

CS230: Lecture 5

Advanced topics in Object Detection

Kian Katanforoosh, Andrew Ng
Stanford University

Today's outline

We will learn:

- the **Intuition** behind object detection methods
- a series of computer vision **papers**

I. Object detection Intuition

II. R-CNN

III. Fast R-CNN

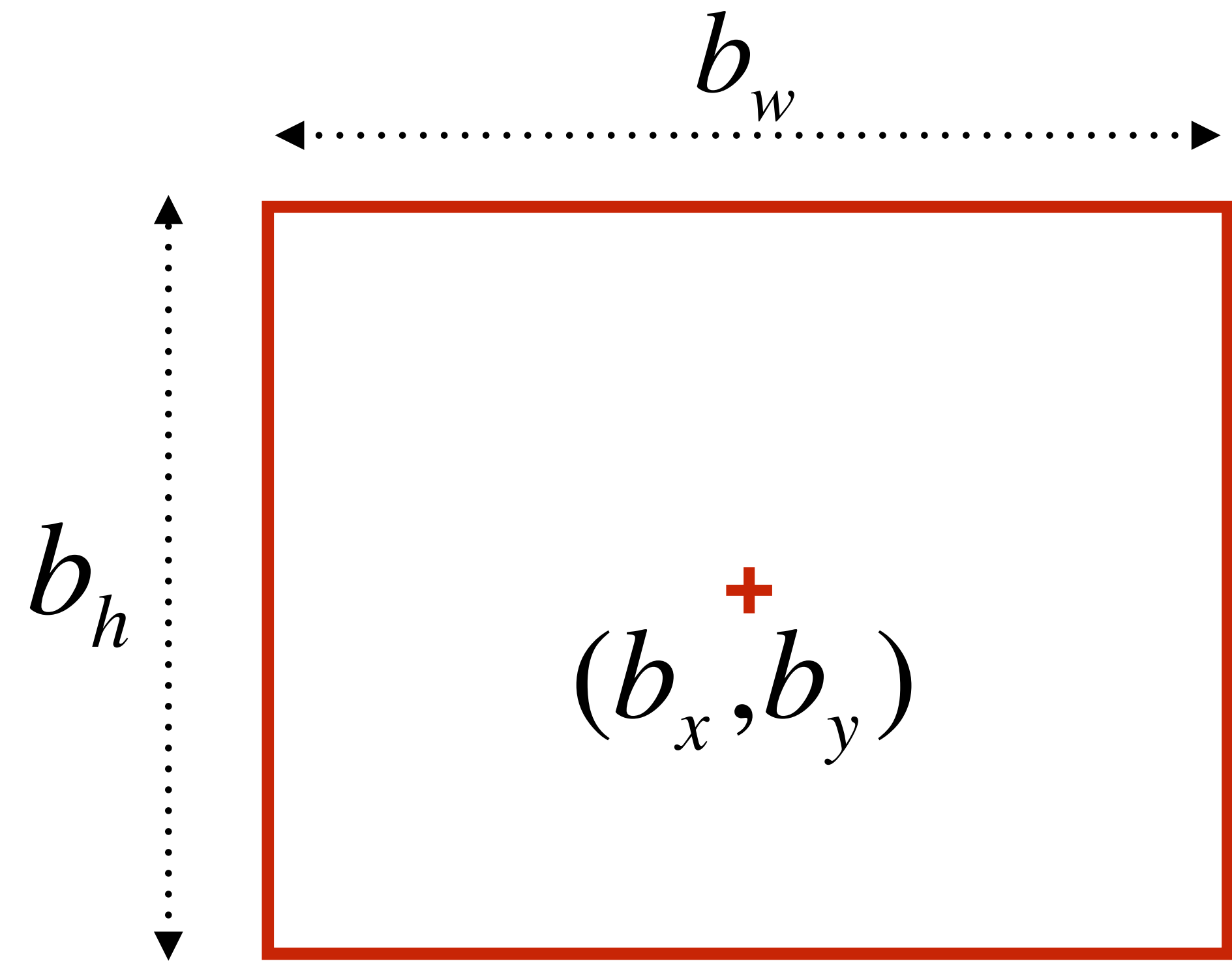
IV. Faster R-CNN

V. YOLO

Let's go

Object Detection Intuition

$$y = (p_c, b_x, b_y, b_h, b_w, c)$$

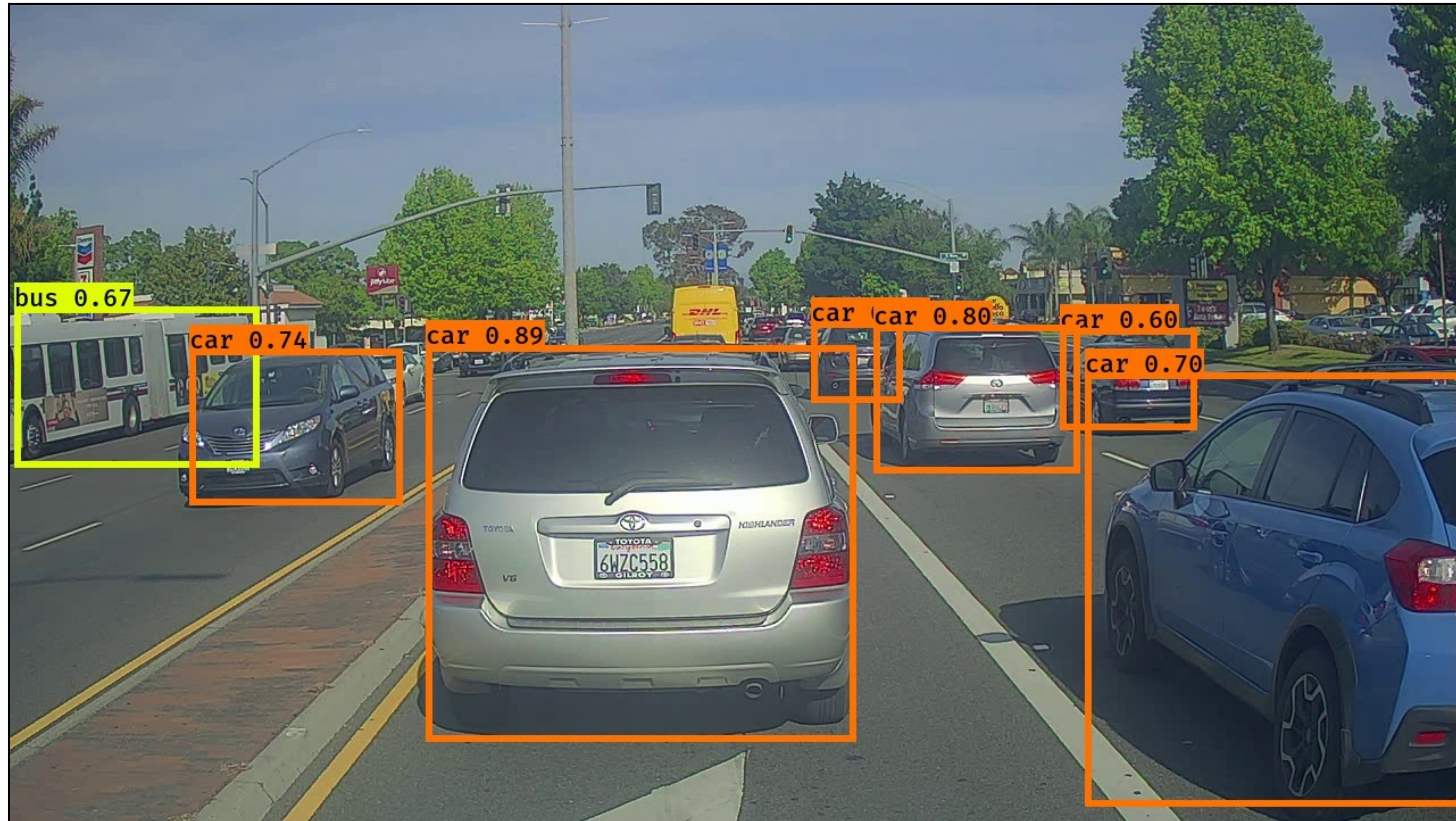


$p_c = 1$: confidence of an object being present in the bounding box

