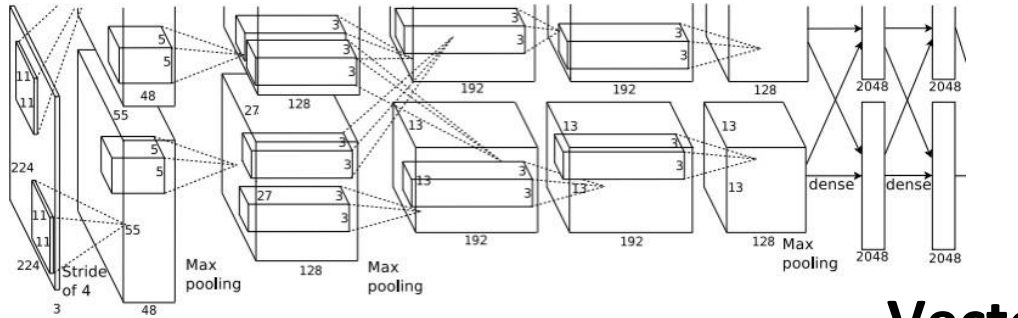# Deep Learning

# So far: Image Classification



**Vector:**
4096

**Fully-Connected:**
4096 to 1000

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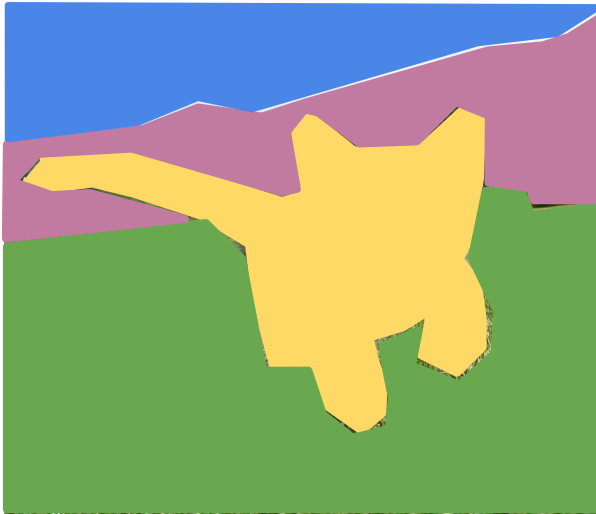...

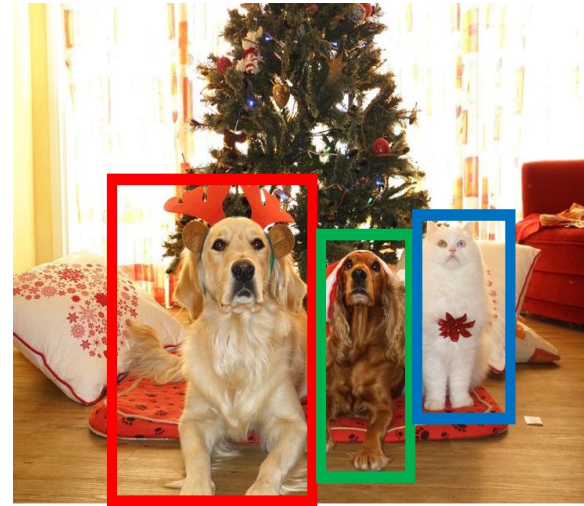# Computer Vision Tasks

**Classification**

**Semantic Segmentation**

**Object Detection**

**Instance Segmentation**

**CAT**

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

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

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

No spatial extent

Multiple Objects

# Today: Object Detection

**Classification**



CAT

No spatial extent

**Semantic Segmentation**



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

**Object Detection**



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

**Instance Segmentation**



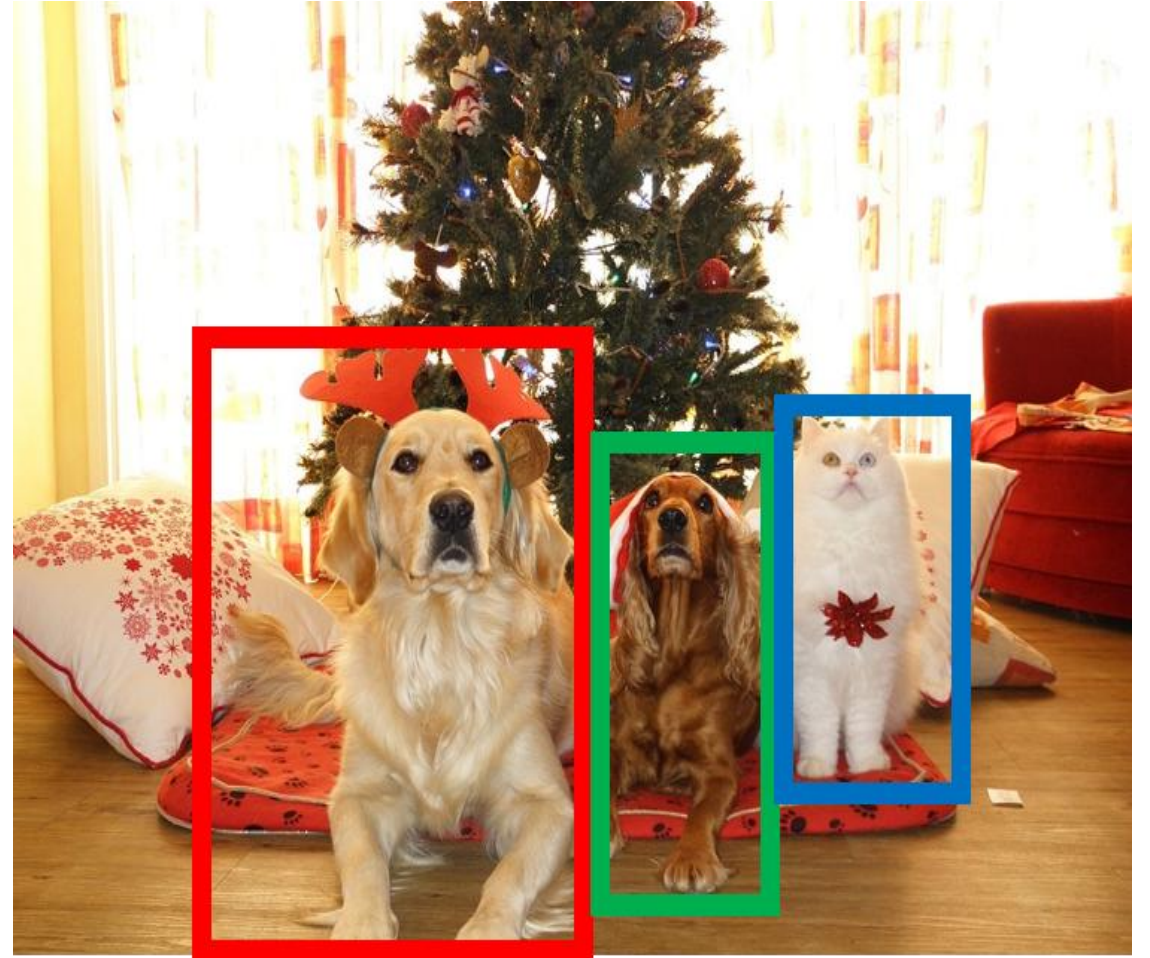**DOG**, **DOG**, **CAT**

Multiple Objects

# Object Detection: Task Definition

**Input**: Single RGB Image

**Output**: A set 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 or higher

# Detecting a single object

**Vector:**
4096

# Detecting a single object

"What"

**Correct label:**
Cat



**Fully Connected:**
4096 to 1000

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

**Softmax Loss**

**Vector:**
4096

# Detecting a single object

"What"

Correct label: Cat

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

**Fully Connected**: 4096 to 1000

**Softmax Loss**

**Vector:** 4096

Treat localization as a regression problem

**Fully Connected**: 4096 to 4

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

"Where"

**L2 Loss**

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

# Detecting a single object

"What"

**Correct label:**
Cat

**Fully Connected**: 4096 to 1000

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

**Softmax Loss**



**Vector:** 4096

Treat localization as a regression problem

**Fully Connected**: 4096 to 4

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

"Where"

**Weighted Sum**

**Loss**

**L2 Loss**

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

# 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**

**Vector:** 4096

Treat localization as a regression problem

**Fully Connected:** 4096 to 4

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

"Where"

**L2 Loss**

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

Multitask Loss

Weighted Sum

**Loss**

# Detecting a single object

"What"

Often pretrained on ImageNet (Transfer learning)

**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

Treat localization as a regression problem

**Fully Connected**: 4096 to 4

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

"Where"

**L2 Loss**

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

# Detecting Multiple Objects



CAT: (x, y, w, h)

4 numbers

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

12 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

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

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



Dog? NO
Cat? YES
Background? NO

# Detecting Multiple Objects: Sliding Window

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



Consider a box of size h x w in an image 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)

# Detecting Multiple Objects: Sliding Window



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

800 x 600 image
has ~58M boxes

Consider a box of size h x w in an image of size H x W:

