# تشخیص اشیاء
# Object Detection

Alireza AkhavanPour

Akhavanpour.ir
CLASS.VISION

# So far: Image Classification



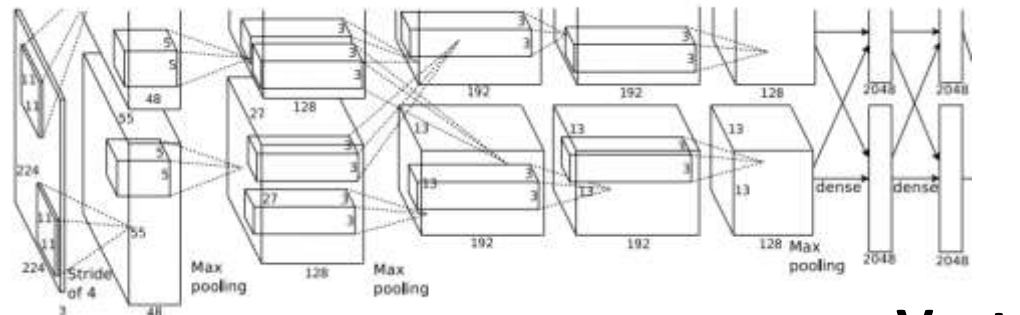This image is CC0 public domain

Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

**Vector:**
4096

**Fully-Connected**:
4096 to 1000

**Class Scores**
Cat: 0.9
Dog: 0.05
Car: 0.01
...

# Where is the object?



Let's code…

**1-simple-regression-train.ipynb**

# Where is the object?



Let's code…

**2-simple-regression-inference.ipynb**

**Object Detection**

تشخیص اشیاء

**Alireza Akhavanpour**

علیرضا اخوان پور

CLASS.
VISION

# Where is the object?
# What about class names?



Let's code...

**3-object-classification-and-localization.ipynb**

**Object Detection**
**Alireza Akhavanpour**

تشخیص اشیاء
علیرضا اخوان پور

CLASS.
VISION

# Where is the object?
# What about class names?



Let's code…

**4-object-classification-and-localization-inference.ipynb**

**Object Detection**

**Alireza Akhavanpour**

تشخیص اشیاء

علیرضا اخوان پور

CLASS.
VISION

# Computer Vision Tasks

| Classification | Semantic Segmentation | Object Detection | Instance Segmentation |
|---|---|---|---|



**CAT**

**GRASS**, **CAT**, **TREE**, **SKY**

**DOG**, **DOG**, **CAT**

**DOG**, **DOG**, **CAT**

No spatial extent    No objects, just pixels

Multiple Objects

This image is CC0 public domain
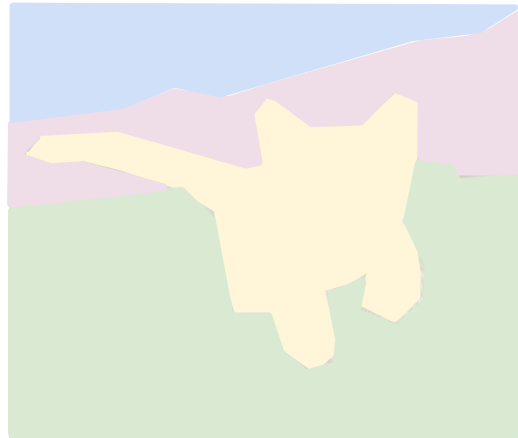
# Today: Object Detection



