

25 YEARS ANNIVERSARY
SOKT

ĐẠI HỌC BÁCH KHOA HÀ NỘI
VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG



HA NOI UNIVERSITY OF SCIENCE AND TECHNOLOGY
SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

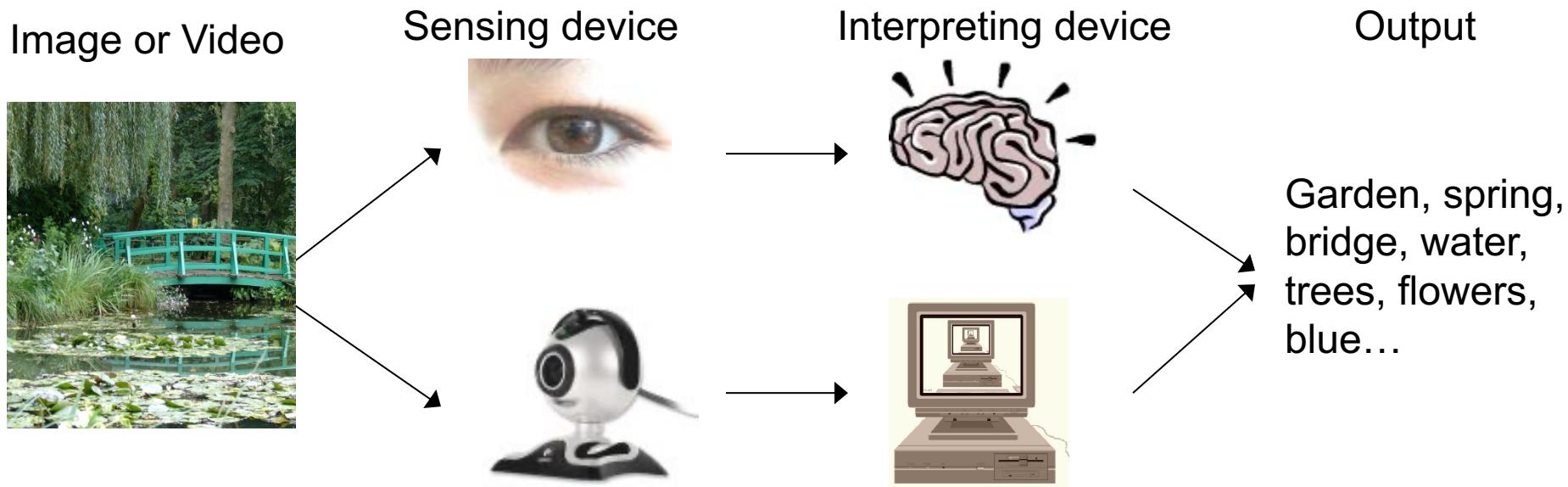
Lesson 6: Object detection

Outline

1. Overview of computer vision and applications
2. Introduction to object detection
3. Regional proposal networks: R-CNN, Fast R-CNN, Faster R-CNN...
4. None regional proposal networks: SSD, Yolo...

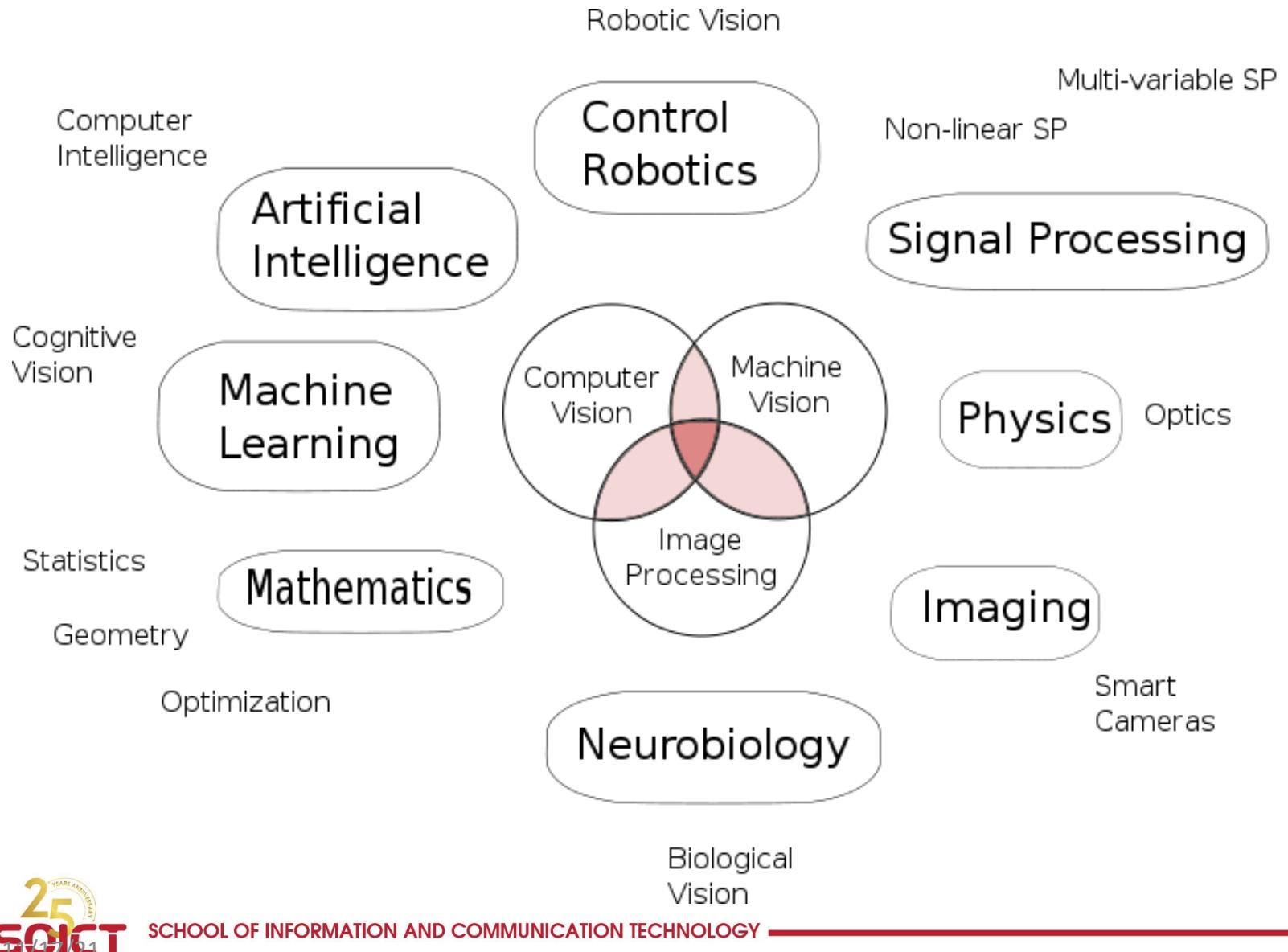
Overview of computer vision and applications

What is computer vision?



- CV is a scientific field that extracts information out of digital images.
- CV Is building algorithms that can understand the content of images and use it for other applications

Overview of relations between computer vision and other fields



Human eyes are very sensitive



Can u detect a person?

100 ms per frame,

Never seen the picture, do not know the
person

Can do effortlessly

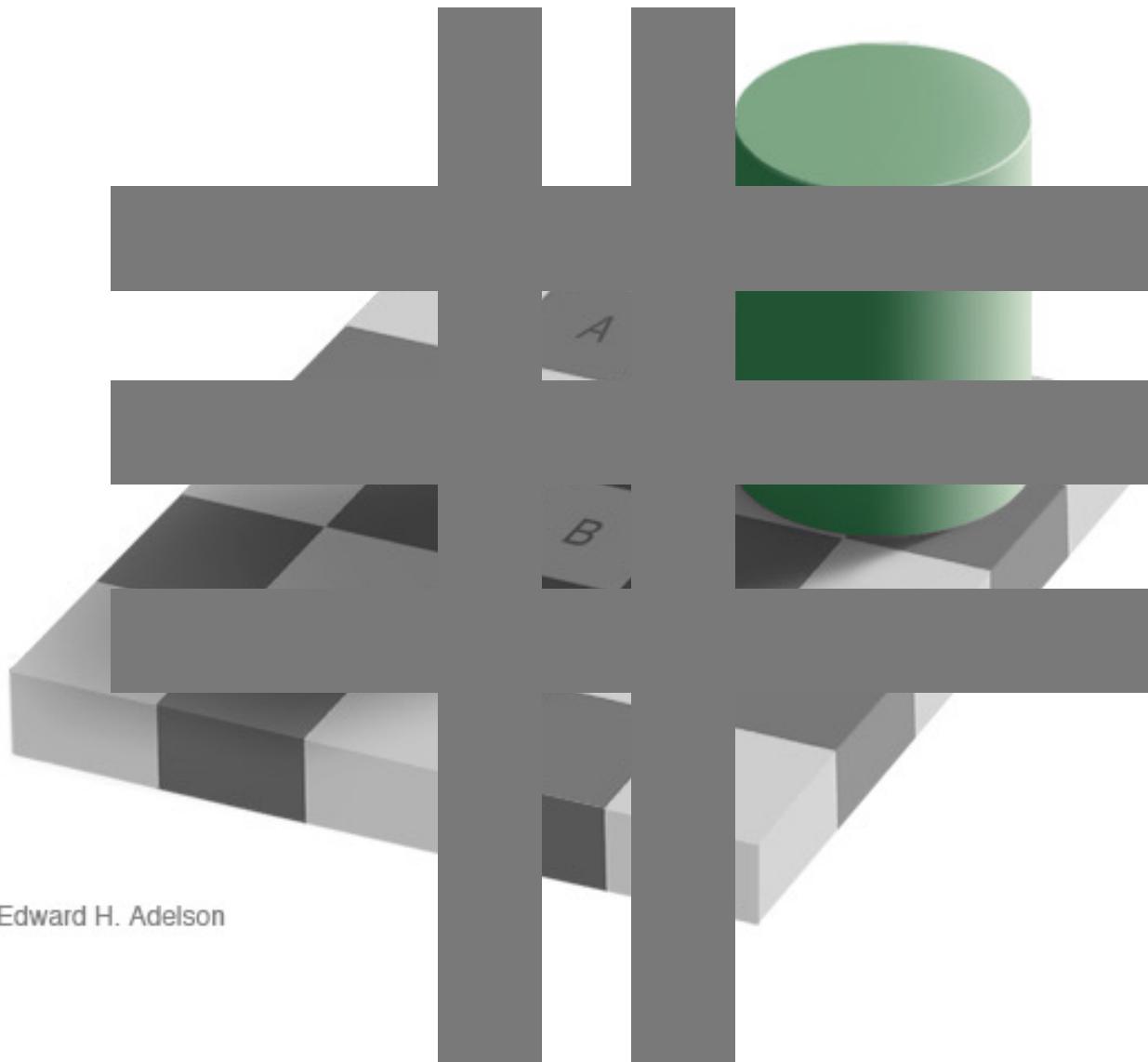
Potter, Biederman, etc. 1970s

However, this speed is obtained at the price of some drawbacks

- Can you recognize who is it?



Do you think the colors in boxes A and B are different?



Edward H. Adelson

Objectives of computer vision

- Bridges between numerically represented pixels with semantics



La Gare Montparnasse, 1895

0	3	2	5	4	7	6	9	8
3	0	1	2	3	4	5	6	7
2	1	0	3	2	5	4	7	6
5	2	3	0	1	2	3	4	5
4	3	2	1	0	3	2	5	4
7	4	5	2	3	0	1	2	3
6	5	4	3	2	1	0	3	2
9	6	7	4	5	2	3	0	1
8	7	6	5	4	3	2	1	0

Why learn computer vision?

- Helpful: Photos and videos are everywhere!



