# Object Detection



CAT, DOG, DUCK

The task of assigning a **label** and a **bounding box** to all objects in the image

# Object Detection: Datasets



20 categories
6k training images
6k validation images
10k test images

80 categories
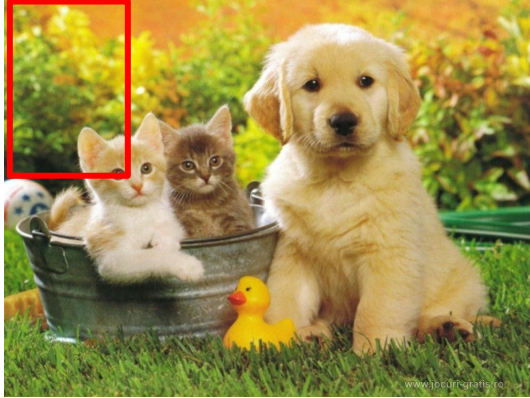200k training images
60k val + test images

200 categories
456k training images
60k validation + test images

# Outline

**Two-stage methods**
One-stage methods

# Object Detection as Classification



Classes = [cat, dog, duck]

Cat ? NO

Dog ? NO

Duck? NO

# Object Detection as Classification



Classes = [cat, dog, duck]

Cat ? NO

Dog ? NO

Duck? NO

# Object Detection as Classification
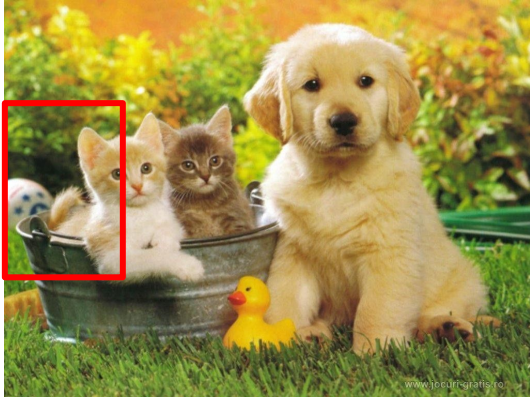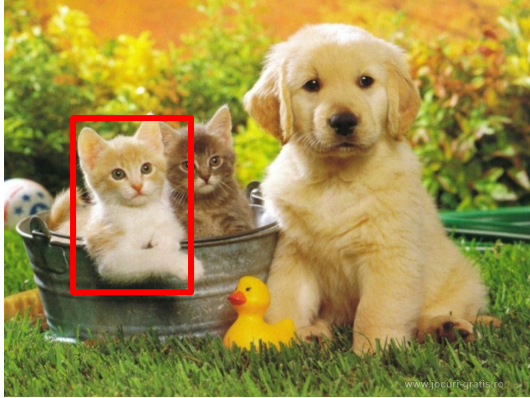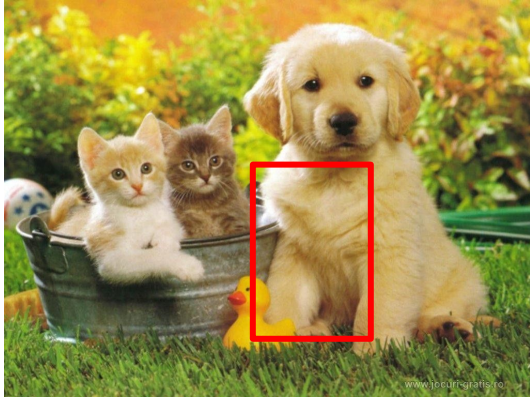
Classes = [cat, dog, duck]

Cat ? YES

Dog ? NO

Duck? NO

# Object Detection as Classification
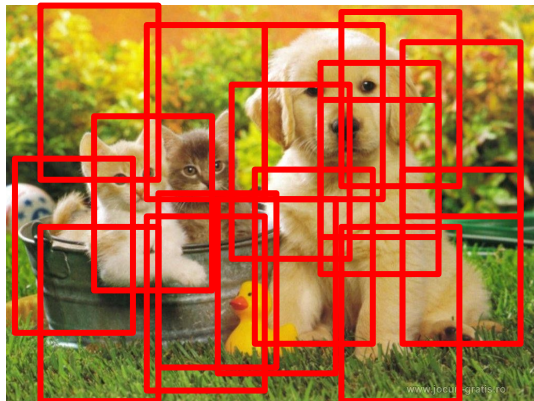


Classes = [cat, dog, duck]

Cat ? NO

Dog ? NO

Duck? NO

# Object Detection as Classification



Problem:
Too many positions & scales to test

Solution: If your classifier is fast enough, go for it

# Object Detection with ConvNets?



Convnets are computationally demanding. We can't test all positions & scales !

Solution: Look at a tiny subset of positions & choose them wisely

# Region Proposals

- Find "blobby" image regions that are likely to contain objects
- "Class-agnostic" object detector
- Look for "blob-like" regions

# Region Proposals



Selective Search (SS)



Multiscale Combinatorial Grouping (MCG)

[SS] Uijlings et al. Selective search for object recognition. IJCV 2013

[MCG]  Arbeláez, Pont-Tuset et al. Multiscale combinatorial grouping. CVPR 2014

# Object Detection with Convnets: R-CNN



warped region

**1**. Input image

**2**. Extract region proposals (~2k)

**3**. Compute CNN features

**4**. Classify regions

aeroplane? no.

person? yes.

tvmonitor? no.

CNN

Girshick et al. Rich feature hierarchies for accurate object detection and semantic segmentation. CVPR 2014

# R-CNN

1. Train network on proposals



**1.** Input image    **2.** Extract region proposals (~2k)    **3.** Compute CNN features    **4.** Classify regions

warped region

aeroplane? no.

person? yes.

tvmonitor? no.

