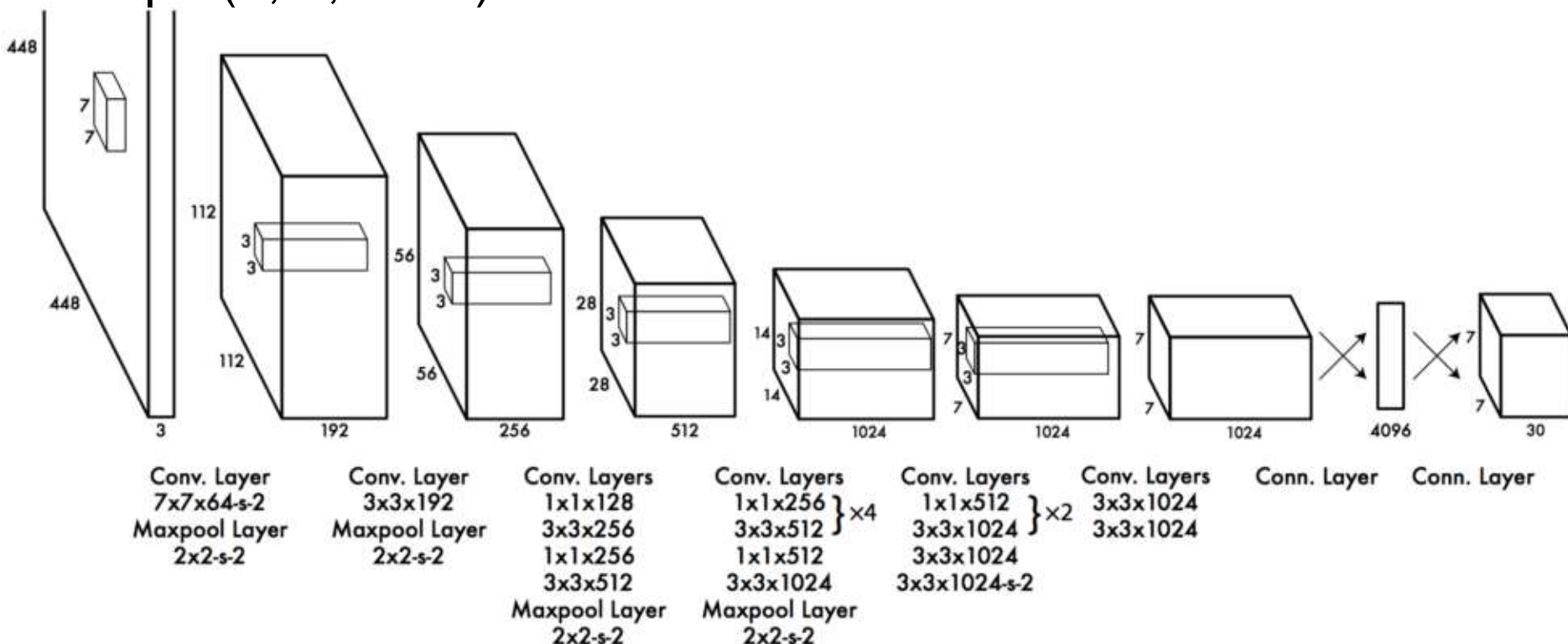
YOLO

1. Take conv feature maps at 7x7 resolution
2. Add two FC layers to predict, at each location, a score for each class and 2 bboxes w/ confidences
 - For PASCAL, output is 7x7x30 ($30 = 20 + 2*(4+1)$)



Network design

YOLO has 24 convolutional layers followed by 2 fully connected layers (FC). The last convolution layer outputs a tensor with shape (7, 7, 1024)



Loss function

YOLO predicts multiple bounding boxes per grid cell. To compute the loss for the true positive, we only want one of them to be **responsible** for the object. For this purpose, we select the one with the highest IoU (intersection over union) with the ground truth.