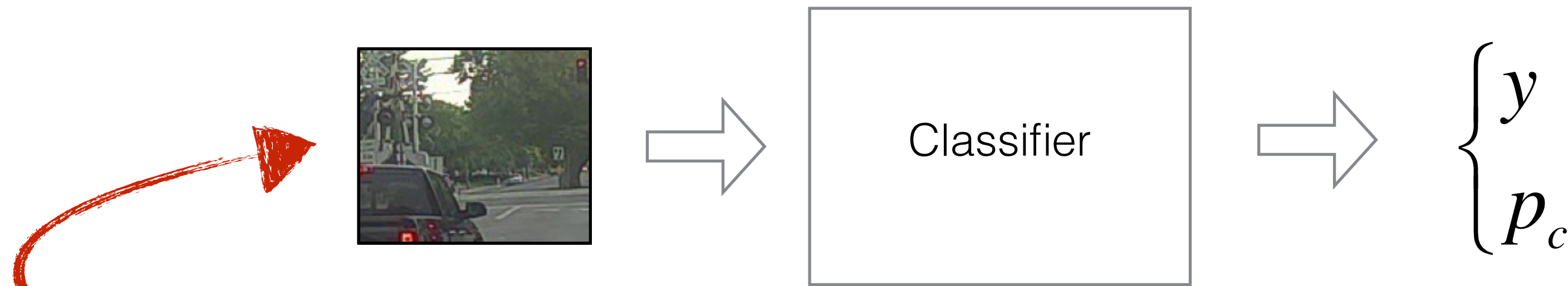
$c = 0$: class of the object being detected (here 0 for “car”)

Object Detection Intuition

$$y = (p_c, b_x, b_y, b_h, b_w, c)$$



Object Detection Intuition



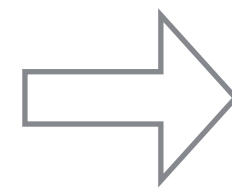
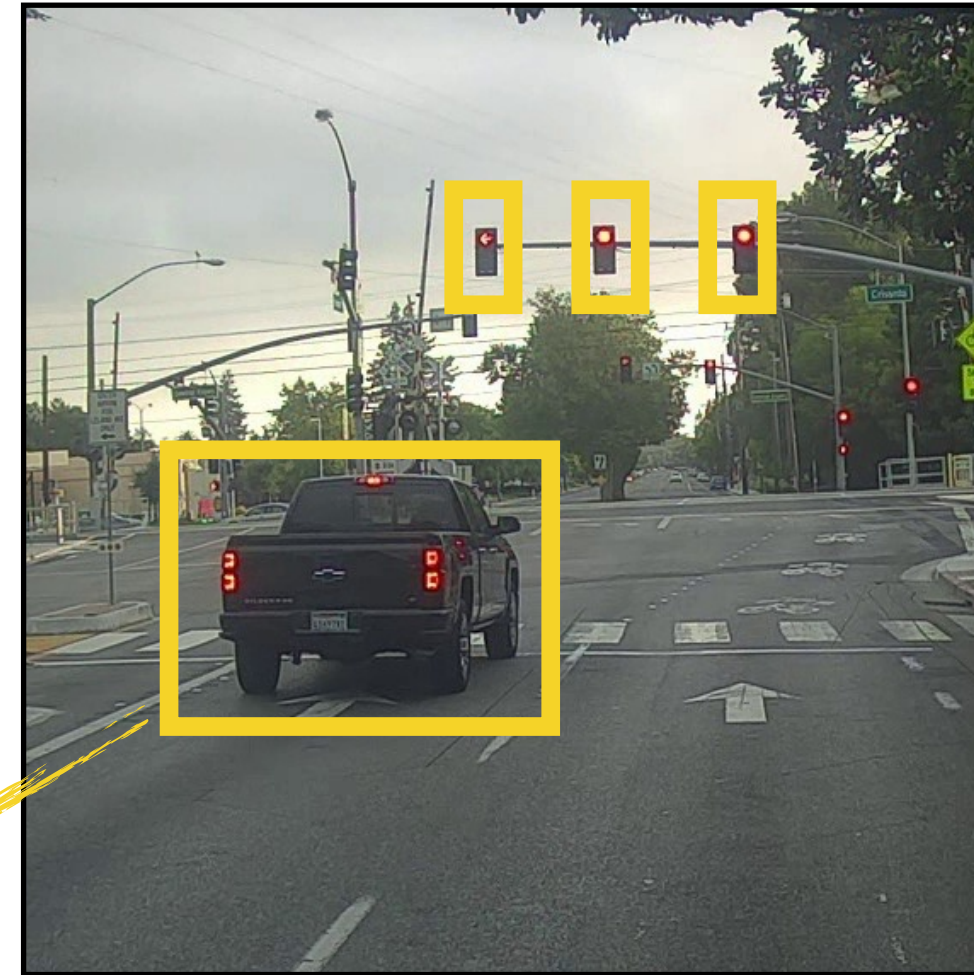
Far too many computations

R-CNN (Regions with CNN features)

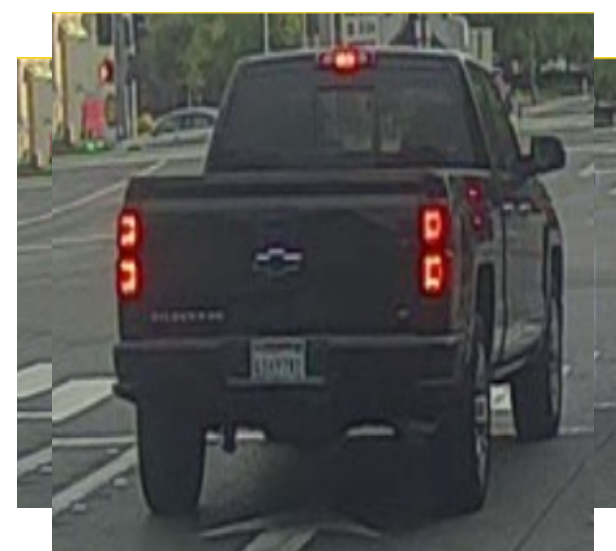
input image



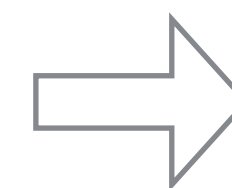
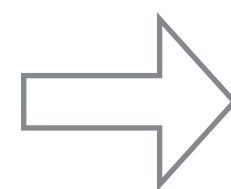
region proposal



How do you find these regions?