CNN

2. Post-hoc training of SVM classifiers & bounding box regressors on fc7 features



bounding box regressor predicts coordinate offsets

Girshick et al. Rich feature hierarchies for accurate object detection and semantic segmentation. CVPR 2014

# R-CNN

We expect:

We get:

# R-CNN

1. Train network on proposals



**1**. Input image **2**. Extract region proposals (~2k) **3**. Compute CNN features **4**. Classify regions

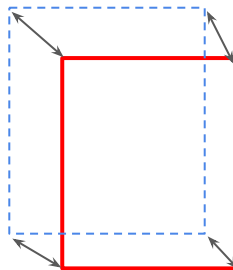warped region

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

CNN

2. Post-hoc training of SVMs & Box regressors on fc7 features

3. **Non Maximum Suppression + score threshold**

Girshick et al. Rich feature hierarchies for accurate object detection and semantic segmentation. CVPR 2014

# R-CNN



Girshick et al. Rich feature hierarchies for accurate object detection and semantic segmentation. CVPR 2014

# R-CNN: Problems

1. Slow at test-time: need to run full forward pass of CNN for each region proposal

2. SVMs and regressors are post-hoc: CNN features not updated in response to SVMs and regressors

3. Complex multistage training pipeline

Slide Credit: CS231n

# Fast R-CNN

R-CNN Problem #1: Slow at test-time: need to run full forward pass of CNN for each region proposal



Solution: Share computation of convolutional layers between region proposals for an image

Girshick Fast R-CNN. ICCV 2015

# Fast R-CNN: Sharing features



Convolution and Pooling

Max-pool within each grid cell

Fully-connected layers

Hi-res input image:
3 x 800 x 600
with region proposal

Hi-res conv features:
C x H x W
with region proposal

RoI conv features:
C x h x w
for region proposal

Fully-connected layers expect low-res conv features:
C x h x w

# Fast R-CNN

R-CNN Problem #2&3: SVMs and regressors are post-hoc. Complex training.



Solution: Train it all at together E2E

Girshick Fast R-CNN. ICCV 2015

# Fast R-CNN

|  |  | R-CNN | Fast R-CNN |
|---|---|---|---|
| Faster! | Training Time: | 84 hours | **9.5 hours** |
|  | (Speedup) | 1x | **8.8x** |
| FASTER! | Test time per image | 47 seconds | **0.32 seconds** |
|  | (Speedup) | 1x | **146x** |
| Better! | mAP (VOC 2007) | 66.0 | **66.9** |

Using VGG-16 CNN on Pascal VOC 2007 dataset

# Fast R-CNN: Problem

Test-time speeds don't include region proposals

|  | R-CNN | Fast R-CNN |
|---|---|---|
| Test time per image | 47 seconds | **0.32 seconds** |
| (Speedup) | 1x | **146x** |
| Test time per image with Selective Search | 50 seconds | **2 seconds** |
| (Speedup) | 1x | **25x** |

# Faster R-CNN

Learn proposals end-to-end sharing parameters with the classification network



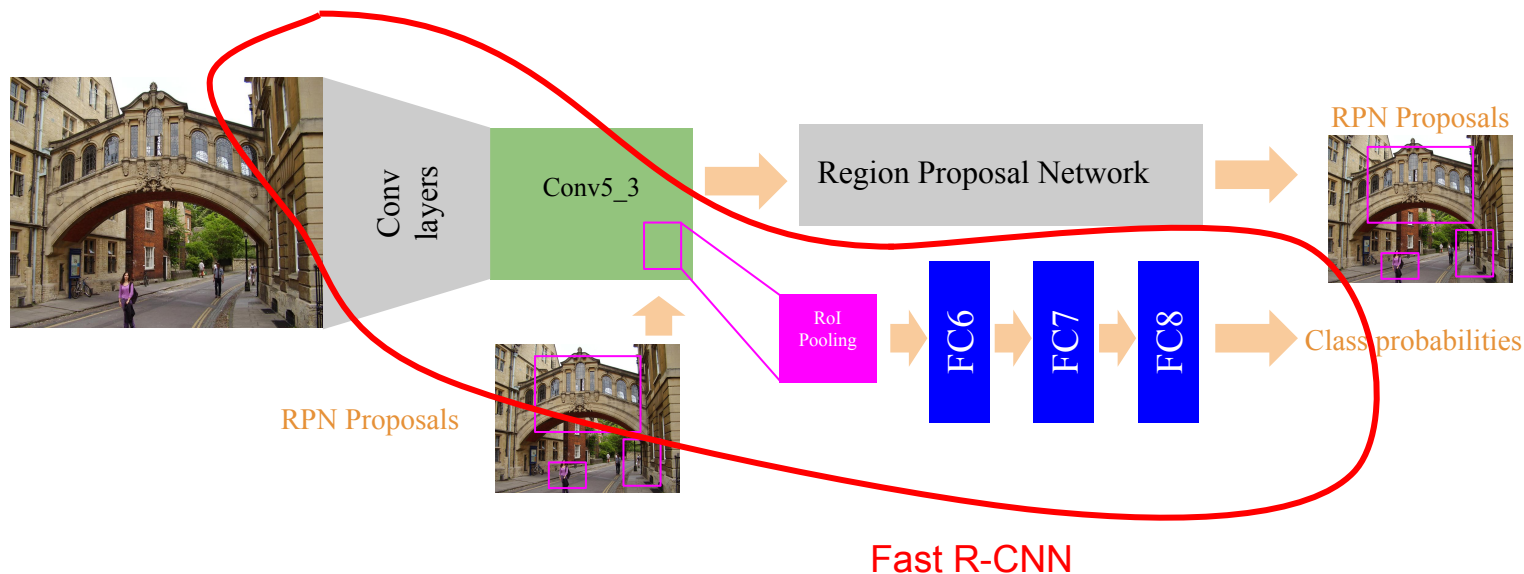Ren et al. Faster R-CNN: Towards real-time object detection with region proposal networks. NIPS 2015

# Faster R-CNN

Learn proposals end-to-end sharing parameters with the classification network



Fast R-CNN

Ren et al. Faster R-CNN: Towards real-time object detection with region proposal networks. NIPS 2015

# Region Proposal Network



Bounding Box Regression

Objectness scores (object/no object)

2$k$ scores

$cls$ layer

4$k$ coordinates

$reg$ layer

$k$ anchor boxes

256-d

intermediate layer

sliding window

conv feature map

In practice, k = 9 (3 different scales and 3 aspect ratios) → 18k boxes for a 40x50 input feature map