YOLO uses sum-squared error between the predictions and the ground truth to calculate loss. The loss function composes of:

- the **classification loss**.
- the **localization loss** (errors between the predicted boundary box and the ground truth).
- the **confidence loss** (the objectness of the box).

Localization loss

The localization loss measures the errors in the predicted boundary box **locations** and **sizes**.

We do not want to weight absolute errors in large boxes and small boxes equally. i.e. a 2-pixel error in a large box is the same for a small box. To partially address this, YOLO predicts the square root of the bounding box width and height instead of the width and height.

$$\lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} \left[(x_i - \hat{x}_i)^2 + (y_i - \hat{y}_i)^2 \right] \\ + \lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} \left[\left(\sqrt{w_i} - \sqrt{\hat{w}_i} \right)^2 + \left(\sqrt{h_i} - \sqrt{\hat{h}_i} \right)^2 \right]$$

where

$\mathbb{1}_{ij}^{\text{obj}} = 1$ if the j th boundary box in cell i is responsible for detecting the object, otherwise 0.

λ_{coord} increase the weight for the loss in the boundary box coordinates.

Confidence loss

$$\sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{obj} (C_i - \hat{C}_i)^2$$

If an object is detected in the box, the confidence loss (measuring the objectness of the box)

where

\hat{C}_i is the box confidence score of the box j in cell i .

$\mathbb{1}_{ij}^{obj} = 1$ if the j th boundary box in cell i is responsible for detecting the object, otherwise 0.

If an object is not detected in the box, the confidence loss is:

$$\lambda_{noobj} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{noobj} (C_i - \hat{C}_i)^2$$

where

$\mathbb{1}_{ij}^{noobj}$ is the complement of $\mathbb{1}_{ij}^{obj}$.

\hat{C}_i is the box confidence score of the box j in cell i .

λ_{noobj} weights down the loss when detecting background.

YOLO

- Objective function:

$$\begin{aligned} & \lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} \left[(x_i - \hat{x}_i)^2 + (y_i - \hat{y}_i)^2 \right] \\ & + \lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} \left[\left(\sqrt{w_i} - \sqrt{\hat{w}_i} \right)^2 + \left(\sqrt{h_i} - \sqrt{\hat{h}_i} \right)^2 \right] \\ & + \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} (C_i - \hat{C}_i)^2 \\ & + \lambda_{\text{noobj}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{noobj}} (C_i - \hat{C}_i)^2 \\ & + \sum_{i=0}^{S^2} \mathbb{1}_i^{\text{obj}} \sum_{c \in \text{classes}} (p_i(c) - \hat{p}_i(c))^2 \end{aligned}$$

Regression

Object/no object confidence

Class prediction

YOLO

- Objective function:

$$\begin{aligned}
 & \lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} \left[(x_i - \hat{x}_i)^2 + (y_i - \hat{y}_i)^2 \right] \\
 & + \lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} \left[\left(\sqrt{w_i} - \sqrt{\hat{w}_i} \right)^2 + \left(\sqrt{h_i} - \sqrt{\hat{h}_i} \right)^2 \right] \\
 & + \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} \left(C_i - \hat{C}_i \right)^2 \\
 & + \lambda_{\text{noobj}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{noobj}} \left(C_i - \hat{C}_i \right)^2 \\
 & + \sum_{i=0}^{S^2} \mathbb{1}_i^{\text{obj}} \sum_{c \in \text{classes}} \left(p_i(c) - \hat{p}_i(c) \right)^2
 \end{aligned}$$