Google
Image Search™

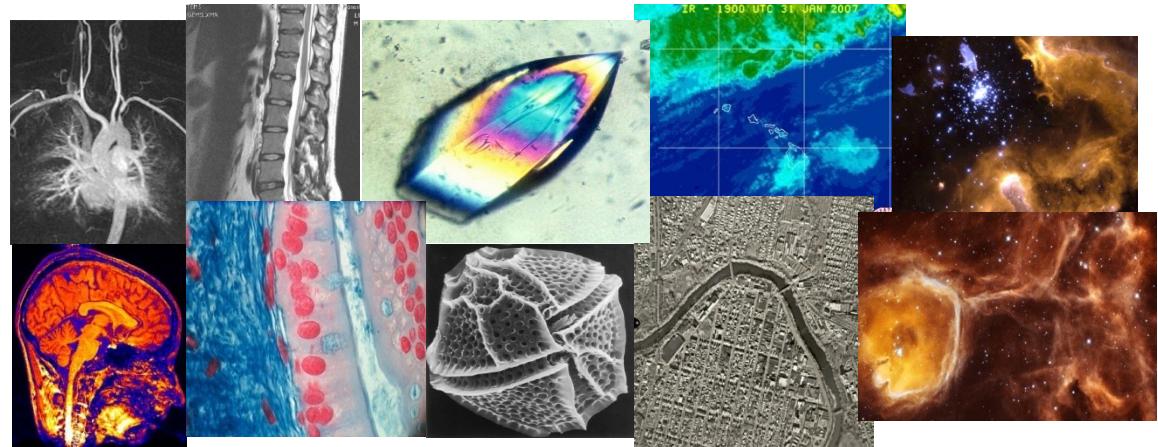
Google Photos

flickr GAMMA™

webshots beta

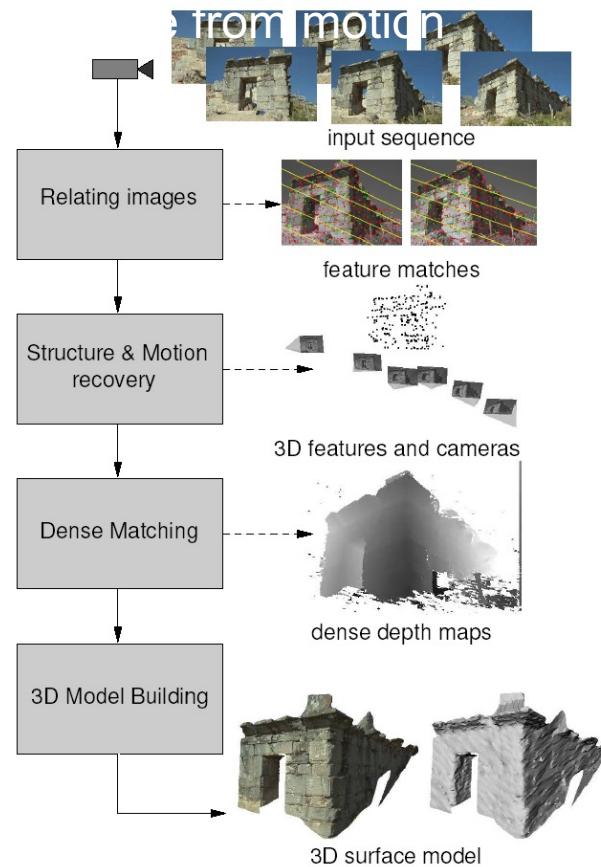
picsearch™

YouTube
Broadcast Yourself™

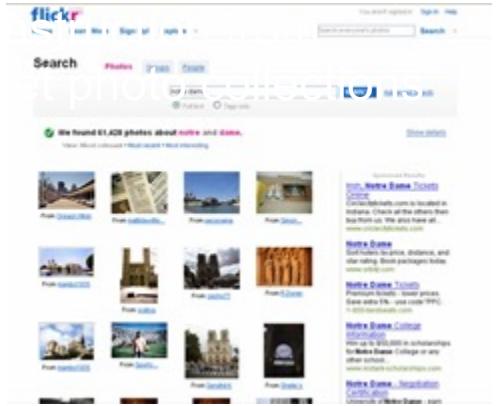


Machine vision can be used as instrumentation

Real-time stereo

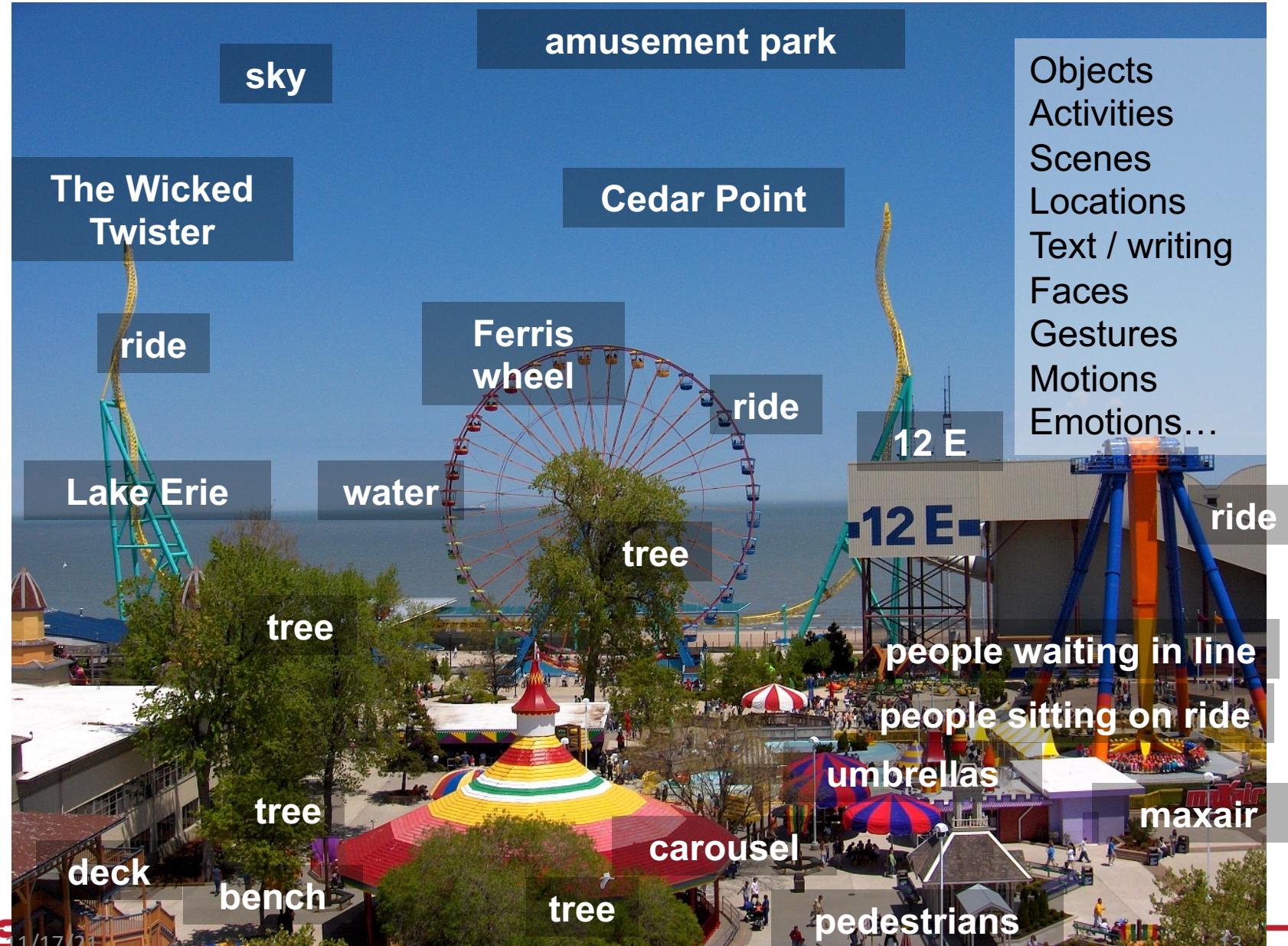


Pollefeys et al.

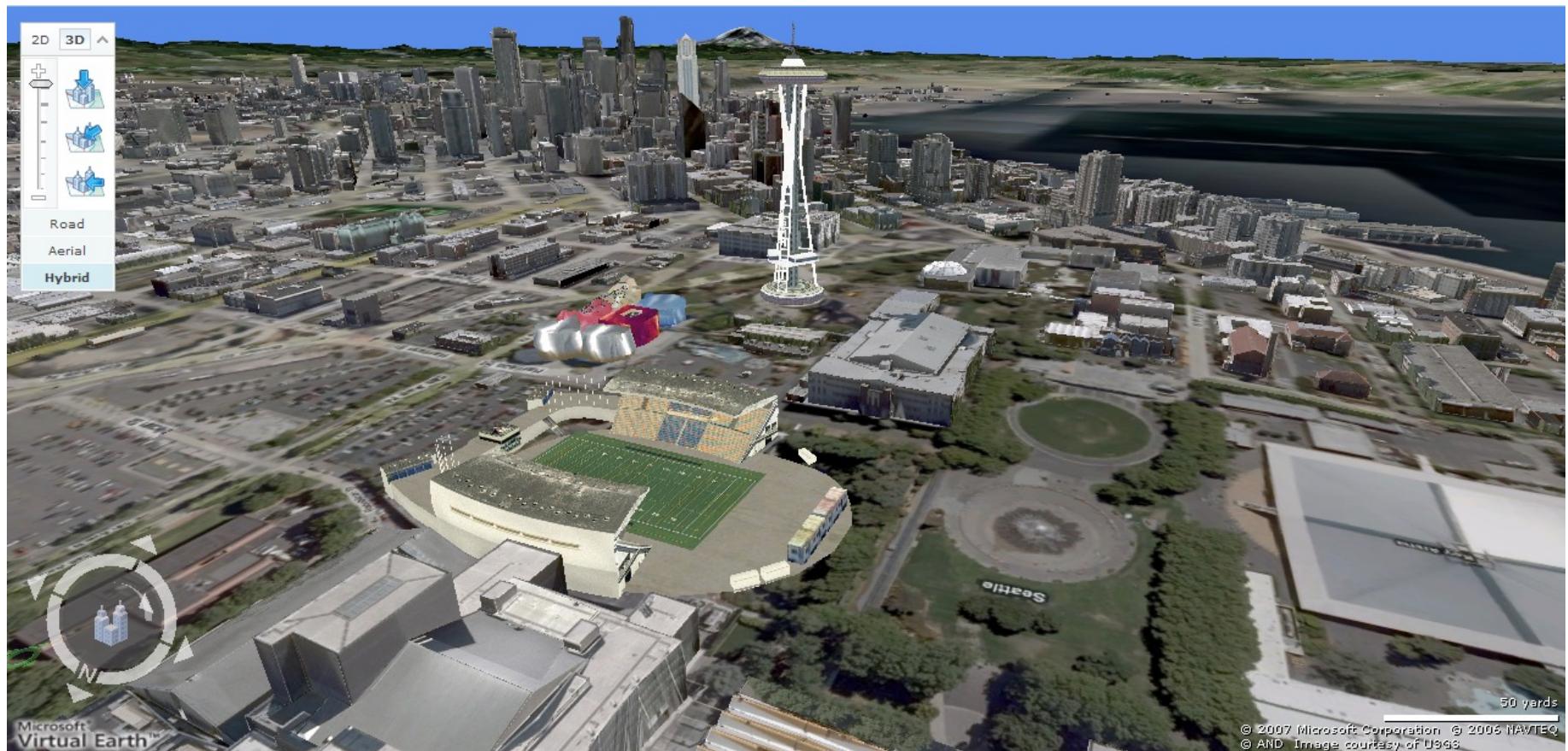


Goesele et al.

Machine vision is a source of semantic information



City 3D modeling



Bing maps, Google Streetview

Source: S. Seitz

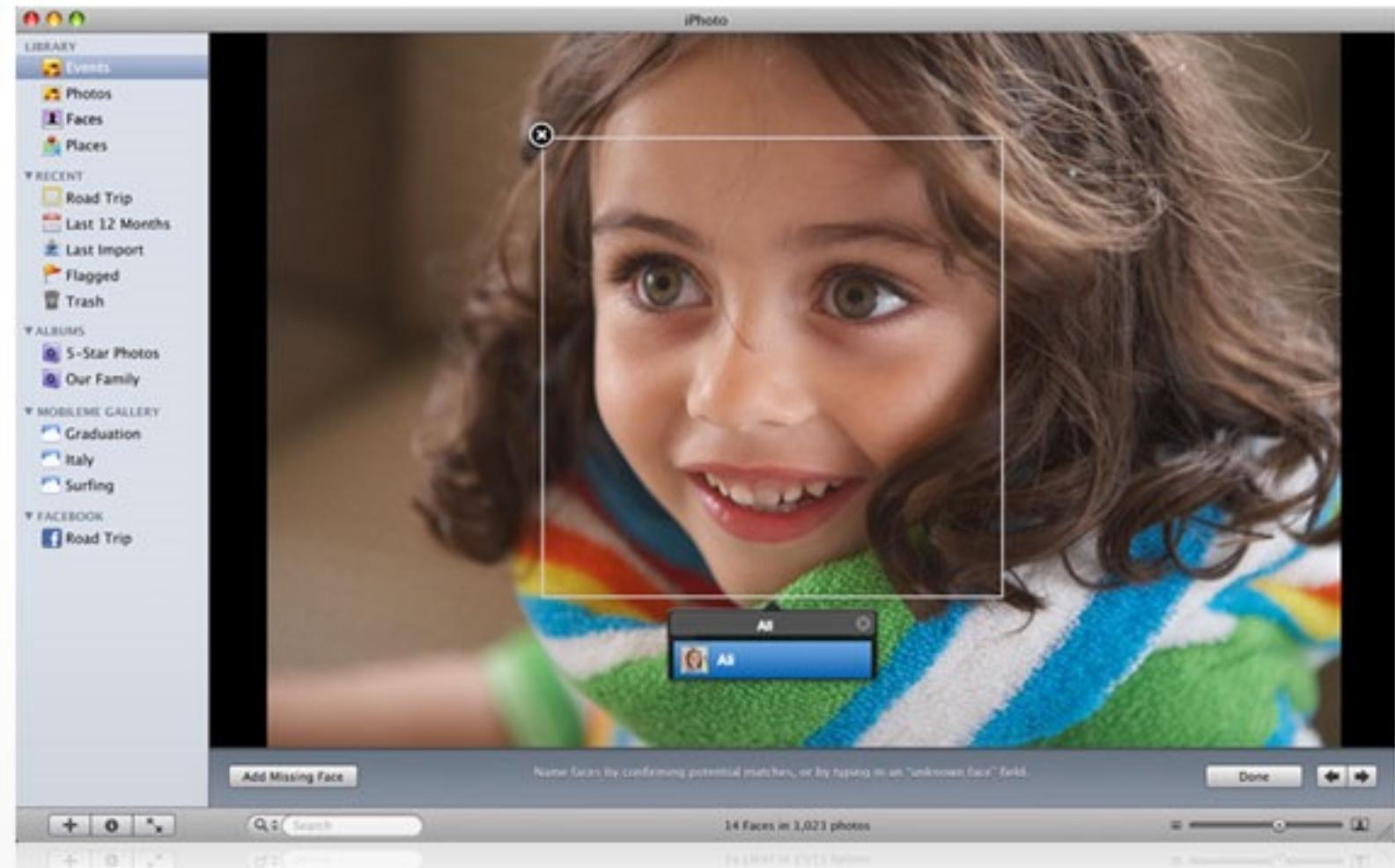
Face detection



- Many digital cameras can automatically detect faces
 - Canon, Sony, Fuji, ...

Source: S. Seitz

Face recognition: Apple iPhoto

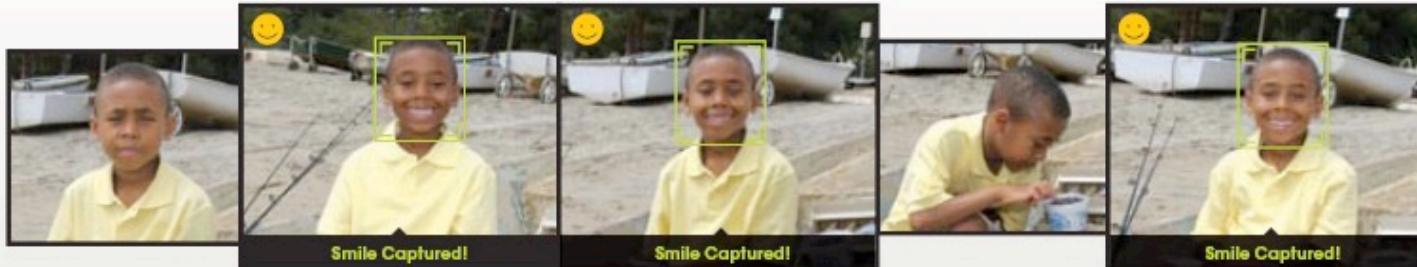


<http://www.apple.com/ilife/iphoto/>

Smile detection

The Smile Shutter flow

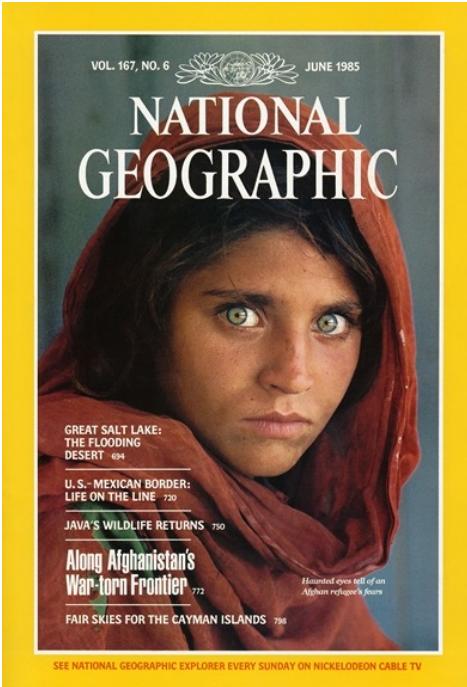
Imagine a camera smart enough to catch every smile! In Smile Shutter Mode, your Cyber-shot® camera can automatically trip the shutter at just the right instant to catch the perfect expression.



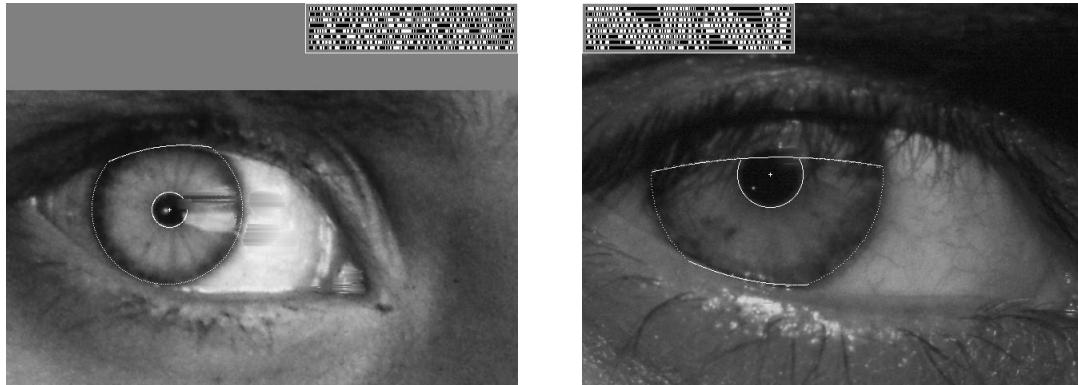
[Sony Cyber-shot® T70 Digital Still Camera](#)

Source: S. Seitz

Biometrics



How the Afghan Girl was Identified by Her Iris Patterns



Source: S. Seitz

Biometrics

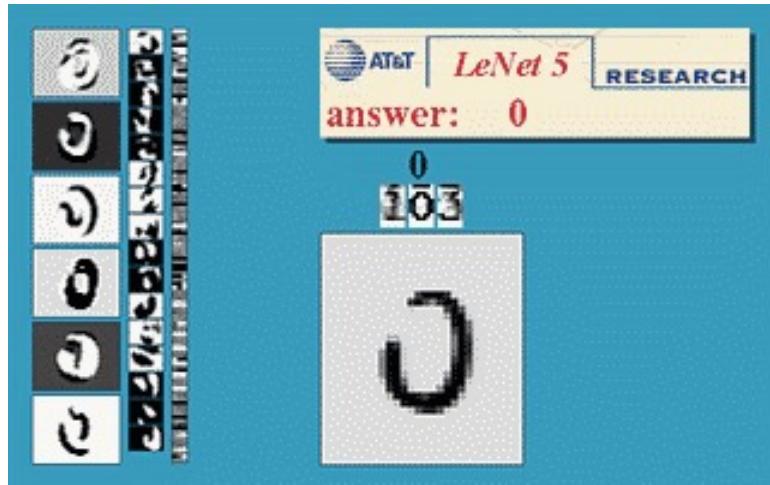


Fingerprint scanners on many laptops and devices



Facial recognition systems appear more and more
Example: iPhone X just introduced FaceID