Ren et al. Faster R-CNN: Towards real-time object detection with region proposal networks. NIPS 2015

# Faster R-CNN: Training

RoI Pooling is not differentiable w.r.t box coordinates. Solutions:
- Alternate training
- Ignore gradient of classification branch w.r.t proposal coordinates
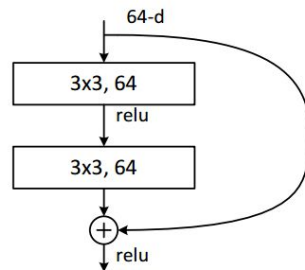- Make pooling function differentiable



Ren et al. Faster R-CNN: Towards real-time object detection with region proposal networks. NIPS 2015

# Faster R-CNN

| | R-CNN | Fast R-CNN | Faster R-CNN |
|---|---|---|---|
| Test time per image (with proposals) | 50 seconds | 2 seconds | **0.2 seconds** |
| (Speedup) | 1x | 25x | **250x** |
| mAP (VOC 2007) | 66.0 | **66.9** | **66.9** |

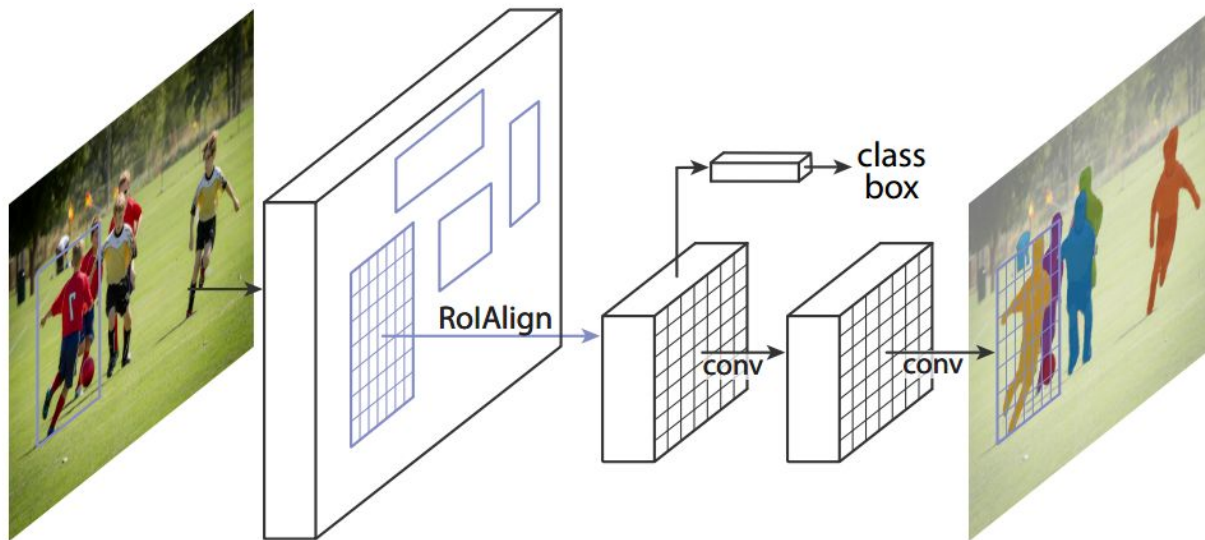Ren et al. Faster R-CNN: Towards real-time object detection with region proposal networks. NIPS 2015

# Better Encoder: ResNet

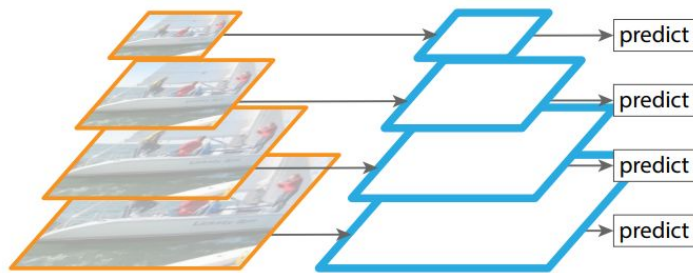- Faster R-CNN was the basis of the winners of COCO and ILSVRC 2015&2016 object detection competitions.



He et al. Deep residual learning for image recognition. CVPR 2016
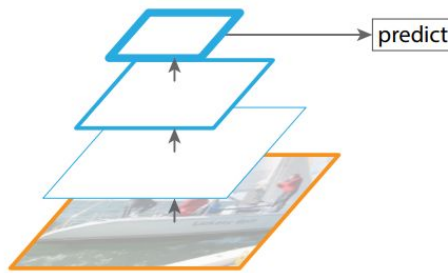
# Better Region Pooling : RoI Align



He et al. Mask R-CNN. ICCV 2017

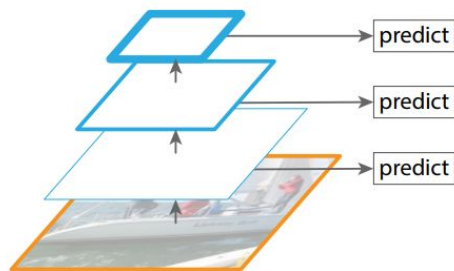# Better Representations: FPN
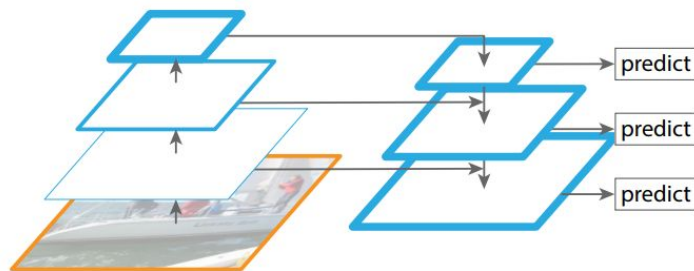