Cell i contains object, predictor j is responsible for it

Small deviations matter less for larger boxes than for smaller boxes

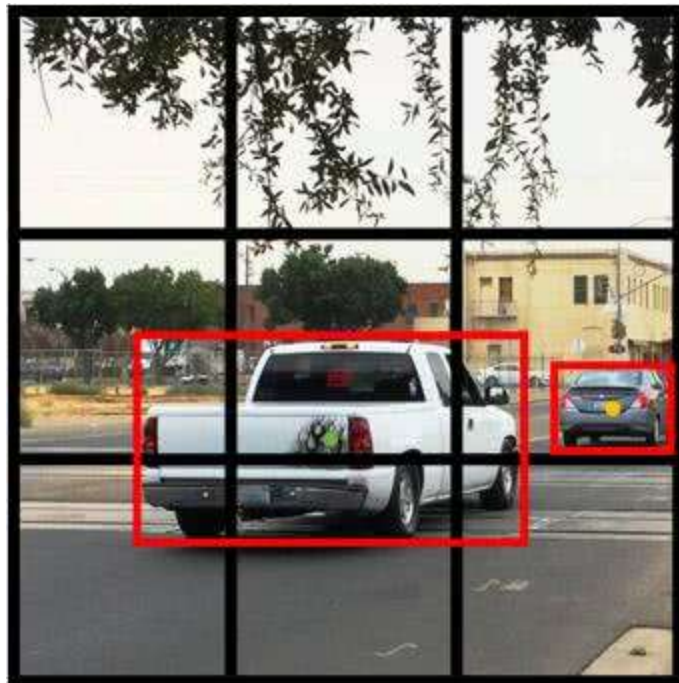
Confidence for object

Confidence for no object

Down-weight loss from boxes that don't contain objects ($\lambda_{\text{noobj}} = 0.5$)