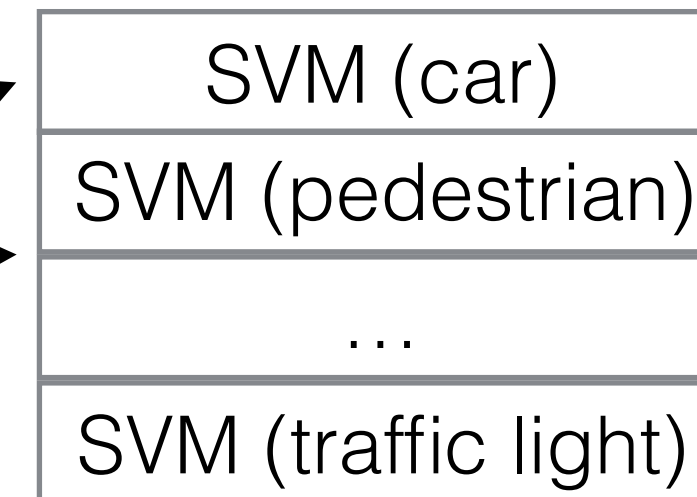
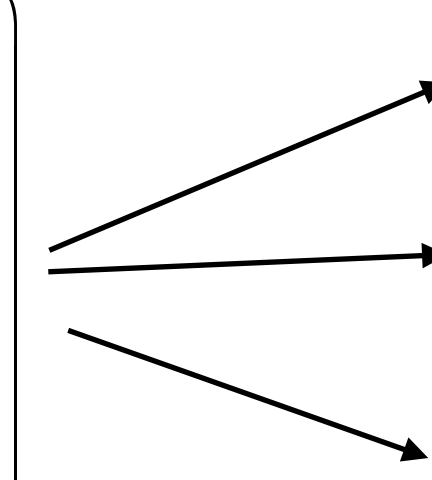


warped proposal



4096-d

(0.931
0.433
0.331
:
0.942
0.158
0.039)



“yes”
“no”
...
“no”

R-CNN (Regions with CNN features)

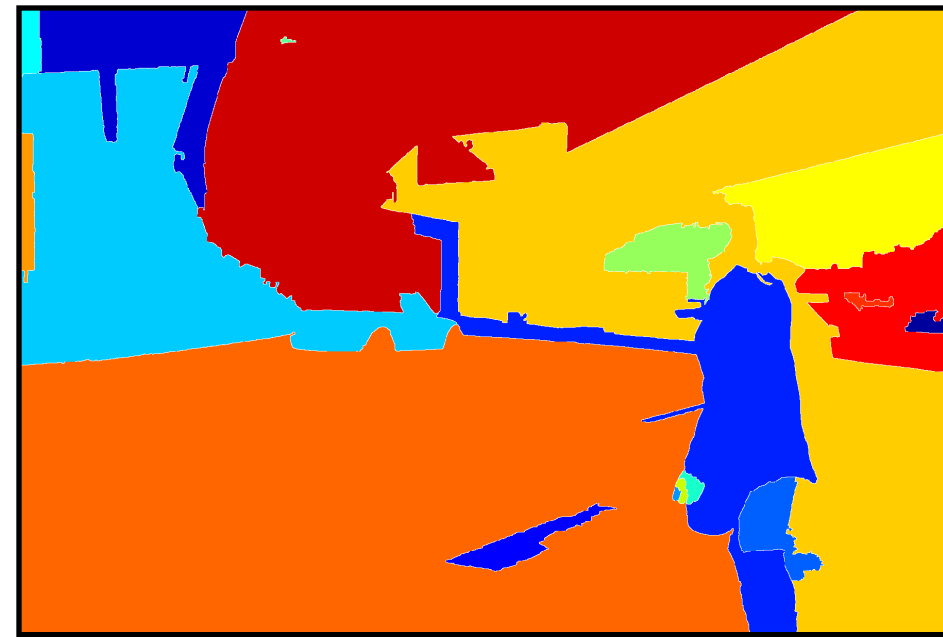
input image



“segmentation”



segmentation labels



“proposals”

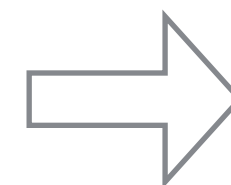
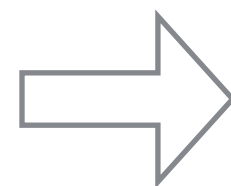


region proposals



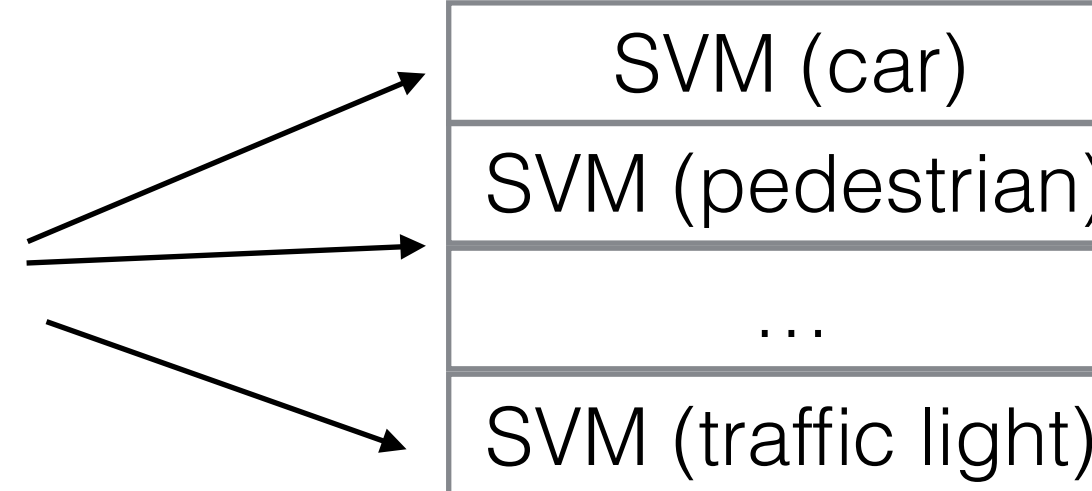
At test time on a new image, ≈ 2000 regions proposals!

region



4096-d

$\begin{pmatrix} 0.931 \\ 0.433 \\ 0.331 \\ \vdots \\ 0.942 \\ 0.158 \\ 0.039 \end{pmatrix}$



“yes”

“no”

...