(a) Featurized image pyramid

(b) Single feature map

(c) Pyramidal feature hierarchy

(d) Feature Pyramid Network

Lin et al. Feature Pyramid Networks for Object Detection. CVPR 2017
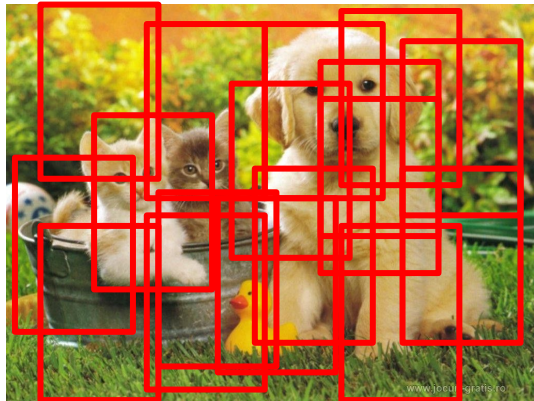
# Outline

Two-stage methods
**One-stage methods**

# One-stage methods

Previously… :



Problem:
Too many positions & scales to test

Solution: If your classifier is fast enough, go for it

# One-stage methods

Previously… :



Problem:
Too many positions & scales to test
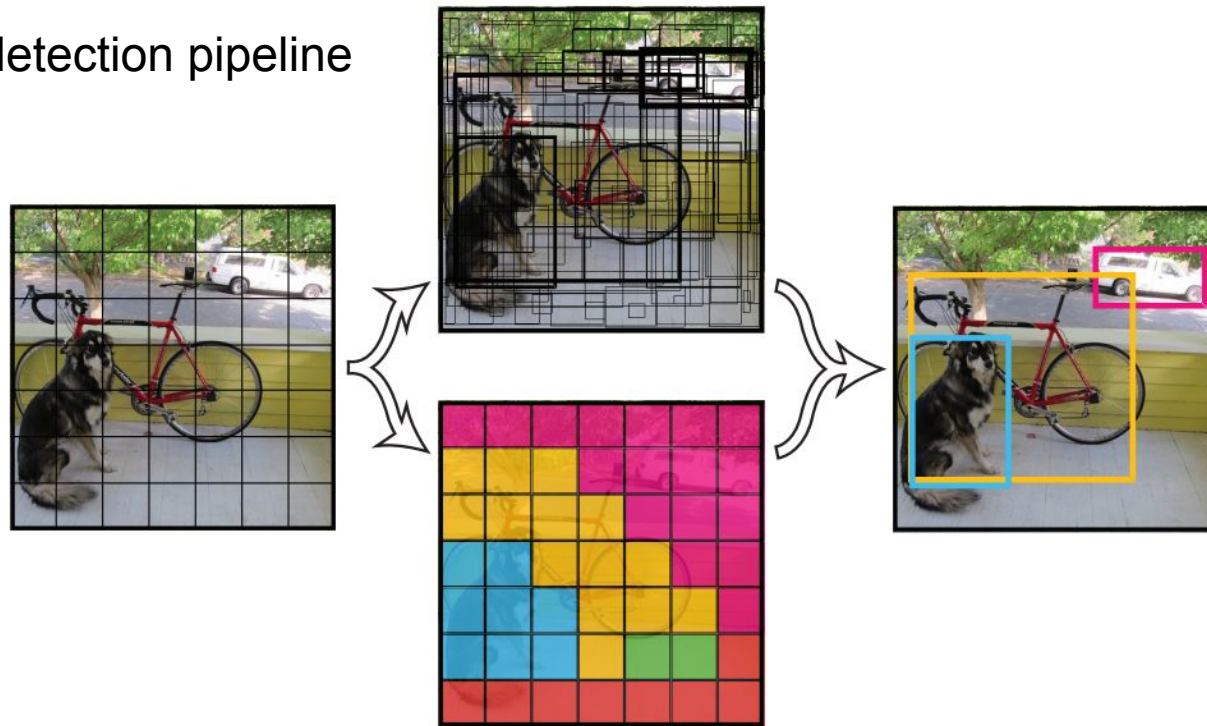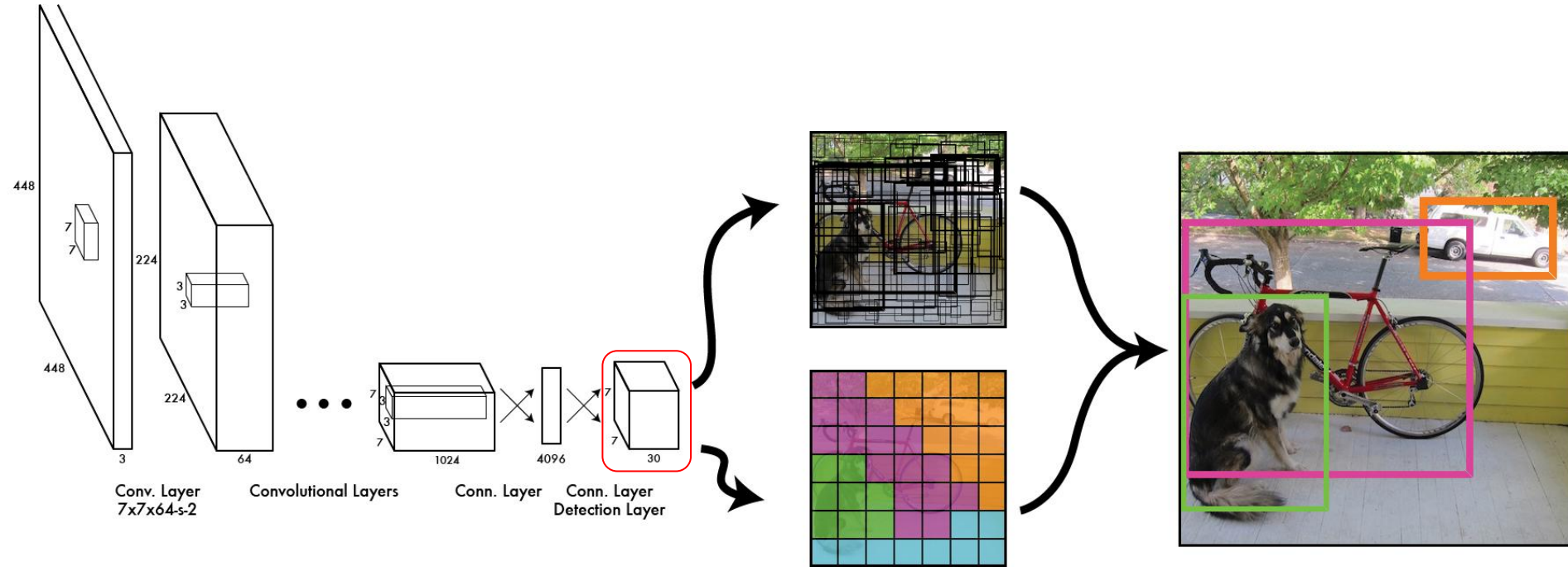
**Modern detectors parallelize feature extraction across all locations.
Region classification is not slow anymore!**