**Classification**

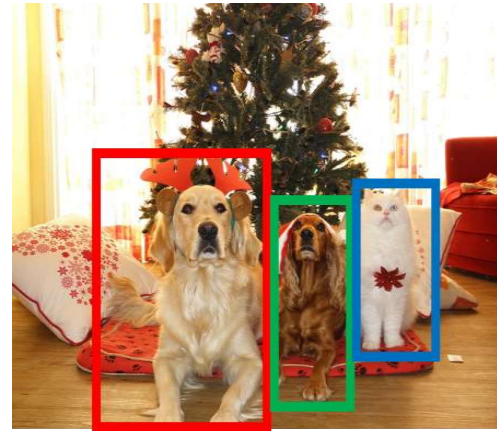CAT

No spatial extent

**Semantic Segmentation**

GRASS, CAT, TREE, SKY

No objects, just pixels

**Object Detection**

DOG, DOG, CAT

**Instance Segmentation**
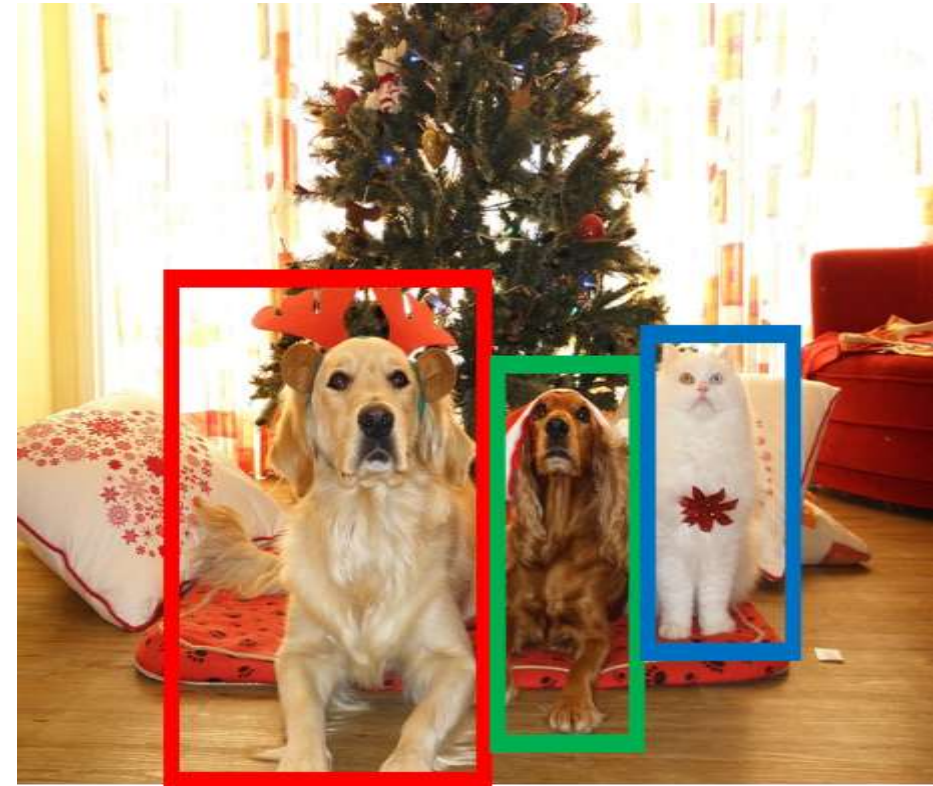
DOG, DOG, CAT

Multiple Objects

Object Detection

Alireza Akhavanpour

تشخیص اشیاء

علیرضا اخوان پور

CLASS. vision

# Object Detection: Task Definition
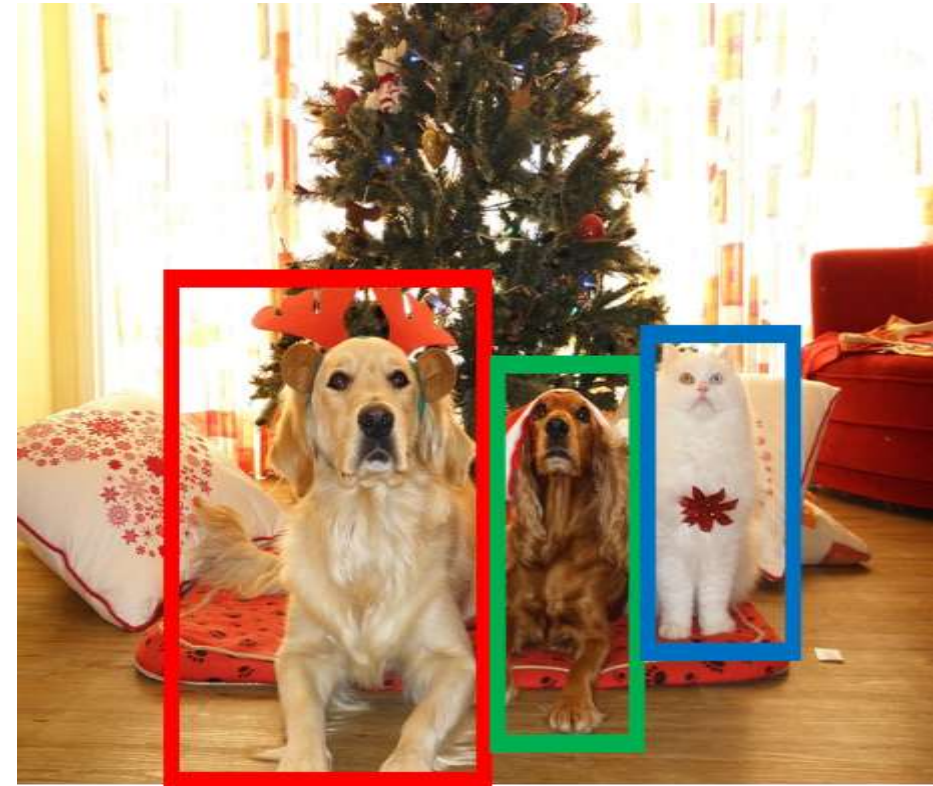
**Input**: Single RGB Image

**Output**: A <u>set</u> of detected objects;
For each object predict:

1. Category label (from fixed, known set of categories)
2. Bounding box (four numbers: x, y, width, height)

# Object Detection: Challenges

- **Multiple outputs**:
Need to output variable numbers of objects per image

- **Multiple types of output**:
Need to predict "what" (category label) as well as "where" (bounding box)

- **Large images**:
Classification works at 224x224; need higher resolution for detection, often ~800x600
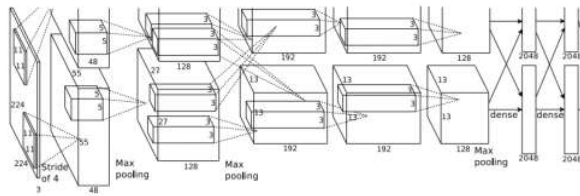
**Object Detection**

**Alireza Akhavanpour**

تشخیص اشیاء

علیرضا اخوان پور

CLASS.
VISION

# Detecting a single object



This image is CC0 public domain

**Vector:**
4096

# Detecting a single object

**Correct label:** Cat

**Class Scores**
Cat: 0.9
Dog: 0.05
Car: 0.01
...

**Softmax Loss**

**Fully Connected:**
4096 to 1000



This image is CC0 public domain

**Vector:**
4096

# Detecting a single object

**Vector:** 4096

"What"

**Fully Connected:** 4096 to 1000

**Class Scores**
Cat: 0.9
Dog: 0.05
Car: 0.01
...

**Correct label:** Cat

**Softmax Loss**

# Detecting a single object



"What"

Fully Connected: 4096 to 1000

**Class Scores**
Cat: 0.9
Dog: 0.05
Car: 0.01
...

Correct label: Cat

**Softmax Loss**

**Multitask Loss**

**Weighted Sum**

**Loss**

Vector: 4096

Fully Connected: 4096 to 4

**Box Coordinates**
(x, y, w, h)

"Where"

Correct box: (x', y', w', h')

**L2 Loss**

This image is CCD public domain

# Detecting a single object

Often pretrained on ImageNet (Transfer learning)



This image is CCD public domain

"What"

Fully Connected: 4096 to 1000

**Class Scores**
Cat: 0.9
Dog: 0.05
Car: 0.01
...

**Correct label:** Cat

**Softmax Loss**

Vector: 4096

Fully Connected: 4096 to 4

**Box Coordinates** (x, y, w, h)
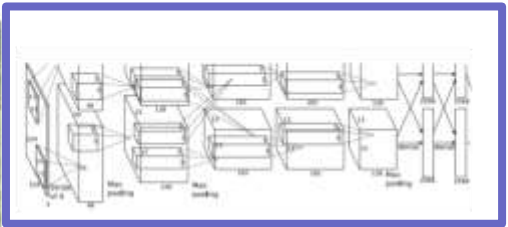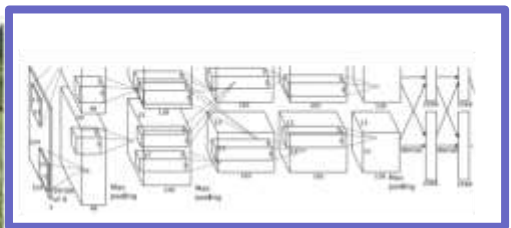
"Where"

**Correct box:** (x', y', w', h')

**L2 Loss**

**Multitask Loss**

**Weighted Sum**

**Loss**

**Object Detection**
**Alireza Akhavanpour**

تشخیص اشیاء
علیرضا اخوان پور

CLASS.
VISION

# Detecting a single object

Often pretrained on
ImageNet (Transfer learning)



This image is CCD public domain

Treat localization as a
regression problem!

"What"

Fully Connected:
4096 to 1000

**Class Scores**
Cat: 0.9
Dog: 0.05
Car: 0.01
...

**Correct label:** Cat

**Softmax Loss**

**Vector:** 4096

Fully Connected:
4096 to 4

**Box Coordinates**
(x, y, w, h)

"Where"

**Correct box:**
(x', y', w', h')

**L2 Loss**

**Weighted Sum**

Multitask Loss

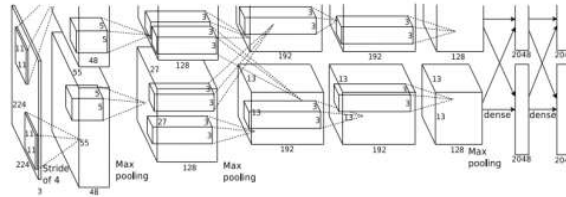**Loss**

**Problem**: Images can have
more than one object!

# Detecting Multiple Objects

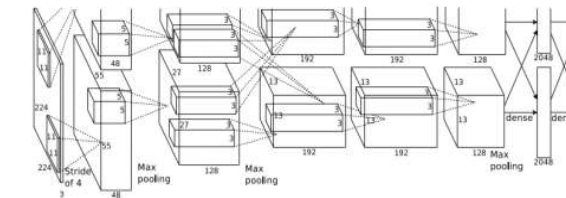**Need different numbers of outputs per image**



CAT: (x, y, w, h)

4 numbers

DOG: (x, y, w, h)
DOG: (x, y, w, h)
CAT: (x, y, w, h)

16 numbers
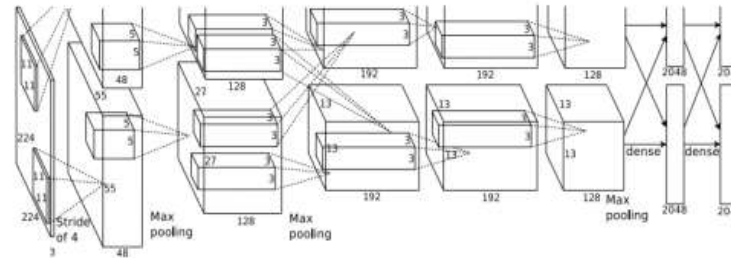
DUCK: (x, y, w, h)
DUCK: (x, y, w, h)
....

Many numbers!