Possible x positions: W – w + 1

Possible y positions: H – h + 1

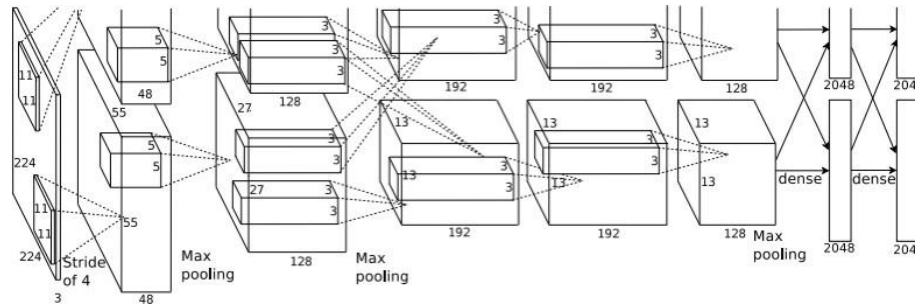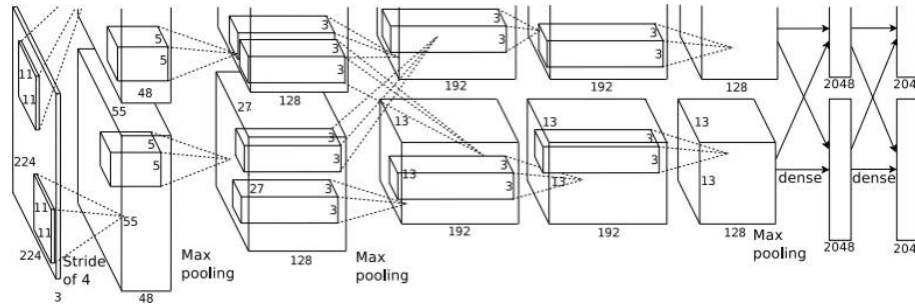Possible positions: (W – w + 1) * (H – h + 1)

# 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

# R-CNN: Region-Based CNN

Input
image



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

# 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.

# R-CNN: Region-Based CNN



Warped image regions (224x224)

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.

# R-CNN: Region-Based CNN



Forward each region through ConvNet

Warped image regions (224x224)

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..

# R-CNN: Region-Based CNN

Classify each region



Forward each region through ConvNet

Class

Warped image regions (224x224)

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.

# R-CNN: Region-Based CNN



Classify each region

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

Forward each
region through
ConvNet

Warped image
regions (224x224)

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

Input
image

Girshick et al, "Rich feature hierarchies for accurate object detection and
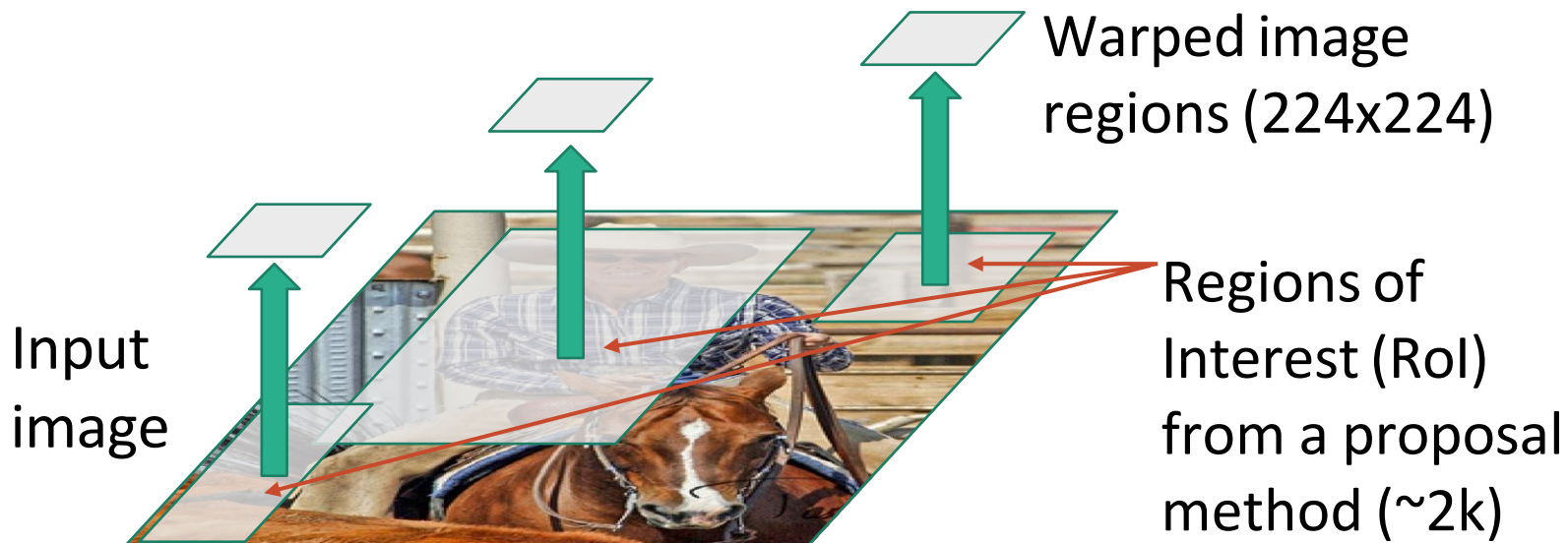semantic segmentation", CVPR 2014..
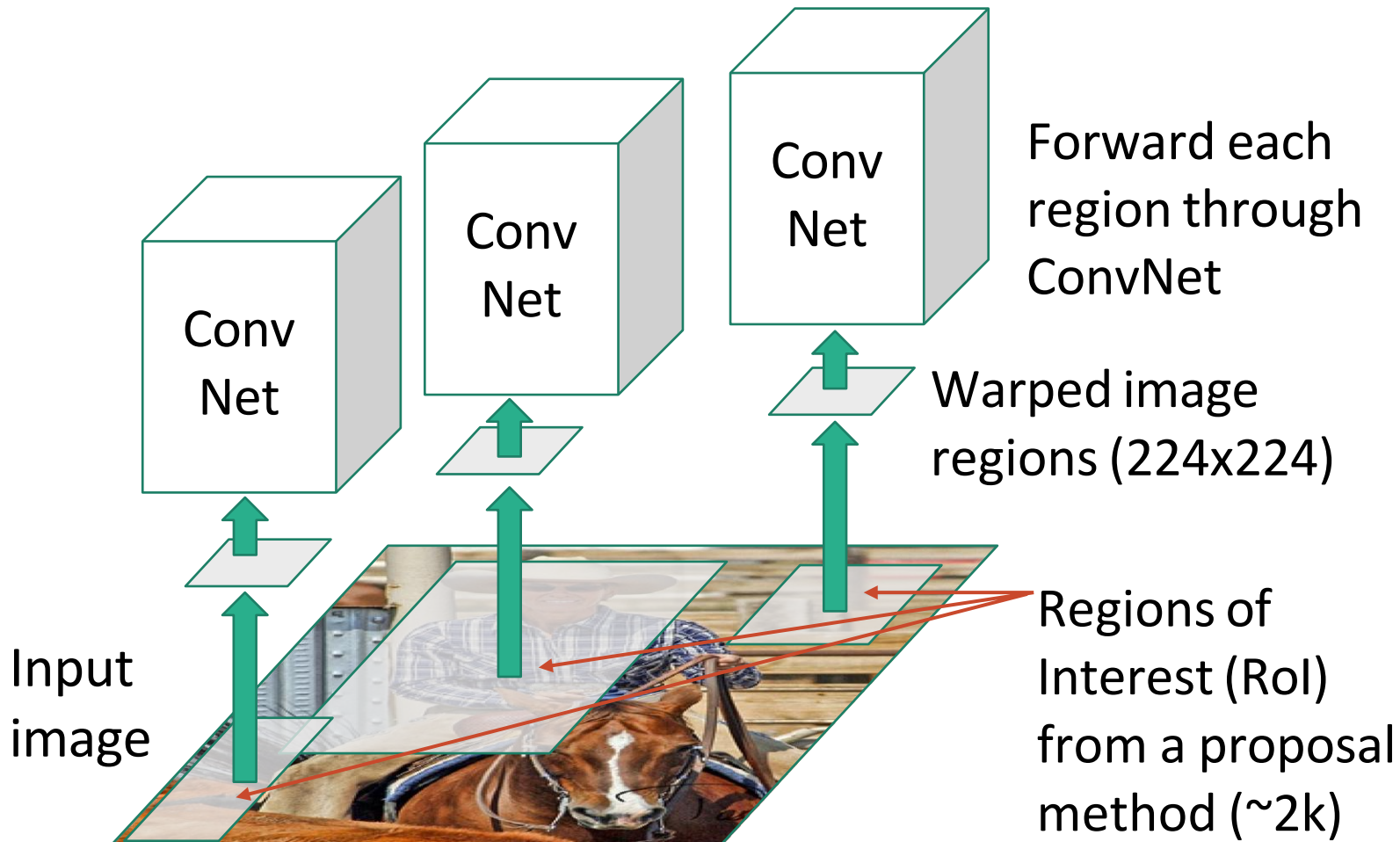
# R-CNN: Region-Based CNN



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

Conv Net

Conv Net

Conv Net

Forward each region through ConvNet

Warped image regions (224x224)