# YOLO: You Only Look Once

Proposal-free object detection pipeline

Redmon et al. You Only Look Once: Unified, Real-Time Object Detection, CVPR 2016

# YOLO: You Only Look Once



Redmon et al. You Only Look Once: Unified, Real-Time Object Detection, CVPR 2016

# YOLO: You Only Look Once

Each cell predicts:

- For each bounding box:
  - 4 coordinates (x, y, w, h)
  - 1 confidence value
- Some number of class probabilities

For Pascal VOC:

- 7x7 grid
- 2 bounding boxes / cell
- 20 classes



$$7 \times 7 \times (2 \times 5 + 20) = 7 \times 7 \times 30 \text{ tensor} = \textbf{1470 outputs}$$

# YOLO: You Only Look Once

Predict class probability for each cell



Redmon et al. You Only Look Once: Unified, Real-Time Object Detection, CVPR 2016

38

# YOLO: You Only Look Once



+ NMS
+ Score threshold

Redmon et al. You Only Look Once: Unified, Real-Time Object Detection, CVPR 2016

# SSD: Single Shot MultiBox Detector

Same idea as YOLO, + several predictors at different stages in the network & uses anchor boxes



Liu et al. SSD: Single Shot MultiBox Detector, ECCV 2016

# YOLOv2

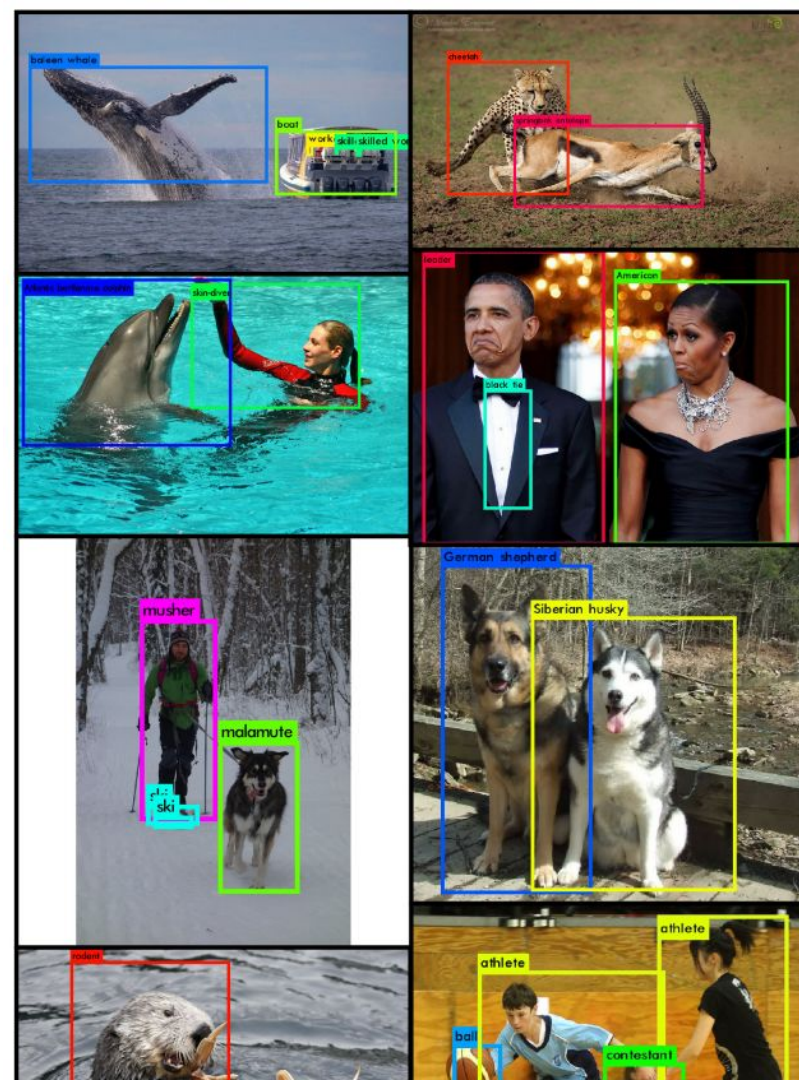| | YOLO | | | | | | | | YOLOv2 |
|---|---|---|---|---|---|---|---|---|---|
| batch norm? | | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| hi-res classifier? | | | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| convolutional? | | | | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| anchor boxes? | | | | ✓ | ✓ | | | | |
| new network? | | | | | ✓ | ✓ | ✓ | ✓ | ✓ |
| dimension priors? | | | | | | ✓ | ✓ | ✓ | ✓ |
| location prediction? | | | | | | ✓ | ✓ | ✓ | ✓ |
| passthrough? | | | | | | | ✓ | ✓ | ✓ |
| multi-scale? | | | | | | | | ✓ | ✓ |
| hi-res detector? | | | | | | | | | ✓ |
| VOC2007 mAP | 63.4 | 65.8 | 69.5 | 69.2 | 69.6 | 74.4 | 75.4 | 76.8 | **78.6** |



Redmon & Farhadi. YOLO900: Better, Faster, Stronger. CVPR 2017

# YOLOv3



YOLO v3 network Architecture

Legend:
- (*) Concatenation
- (+) Addition
- Residual Block
- Detection Layer
- Upsampling Layer
- • Further Layers

Scale 1 Stride: 32

Scale 2 Stride: 16

Scale 3 Stride: 8

# RetinaNet

**Matching proposal-based performance with a one-stage approach**



(a) ResNet    (b) feature pyramid net    (c) class subnet (top)    (d) box subnet (bottom)