Optical character recognition (OCR)



Handwritten digit recognition, AT&T labs

4YCH428

4YCH428

4YCH428

Plate reader

http://en.wikipedia.org/wiki/Automatic_number_plate_recognition

Human computer interactions and games



Microsoft's Kinect



Sony EyeToy



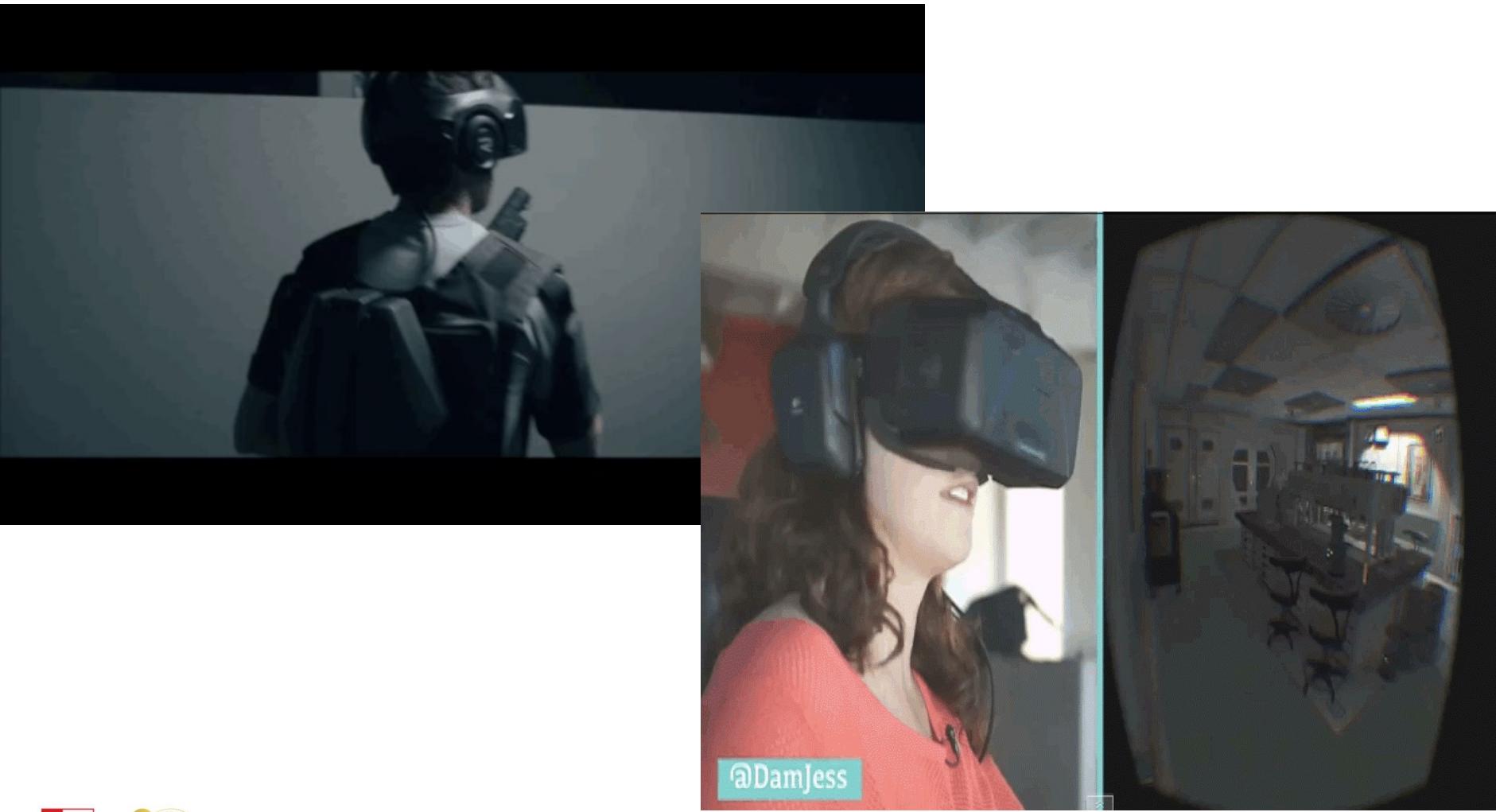
Assistive technologies

Source: S. Seitz

Augmented reality



Virtual reality



Applications in robotics and space exploration



[NASA'S Mars Exploration Rover Spirit](#) captured this westward view from atop a low plateau where Spirit spent the closing months of 2007.

Vision used for many different tasks:

- Panorama
- 3D modeling of the Martian surface
- Obstacle detection, location tracking

For details see “Computer Vision on Mars” by Matthies et al.

Source: S. Seitz

Introduction to object detection

Computer vision problems

Semantic Segmentation



GRASS, CAT,
TREE, SKY

No objects, just pixels

Classification + Localization



CAT

Single Object

Object Detection



DOG, DOG, CAT

Multiple Object

Instance Segmentation

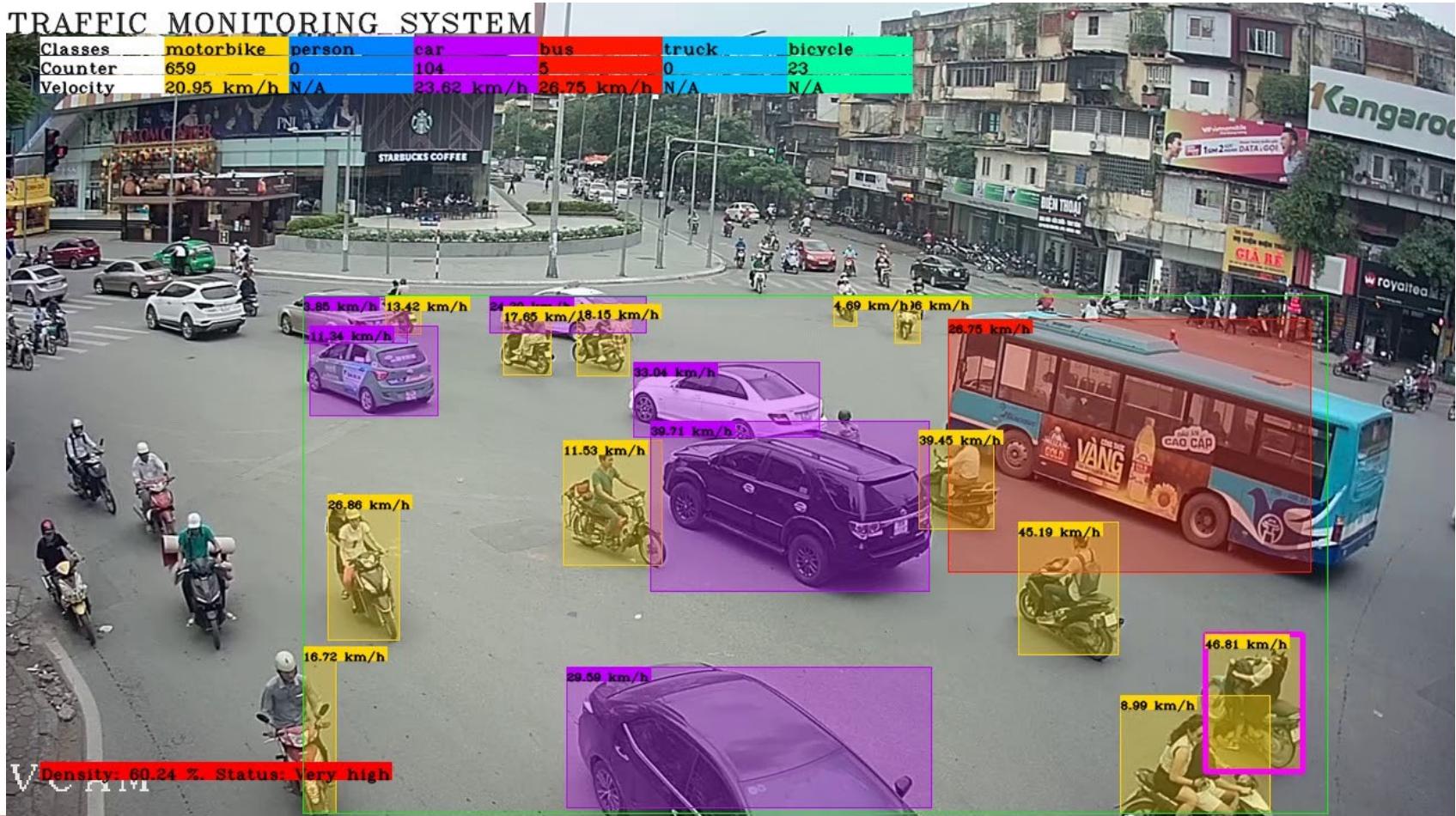


DOG, DOG, CAT

[This image is CC0 public domain](#)