Input image

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

Region proposal: ($p_x$, $p_y$, $p_h$, $p_w$)
Transform: ($t_x$, $t_y$, $t_h$, $t_w$)
Output box: ($b_x$, $b_y$, $b_h$, $b_w$)

Translate relative to box size:
$b_x = p_x + p_w t_x$        $b_y = p_y + p_h t_y$

Log-space scale transform:
$b_w = p_w exp(t_w)$      $b_h = p_h exp(t_h)$

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

# R-CNN: Test-time



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
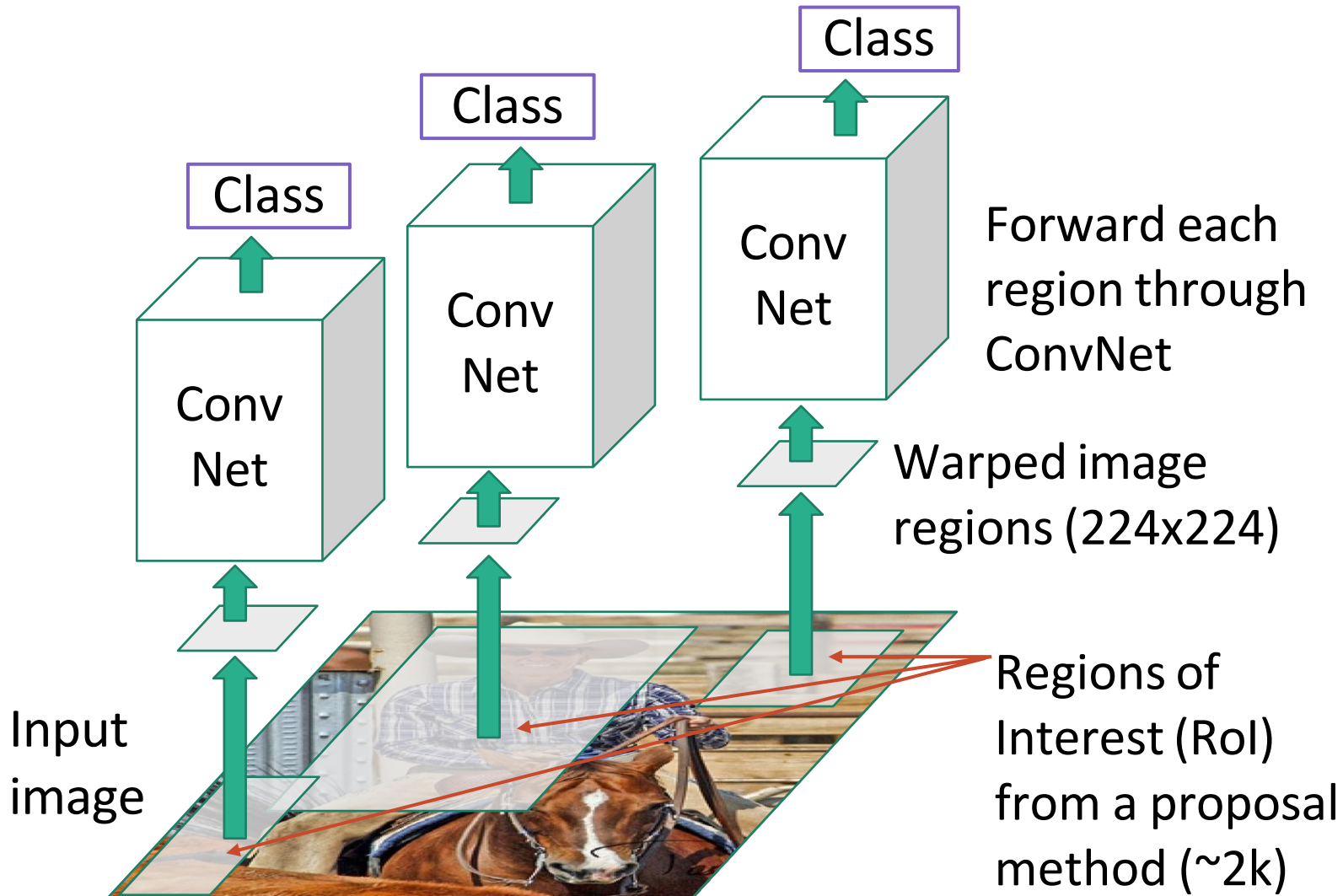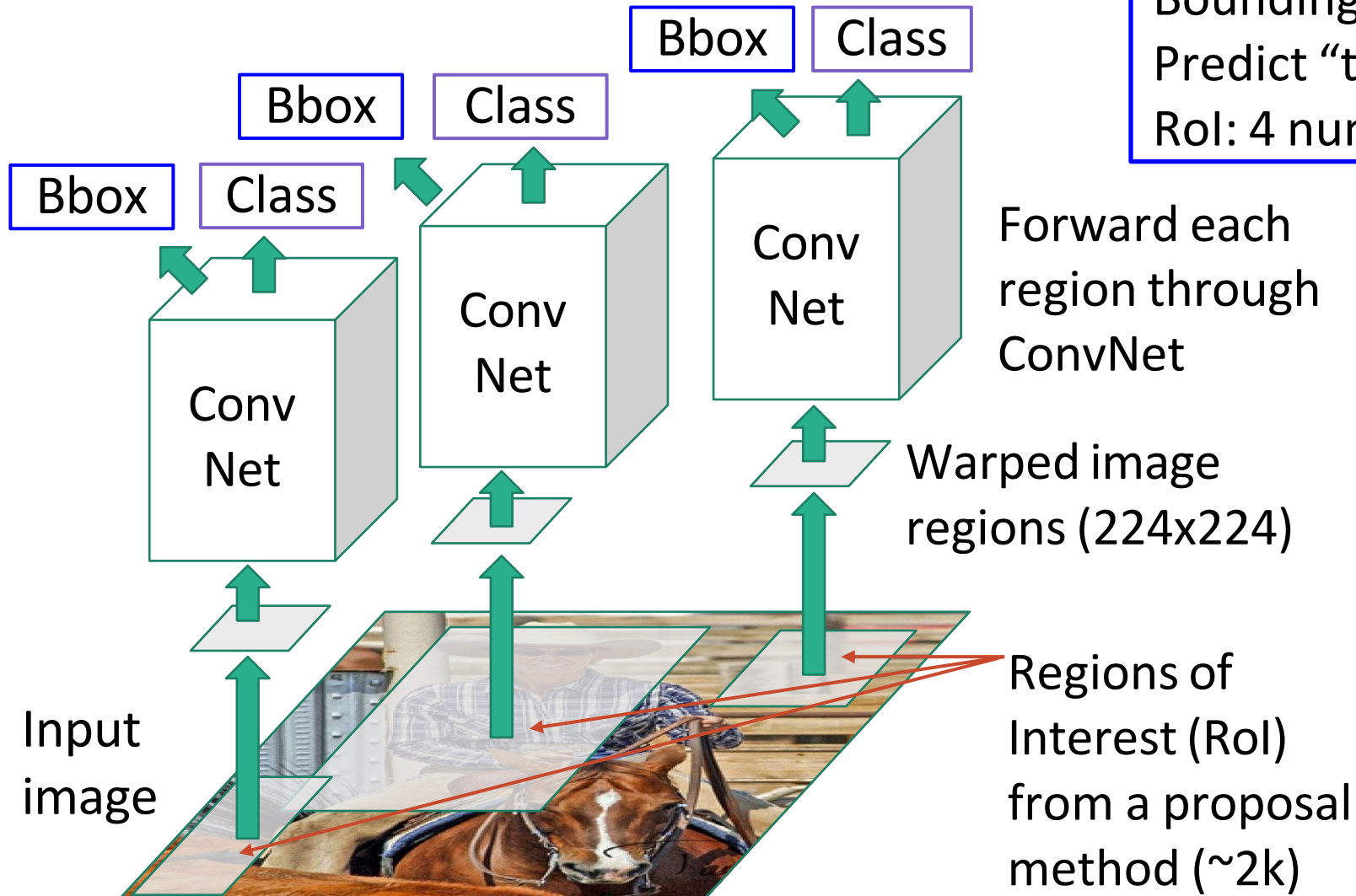   (Many choices here: threshold on background, per-category, or take top K proposals per image)
4. Compare with ground-truth boxes