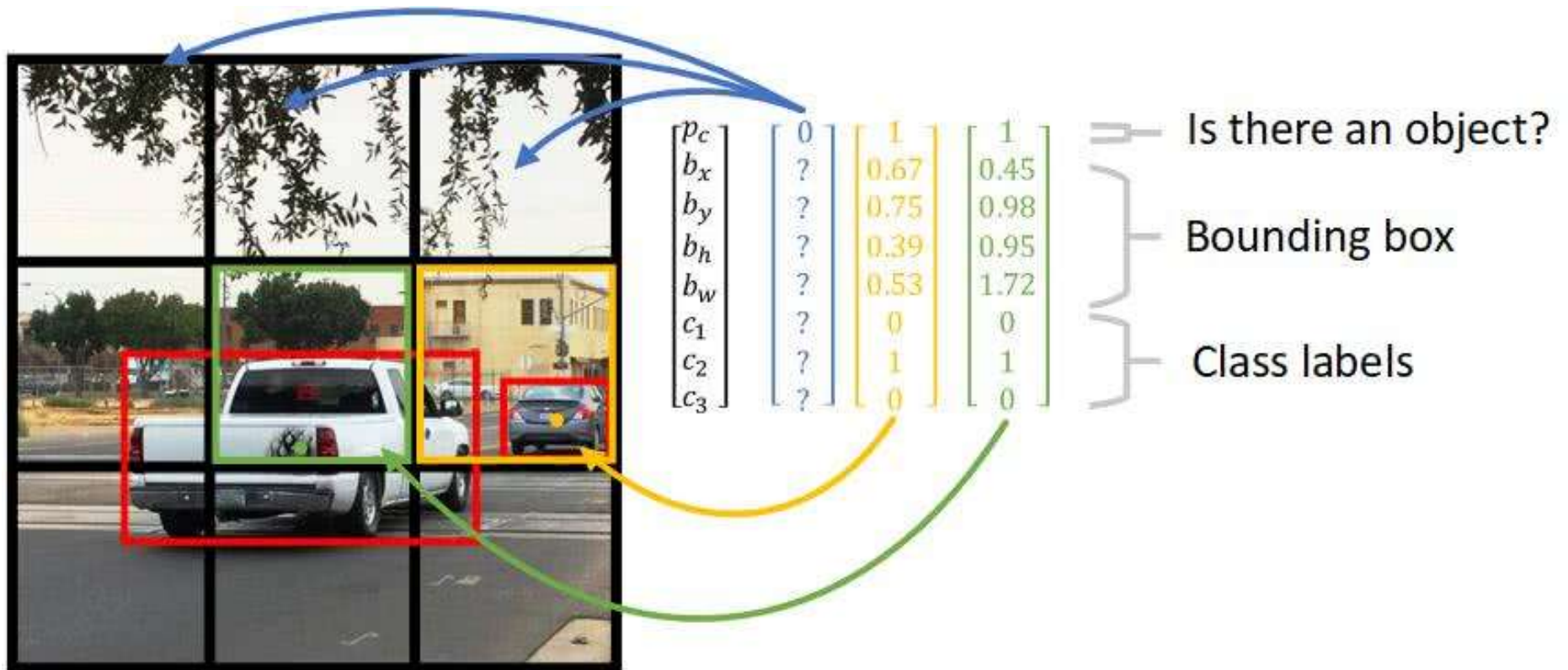
Class probability

YOLO Training



$$y = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \\ c_4 \end{bmatrix}$$

Training



Training

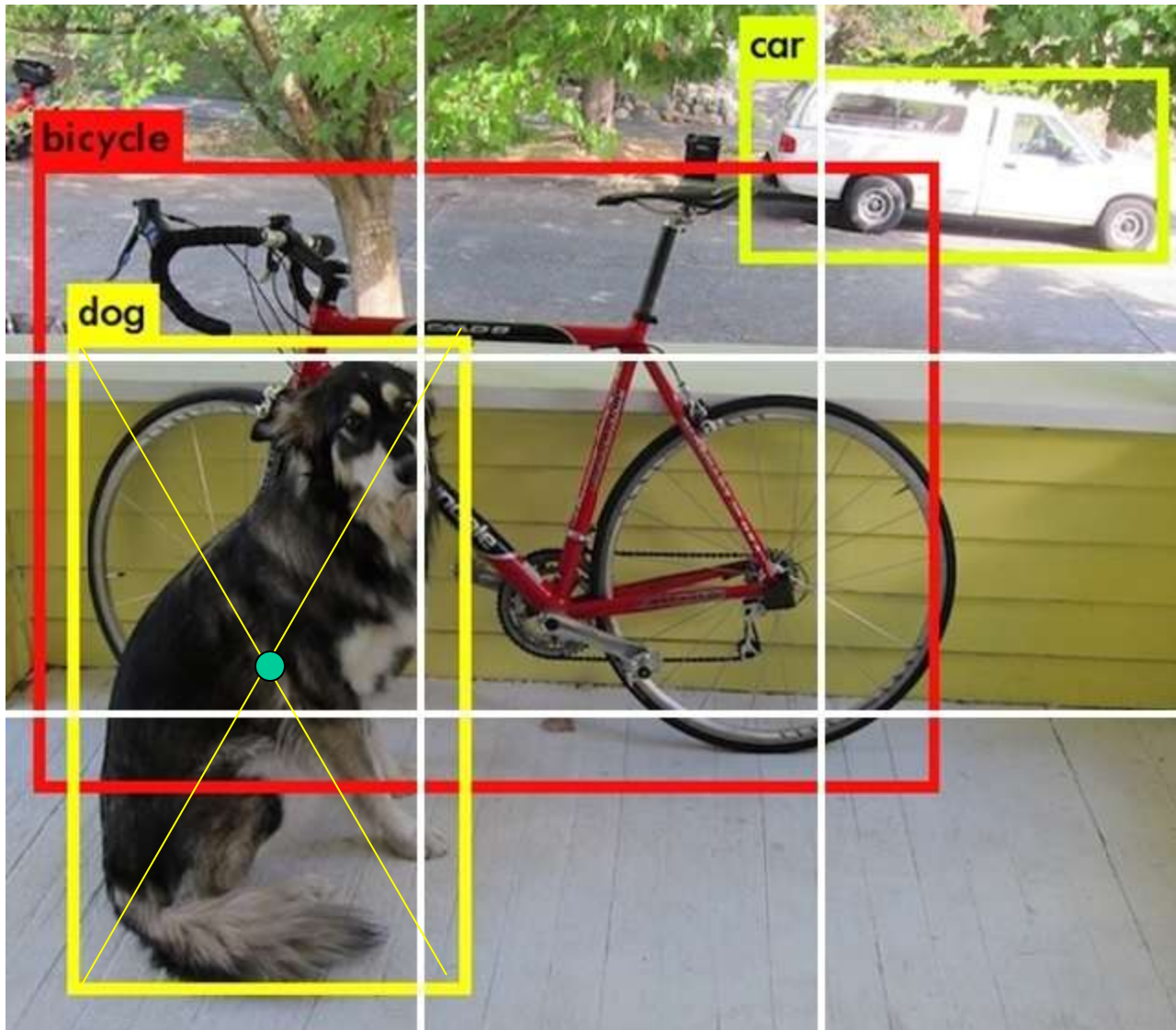
$(b_x \ b_y)$: center of the object corresponding to the upper left corner of the grid cell