# Detecting Multiple Objects: **Sliding Window**

Apply a CNN to many different crops of the image, CNN classifies each crop as object or background



Dog? NO
Cat? NO
Background? YES

# Detecting Multiple Objects: **Sliding Window**

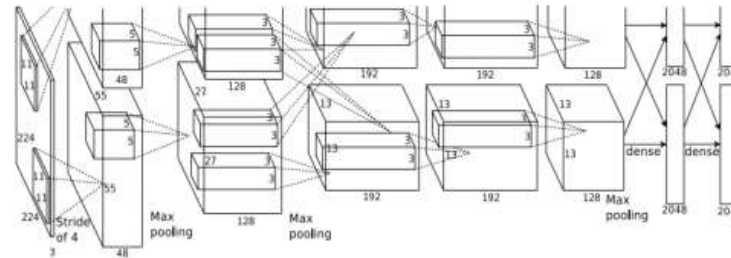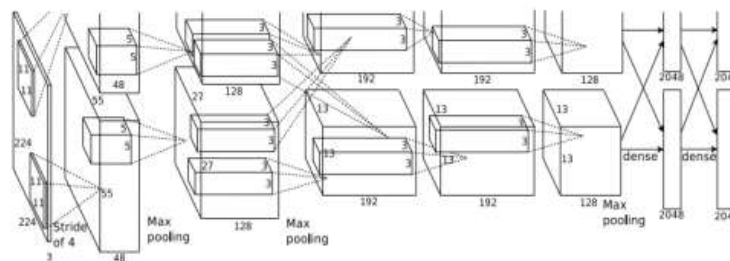Apply a CNN to many different crops of the image, CNN classifies each crop as object or background



Dog? YES
Cat? NO
Background? NO

# Detecting Multiple Objects: **Sliding Window**


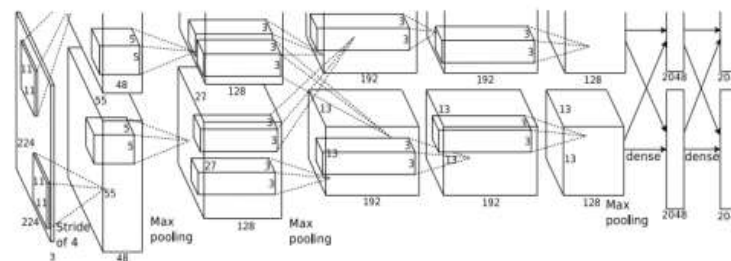
Dog? YES
Cat? NO
Background? NO

# Detecting Multiple Objects: **Sliding Window**



Dog? NO
Cat? YES
Background? NO

# Detecting Multiple Objects: **Sliding Window**

**Question**: How many possible boxes are there in an image of size H x W?

# Detecting Multiple Objects: **Sliding Window**

**Question**: How many possible boxes are there in an image of size H x W?



Consider a box of size h x w:

Possible x positions: W – w + 1

Possible y positions: H – h + 1

Possible positions:     (W – w + 1) * (H – h + 1)

**Object Detection**

**Alireza Akhavanpour**

تشخیص اشیاء

علیرضا اخوان پور

CLASS.
VISION

# Detecting Multiple Objects: **Sliding Window**

**Question**: How many possible boxes are there in an image of size H x W?



Consider a box of size h x w:

Possible x positions: W − w + 1

Possible y positions: H − h + 1

Possible positions:     (W − w + 1) * (H − h + 1)

Total possible boxes:

$$\sum_{h=1}^{H}\sum_{w=1}^{W}(W - w + 1)(H - h + 1)$$

$$= \frac{H(H + 1)}{2}\frac{W(W + 1)}{2}$$

# Detecting Multiple Objects: **Sliding Window**

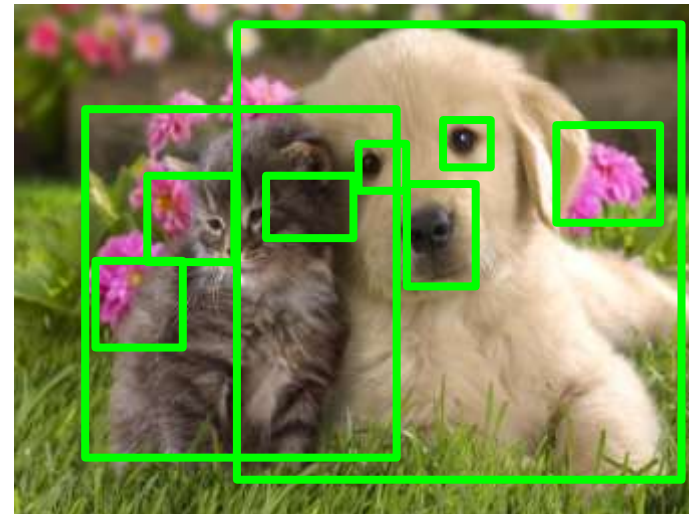**Question**: How many possible boxes are there in an image of size H x W?

Total possible boxes:

$$\sum_{h=1}^{H}\sum_{w=1}^{W}(W - w + 1)(H - h + 1)$$

$$= \frac{H(H + 1)}{2}\frac{W(W + 1)}{2}$$

800 x 600 image has ~58M boxes! No way we can evaluate them all

**Object Detection**

**Alireza Akhavanpour**

تشخیص اشیاء

علیرضا اخوان پور

CLASS.
VISION

# Region Proposals

- Find a small set of boxes that are likely to cover all objects
- Often based on heuristics: e.g. look for "blob-like" image regions
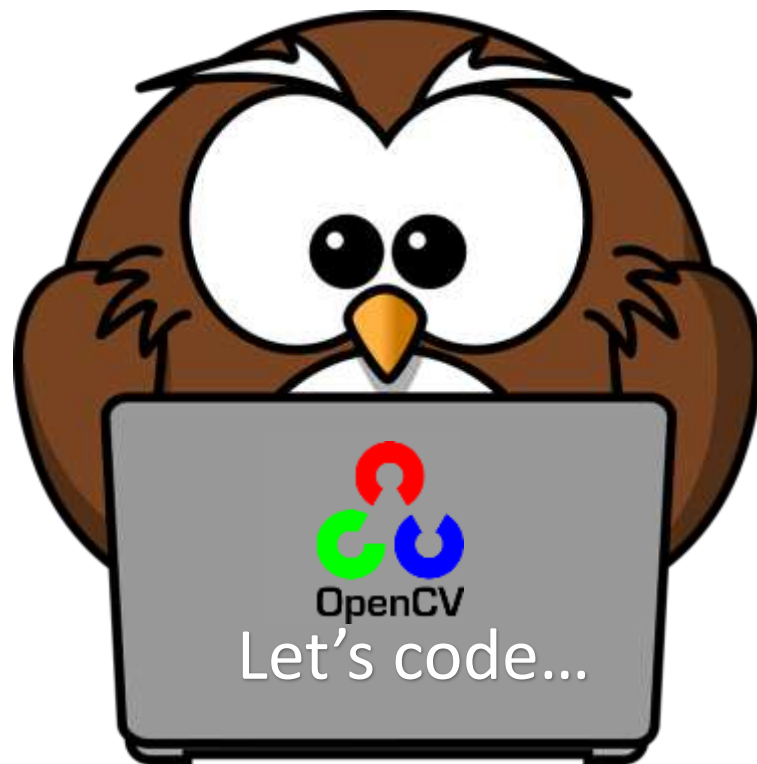- Relatively fast to run; e.g. Selective Search gives 2000 region proposals in a few seconds on CPU



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

**Object Detection**

**Alireza Akhavanpour**

تشخیص اشیاء

علیرضا اخوان پور

CLASS.
VISION

# Region Proposals



Let's code...

**5-selective-search.ipynb**

**Object Detection**

**Alireza Akhavanpour**

تشخیص اشیاء

علیرضا اخوان پور

CLASS.
VISION

# R-CNN: Region-Based CNN



Input image

Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.
Figure copyright Ross Girshick, 2015; source. Reproduced with permission.

# R-CNN: Region-Based CNN



Input image

Regions of Interest (RoI) from a proposal method (~2k)

Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.
Figure copyright Ross Girshick, 2015; source. Reproduced with permission.