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

# Comparing Boxes: Intersection over Union (IoU)
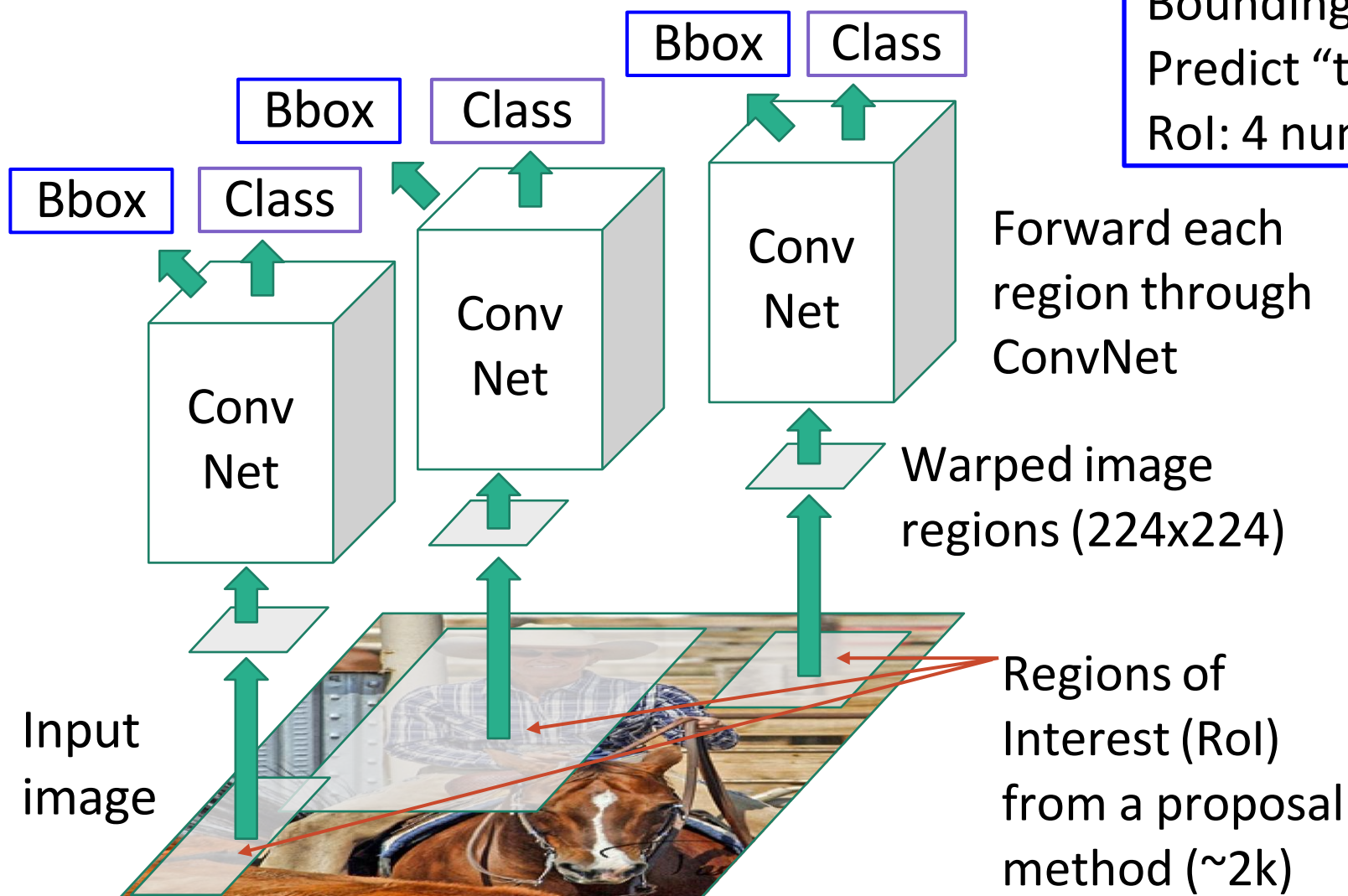
How can we compare our
prediction to the ground-truth box?

# Comparing Boxes: Intersection over Union (IoU)

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

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

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

# Comparing Boxes: Intersection over Union (IoU)

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

**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"

# Comparing Boxes: Intersection over Union (IoU)

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

**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",

# Comparing Boxes: Intersection over Union (IoU)

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

**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",
IoU > 0.9 is "almost perfect"

# Overlapping Boxes: Non-Max Suppression (NMS)

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

# 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

# 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

IoU( ■ , ■ ) = **0.74**

# 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
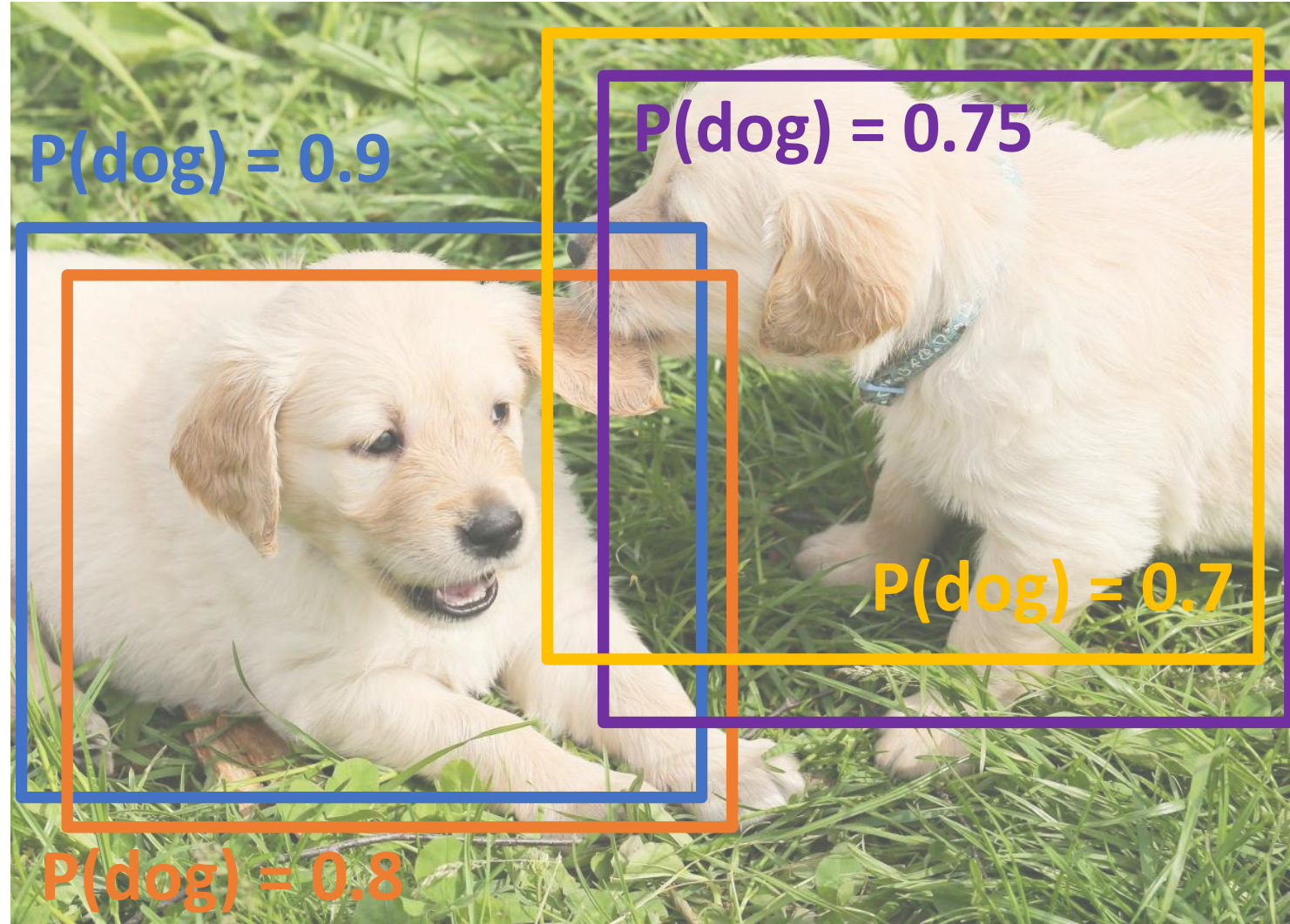


P(dog) = 0.9

P(dog) = 0.75

# 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

**Problem**: NMS may eliminate "good" boxes when objects are highly overlapping → Soft-NMS

# Evaluating Object Detectors:
# Mean Average Precision (mAP)

1. Run object detector on all test images (with NMS)
2. For each category, compute Average Precision (AP) =
   area under Precision vs Recall Curve

# Evaluating Object Detectors: Mean Average Precision (mAP)

1. Run object detector on all test images (with NMS)
2. For each category, compute Average Precision (AP) = area under Precision vs Recall Curve
   1. For each detection (highest score to lowest score)

All dog detections sorted by score

| 0.99 | 0.95 | 0.90 | 0.5 | 0.50 |

All ground-truth dog boxes

# Evaluating Object Detectors: Mean Average Precision (mAP)

1. Run object detector on all test images (with NMS)
2. For each category, compute Average Precision (AP) = area under Precision vs Recall Curve
   1. For each detection (highest score to lowest score)
      1. If it matches some GT box with IoU > 0.5, mark it as positive and eliminate the GT
      2. Otherwise mark it as negative

All dog detections sorted by score

| 0.99 | 0.95 | 0.90 | 0.5 | 0.50 |
|------|------|------|-----|------|

Match: IoU > 0.5

All ground-truth dog boxes

# Evaluating Object Detectors: Mean Average Precision (mAP)

1. Run object detector on all test images (with NMS)
2. For each category, compute Average Precision (AP) = area under Precision vs Recall Curve
   1. For each detection (highest score to lowest score)
      1. If it matches some GT box with IoU > 0.5, mark it as positive and eliminate the GT
      2. Otherwise mark it as negative
      3. Plot a point on PR Curve