“no”

x 2000

But we only needed two boxes: 1 car + 1 pedestrian ...

Non-max suppression

input image



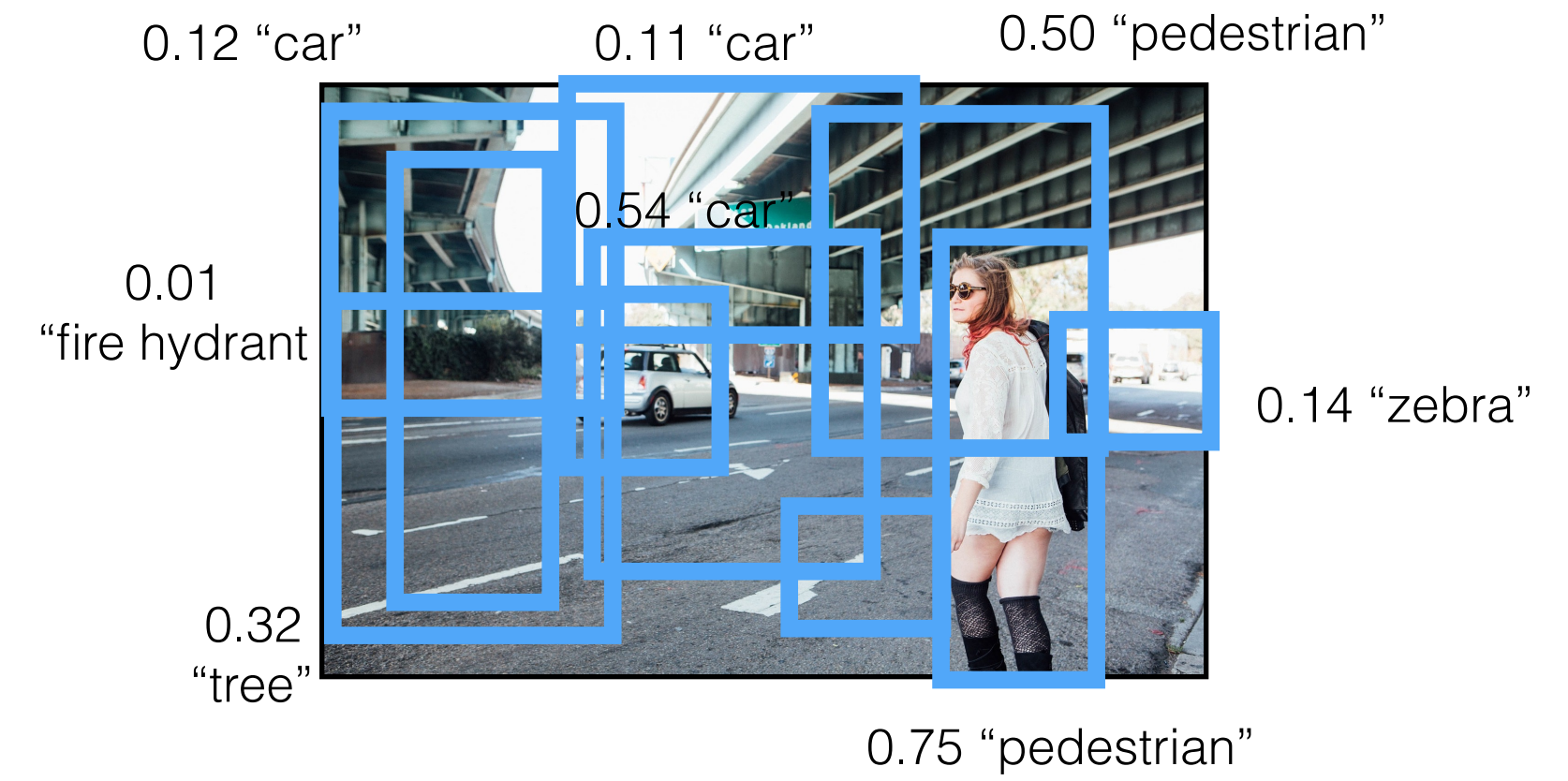
selective search
→

region proposals

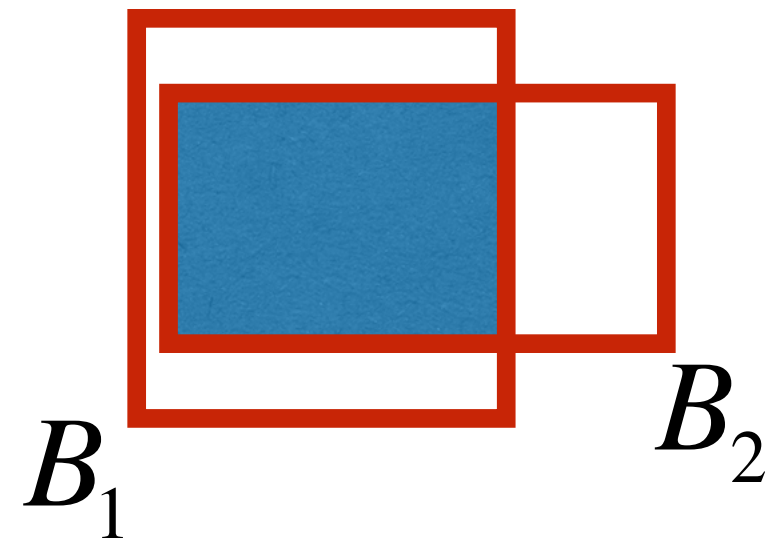


CNN + SVMs
→

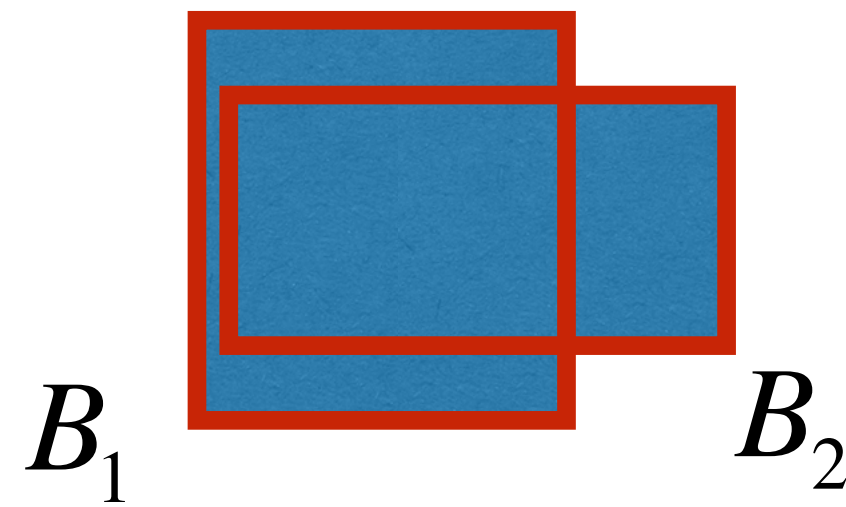
region + scores + classes



Intersection



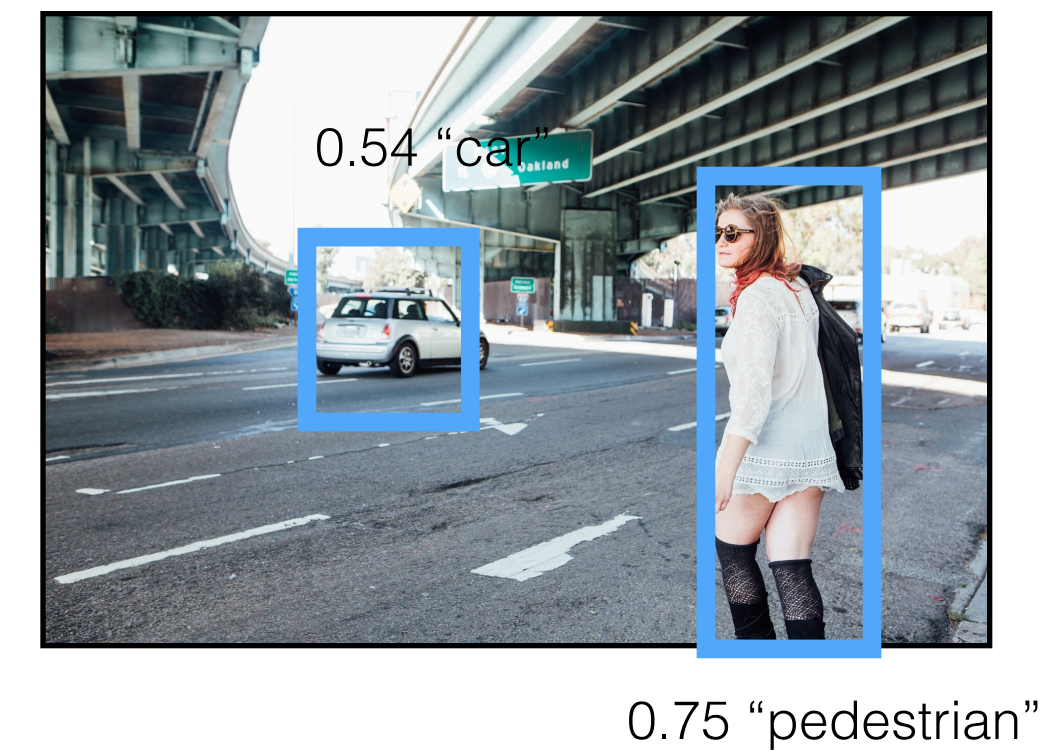
Union



Intersection over Union

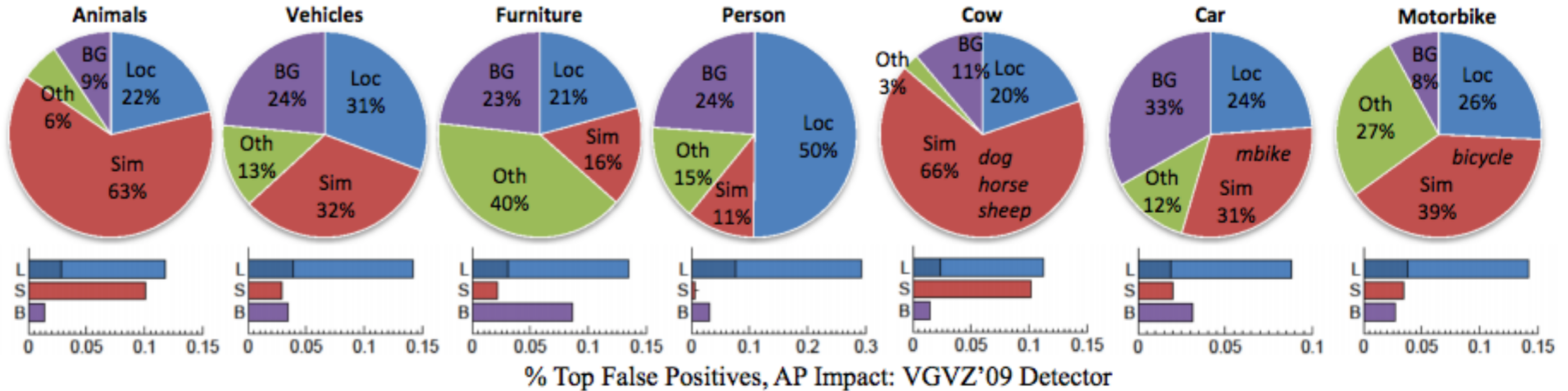
$$IoU = \frac{B_1 \cap B_2}{B_1 \cup B_2} = \frac{\text{Intersection Area}}{\text{Union Area}}$$

NMS + score threshold
↓



Still slow...
And training is not unified: CNN and each SVM have to be trained separately!

Evaluate the performance of your model



- BG: background and non-labelled objects
- Loc: localization error (IoU)
- Sim: Confusion with similar objects
- Oth: Confusion with dissimilar objects

Fast R-CNN

R-CNN drawbacks:

- Training is a multi-stage pipeline
- Training is expensive in space and time
- Testing is slow

Fast R-CNN solutions:

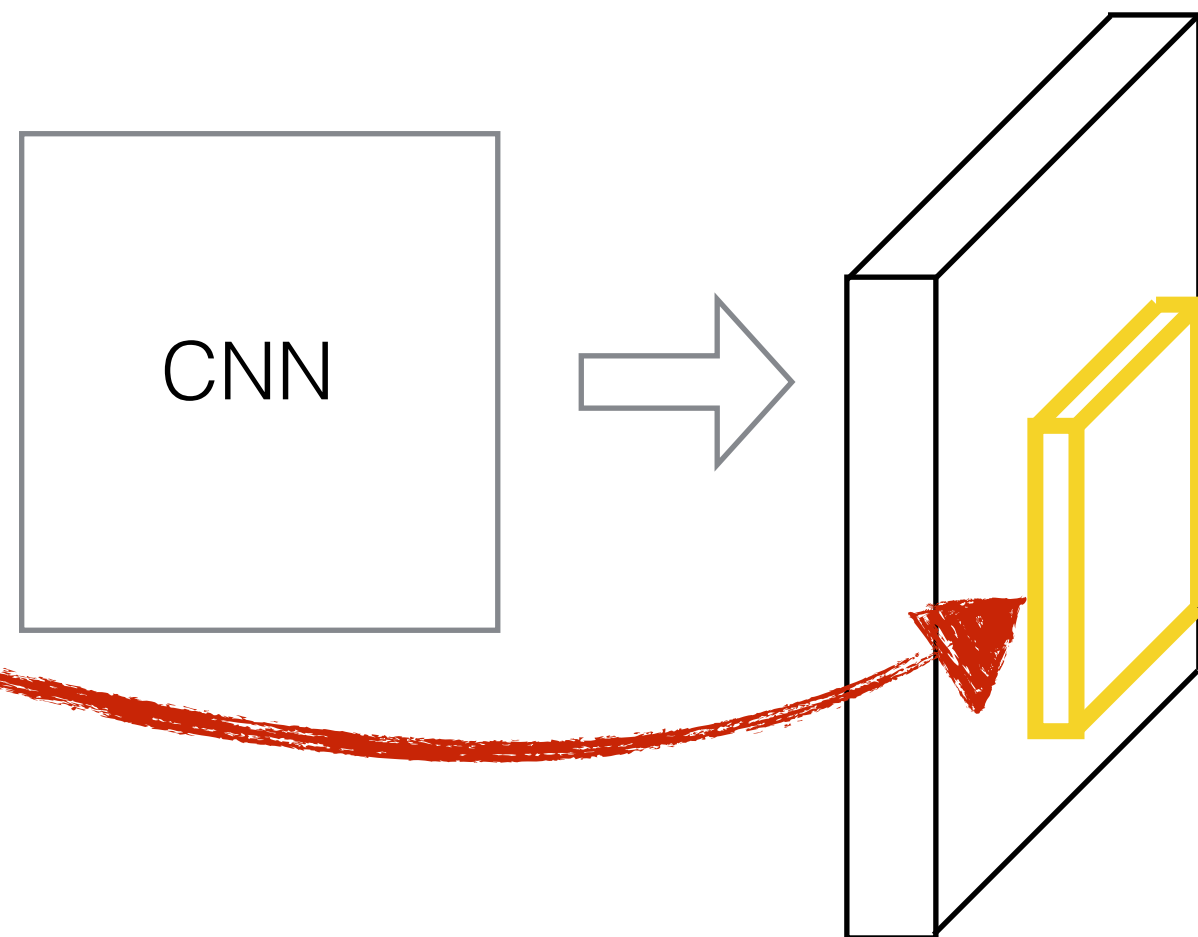
- Share computations for the CNN (not per-region anymore)
- Make training single-stage with multi-task loss
- No disk storage

Fast R-CNN

input image + Rols



feature map



RoI Pooling layer



FCs

4096-d

$\begin{pmatrix} 0.931 \\ 0.433 \\ 0.331 \\ \vdots \\ 0.942 \\ 0.158 \\ 0.039 \end{pmatrix}$

#classes

$\begin{pmatrix} 0.31 \\ 0.13 \\ 0.47 \\ \vdots \\ 0.43 \\ 0.81 \\ 0.09 \end{pmatrix}$

FCs (Softmax)

How do you train this?

$\begin{pmatrix} 0.24 \\ 0.32 \\ 0.18 \\ 0.14 \end{pmatrix} \dots \begin{pmatrix} 0.14 \\ 0.54 \\ 0.19 \\ 0.34 \end{pmatrix}$

\longleftrightarrow
 #classes

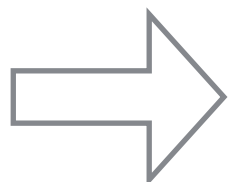
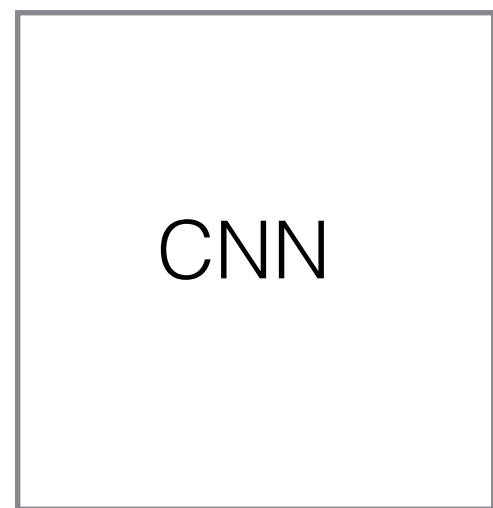
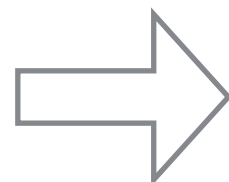
FCs (Bbox regressor)

$$L(p, u, t^u, v) = L_{cls}(p, u) + \lambda \mathbf{1}_{[u \geq 1]} L_{loc}(t^u, v)$$

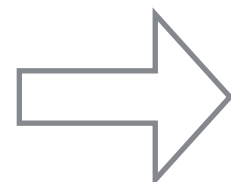
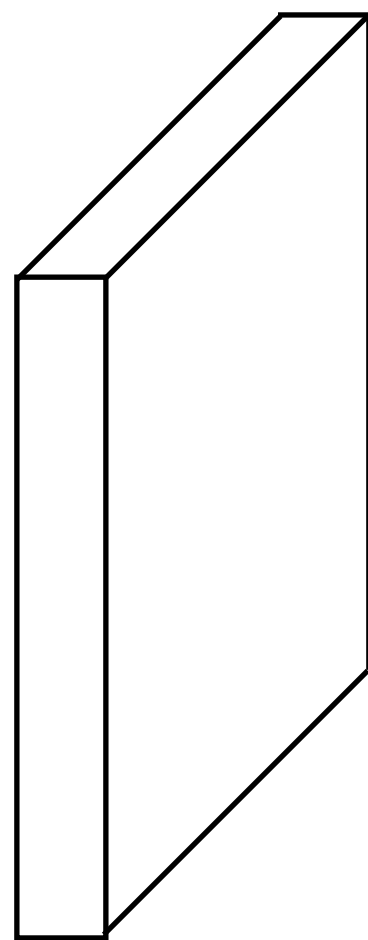
Let's make it even faster

Faster R-CNN

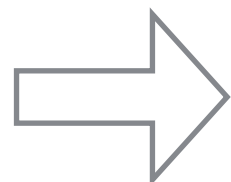
input image (no RoI)



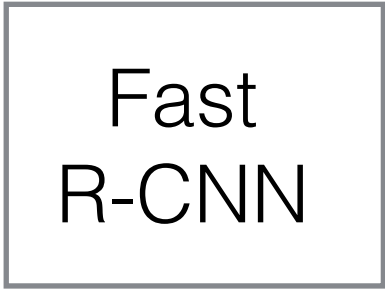
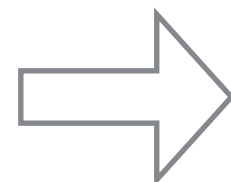
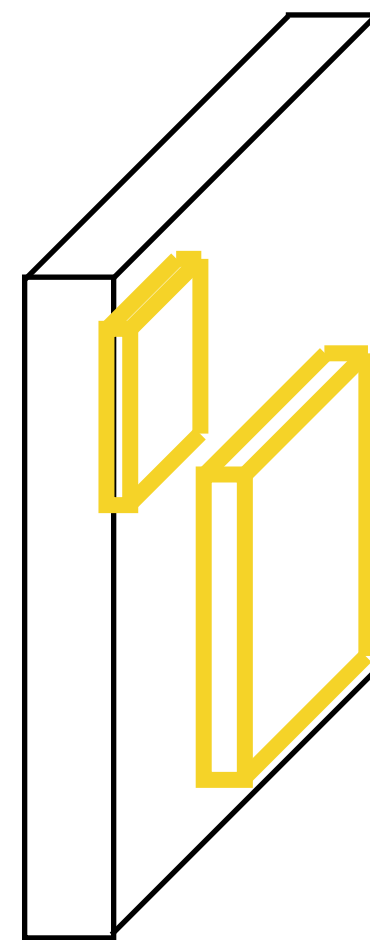
feature map