Key idea is to lower loss weight for well classified samples, increase it for difficult ones

Lin et al. Focal Loss for Dense Object Detection. ICCV 2017

# Overview



| Method | mAP | time |
|---|---|---|
| [B] SSD321 | 28.0 | 61 |
| [C] DSSD321 | 28.0 | 85 |
| [D] R-FCN | 29.9 | 85 |
| [E] SSD513 | 31.2 | 125 |
| [F] DSSD513 | 33.2 | 156 |
| [G] FPN FRCN | 36.2 | 172 |
| RetinaNet-50-500 | 32.5 | 73 |
| RetinaNet-101-500 | 34.4 | 90 |
| RetinaNet-101-800 | **37.8** | 198 |
| **YOLOv3-320** | 28.2 | **22** |
| **YOLOv3-416** | 31.0 | 29 |
| **YOLOv3-608** | 33.0 | 51 |

# Summary

## Two-stage methods
- R-CNN
- Fast R-CNN
- Faster R-CNN

## One-stage methods
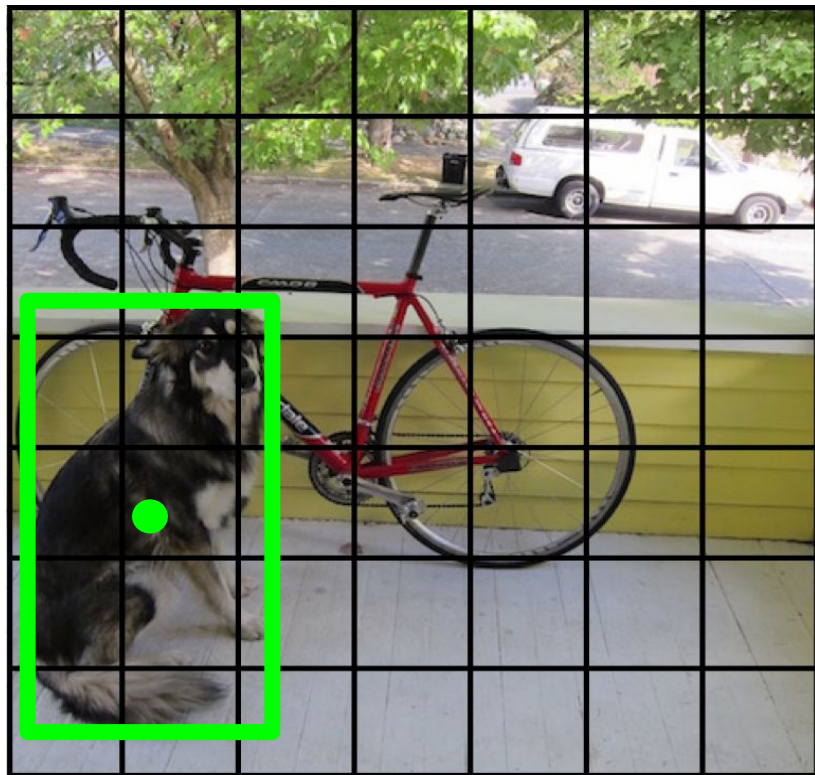- YOLO
- SSD
- RetinaNet

# Questions?

# Resources

APIs including implementations to most popular detectors:

- [Detectron: Facebook Object Detection API](#) (Caffe2)
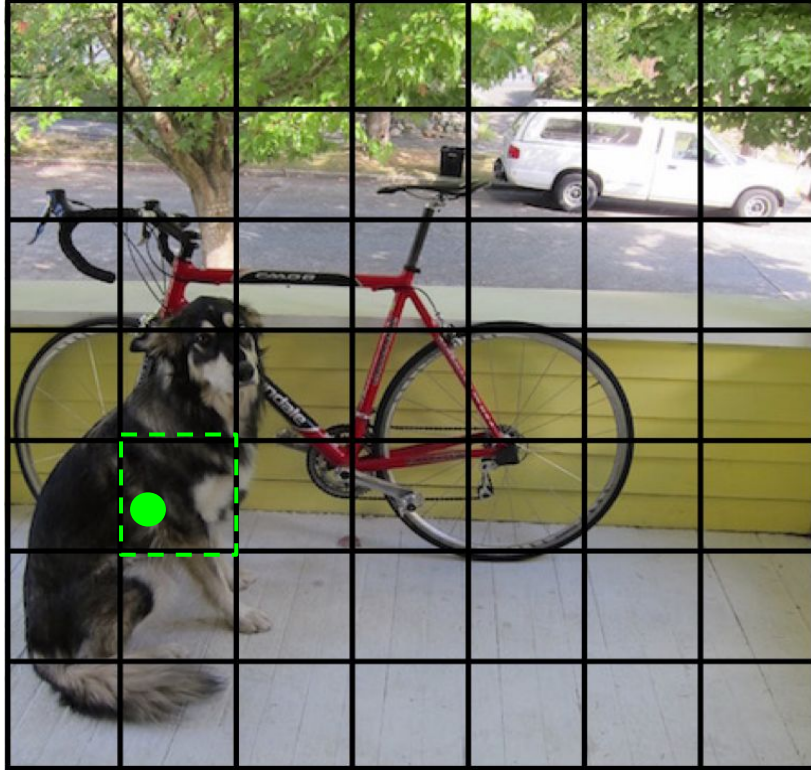- [Google's Tensorflow Object Detection API](#)