All dog detections sorted by score

| 0.99 | 0.95 | 0.90 | 0.5 | 0.50 |

Match: IoU > 0.5

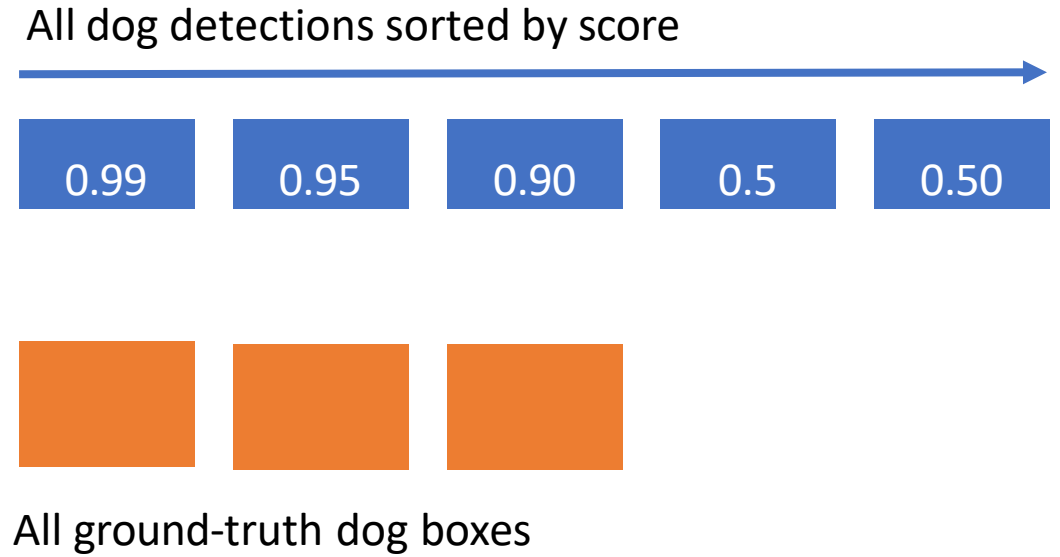All ground-truth dog boxes

Precision = 1/1 = 1.0
Recall = 1/3 = 0.33

# Evaluating Object Detectors: Mean Average Precision (mAP)

1. Run object detector on all test images (with NMS)
2. For each category, compute Average Precision (AP) = area under Precision vs Recall Curve
   1. For each detection (highest score to lowest score)
      1. If it matches some GT box with IoU > 0.5, mark it as positive and eliminate the GT
      2. Otherwise mark it as negative
      3. Plot a point on PR Curve

All dog detections sorted by score

| 0.99 | 0.95 | 0.90 | 0.5 | 0.50 |

Match: IoU > 0.5

All ground-truth dog boxes

Precision = 2/2 = 1.0
Recall = 2/3 = 0.67

Precision

Recall            1.0

# Evaluating Object Detectors: Mean Average Precision (mAP)

All dog detections sorted by score

| 0.99 | 0.95 | 0.90 | 0.5 | 0.50 |

No match > 0.5 IoU with GT

All ground-truth dog boxes

1. Run object detector on all test images (with NMS)
2. For each category, compute Average Precision (AP) = area under Precision vs Recall Curve
   1. For each detection (highest score to lowest score)
      1. If it matches some GT box with IoU > 0.5, mark it as positive and eliminate the GT
      2. Otherwise mark it as negative
      3. Plot a point on PR Curve

Precision = 2/3 = 0.67
Recall = 2/3 = 0.67
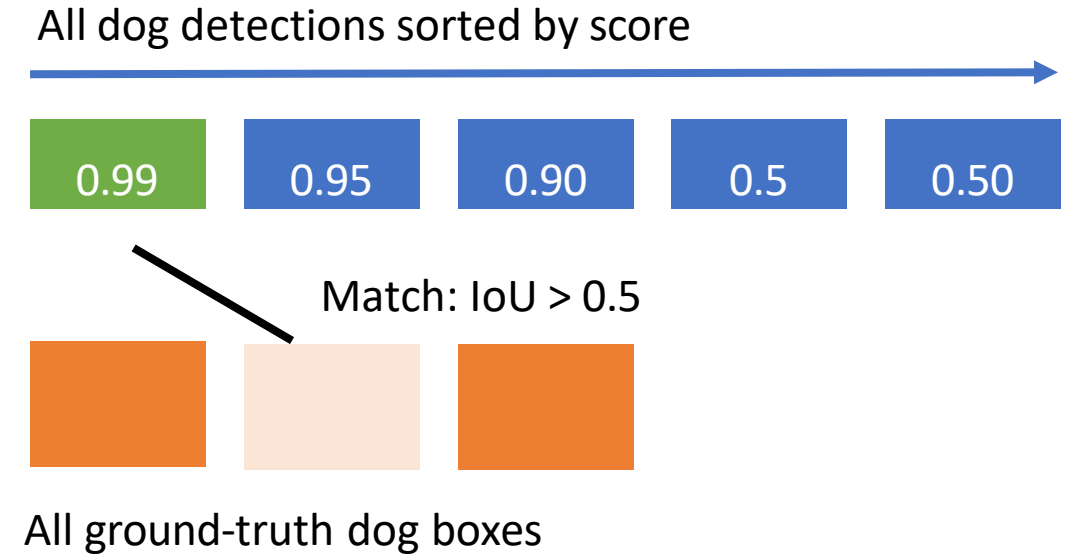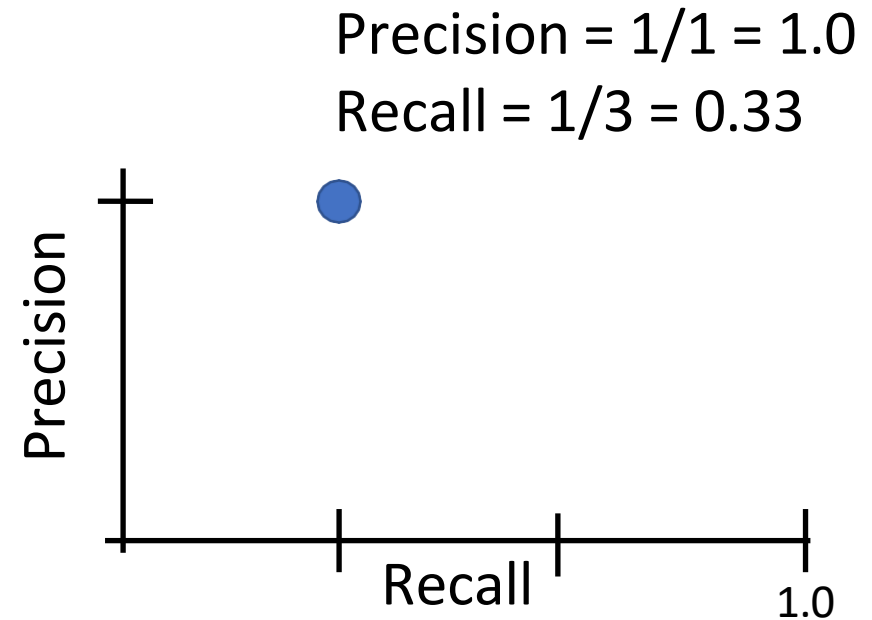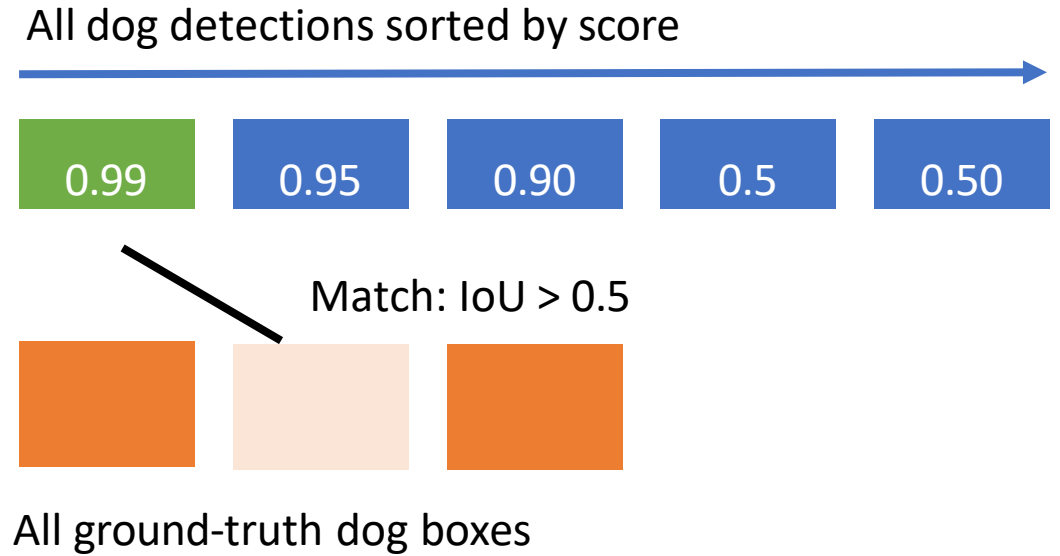
Precision