Object Detection

Alireza Akhavanpour
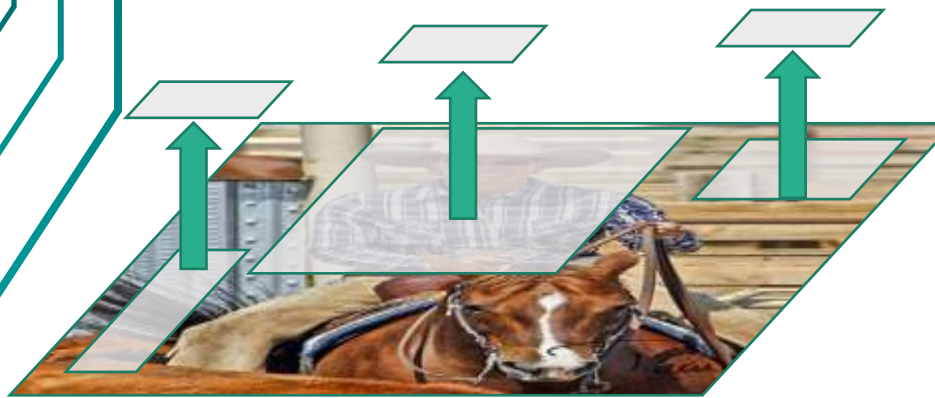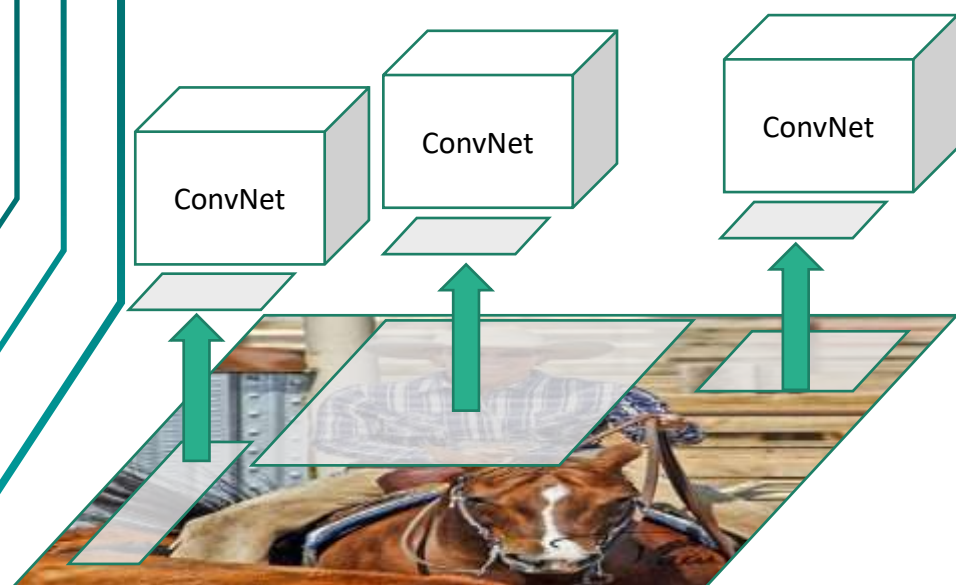
تشخیص اشیاء

علیرضا اخوان پور

CLASS.
VISION

# R-CNN: Region-Based CNN



Warped image regions (224x224)

Input image

Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.
Figure copyright Ross Girshick, 2015; source. Reproduced with permission.

# R-CNN: Region-Based CNN



Forward each region through Convolutional network

Input image

Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.
Figure copyright Ross Girshick, 2015; source. Reproduced with permission.

# R-CNN: Region-Based CNN

Class

Class

Class

ConvNet

ConvNet

ConvNet

Input image

Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.
Figure copyright Ross Girshick, 2015; source. Reproduced with permission.

# R-CNN: Region-Based CNN

Classify each region

Bounding box regression:
Predict "transform" to correct the RoI:
4 numbers ($t_x$, $t_y$, $t_h$, $t_w$)

Bbox    Class        Bbox    Class

Bbox    Class

ConvNet

ConvNet

ConvNet

ConvNet

Input image

Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.
Figure copyright Ross Girshick, 2015; source. Reproduced with permission.

CLASS.
VISION

# R-CNN: Test-time
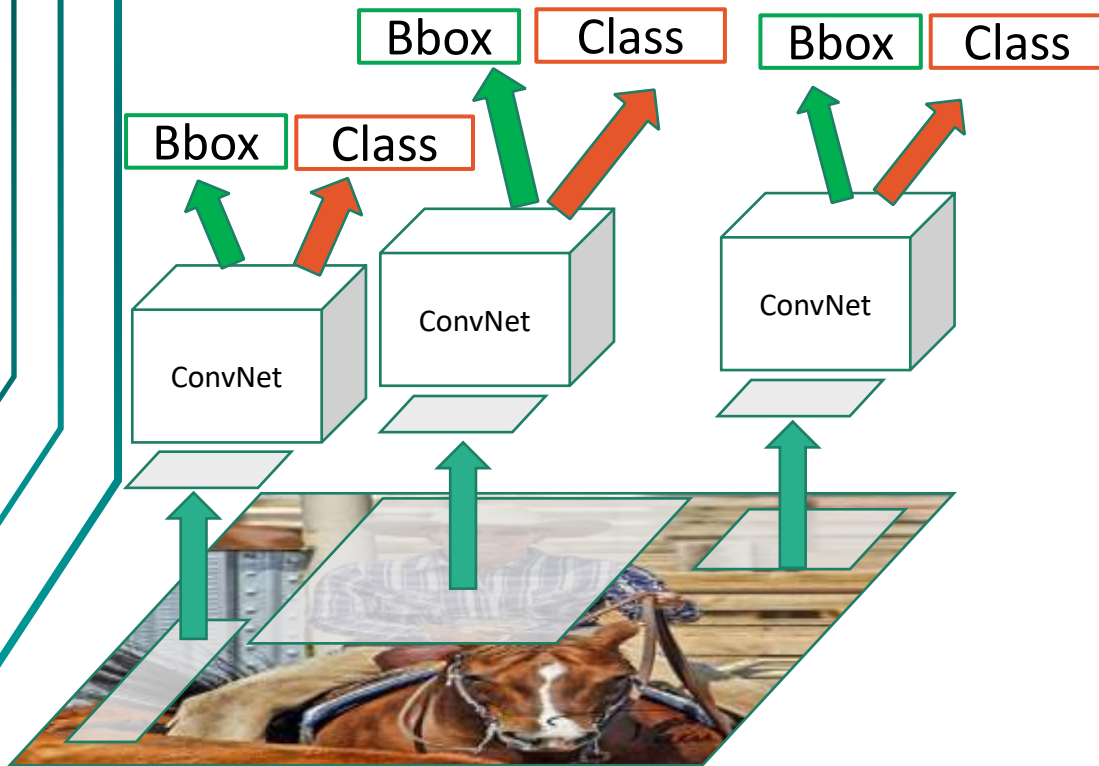


Input image

Input: Single RGB Image

1. Run region proposal method to compute ~2000 region proposals
2. Resize each region to 224x224 and run independently through CNN to predict class scores and bbox transform
3. Use scores to select a subset of region proposals to output
(Many choices here: threshold on background, or per-category? Or take top K proposals per image?)
4. Compare with ground-truth boxes

Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.
Figure copyright Ross Girshick, 2015; source. Reproduced with permission.