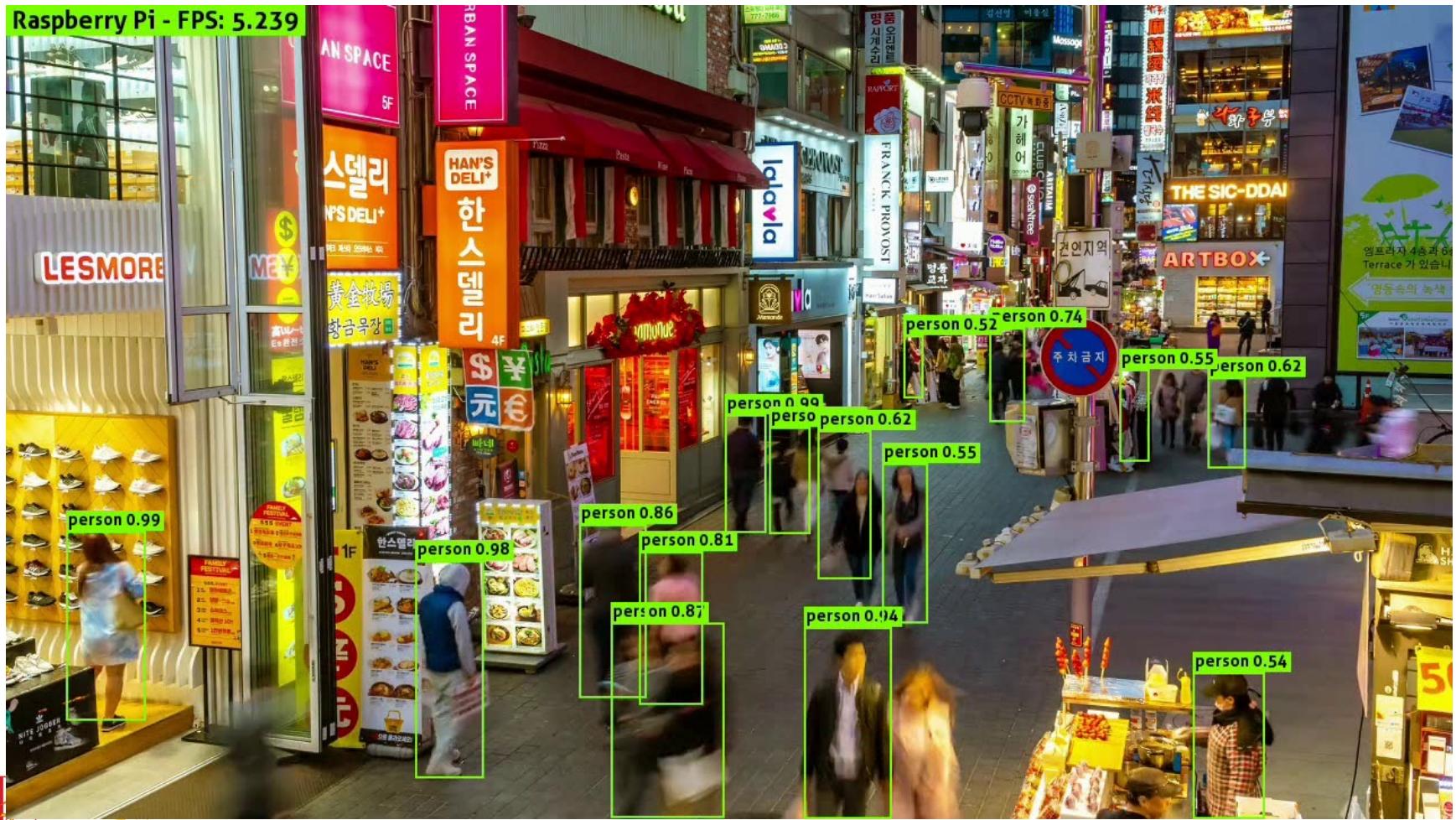
Some applications of object detection

- Smart traffic



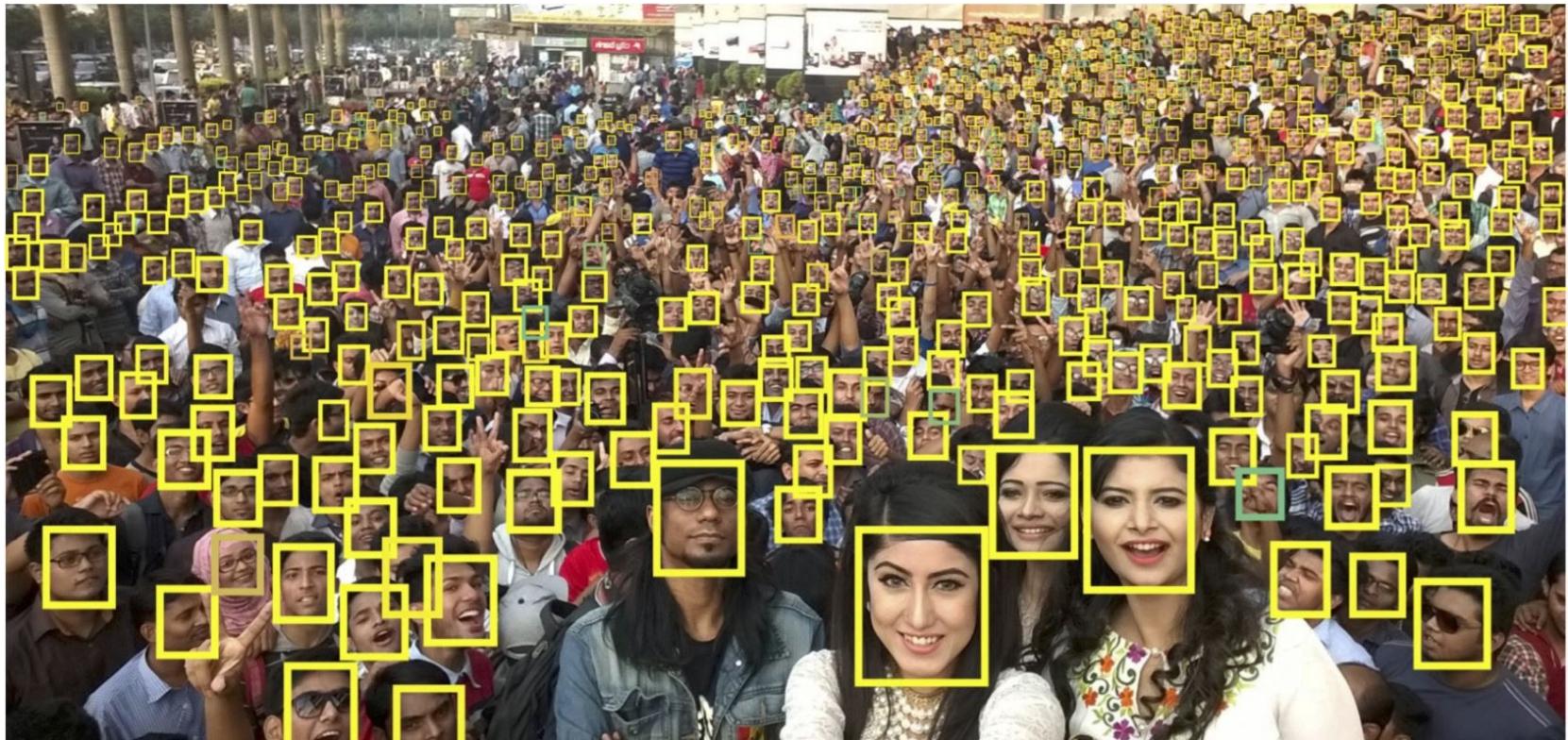
Some applications of object detection

- Human detection



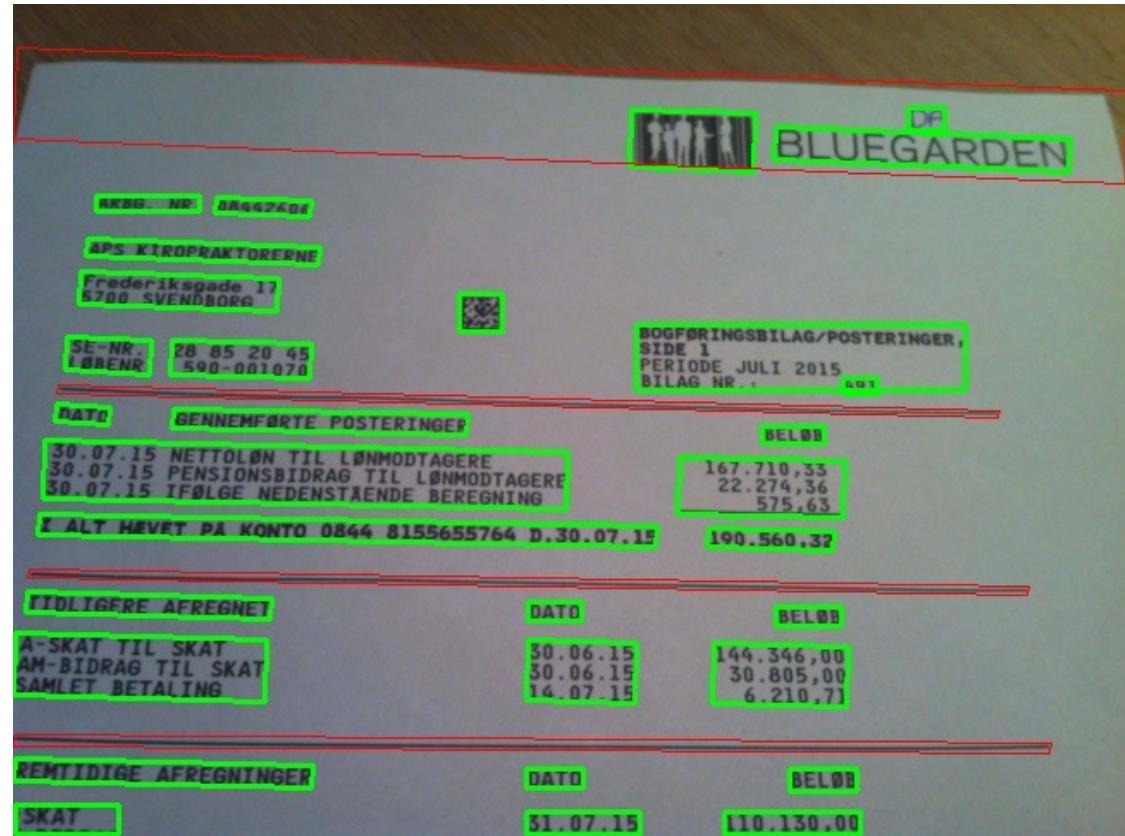
Some applications of object detection

- Face detection



Some applications of object detection

- Text detection



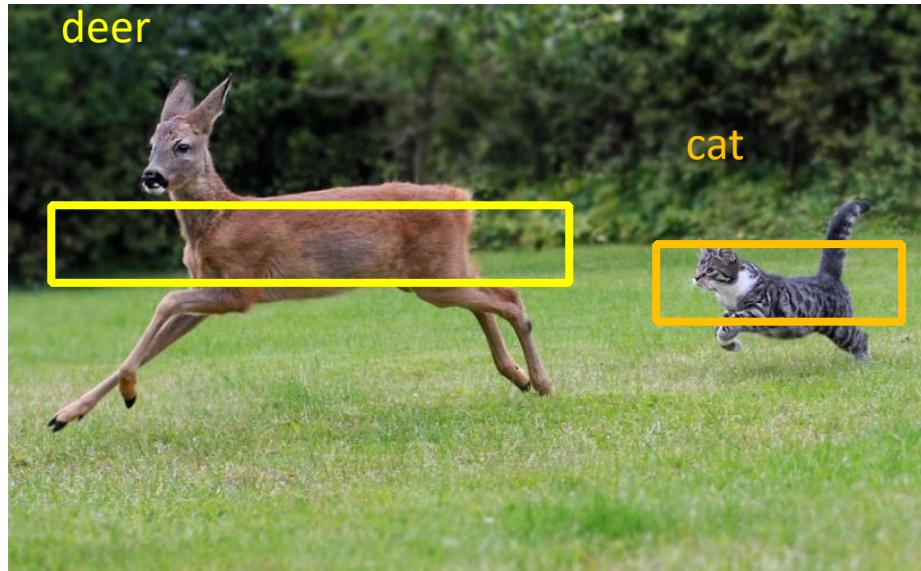
Some applications of object detection

- Automatic robot picking strawberries



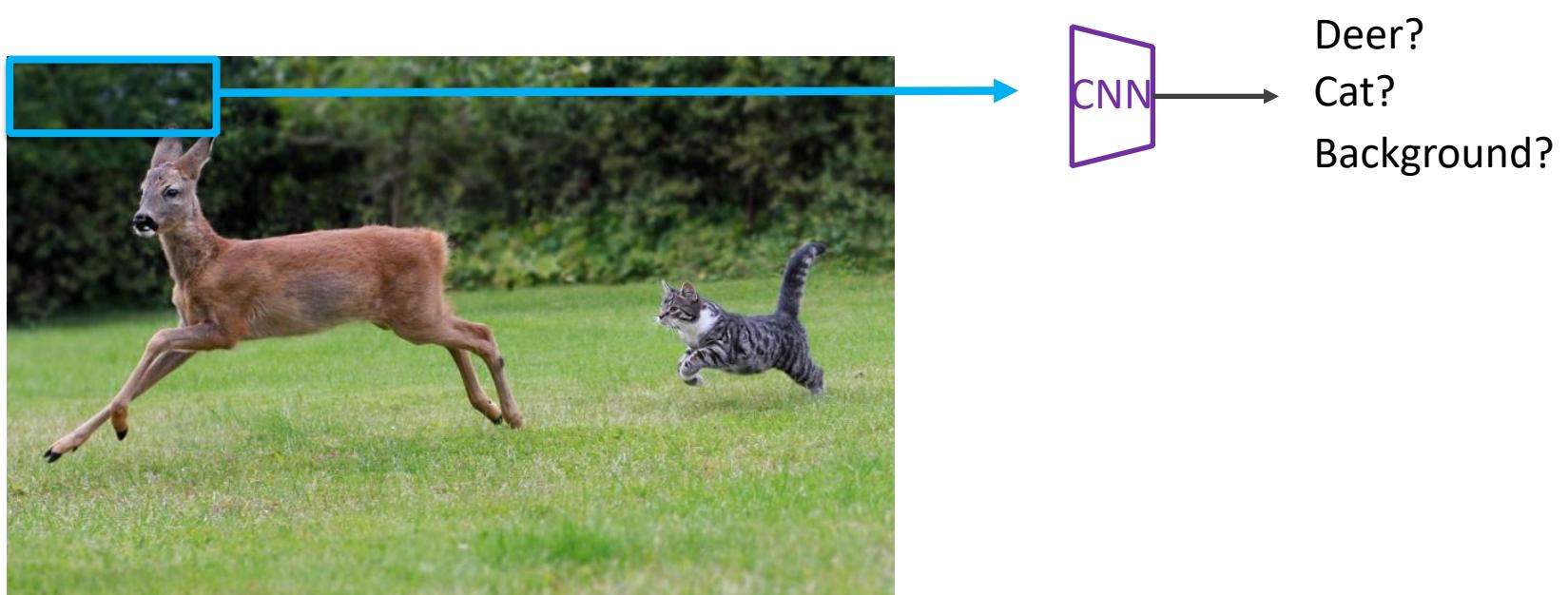
Regional proposal networks (two-stage object detectors)

Sliding window approach



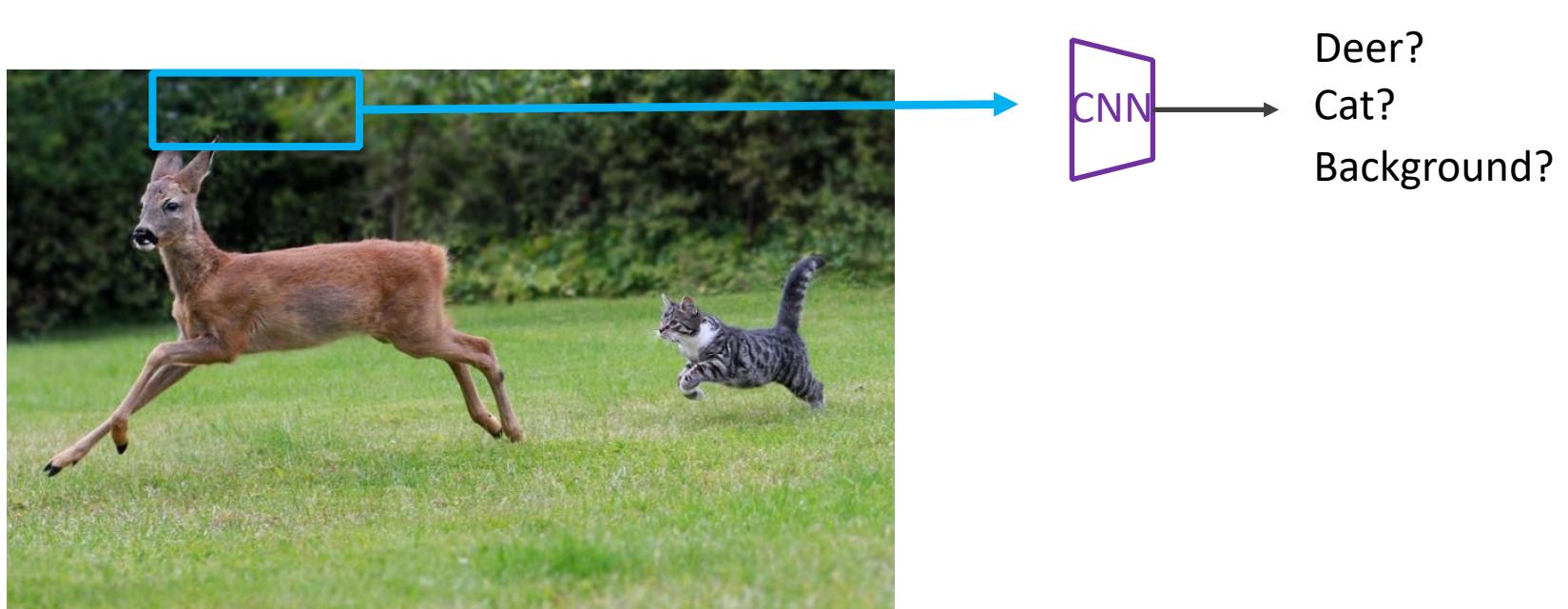
Sliding window

- Scan the window from left to right, top to bottom. At each position, the current window area is classified into different classes plus the background layer.



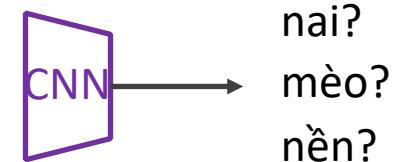
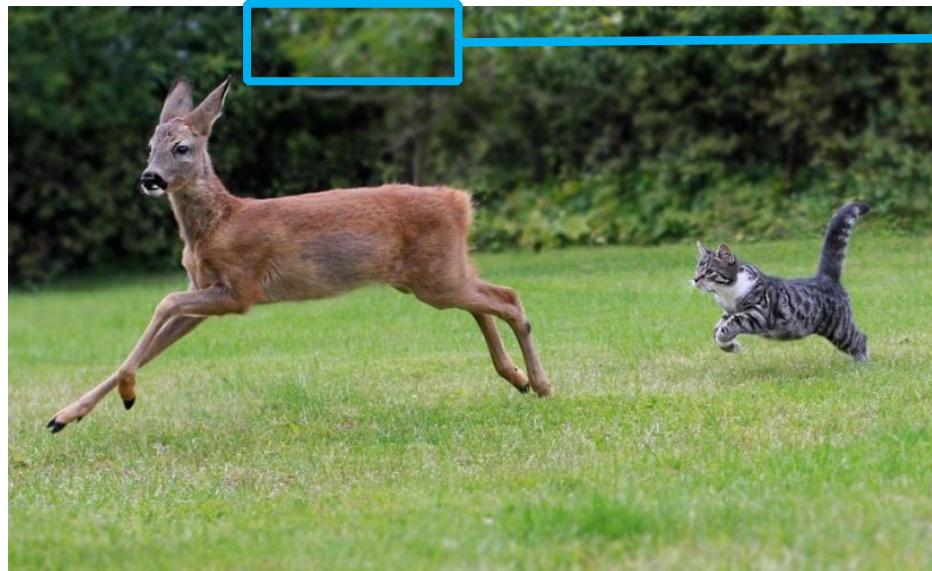
Sliding window (2)

- Scan the window from left to right, top to bottom. At each position, the current window area is classified into different classes plus the background layer.



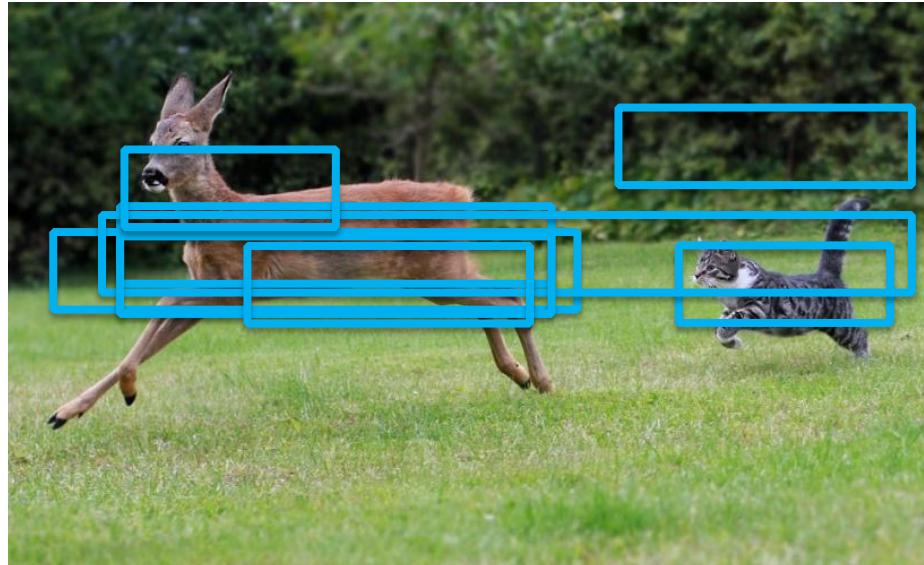
Sliding window (3)

- Scan the window from left to right, top to bottom. At each position, the current window area is classified into different classes plus the background layer.