b_h : height of the bounding box

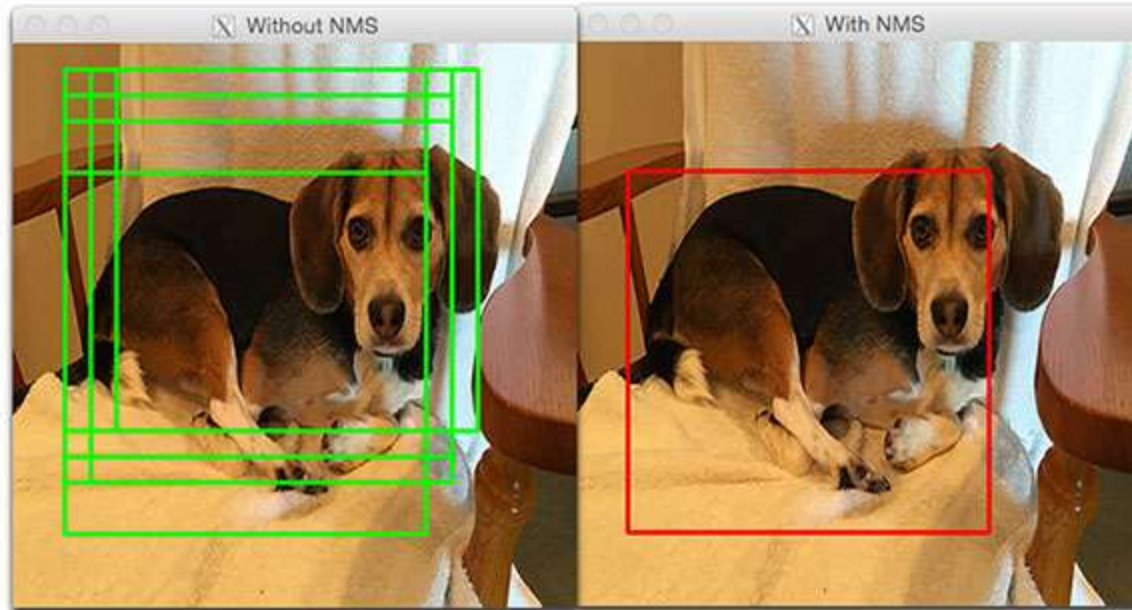
b_w : width of the bounding box

value c is corresponding to a class of an object (f.e. we have 3 types of targets to detect: car, bicycle, dog) Respectively.

$$y = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix}$$



Non-max suppression



Even though we ignored weak detections, YOLO can make duplicate detections for the same object. To fix this, YOLO applies non-maximal suppression to remove duplications with lower confidence. Non-maximal suppression adds 2-3% in mAP.