

## CSCE 496/896 Lecture 9: Object Detection

Stephen Scott

sscott@cse.unl.edu

1/21

Navigation icons

## Introduction

- We know that CNNs are useful in image classification
- Now consider **object detection**
  - Given an input image, identify what objects (plural) are in it and where they are
  - Output **bounding box** of each object

Navigation icons

## Outline

- Performance measures
- RCNN
- SPP-net
- Fast RCNN
- YOLO

3/21

Navigation icons

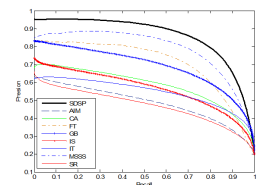
## Performance Measures

### Mean Average Precision

**Mean average precision (mAP)** to measure how well objects are identified

- Recall from Lecture 3

- **Precision** is fraction of those labeled positive that are positive
- **Recall** is fraction of the true positives that are labeled positive
- **Precision-recall curve** plots precision vs recall



Navigation icons

## Performance Measures

### Mean Average Precision (2)

- Given a ranking (by confidence values) of  $n$  items, **average precision at  $n$  (AP@ $n$ )** is average of precision values at each position in the ranking:

$$AP = \sum_{k=1}^n P(k) \Delta r(k) ,$$

where  $P(k)$  is precision at position  $k$  and  $\Delta r(k)$  is change in recall:  $r(k) - r(k-1)$  ( $= 0$  if instance  $k$  is negative,  $= 1/N_p$  if  $k$  is one of  $N_p$  positives)

- E.g., if ranking =  $\langle +, +, -, +, - \rangle$ ,  $AP@5$   
 $= (1)(1/3) + (1)(2/3) + (2/3)(2/3) + (3/4)(1) + (3/5)(1)$
- Larger as more positives ranked above negatives
- mAP is mean of average precision across all classes

5/21

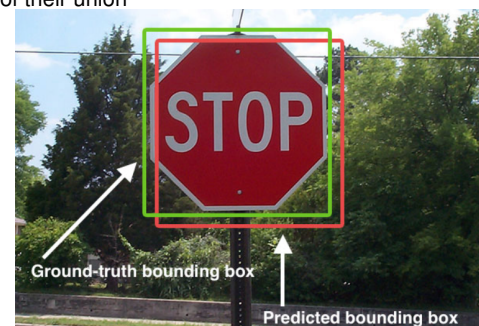
Navigation icons

## Performance Measures

### Intersection Over Union

**Intersection over union (IoU)** to measure quality of bounding boxes

- Divide the size of the two boxes' intersection by the size of their union



6/21

Navigation icons

## Basic Idea of Object Detection

Split input image into **regions** and classify each region with a CNN and other machinery

- Region boundary is bounding box
- Object detected in region is object in BB



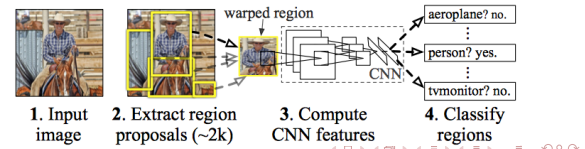
## Issues:

- Limited to bounding boxes of fixed sizes and locations
- An object could span regions

## Region CNN (Girshick et al. 2014)

- R-CNN **proposes** collection of 2000 regions in image
- Warps each region to match input dimensions ( $227 \times 227 \times 3$ ) of CNN to get 4096-dimensional embedded representation
- Classifies each embedded vector with class-specific binary SVMs
- Apply class-specific regressors to fine-tune bounding boxes

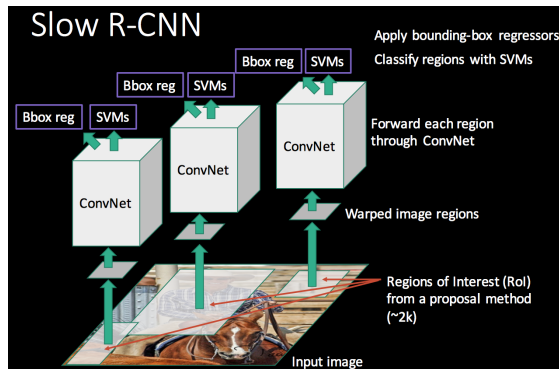
## R-CNN: Regions with CNN features



## Region CNN (Girshick et al. 2014)

Example from Girshick (2015)

## Slow R-CNN

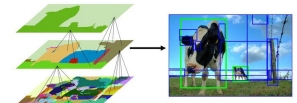


## Region CNN (Girshick et al. 2014)

Selective Search

Popular method to propose Rols: **selective search**

- Segment** the image
- Compute bounding boxes of segments
- Iteratively merge adjacent segments based on similarity
  - Linear combination of similarities of: color, texture, size, shape
- Goto 2