Recall            1.0

# Evaluating Object Detectors: Mean Average Precision (mAP)

1. Run object detector on all test images (with NMS)
2. For each category, compute Average Precision (AP) = area under Precision vs Recall Curve
   1. For each detection (highest score to lowest score)
      1. If it matches some GT box with IoU > 0.5, mark it as positive and eliminate the GT
      2. Otherwise mark it as negative
      3. Plot a point on PR Curve

All dog detections sorted by score

| 0.99 | 0.95 | 0.90 | 0.5 | 0.50 |

No match > 0.5 IoU with GT

All ground-truth dog boxes

Precision = 2/4 = 0.5
Recall = 2/3 = 0.67

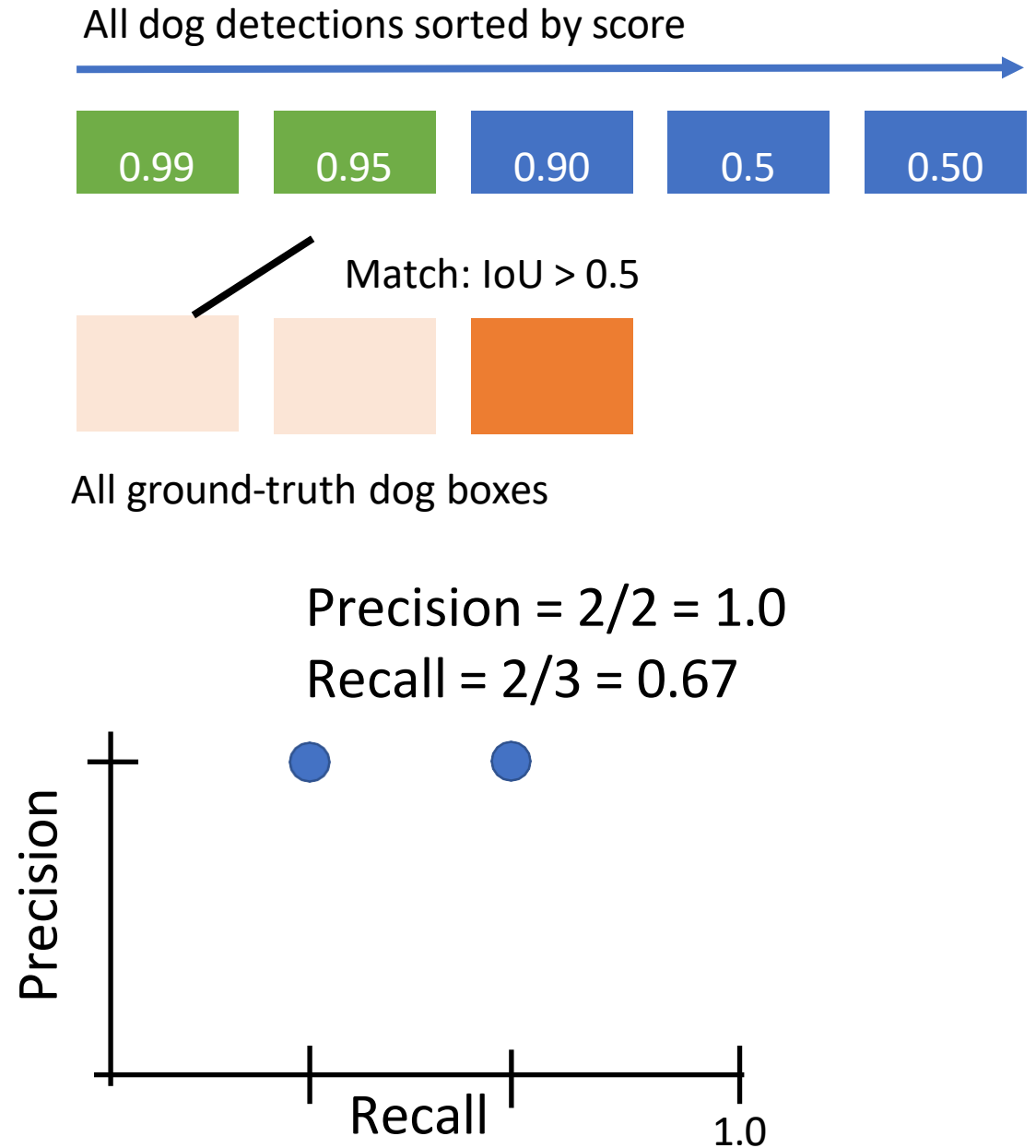# Evaluating Object Detectors: Mean Average Precision (mAP)

1. Run object detector on all test images (with NMS)
2. For each category, compute Average Precision (AP) = area under Precision vs Recall Curve
   1. For each detection (highest score to lowest score)
      1. If it matches some GT box with IoU > 0.5, mark it as positive and eliminate the GT
      2. Otherwise mark it as negative
      3. Plot a point on PR Curve

All dog detections sorted by score

| 0.99 | 0.95 | 0.90 | 0.5 | 0.50 |

Match: > 0.5 IoU

All ground-truth dog boxes

Precision = 3/5 = 0.6
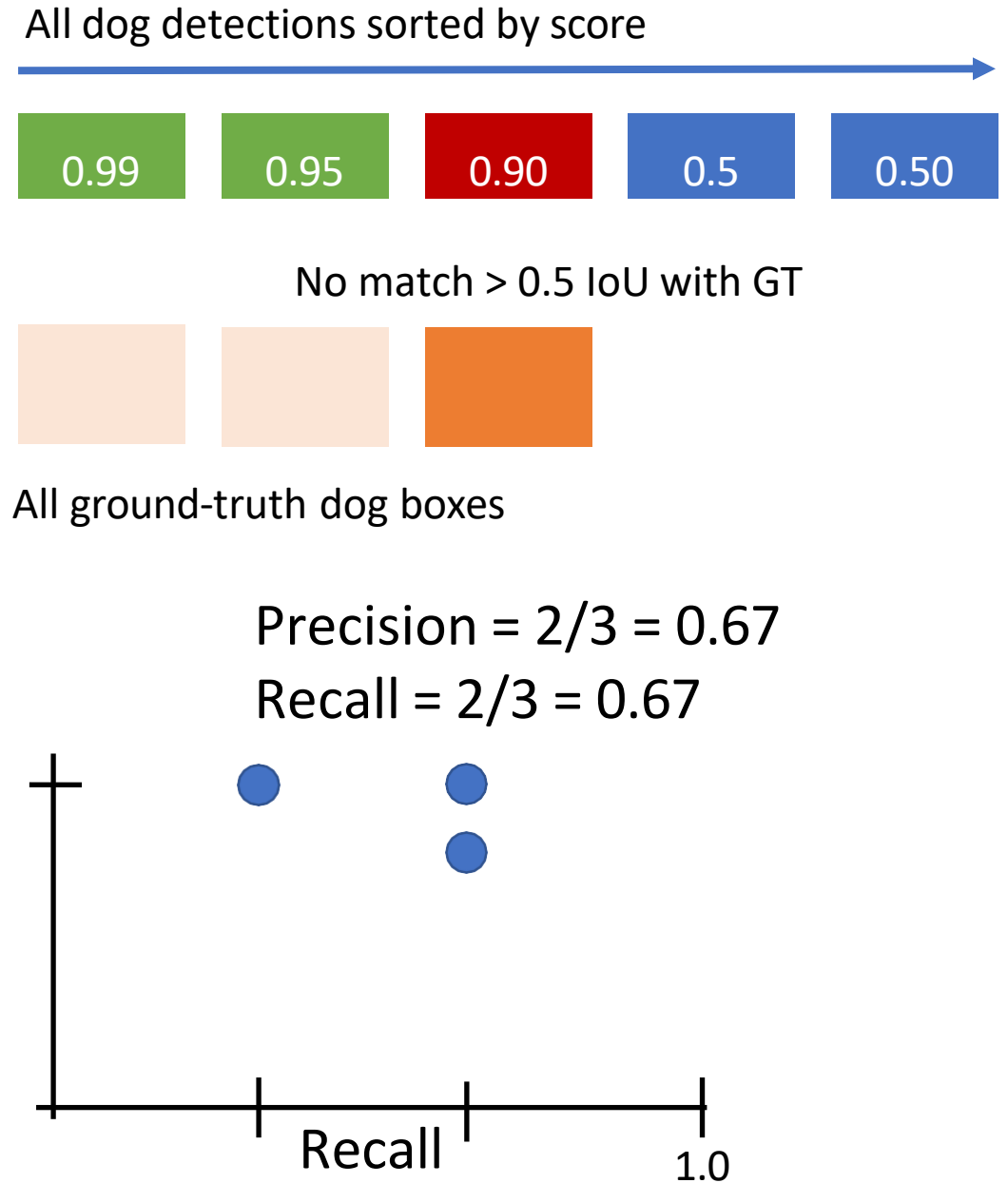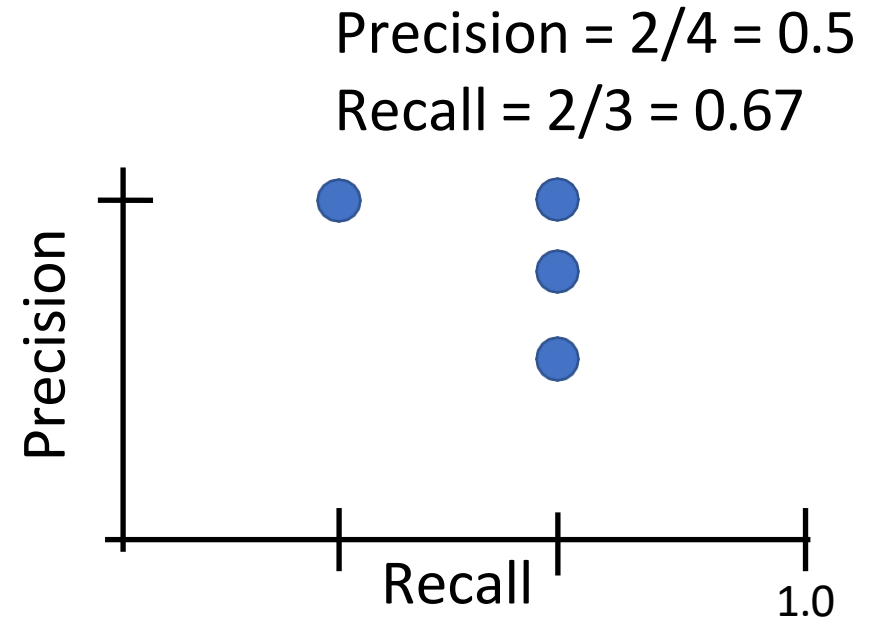Recall = 3/3 = 1.0