Object Detection

Alireza Akhavanpour
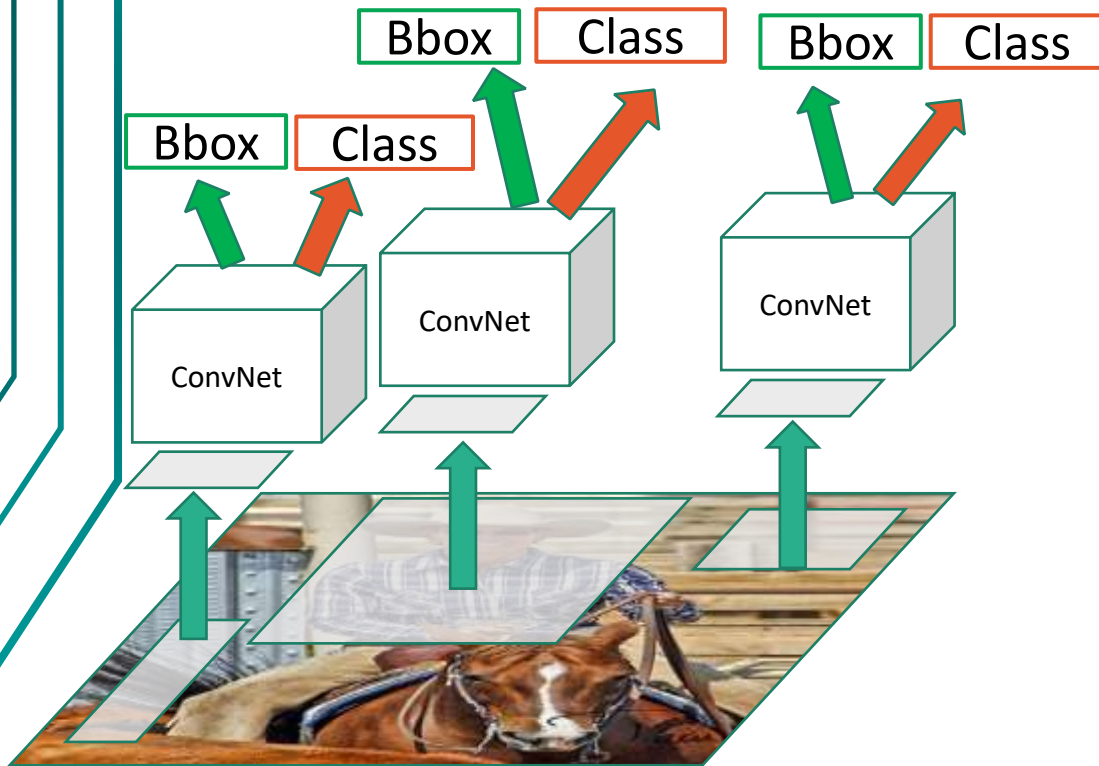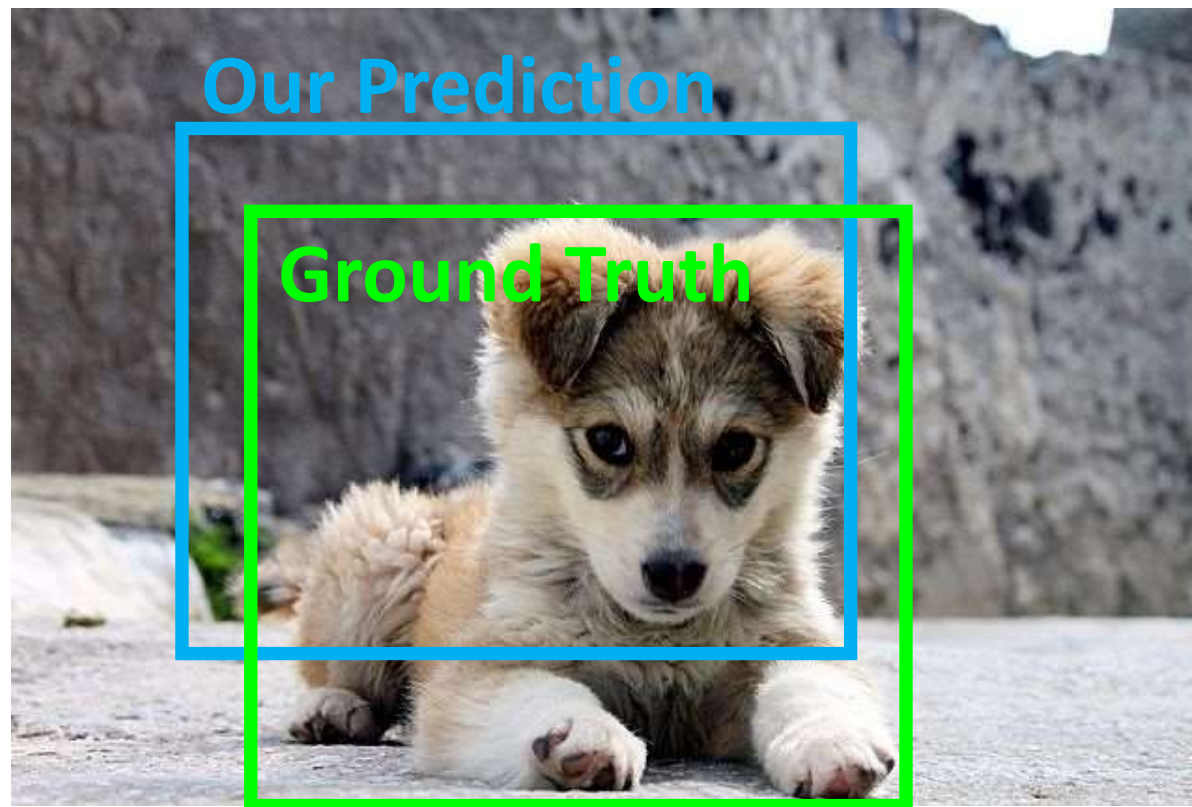
تشخیص اشیاء

علیرضا اخوان پور

CLASS.
VISION

# Comparing Boxes: Intersection over Union (IoU)

How can we compare our prediction to the ground-truth box?

**Object Detection**

**Alireza Akhavanpour**

تشخیص اشیاء

علیرضا اخوان پور

# Comparing Boxes: Intersection over Union (IoU)

How can we compare our prediction to the ground-truth box?

**Our Prediction**

**Ground Truth**

**Intersection over Union** (IoU)
(Also called "Jaccard similarity" or "Jaccard index"):

$$\frac{\textit{Area of Intersection}}{\textit{Area of Union}}$$

Object Detection

Alireza Akhavanpour

تشخیص اشیاء

علیرضا اخوان پور

CLASS.
VISION

# Comparing Boxes: Intersection over Union (IoU)

How can we compare our prediction to the ground-truth box?



**Our Prediction**

**Ground Truth**

IOU = 0.54

**Intersection over Union** (IoU)
(Also called "Jaccard similarity" or "Jaccard index"):

$$\frac{\textit{Area of Intersection}}{\textit{Area of Union}}$$

IOU > 0.5 is "decent"

Object Detection

تشخیص اشیاء

Alireza Akhavanpour

علیرضا اخوان پور

# Comparing Boxes: Intersection over Union (IoU)

How can we compare our prediction to the ground-truth box?



**Our Prediction**

**Ground Truth**

IOU = 0.72

**Intersection over Union** (IoU)
(Also called "Jaccard similarity" or "Jaccard index"):

$$\frac{\textit{Area of Intersection}}{\textit{Area of Union}}$$

IOU > 0.5 is "decent"
IOU > 0.7 is "pretty good"

Object Detection
Alireza Akhavanpour

تشخیص اشیاء
علیرضا اخوان پور

CLASS.
VISION

# Comparing Boxes: Intersection over Union (IoU)

How can we compare our prediction to the ground-truth box?



**Our Prediction**

**Ground Truth**

IOU = 0.93

**Intersection over Union** (IoU)
(Also called "Jaccard similarity" or "Jaccard index"):

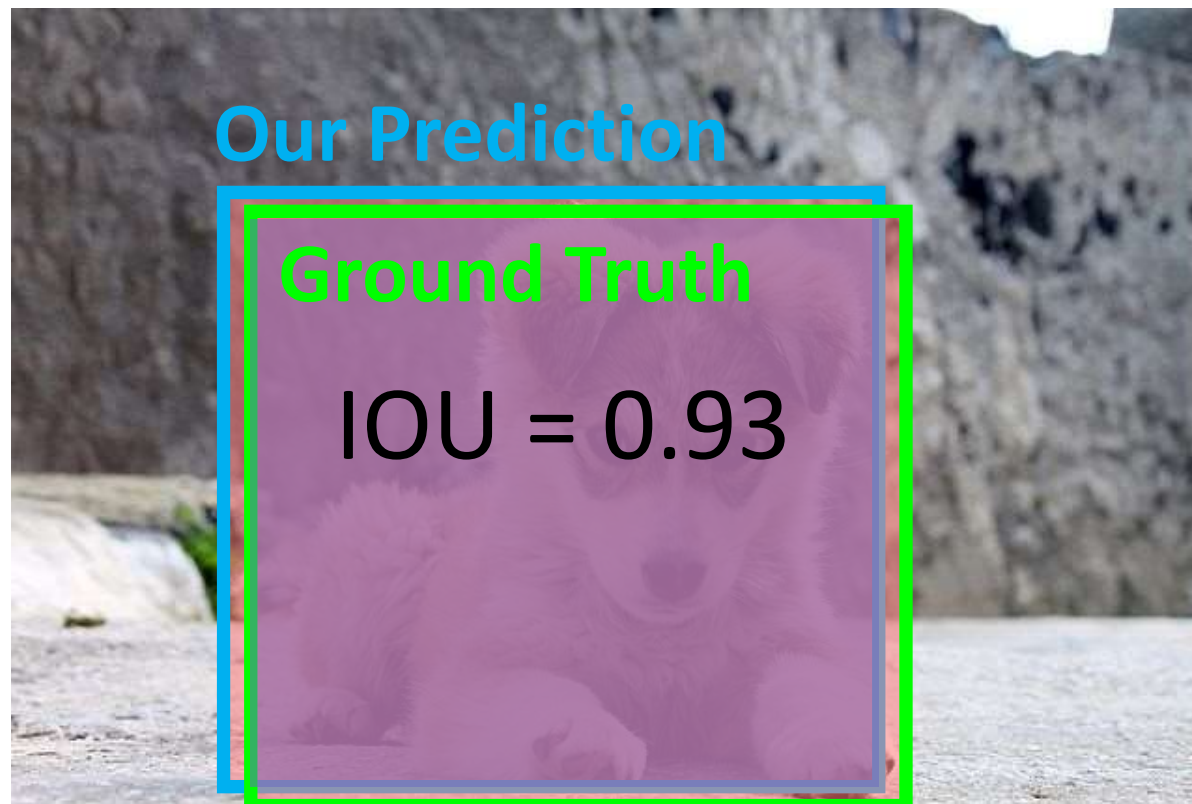$$\frac{\text{Area of Intersection}}{\text{Area of Union}}$$