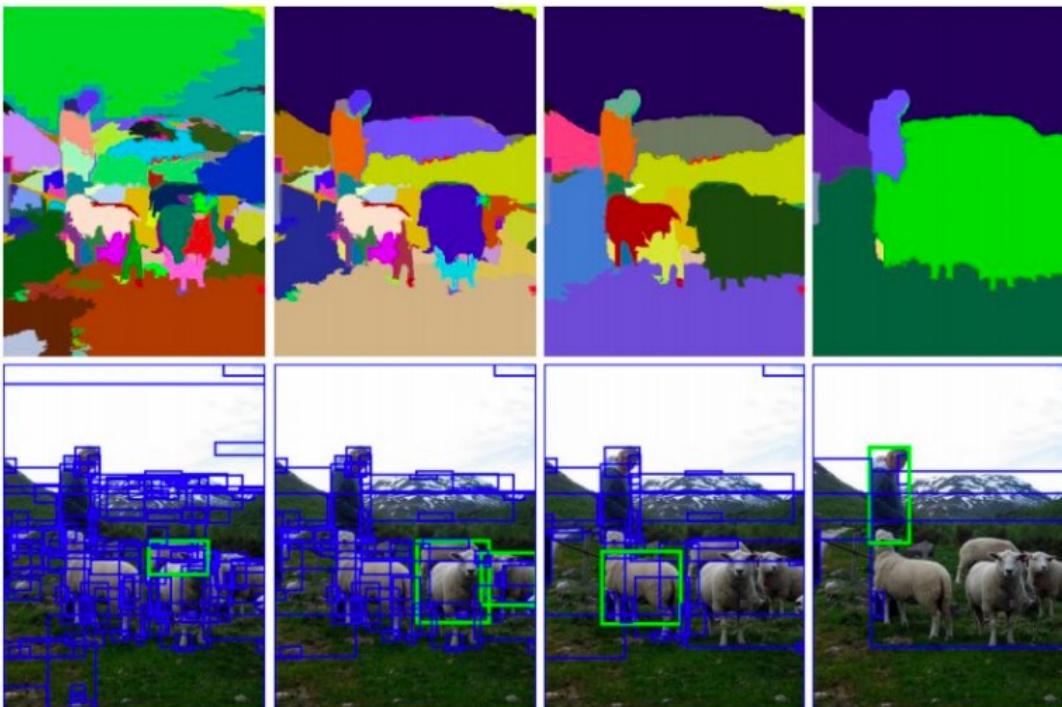
Regional proposal approach

- Instead of scanning all locations (very large numbers!), analyze to only suggest some areas (boxes) with a high probability of containing the object.
- These methods have two stages:
 1. region proposal
 2. process each region to classify and to correct box . coordinates



SS: Selective Search

- Segmentation As Selective Search for Object Recognition. van de Sande et al. ICCV 2011

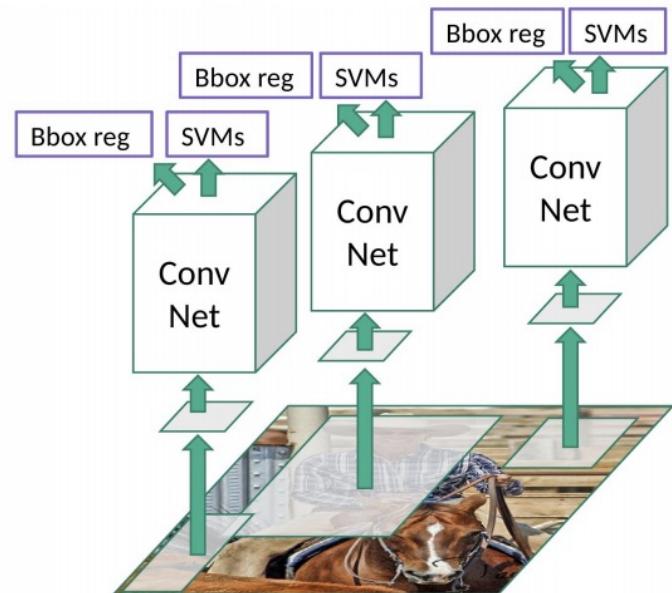


Fast with very high recall
Over-segmenting the image based on intensity of the pixels

1. Add all bounding boxes corresponding to segmented parts to the list of regional proposals
2. Group adjacent segments based on similarity
3. Go to step 1

R-CNN (Region-based ConvNet)

- Suggest some potential regions using other raw algorithms, such as selective search
- Using CNN to extract features of each region and then classify by SVM



R-CNN: *Regions with CNN features*

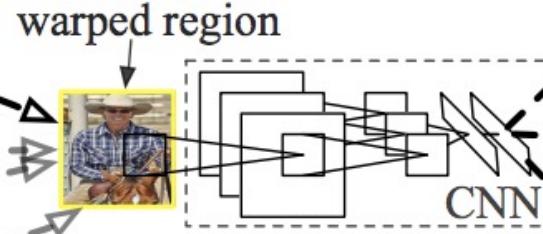


1. Input image



2. Extract region proposals (~2k)

warped region



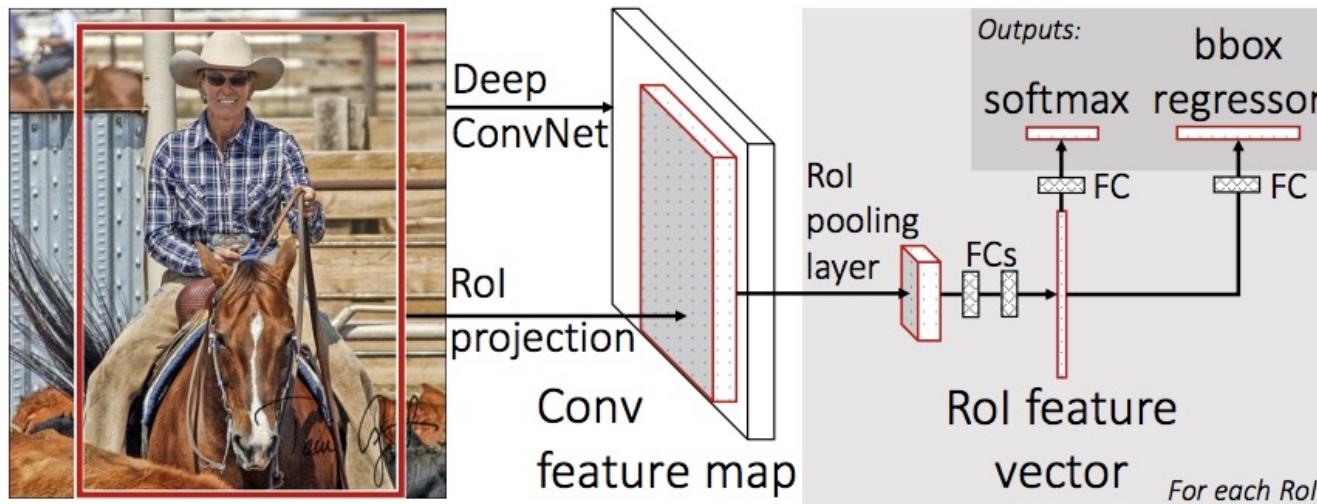
3. Compute CNN features

aeroplane? no.
⋮
person? yes.
⋮
tvmonitor? no.

4. Classify regions

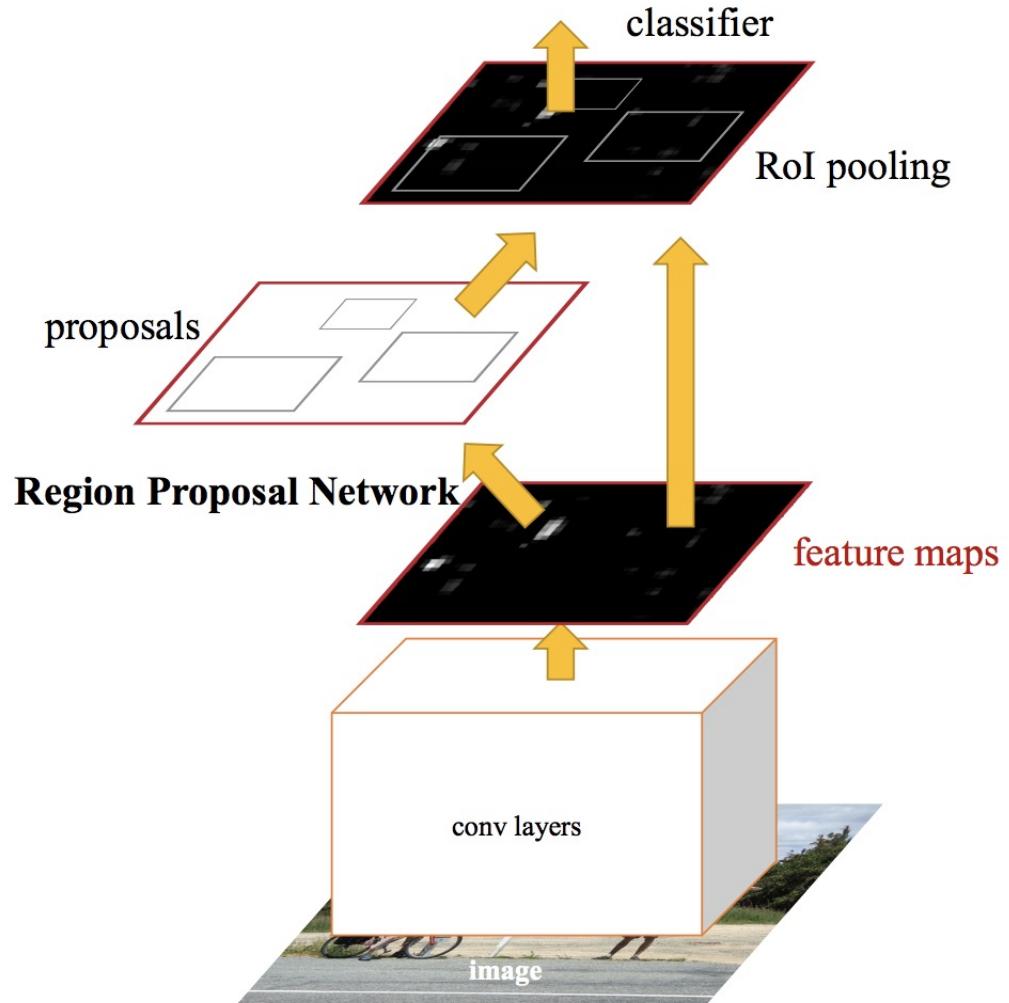
Fast-RCNN

- Push all regions (about 2000) over CNN extractor network at the same time
- Crop the information in the output layer of the CNN instead of cropping the area on the original image like R-CNN
- Push through the classification branch and the bbox regressor branch



Faster-RCNN

- Use a dedicated network to recommend regions instead of selective search
- Also called two-stage object detector method.



None regional proposal networks (one-stage object detectors)

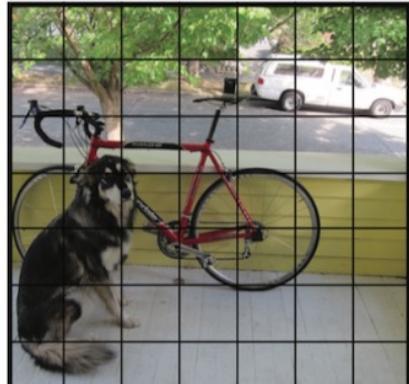
One-Stage Object Detectors

- These networks often suggest a dense box grid on the original image, often with a consistent stride.
- Each of these boxes will be classified and corrected for coordinates (if the box contains objects) using the CNN network
- One-stage networks are generally faster and simpler than two-stage networks, but the accuracy may not be as high.

YOLO- You Only Look Once

$S \times S \times B$ bounding boxes

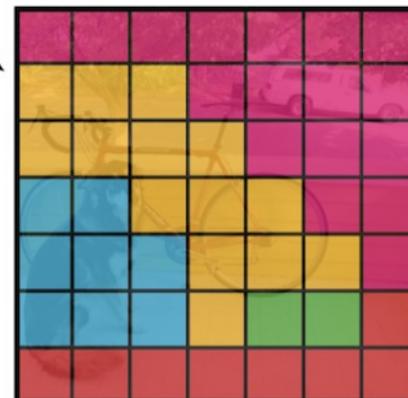
confidence = $Pr(\text{object}) \times \text{IoU}(\text{pred}, \text{truth})$



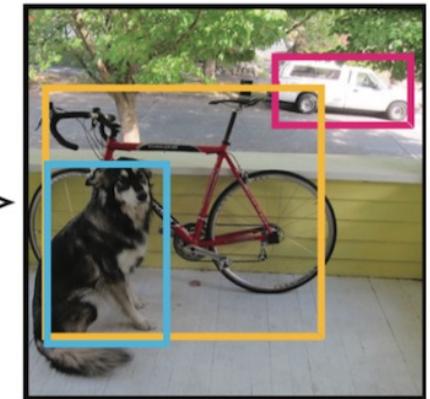
$S \times S$ grid on input



Bounding boxes + confidence



Class probability map

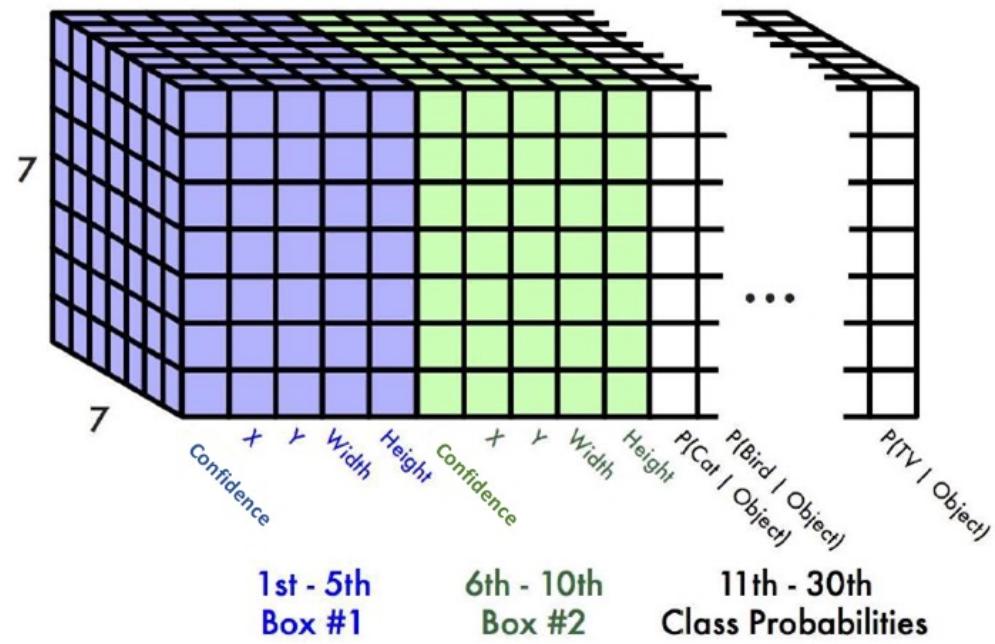
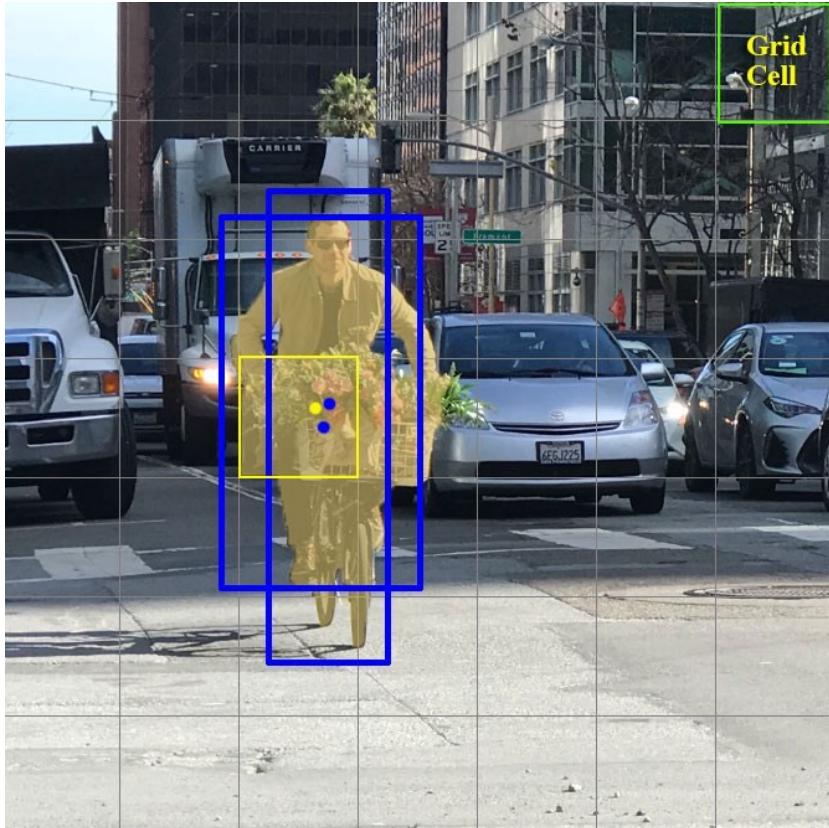