# Evaluating Object Detectors: Mean Average Precision (mAP)

1. Run object detector on all test images (with NMS)
2. For each category, compute Average Precision (AP) = area under Precision vs Recall Curve
   1. For each detection (highest score to lowest score)
      1. If it matches some GT box with IoU > 0.5, mark it as positive and eliminate the GT
      2. Otherwise mark it as negative
      3. Plot a point on PR Curve
   2. Average Precision (AP) = area under PR curve

All dog detections sorted by score

| 0.99 | 0.95 | 0.90 | 0.5 | 0.50 |
|------|------|------|-----|------|

All ground-truth dog boxes

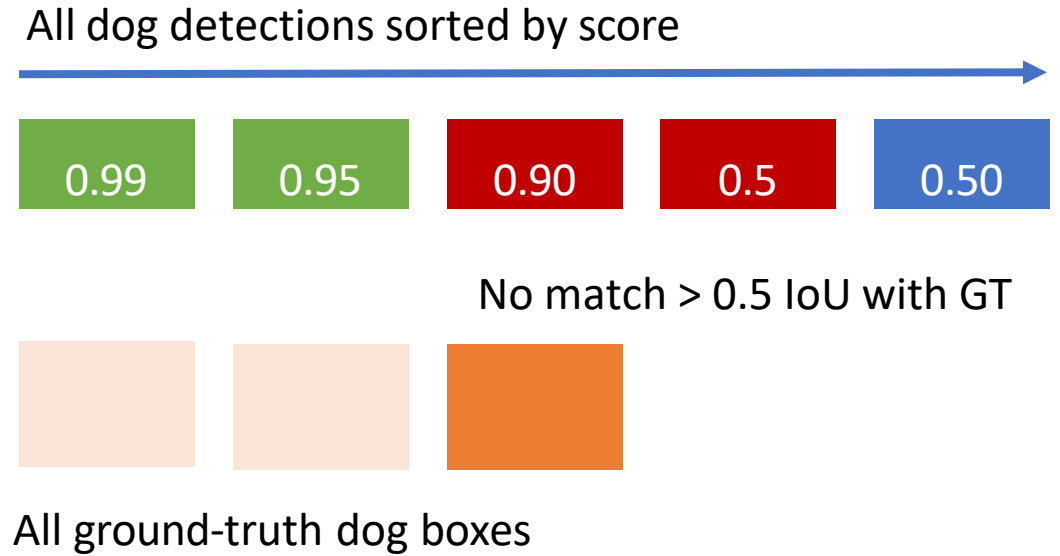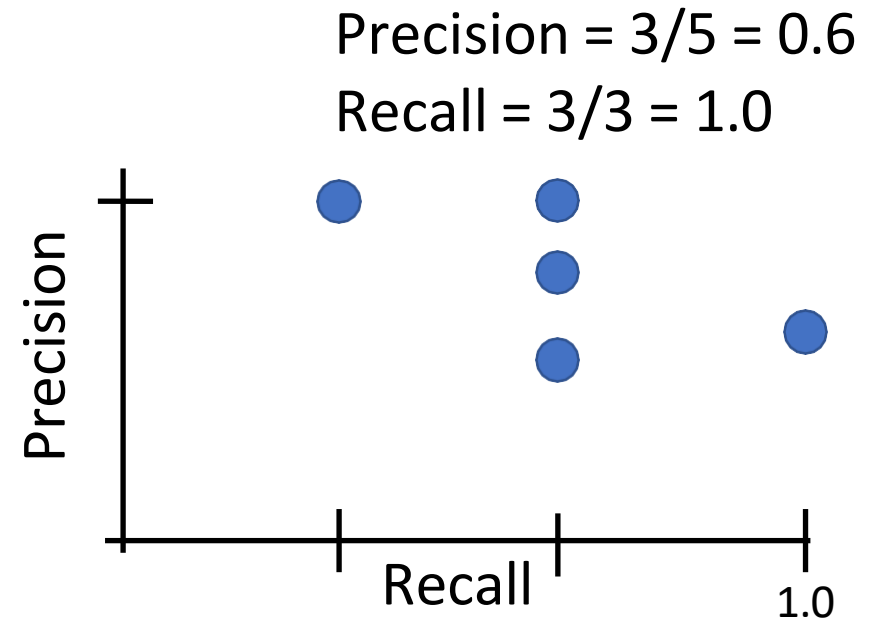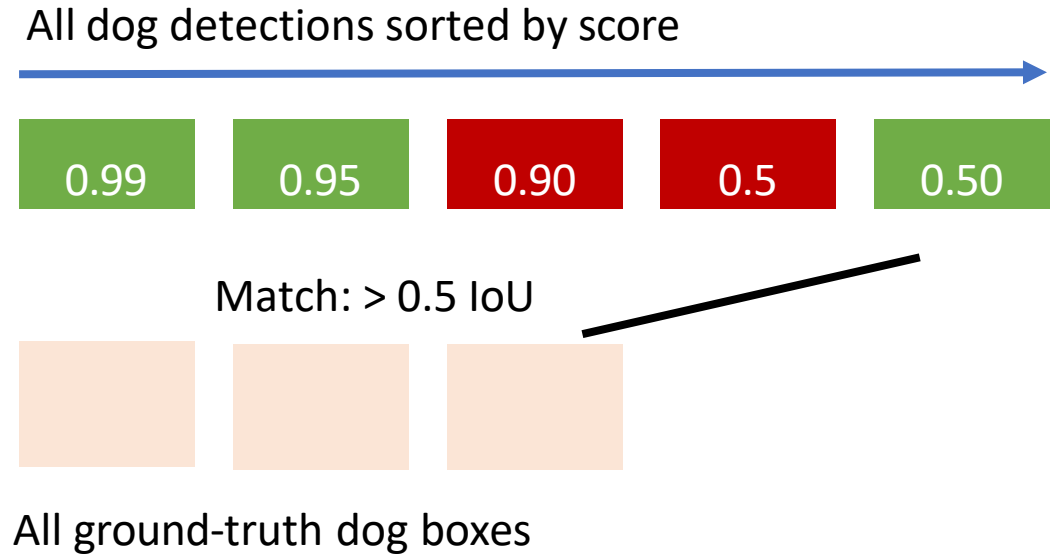Dog AP = 0.86