IOU > 0.5 is "decent"
IOU > 0.7 is "pretty good"
IOU > 0.9 is "almost perfect"

Object Detection

Alireza Akhavanpour

تشخیص اشیاء

علیرضا اخوان پور

**CLASS. VISION**

# What is bBox?



**6-object-detection-and-bounding-boxes.ipynb**

# Intersection over Union (IoU)



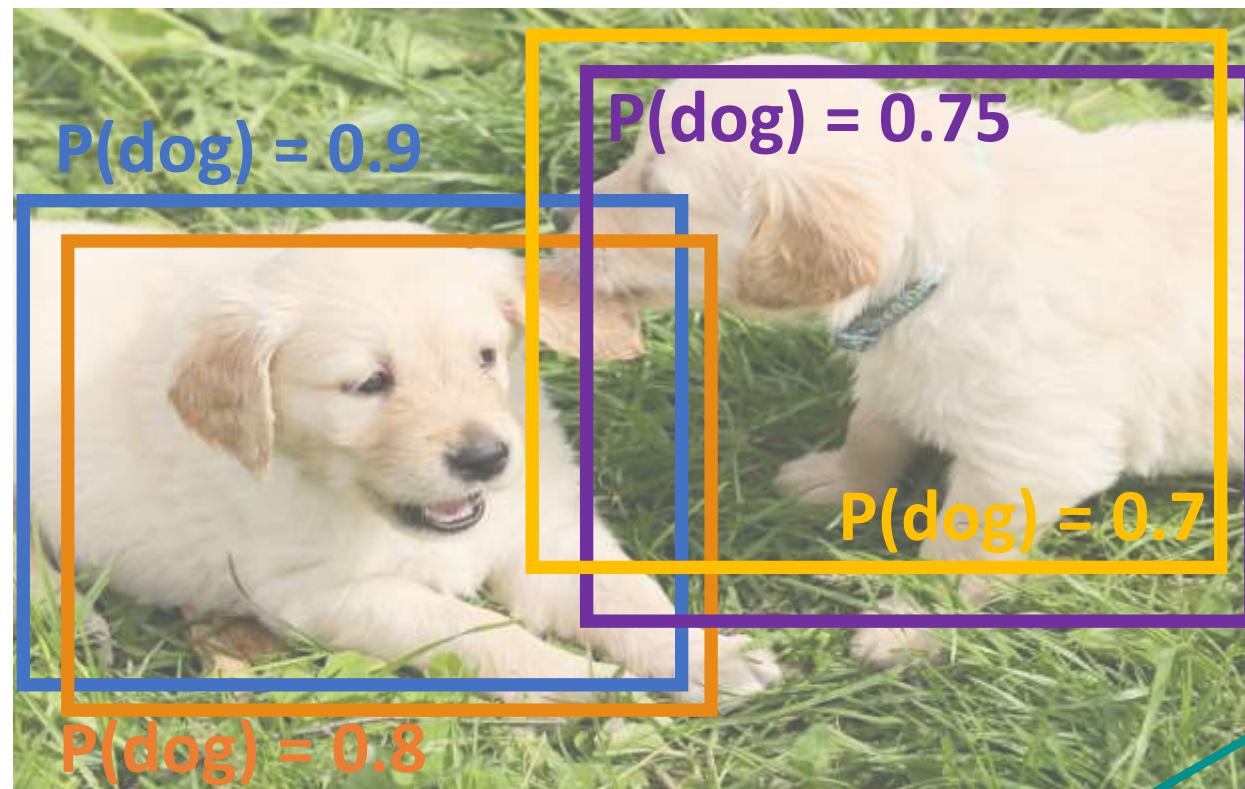Let's code…

**7-Intersection-over-Union(IoU).ipynb**

# Overlapping Boxes

**Problem**: Object detectors often output many overlapping detections:



P(dog) = 0.9
P(dog) = 0.75
P(dog) = 0.7
P(dog) = 0.8

# Overlapping Boxes: **Non-Max Suppression (NMS)**

**Problem**: Object detectors often output many overlapping detections:

**Solution**: Post-process raw detections using **Non-Max Suppression (NMS)**

1. Select next highest-scoring box
2. Eliminate lower-scoring boxes with IoU > threshold (e.g. 0.7)
3. If any boxes remain, GOTO 1



P(dog) = 0.9

P(dog) = 0.75

P(dog) = 0.7

P(dog) = 0.8

**Object Detection**

**Alireza Akhavanpour**

تشخیص اشیاء

علیرضا اخوان پور
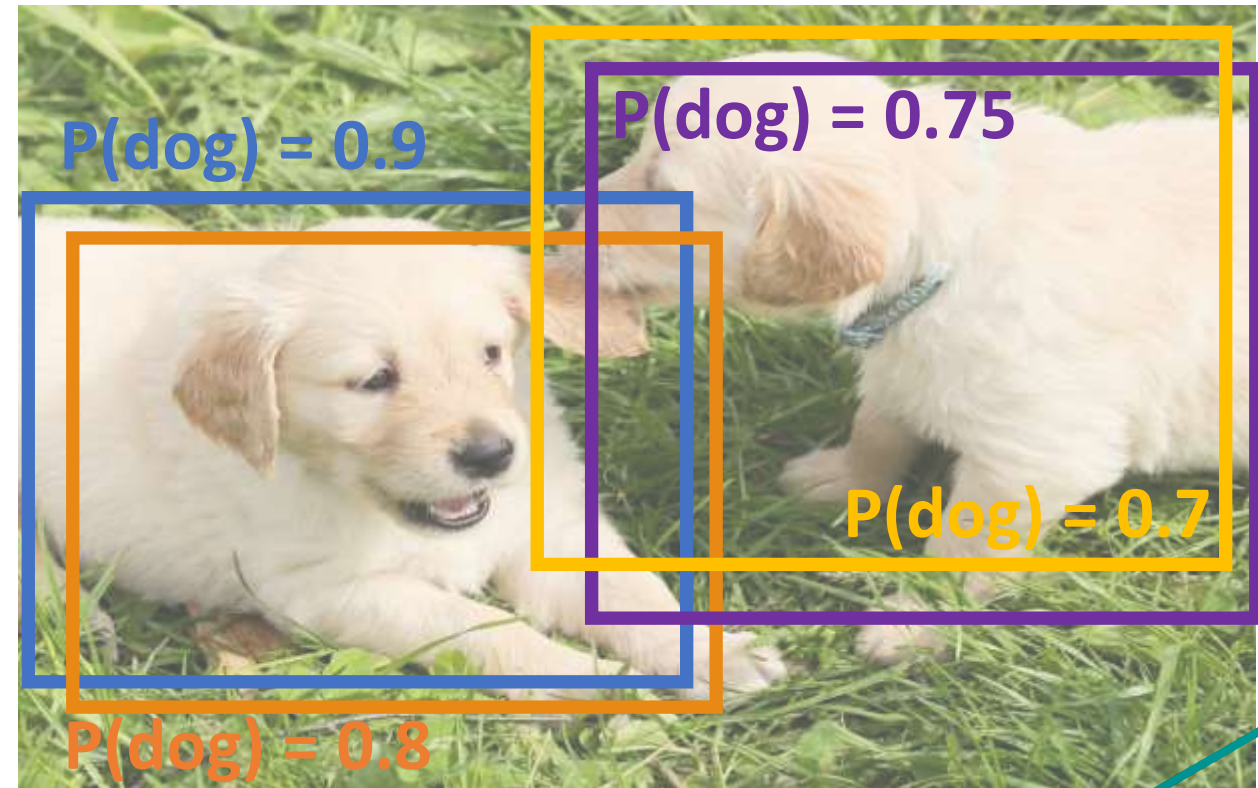
Puppy image is CC0 Public Domain