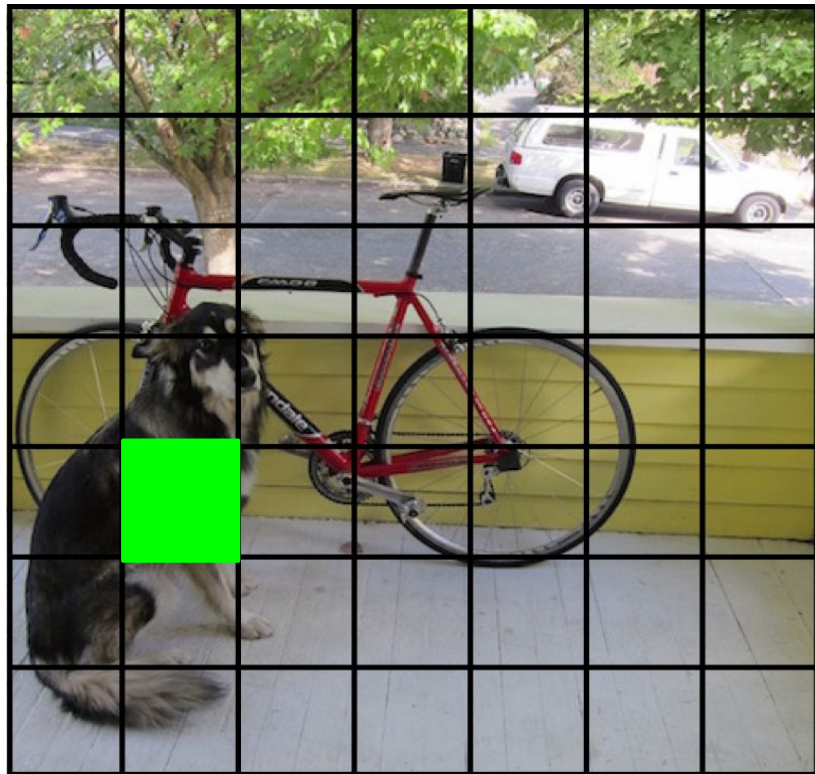Many unofficial ports to other frameworks !

# YOLO: Training



For training, each ground truth bounding box is matched into the right cell
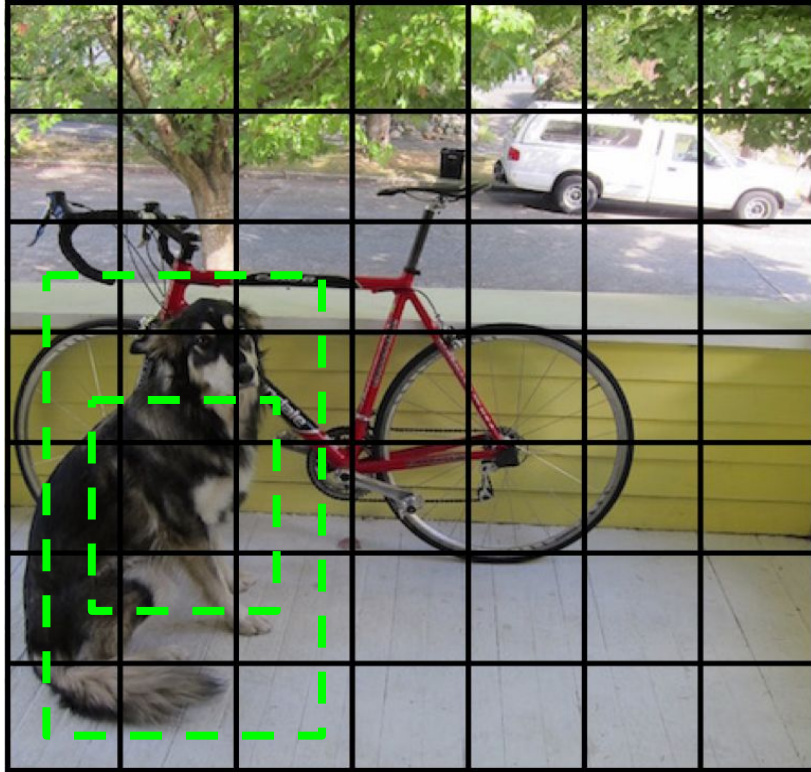
# YOLO: Training



For training, each ground truth bounding box is matched into the right cell
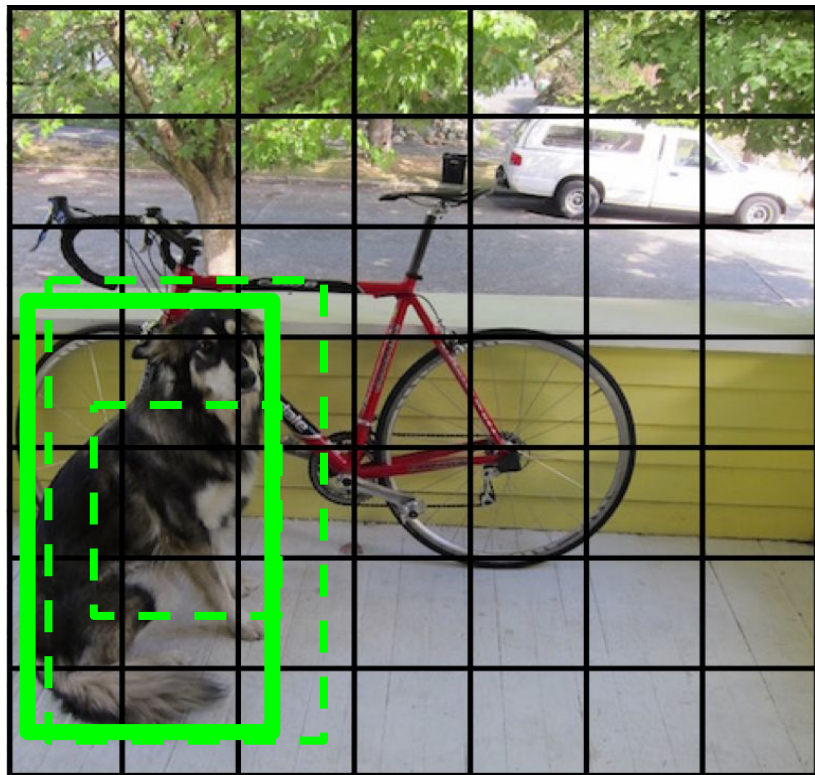
# YOLO: Training



Optimize class prediction in that cell:
dog: 1, cat: 0, bike: 0, ...