Final detections

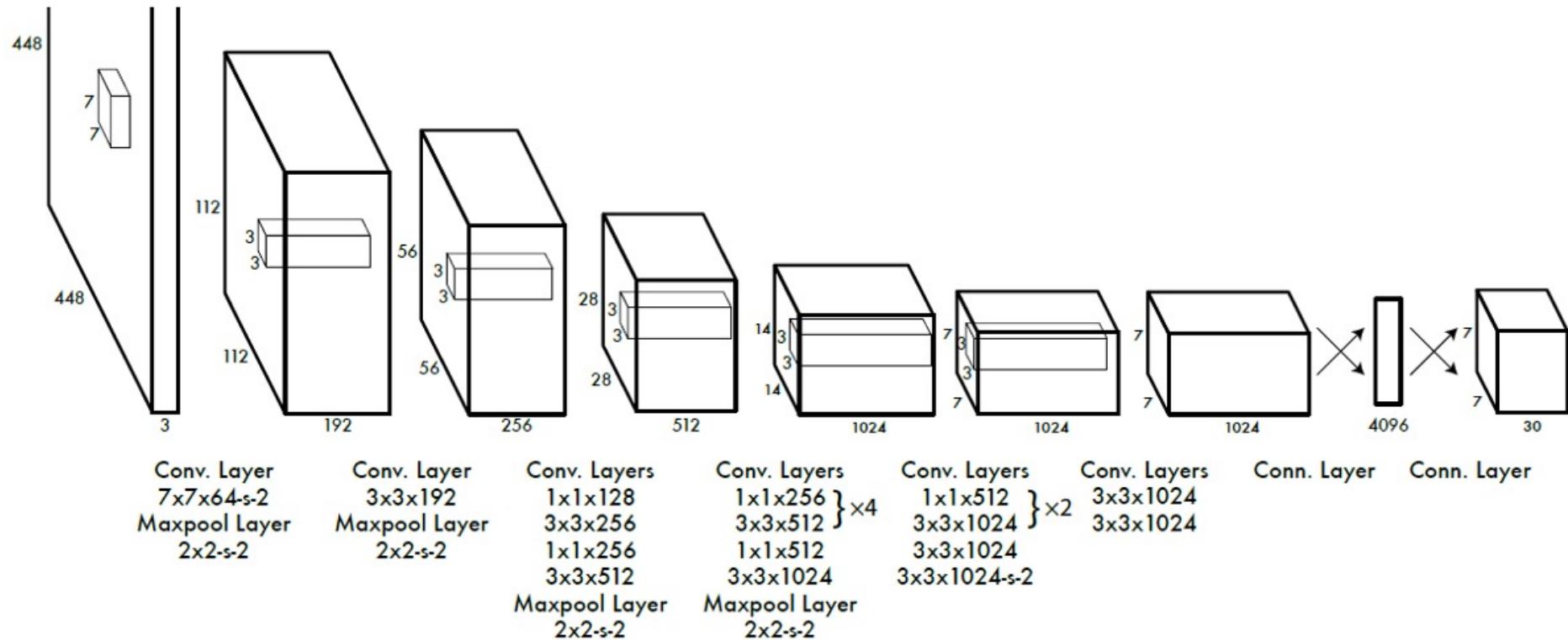
$Pr(\text{Class}_i | \text{object})$

$$\Pr(\text{Class}_i | \text{Object}) * \Pr(\text{Object}) * \text{IOU}_{\text{pred}}^{\text{truth}} = \Pr(\text{Class}_i) * \text{IOU}_{\text{pred}}^{\text{truth}}$$

YOLO- You Only Look Once



YOLO- You Only Look Once

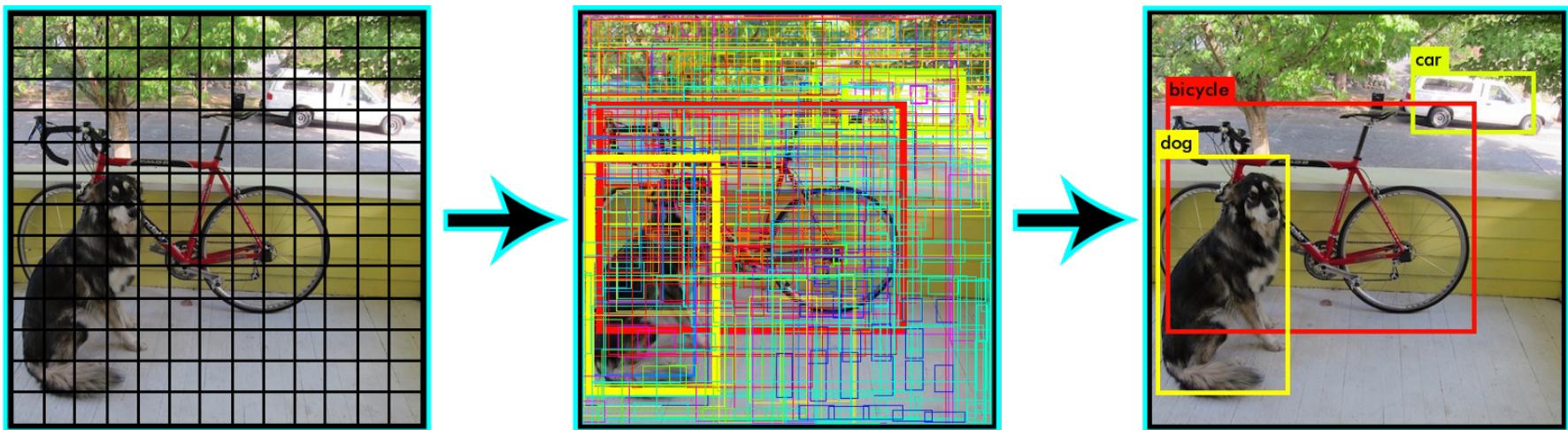


YOLO- You Only Look Once

$$\lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} \left[(x_i - \hat{x}_i)^2 + (y_i - \hat{y}_i)^2 \right] \quad \begin{array}{l} \text{1 when there is object, 0 when there is no object} \\ \text{Bounding Box Location (x, y) when there is object} \end{array}$$
$$+ \lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} \left[(\sqrt{w_i} - \sqrt{\hat{w}_i})^2 + (\sqrt{h_i} - \sqrt{\hat{h}_i})^2 \right] \quad \begin{array}{l} \text{Bounding Box size (w, h)} \\ \text{when there is object} \end{array}$$
$$+ \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} \left(C_i - \hat{C}_i \right)^2 \quad \text{Confidence when there is object}$$
$$+ \lambda_{\text{noobj}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{noobj}} \left(C_i - \hat{C}_i \right)^2 \quad \begin{array}{l} \text{1 when there is no object, 0 when there is object} \\ \text{Confidence when there is no object} \end{array}$$
$$+ \sum_{i=0}^{S^2} \mathbb{1}_i^{\text{obj}} \sum_{c \in \text{classes}} (p_i(c) - \hat{p}_i(c))^2 \quad \text{Class probabilities when there is object}$$

YOLO- You Only Look Once

- Non-maximal suppression: gather the boxes to give the final result

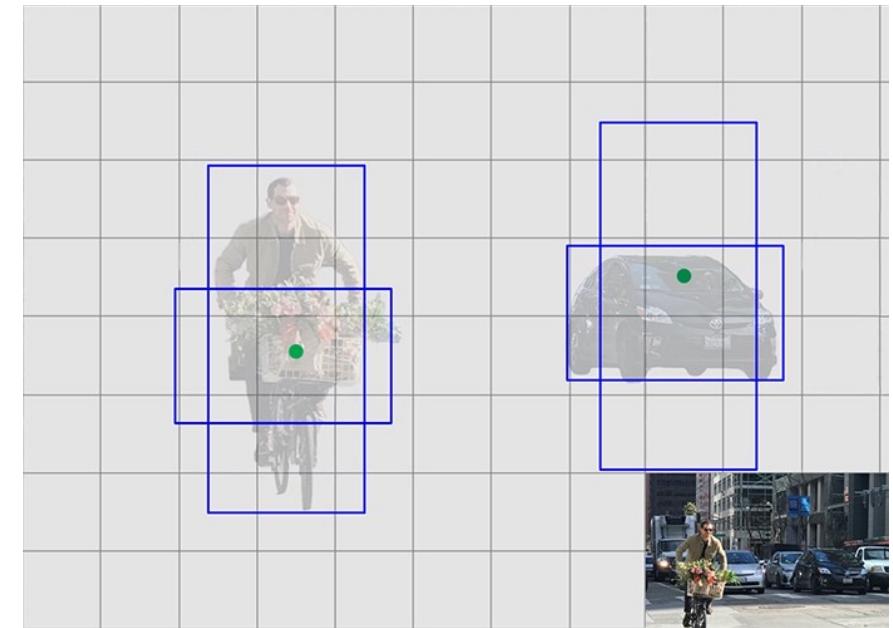
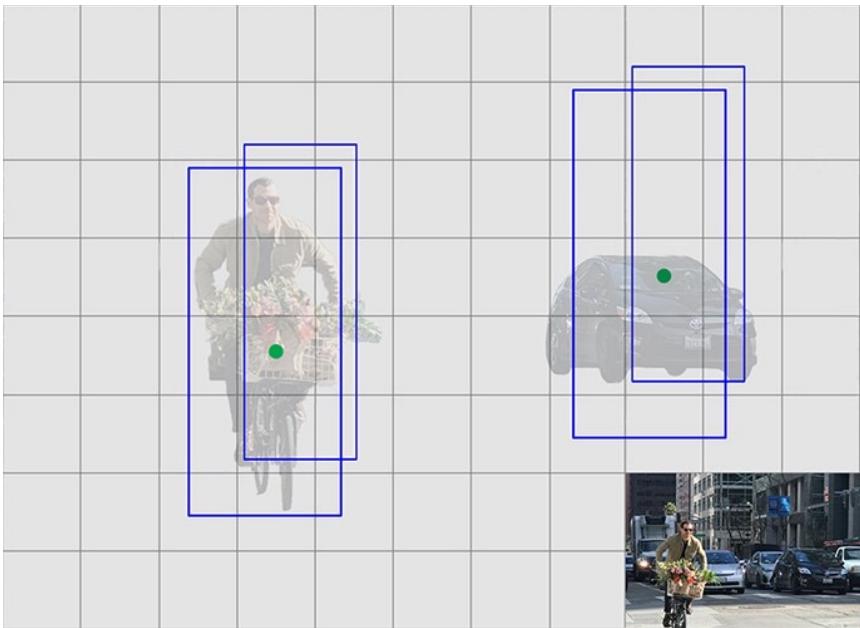


YOLO can make duplicate detections for the same object.

YOLO applies non-maximal suppression to remove duplications with lower confidence

- Sort the predictions by the confidence scores.
- Start from the top scores, ignore any current prediction if we find any previous predictions that have the same class and $\text{IoU} > 0.5$ with the current prediction.
- Repeat step 2 until all predictions are checked.

YOLO v2



Initially, YOLO makes arbitrary guesses on the boundary boxes.

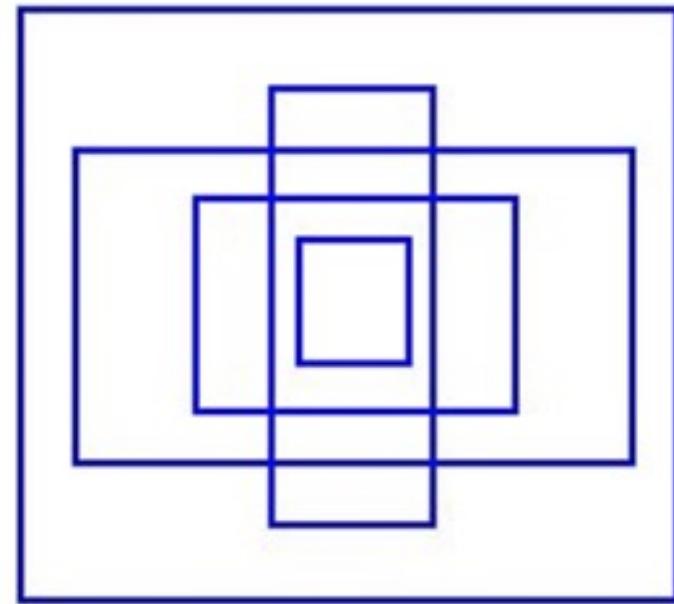
In the real-life domain, the boundary boxes are not arbitrary.

Idea

- We can create 5 **anchor** boxes with the intended shapes.
- We predict offsets to each of the anchor boxes above.
- If we **constrain** the offset values, we can maintain the diversity of the predictions and have each prediction focuses on a specific shape

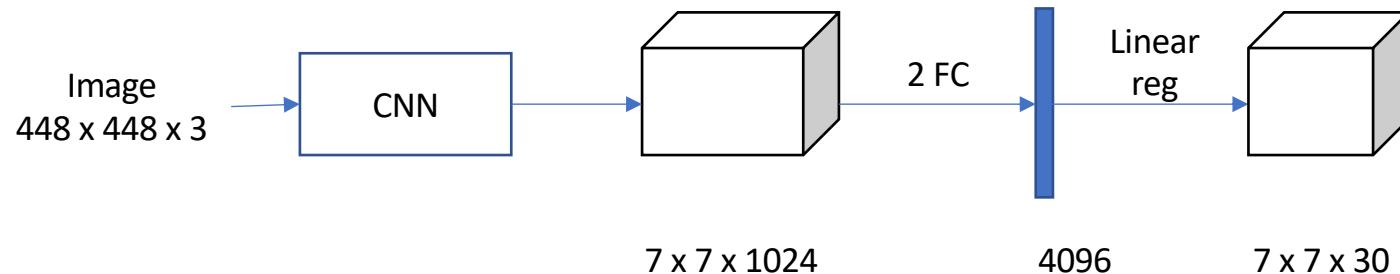
YOLO v2

- Each cell has 5 anchor boxes. For each network anchor will give the following information:
 - box offset: 4 real numbers in the range [0, 1]
 - Confidence that the box is likely to contain the object (objectness score).
 - Probability distribution of objects in that box for different classes (class scores).
- In total, each cell has an output number: $5 * (4 + 1 + 20) = 125$ real numbers