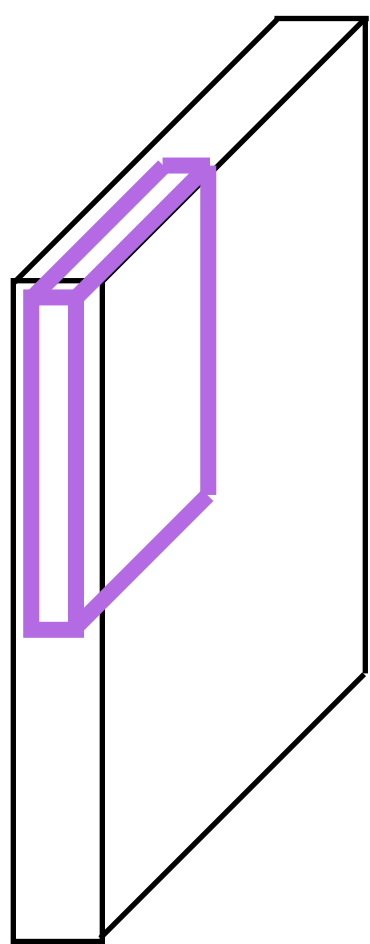
“attention”



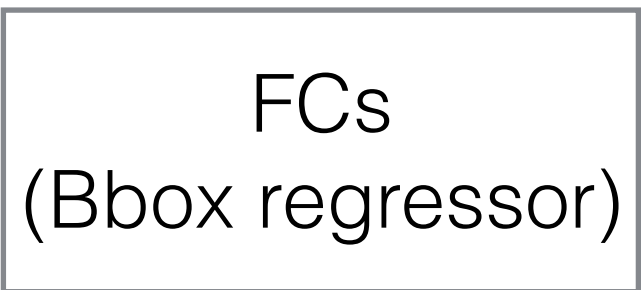
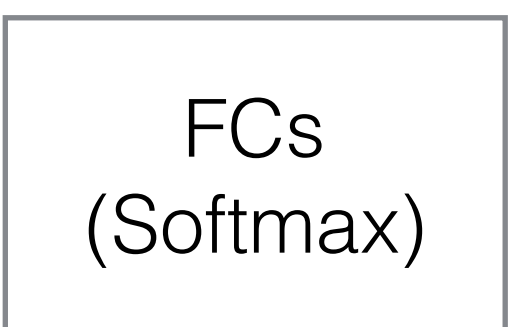
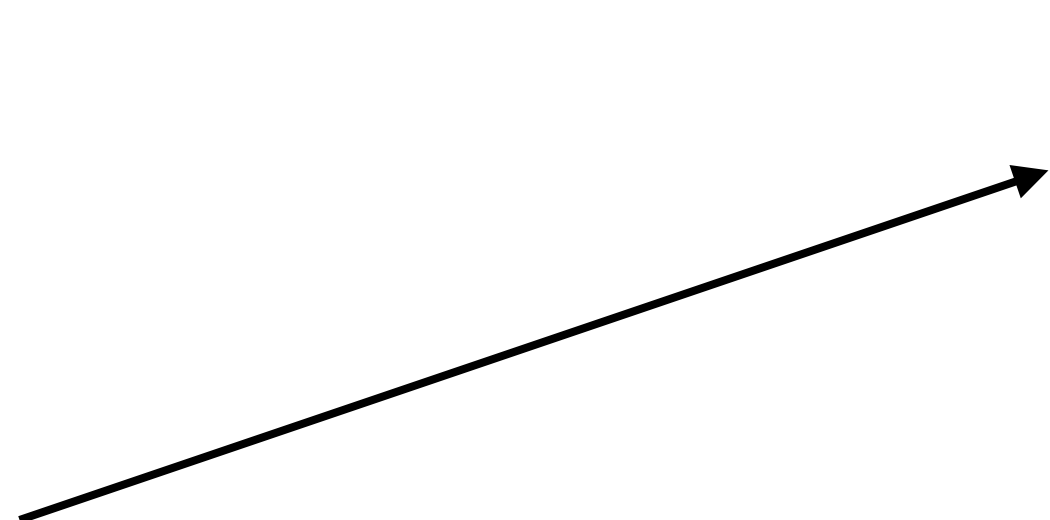
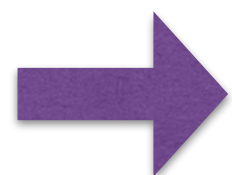
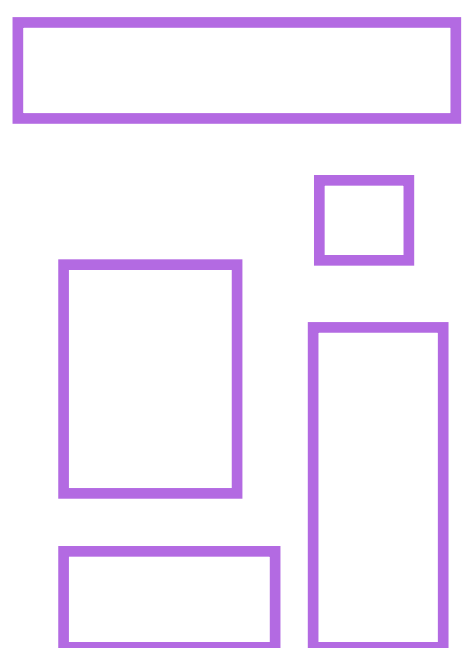
feature map with RoIs