# YOLO: Training



Predicted boxes for this cell

# YOLO: Training



Find the best one wrt ground truth bounding box, optimize it (i.e. adjust its coordinates to be closer to the ground truth's coordinates)

# YOLO: Training



Increase matched box's confidence, decrease non-matched boxes confidence

# YOLO: Training



Increase matched box's confidence, decrease non-matched boxes confidence

# YOLO: Training



For cells with no ground truth detections, confidences of all predicted boxes are decreased

# YOLO: Training



For cells with no ground truth detections:

- Confidences of all predicted boxes are decreased
- Class probabilities are not adjusted

# YOLO: Training, formally

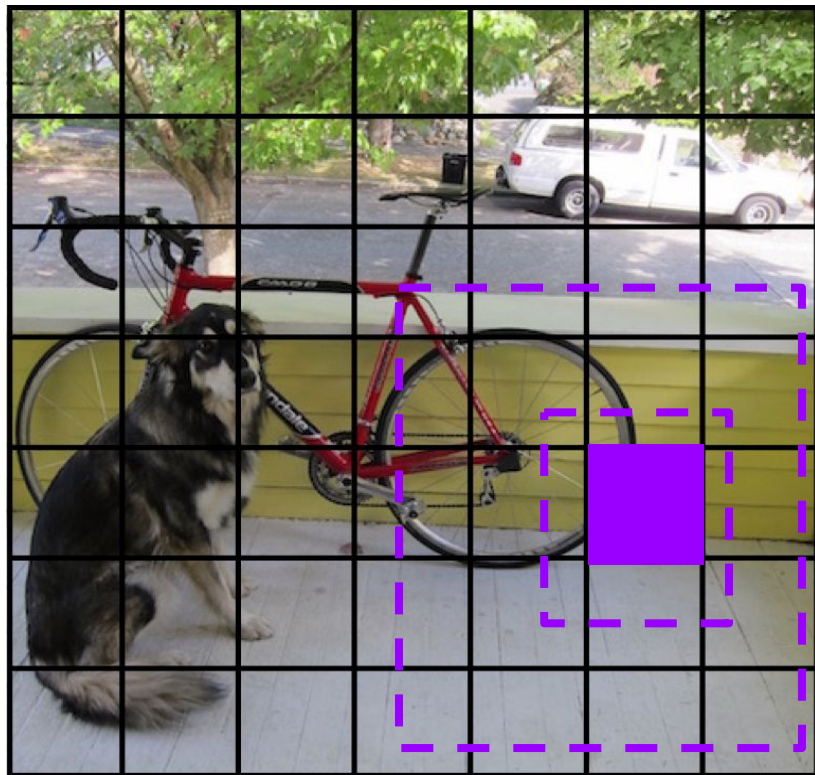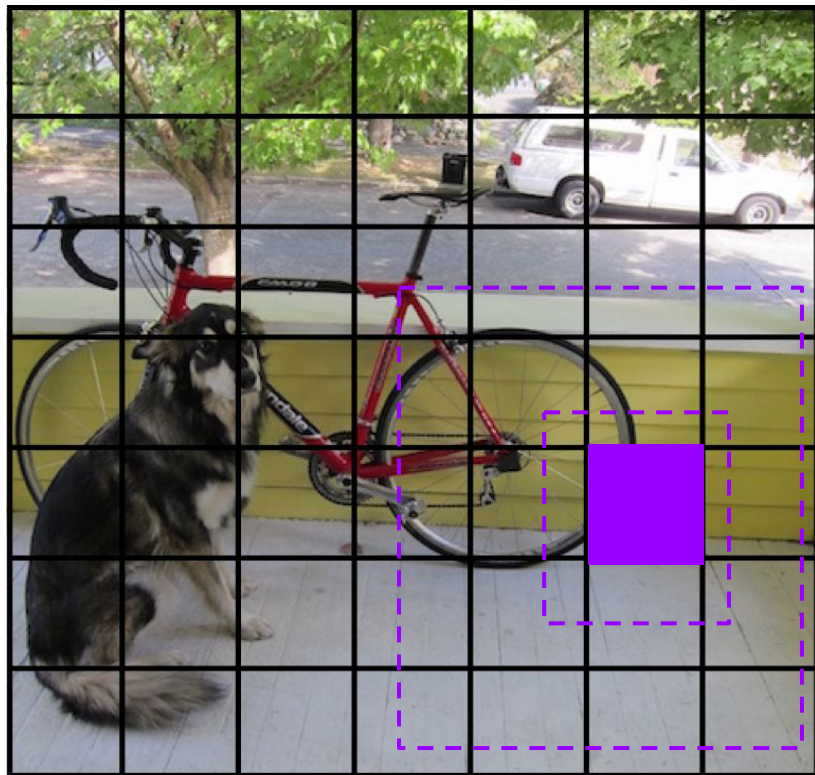= 1 if box $j$ and cell $i$ are matched together, 0 otherwise

Bounding box coordinate regression

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

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

Bounding box score prediction

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

= 1 if box $j$ and cell $i$ are NOT matched together

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

Class score prediction

$$+ \sum_{i=0}^{S^2} \mathbb{1}_{i}^{\text{obj}} \sum_{c \in \text{classes}} (p_i(c) - \hat{p}_i(c))^2$$

= 1 if cell $i$ has an object present

# Fast R-CNN: RoI Pooling



input

| 0.88 | 0.44 | 0.14 | 0.16 | 0.37 | 0.77 | 0.96 | 0.27 |
| 0.19 | 0.45 | 0.57 | 0.16 | 0.63 | 0.29 | 0.71 | 0.70 |
| 0.66 | 0.26 | 0.82 | 0.64 | 0.54 | 0.73 | 0.59 | 0.26 |
| 0.85 | 0.34 | 0.76 | 0.84 | 0.29 | 0.75 | 0.62 | 0.25 |
| 0.32 | 0.74 | 0.21 | 0.39 | 0.34 | 0.03 | 0.33 | 0.48 |
| 0.20 | 0.14 | 0.16 | 0.13 | 0.73 | 0.65 | 0.96 | 0.32 |
| 0.19 | 0.69 | 0.09 | 0.86 | 0.88 | 0.07 | 0.01 | 0.48 |
| 0.83 | 0.24 | 0.97 | 0.04 | 0.24 | 0.35 | 0.50 | 0.91 |

Image resource: deepsense.ai