CLASS.
VISION

# Overlapping Boxes: **Non-Max Suppression (NMS)**

**Problem**: Object detectors often output many overlapping detections:

**Solution**: Post-process raw detections using **Non-Max Suppression (NMS)**

1. Select next highest-scoring box
2. Eliminate lower-scoring boxes with IoU > threshold (e.g. 0.7)
3. If any boxes remain, GOTO 1

IoU(■, ■) = **0.78**
IoU(■, ■) = 0.05
IoU(■, ■) = 0.07



P(dog) = 0.9
P(dog) = 0.75
P(dog) = 0.7
P(dog) = 0.8

# Overlapping Boxes: **Non-Max Suppression (NMS)**

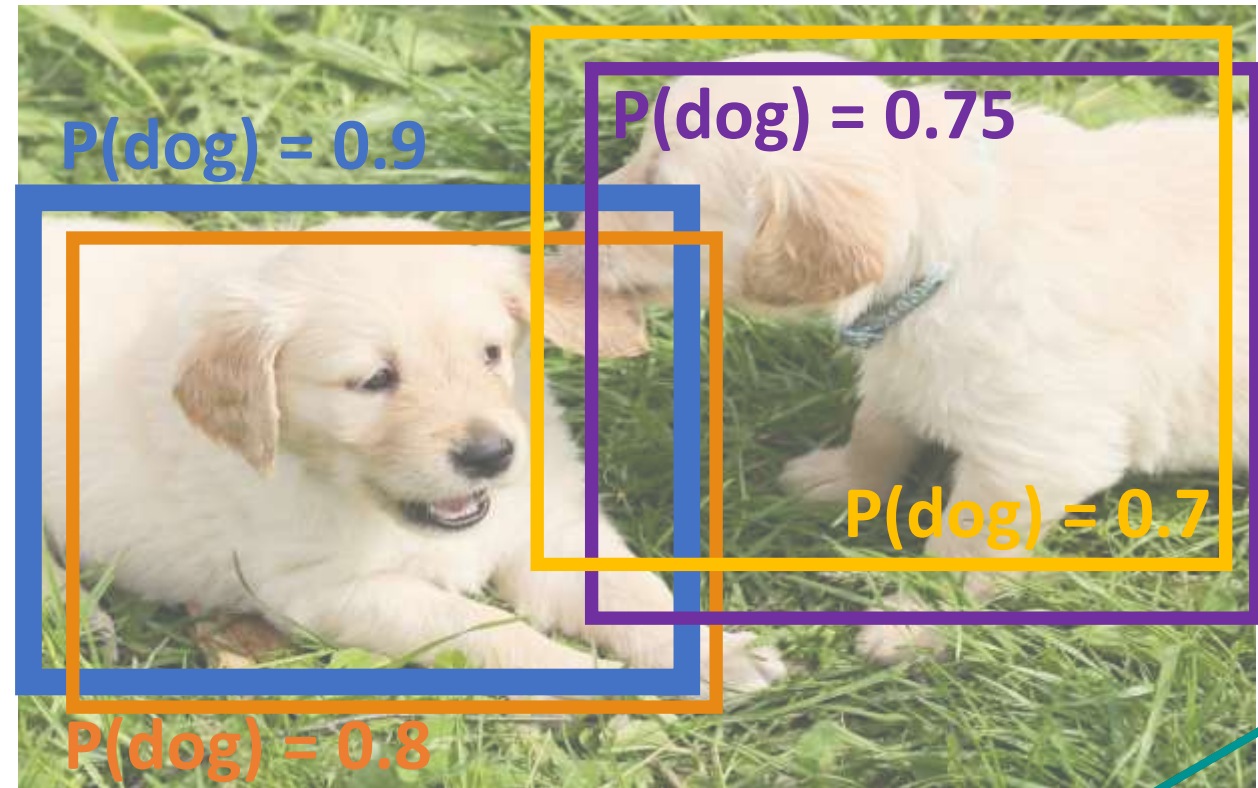**Problem**: Object detectors often output many overlapping detections:

**Solution**: Post-process raw detections using **Non-Max Suppression (NMS)**

1. Select next highest-scoring box
2. Eliminate lower-scoring boxes with IoU > threshold (e.g. 0.7)
3. If any boxes remain, GOTO 1



P(dog) = 0.9

P(dog) = 0.75

P(dog) = 0.7

**Object Detection**

**Alireza Akhavanpour**

تشخیص اشیاء

علیرضا اخوان پور

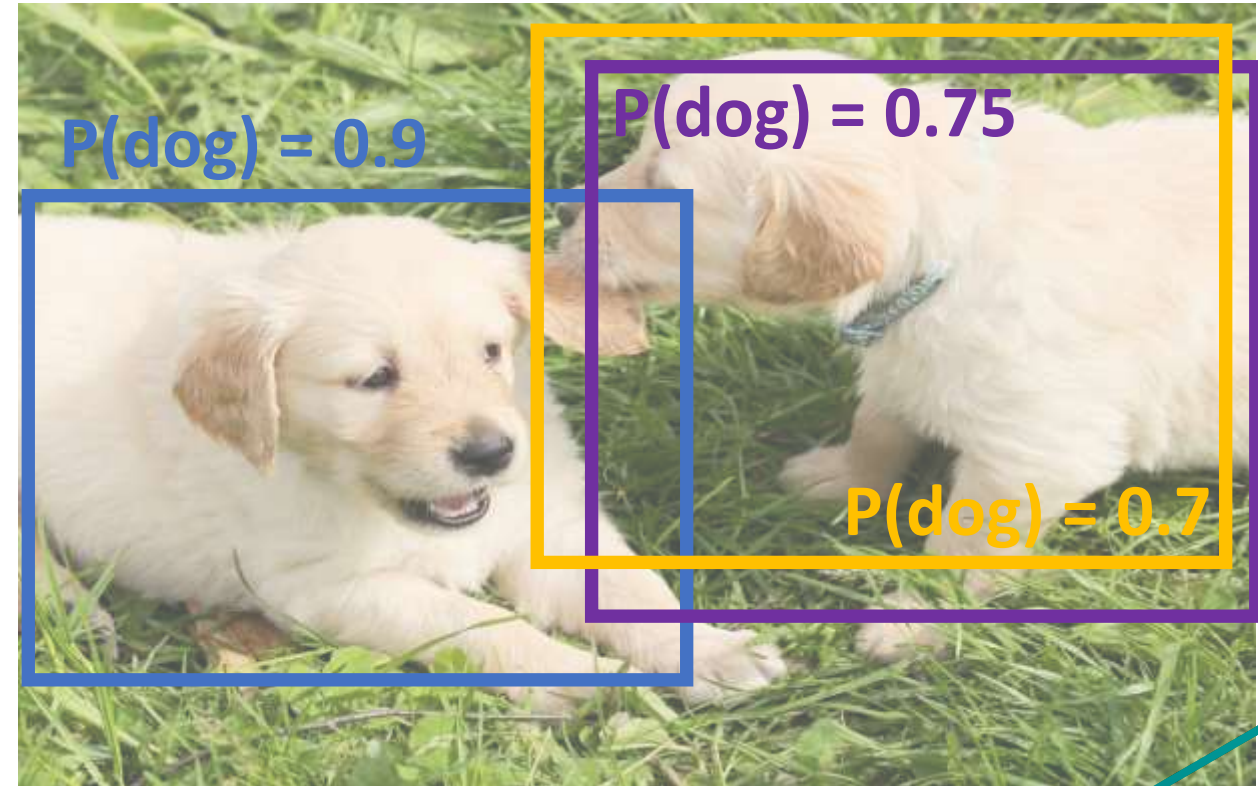Puppy image is CC0 Public Domain

CLASS.
VISION

# Overlapping Boxes: **Non-Max Suppression (NMS)**

**Problem**: Object detectors often output many overlapping detections:

**Solution**: Post-process raw detections using **Non-Max Suppression (NMS)**

1. Select next highest-scoring box
2. Eliminate lower-scoring boxes with IoU > threshold (e.g. 0.7)
3. If any boxes remain, GOTO 1