feature map



k anchors per window



k scores

- 0.31
- 0.13
- 0.47
- ⋮
- 0.43
- 0.81
- 0.09

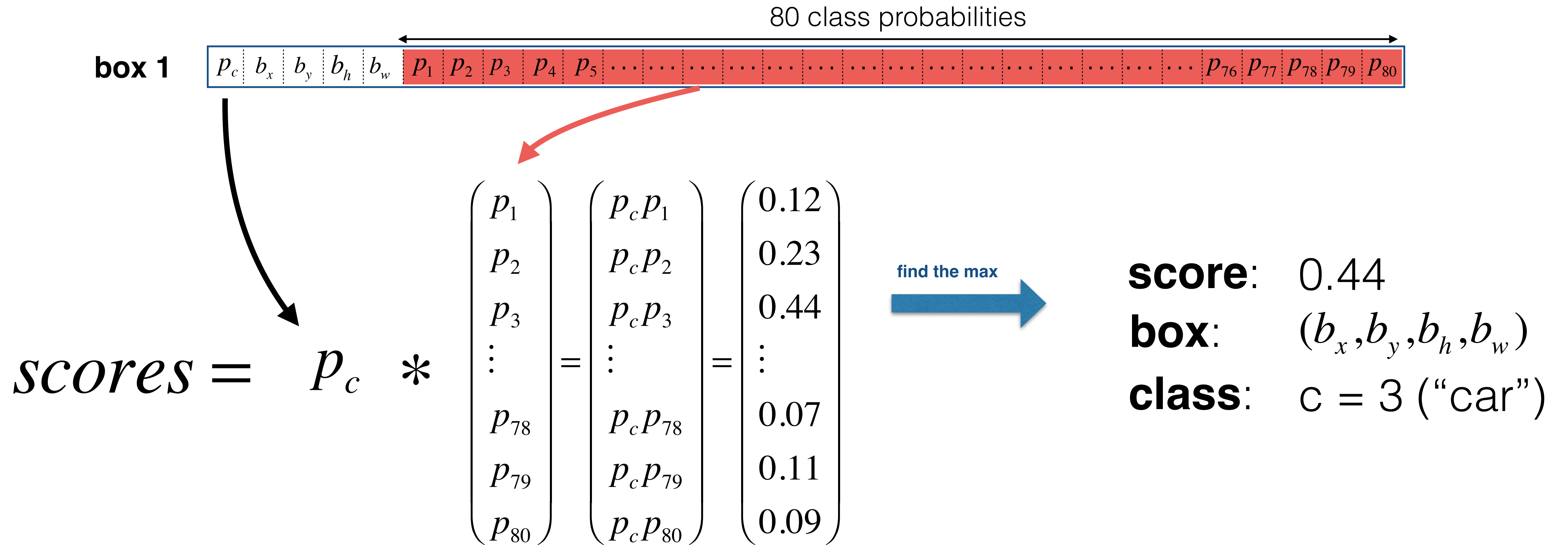
- (0.24)
- (0.32)
- (0.18)
- (0.14)
- ⋯
- (0.14)
- (0.54)
- (0.19)
- (0.34)



$$\begin{aligned}
 t_x &= (x - x_a)/w_a, & t_y &= (y - y_a)/h_a, \\
 t_w &= \log(w/w_a), & t_h &= \log(h/h_a), \\
 t_x^* &= (x^* - x_a)/w_a, & t_y^* &= (y^* - y_a)/h_a, \\
 t_w^* &= \log(w^*/w_a), & t_h^* &= \log(h^*/h_a),
 \end{aligned}$$

$$L(\{p_i\}, \{t_i\}) = \frac{1}{N_{cls}} \sum_i L_{cls}(p_i, p_i^*) + \lambda \frac{1}{N_{reg}} \sum_i p_i^* L_{reg}(t_i, t_i^*)$$

YOLO



the box (b_x, b_y, b_h, b_w) has detected $c = 3$ ("car") with probability score: 0.44

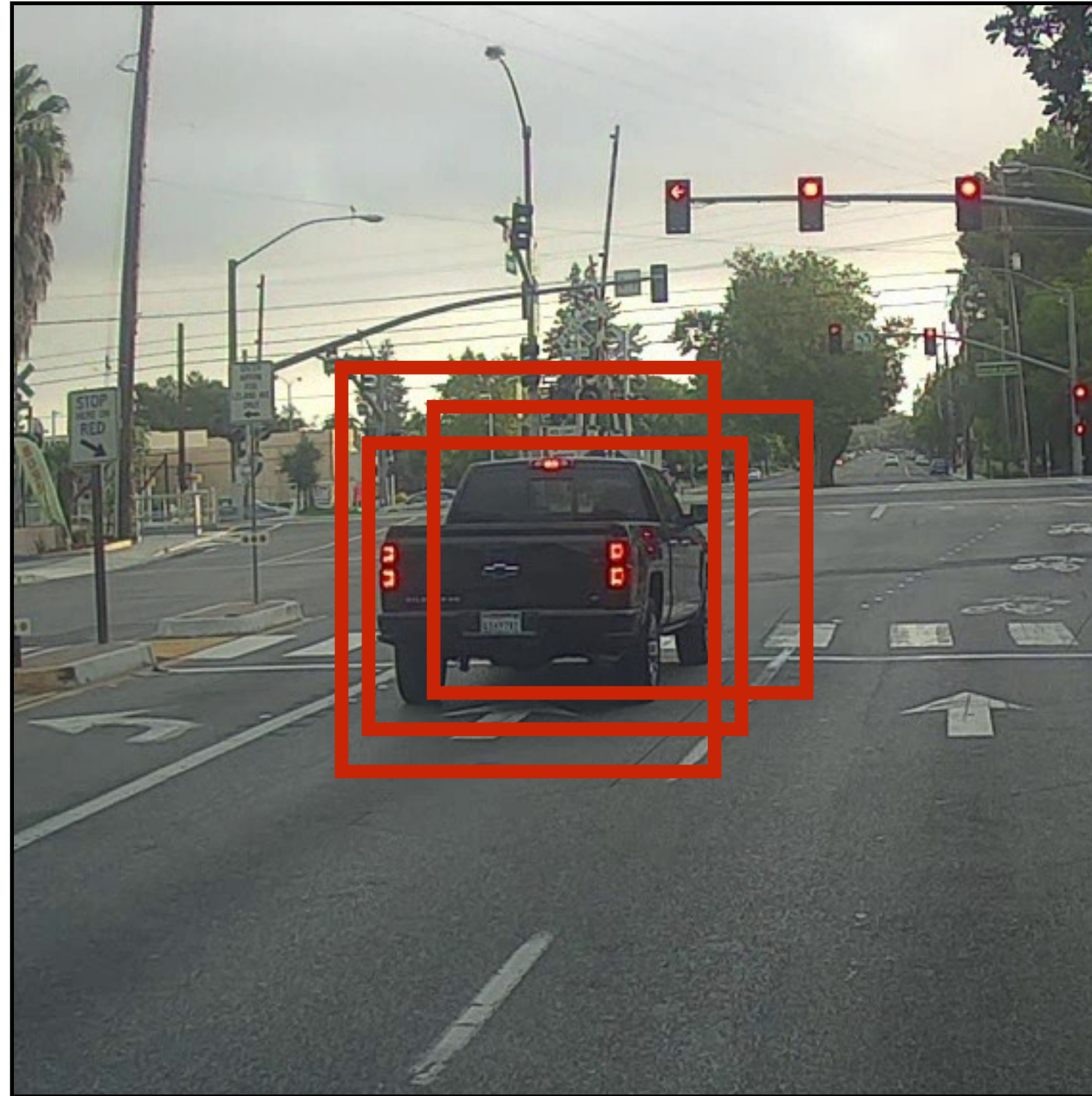
YOLO



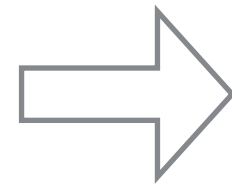
Need to filter:
- Score thresholding
- NMS

YOLO

Before non-max suppression



**Non-Max
Suppression**



After non-max suppression



YOLO

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

demo