## Region CNN (Girshick et al. 2014)

Issues

- Training and detection are slow
- Detection: 13s/image on GPU, 53s/image on CPU
- Due to large number of regions proposed, each run through CNN and classified

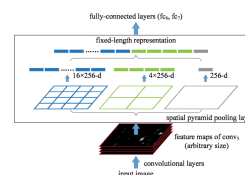
## Spatial Pyramid Pooling (He et al. 2015)

- Part of R-CNN's slowdown at test time is running each RoI through ConvNet separately
- To speed up test time, instead put entire image through **single ConvNet**

- Choose Rols from ConvNet output and run through **spatial pyramid pooling (SPP)** layer

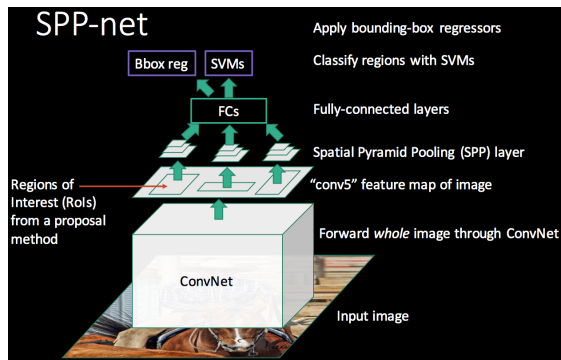
- Max/avg pooling with **fixed** number of bins
- Produces fixed-length vector regardless of input size

- Fixed-length vectors feed to fully connected layers, then SVMs



## Spatial Pyramid Pooling (He et al. 2015)

Example from Girshick (2015)



## Spatial Pyramid Pooling (He et al. 2015)

Drawbacks

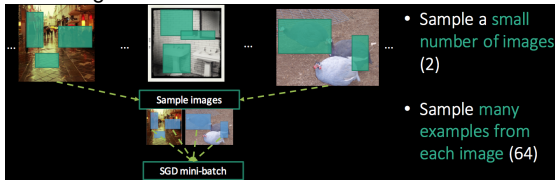
- While training is faster than R-CNN, is still slow and disk-intensive
- Cannot efficiently update ConvNet parameters, so kept frozen
  - Each RoI's receptive field covers most of entire image, so forward pass expensive across all images of mini-batch

## Fast R-CNN (Girshick 2015)

Hierarchical Sampling

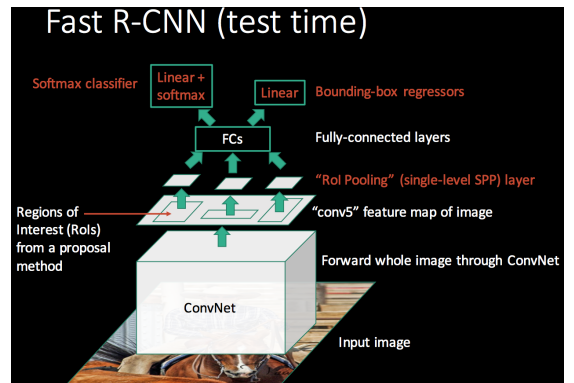
Similar architecture to SPP-net

- Mini-batches constructed via **hierarchical sampling**: Sample a similar number of RoIs over a **smaller number of images**



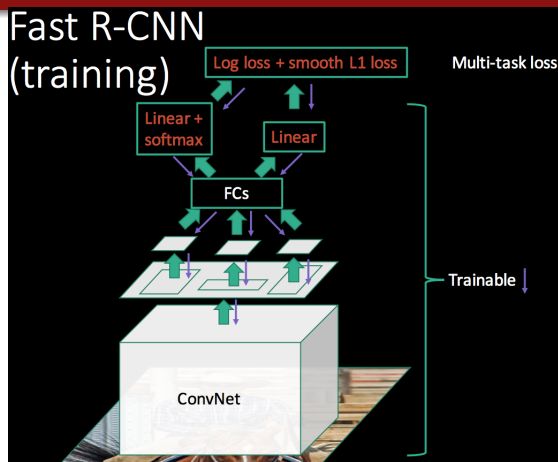
## Fast R-CNN (Girshick 2015)

Example from Girshick (2015)



## Fast R-CNN (Girshick 2015)

Example from Girshick (2015)



## You Only Look Once (Redmon et al. 2016)

- A single, unified network
- Can process 45 frames per second on a GPU (155 fps for Fast YOLO)
- Lower mAP than some R-CNN variants, but much faster
- Highest mAP of real-time detectors ( $\geq 30$  fps)