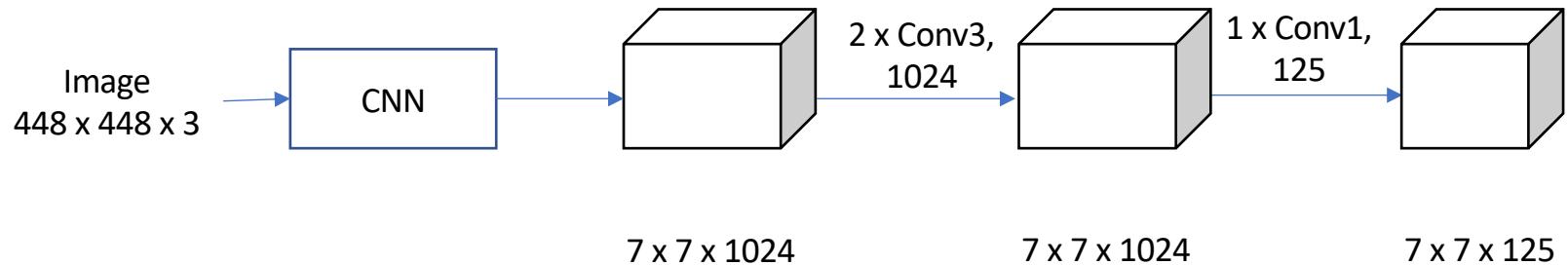


5 anchor box

YOLO v2



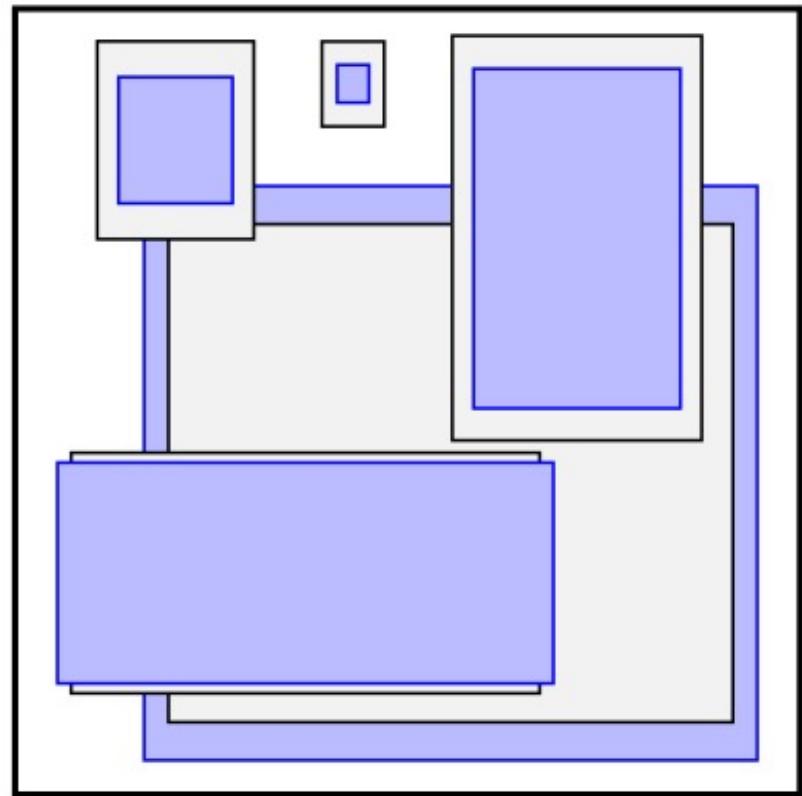
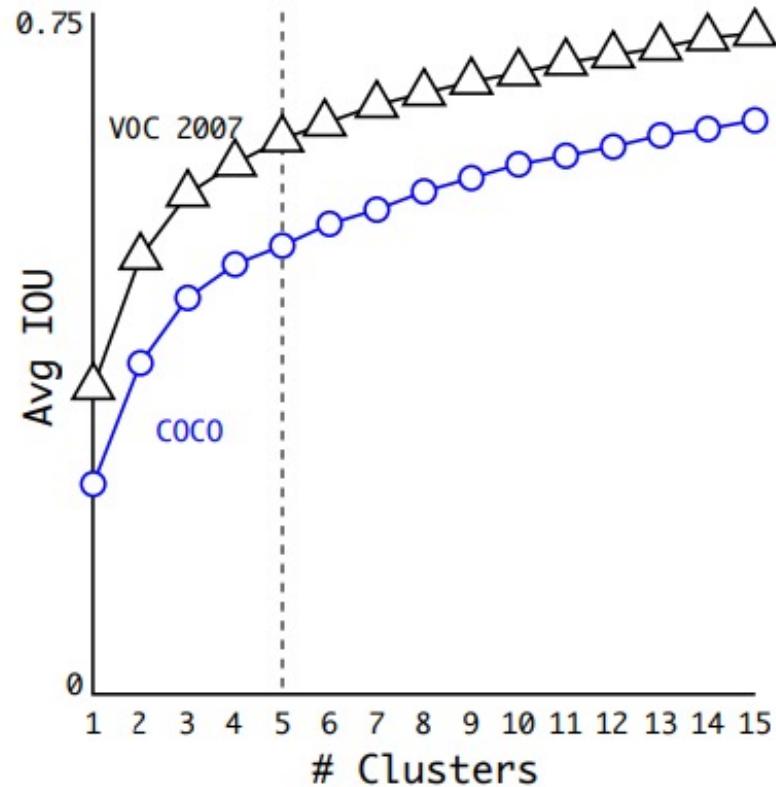
YOLO v1



YOLO v2

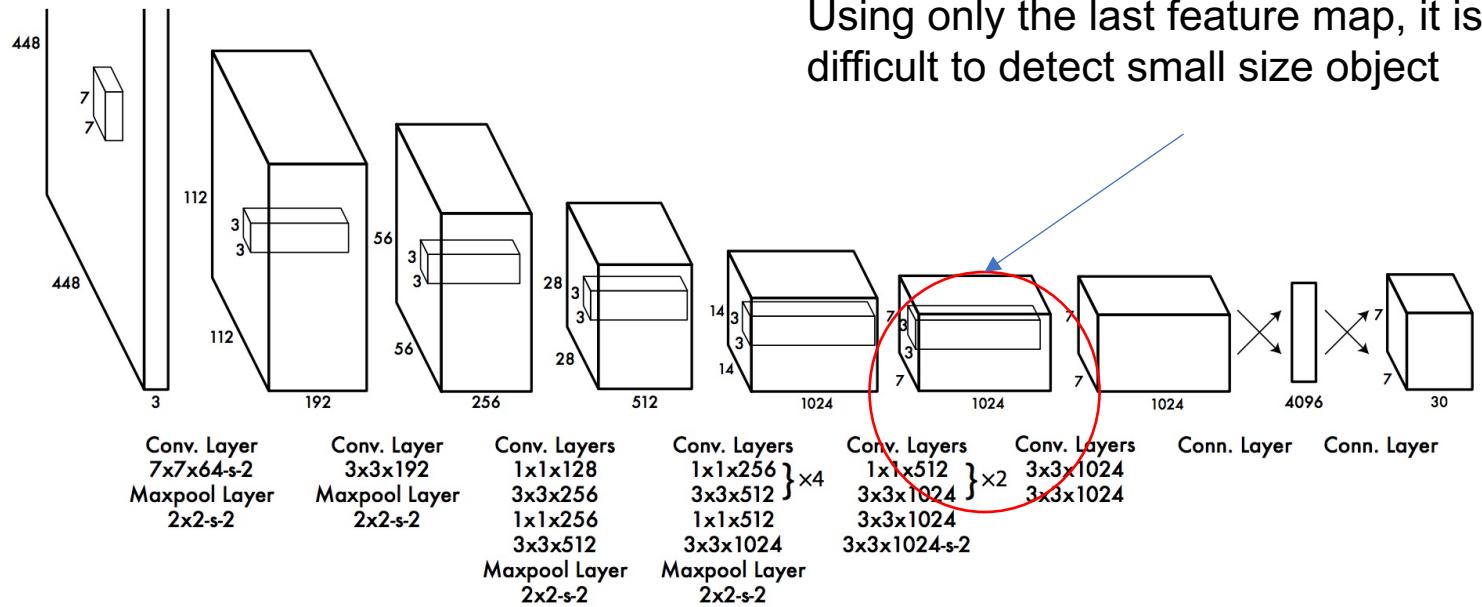
YOLO v2

- Determine the default size of the anchors by applying k-means on the box set of labeled objects in the training set.



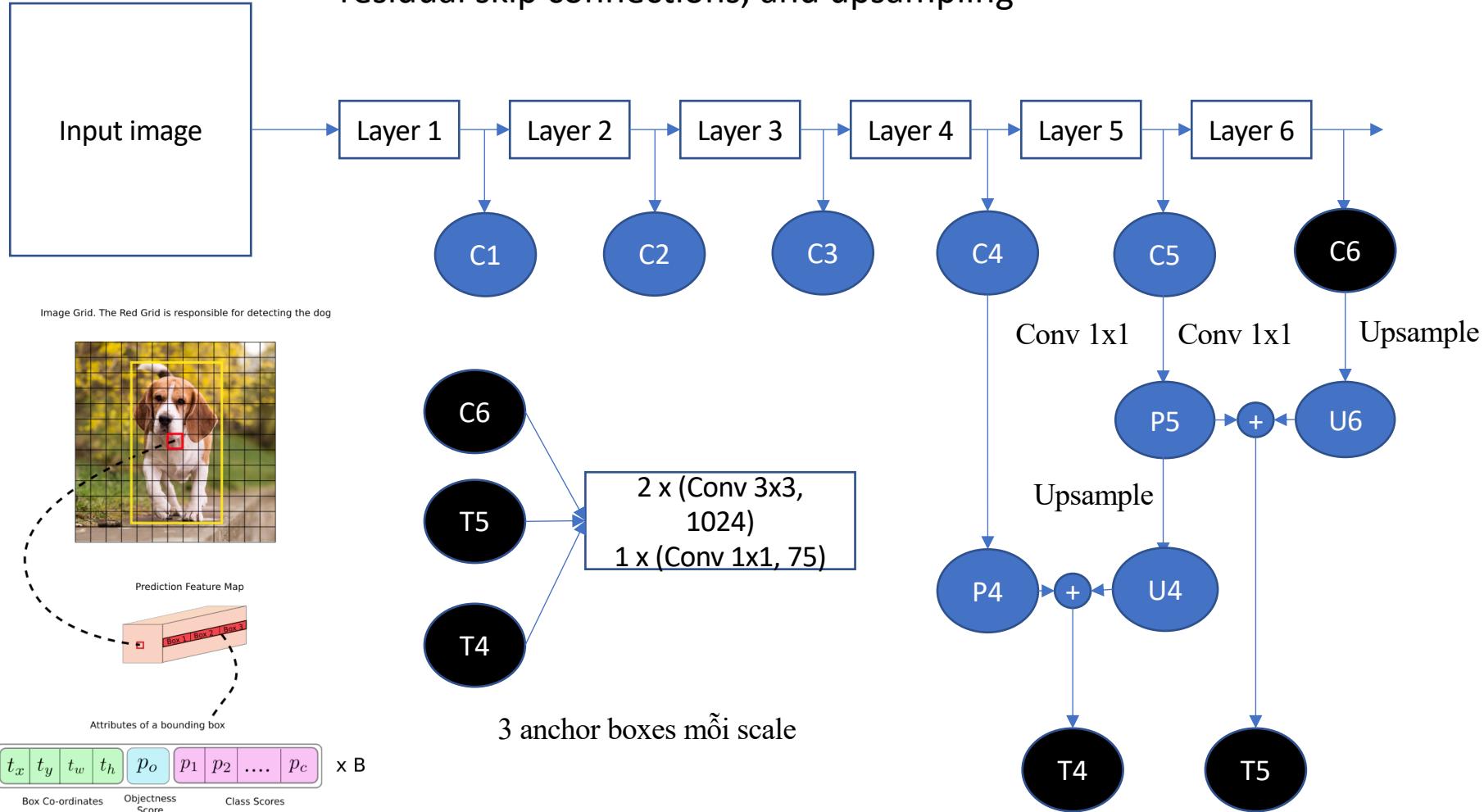
YOLO v2

- Cons of YOLO v1 and v2:



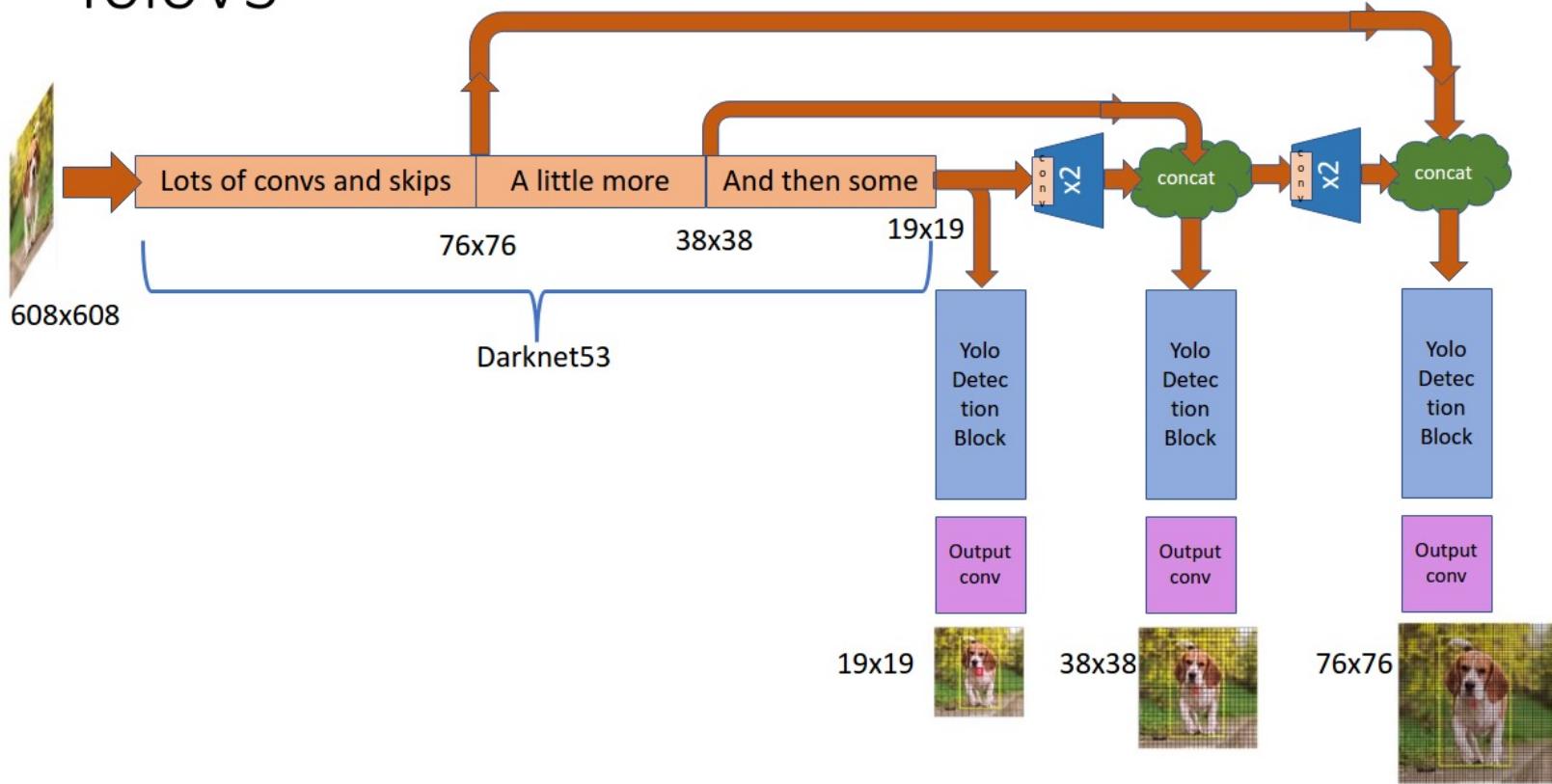
YOLO v3

residual skip connections, and upsampling



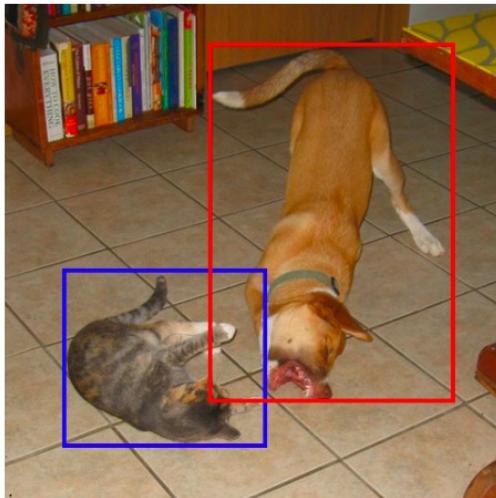
YOLO v3

YoloV3

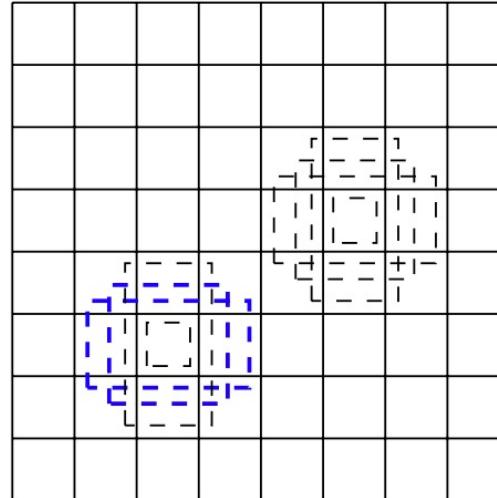


SSD: Single Shot Detector

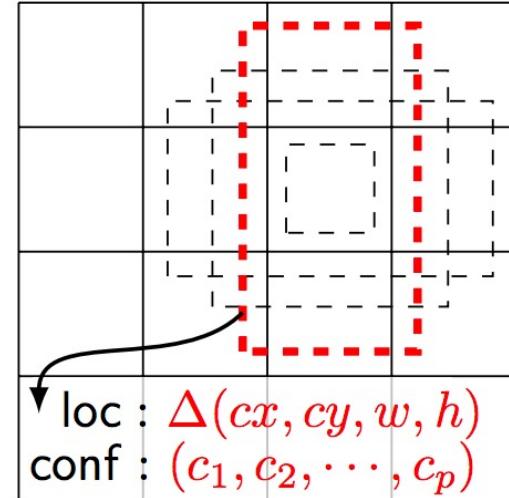
- Similar to YOLO but the box grid is denser, there are many grids with different box sizes
- The backbone network architecture is different from YOLO
- Data augmentation + Hard negative mining