IoU(■, ■) = **0.74**
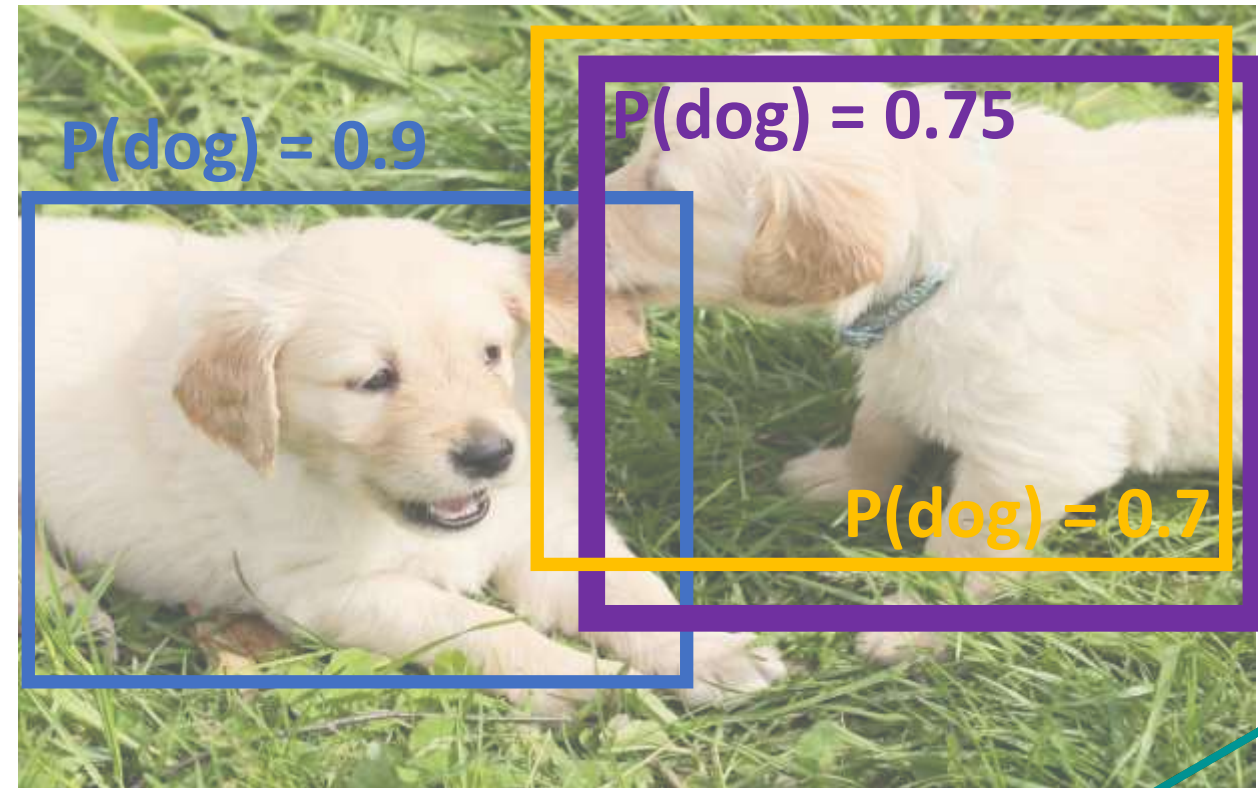


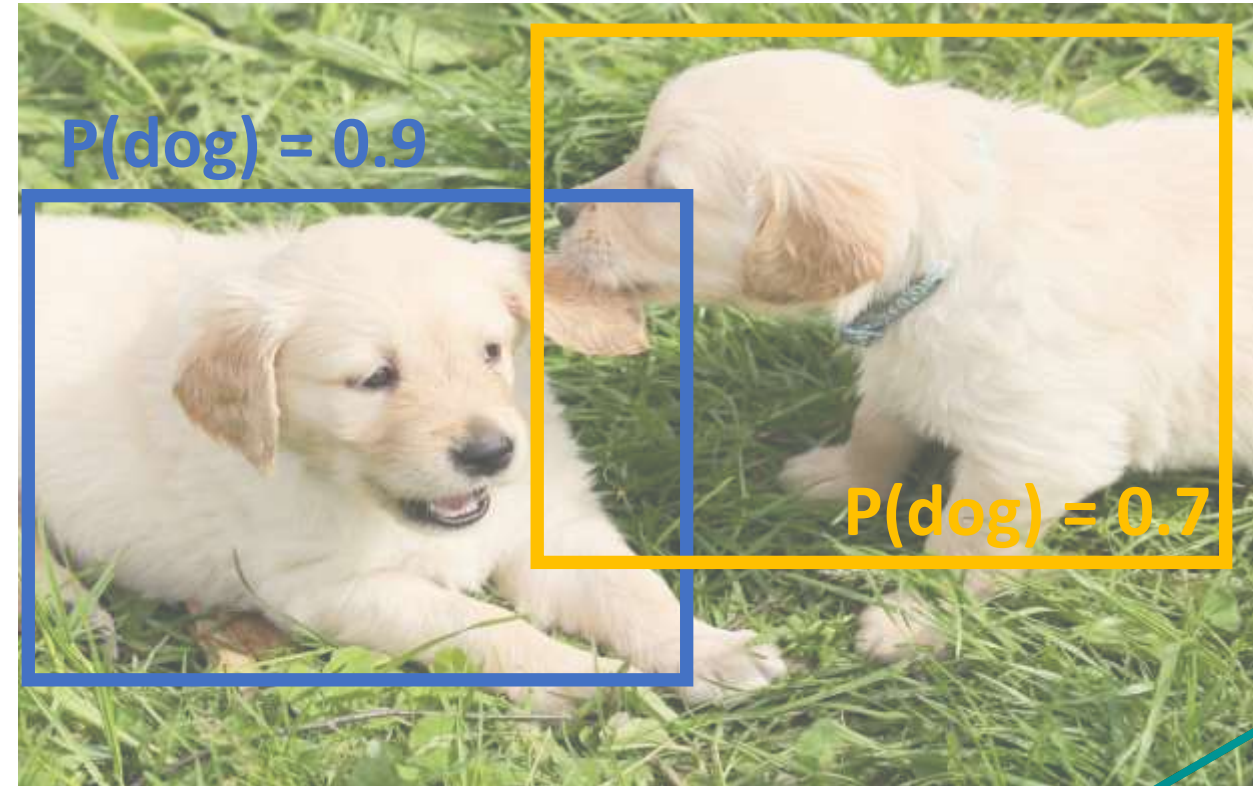P(dog) = 0.9

P(dog) = 0.75

P(dog) = 0.7

# Overlapping Boxes: **Non-Max Suppression (NMS)**

**Problem**: Object detectors often output many overlapping detections:

**Solution**: Post-process raw detections using **Non-Max Suppression (NMS)**

1. Select next highest-scoring box
2. Eliminate lower-scoring boxes with IoU > threshold (e.g. 0.7)
3. If any boxes remain, GOTO 1



P(dog) = 0.9

P(dog) = 0.7

# Overlapping Boxes: **Non-Max Suppression (NMS)**



**Problem**:

NMS may eliminate "good" boxes when objects are highly overlapping!

**Object Detection**

**Alireza Akhavanpour**

تشخیص اشیاء

علیرضا اخوان پور

CLASS.
VISION

# Evaluating Object Detectors: Mean Average Precision (mAP)