(a) Image with GT boxes



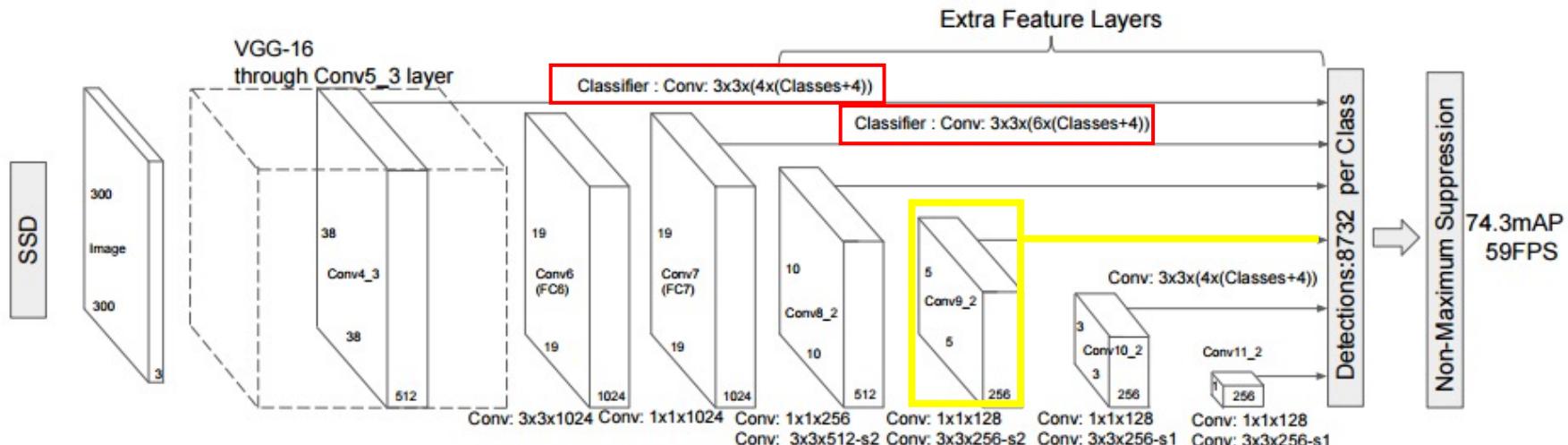
(b) 8×8 feature map



(c) 4×4 feature map

SSD: Single Shot Detector

- Network backbone: VGG-16
- Add additional convolutional layers behind the layers of the backbone network
- Detect objects at different levels in the network (Multi-scale)

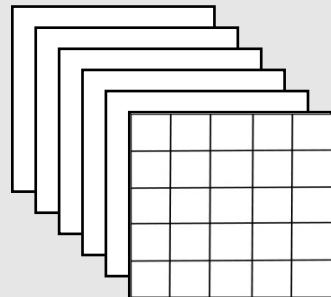


SSD has 8732 bounding boxes which is more than that of YOLO.

Liu et al. ECCV 2016.

SSD: Single Shot Detector

Feature map
đầu vào

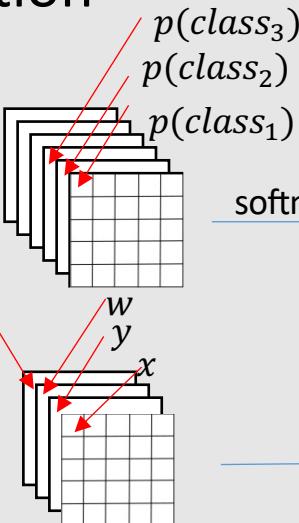


@5x5x256
Feature map

Prediction

5x5x
21classes

5x5x
4 box offset

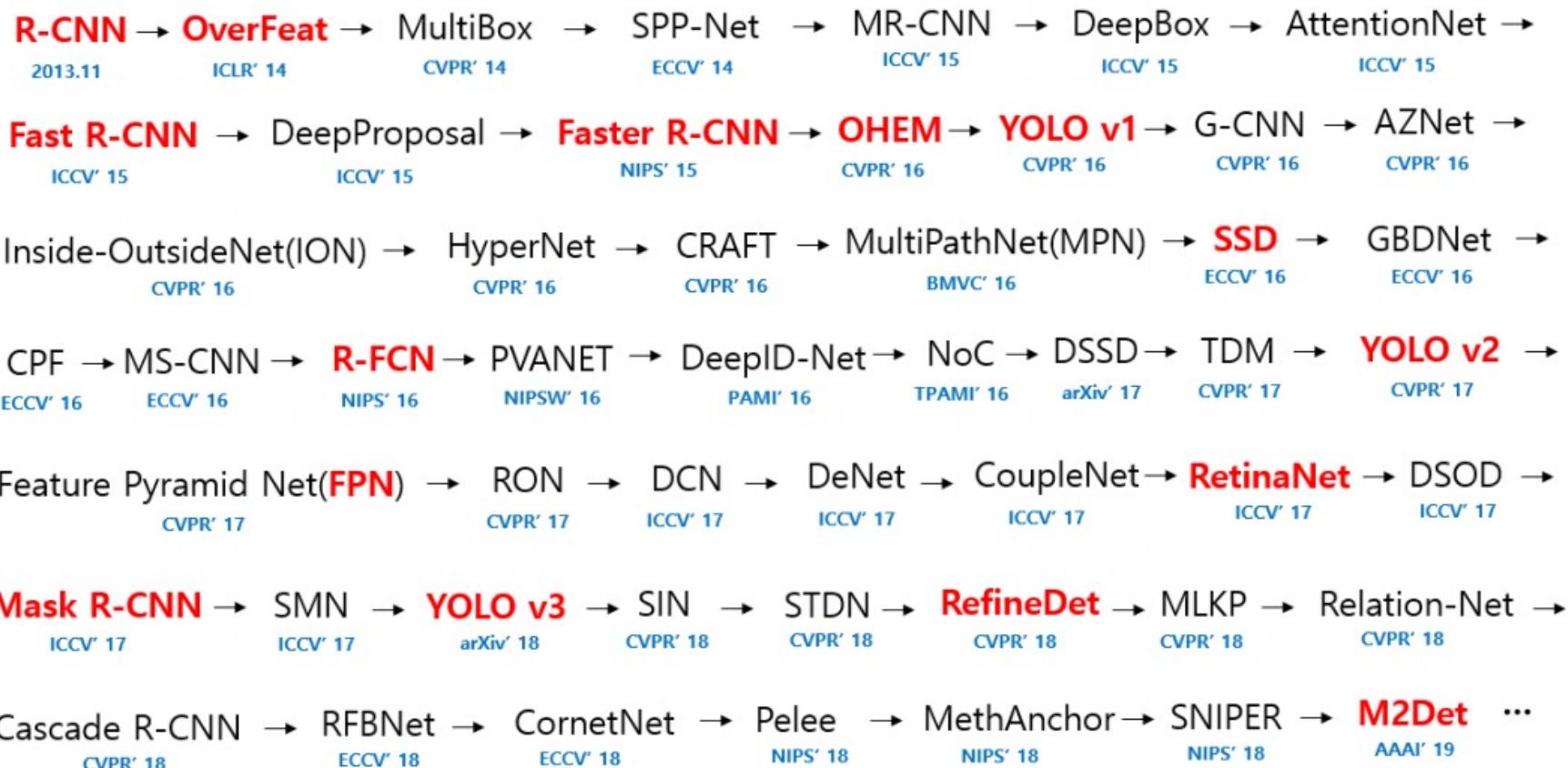


$p(\text{class}_3)$
 $p(\text{class}_2)$
 $p(\text{class}_1)$

softmax

Objective funcs

$$L(x, c, l, g) = \frac{1}{N} (L_{conf}(x, c) + \alpha L_{loc}(x, l, g))$$



One-stage vs two-stage

Faster and
simpler

one-stage object detector

(dense sampling of object
locations, scales, and aspect ratios)

YOLO YOLO-v2 YOLO-v3

SSD

DSSD

MDCN

SqueezeNet

RetinaNet

RedefineDet

CornetNet

CenterNet

EfficientDet



More
accurate

two-stage object detector

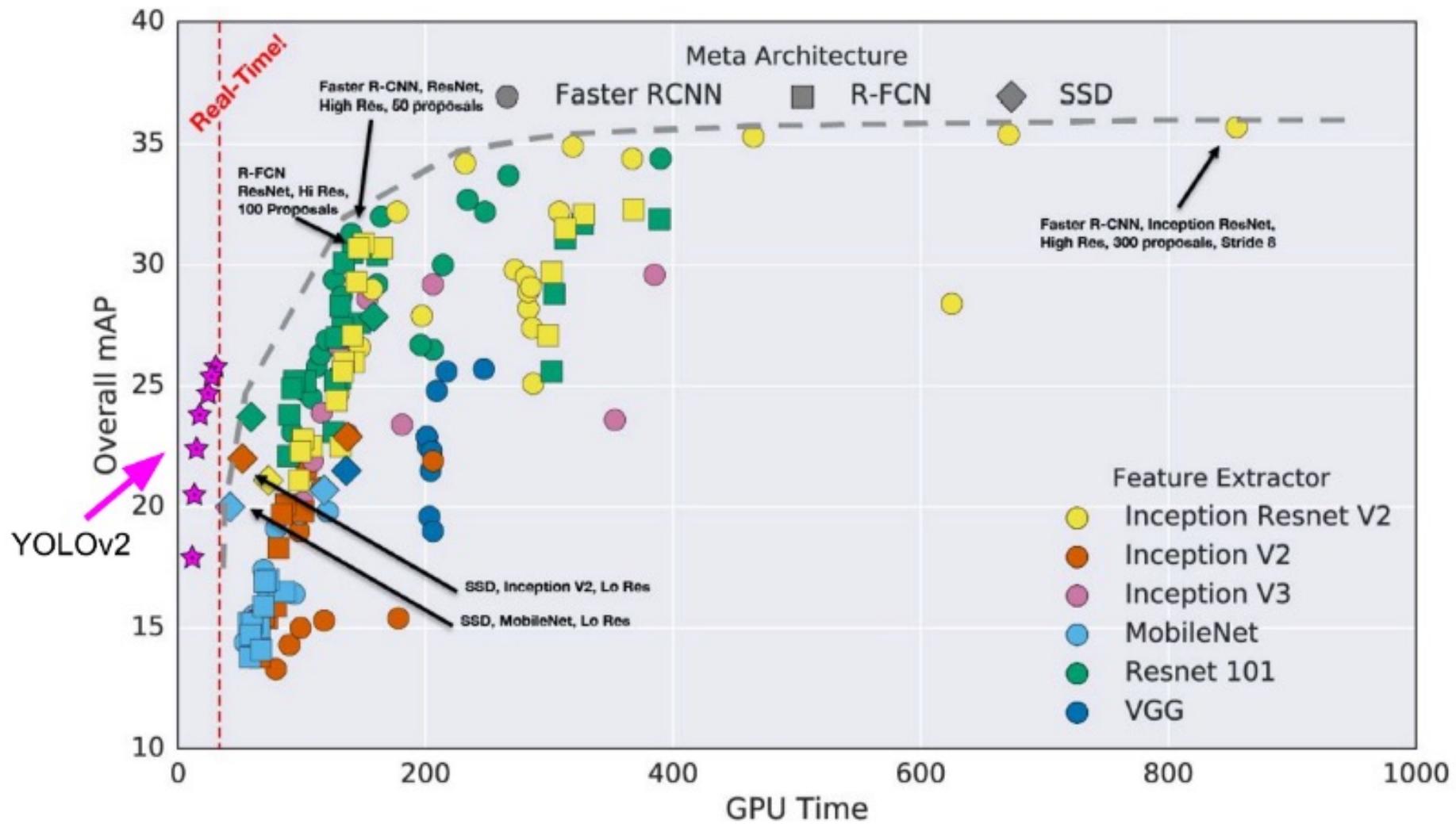
(proposal-driven
mechanism)

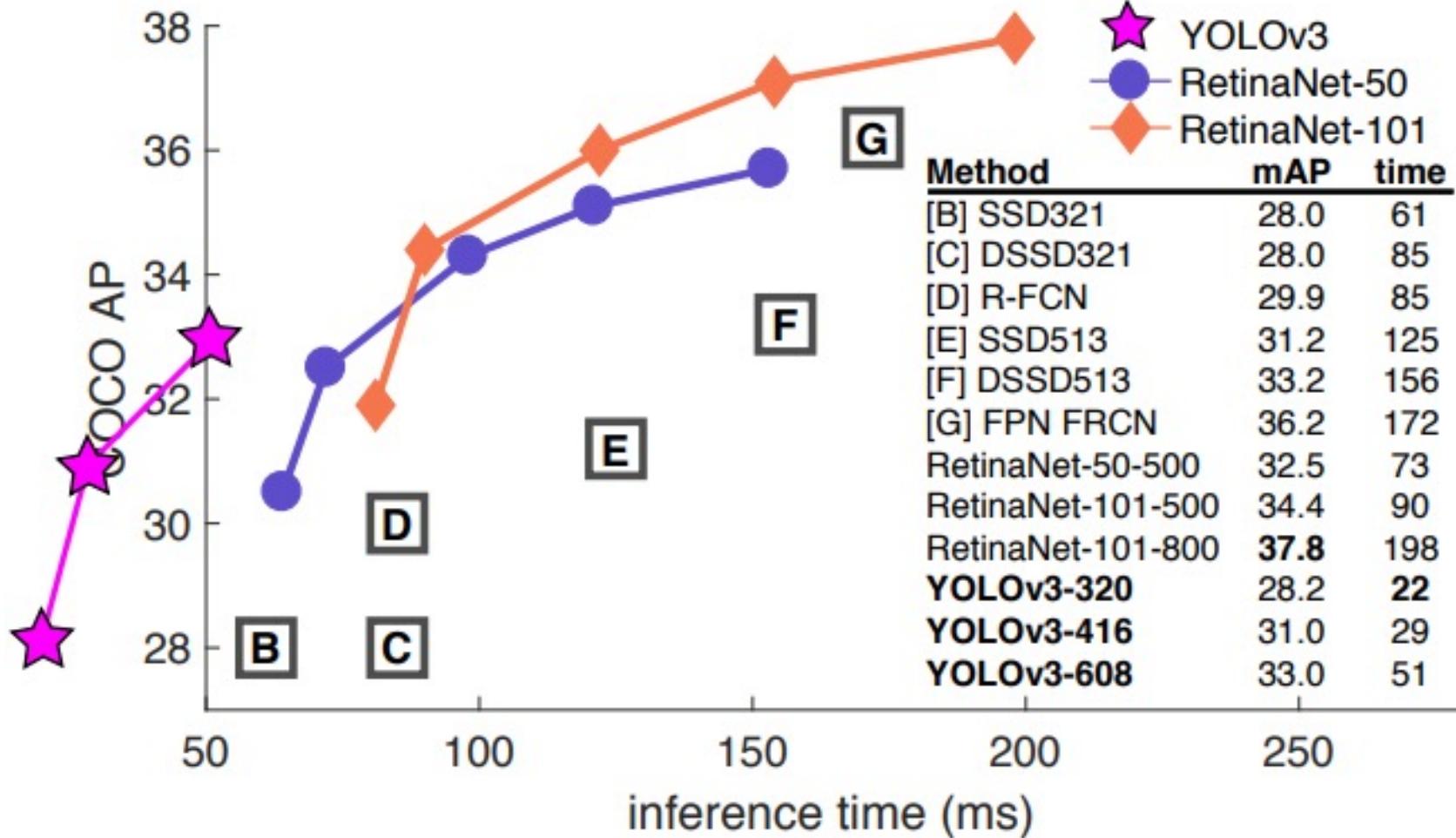
R-CNN

Fast R-CNN

Faster R-
CNN

Feature Pyramid Network
(FPN)
Mask R-
CNN







25
YEARS ANNIVERSARY
SOICT

VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG
SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

**Thank you
for your
attention!**