Precision

Recall 1.0

# Evaluating Object Detectors: Mean Average Precision (mAP)

1. Run object detector on all test images (with NMS)
2. For each category, compute Average Precision (AP) = area under Precision vs Recall Curve
   1. For each detection (highest score to lowest score)
      1. If it matches some GT box with IoU > 0.5, mark it as positive and eliminate the GT
      2. Otherwise mark it as negative
      3. Plot a point on PR Curve
   2. Average Precision (AP) = area under PR curve
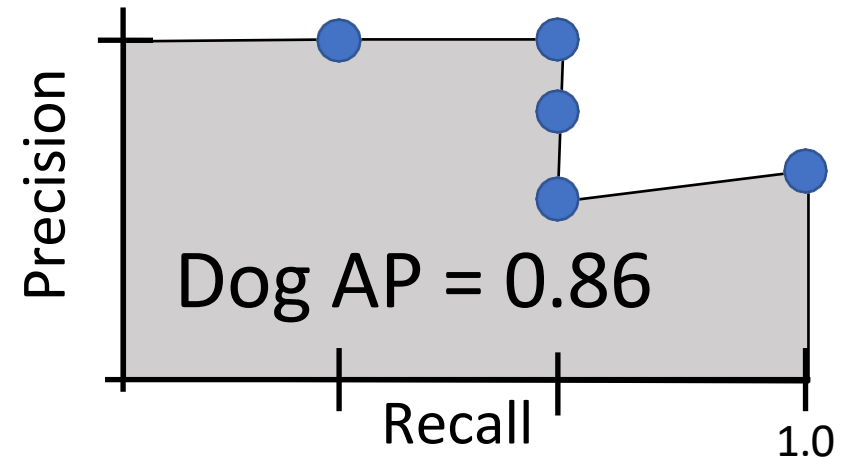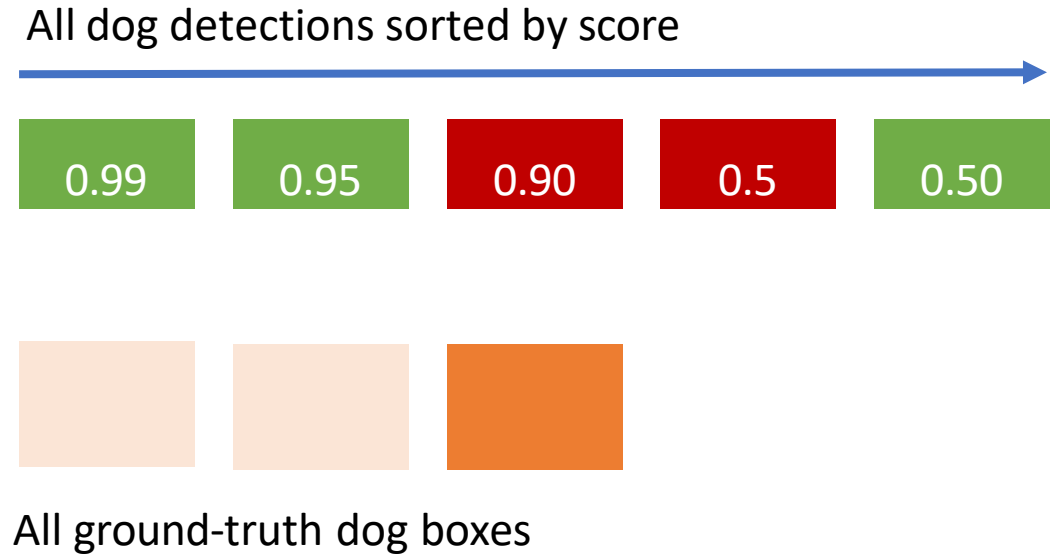3. Mean Average Precision (mAP) = average of AP for each category

Car AP = 0.65
Cat AP = 0.80
Dog AP = 0.86
mAP@0.5 = 0.77

# 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

Conv Net

Conv Net

Conv Net

Forward each region through ConvNet

Warped image regions (224x224)

Input image

Regions of Interest (RoI)

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

# R-CNN: Region-Based CNN



Classify each region

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

**Problem**: Very slow
Need to do ~2k forward
passes for each image

Forward each region through ConvNet

Warped image regions (224x224)

Regions of Interest (RoI)

Input image

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

# 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

Conv Net

Conv Net

Conv Net

Forward each region through ConvNet

Warped image regions (224x224)

Input image

Regions of Interest (RoI)

**Problem**: Very slow
Need to do ~2k forward passes for each image

**Solution**: Run CNN *before* warping

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

"Slow" R-CNN
Process each region independently

# Fast R-CNN

Input image

Bbox  Class

Bbox  Class

Bbox  Class

Conv Net

Conv Net

Conv Net

Input image

# Fast R-CNN



"Backbone"
network:
AlexNet, VGG,
ResNet, etc

ConvNet

Image features

Run whole image
through ConvNet

Input image

"Slow" R-CNN
Process each region
independently

Bbox  Class

Bbox  Class

Bbox  Class

Conv
Net

Conv
Net

Conv
Net

Input
image

# Fast R-CNN

"Slow" R-CNN
Process each region independently

Regions of Interest (RoIs) from a proposal method

Image features

"Backbone" network: AlexNet, VGG, ResNet, etc

ConvNet

Run whole image through ConvNet

Input image

Bbox Class

Bbox Class

Bbox Class

Conv Net

Conv Net

Conv Net

Input image

# Fast R-CNN



"Slow" R-CNN
Process each region independently

Regions of Interest (RoIs) from a proposal method

Crop + Resize features

Image features

"Backbone" network: AlexNet, VGG, ResNet, etc

Run whole image through ConvNet

ConvNet

Input image

Bbox     Class
Bbox     Class
Bbox     Class

Conv Net
Conv Net
Conv Net

Input image

# Fast R-CNN



Regions of Interest (RoIs) from a proposal method

"Backbone" network: AlexNet, VGG, ResNet, etc

Per-Region Network

Crop + Resize features

Image features

Run whole image through ConvNet

ConvNet

Input image

"Slow" R-CNN
Process each region independently

Bbox    Class
Bbox    Class
Bbox    Class

Conv Net
Conv Net
Conv Net

Input image

# Fast R-CNN

**Bbox** **Bbox** **Bbox**

**Class** **Class** **Class**

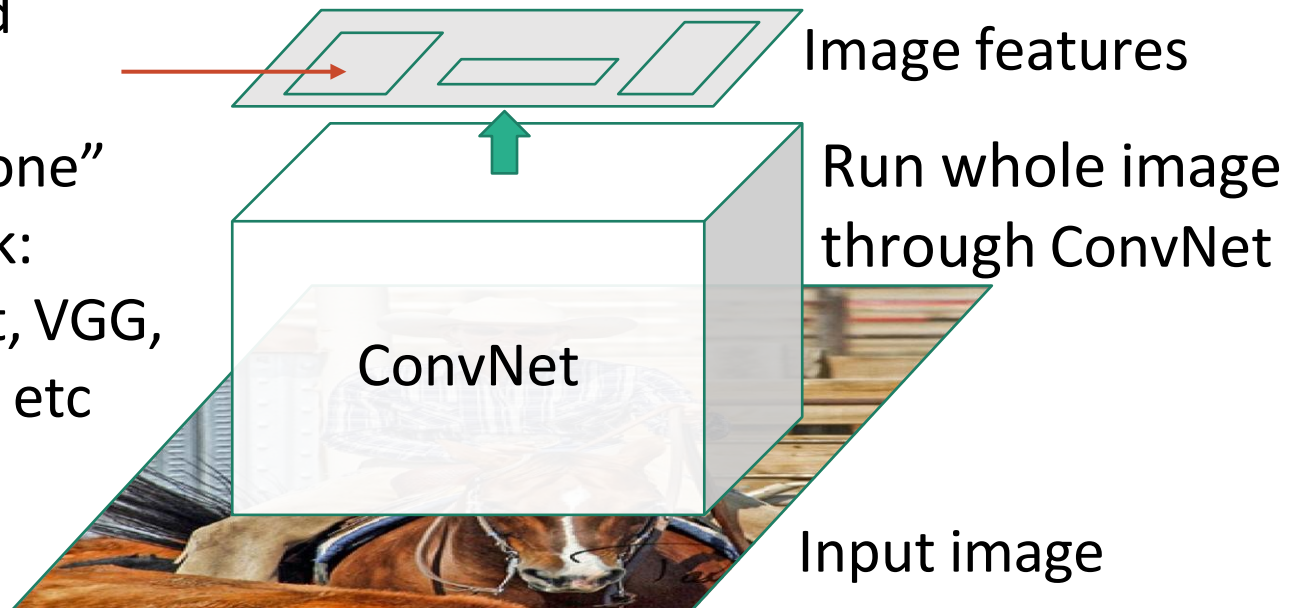Category and box transform per region

Regions of Interest (RoIs) from a proposal method

CNN CNN CNN — Per-Region Network

Crop + Resize features

Image features

"Backbone" network: AlexNet, VGG, ResNet, etc

ConvNet — Run whole image through ConvNet

Input image

## "Slow" R-CNN
Process each region independently

**Bbox** **Class**

**Bbox** **Class**

**Bbox** **Class**

ConvNet

ConvNet

ConvNet

Input image

# Fast R-CNN



Category and box transform per region

Per-Region network is relatively lightweight

Regions of Interest (RoIs) from a proposal method

Per-Region Network

Crop + Resize features

Image features

"Backbone" network: AlexNet, VGG, ResNet, etc

Run whole image through ConvNet

Most of the computation happens in backbone network; this saves work for overlapping region proposals

Input image

# Fast R-CNN



Bbox  Bbox  Bbox

Class  Class  Class

**Category and box transform per region**

Regions of Interest (RoIs) from a proposal method

Per-Region Network

CNN  CNN  CNN

**Crop + Resize features**

Image features

"Backbone" network: AlexNet, VGG, ResNet, etc

Run whole image through ConvNet

ConvNet

Input image

Example: When using AlexNet for detection, five conv layers are used for backbone and two FC layers are used for per-region network

# Fast R-CNN



Category and box transform per region

Regions of Interest (RoIs) from a proposal method

Per-Region Network

Crop + Resize features

Image features

"Backbone" network: AlexNet, VGG, ResNet, etc

Run whole image through ConvNet

Input image

Example: For ResNet, last stage is used as per-region network; the rest of the network is used as backbone

# Fast R-CNN

Bbox  Bbox  Bbox

Class  Class  Class

Category and box transform per region

Regions of Interest (RoIs) from a proposal method

CNN  CNN  CNN

Per-Region Network

Crop + Resize features

Image features

How to crop features?

"Backbone" network: AlexNet, VGG, ResNet, etc

ConvNet

Run whole image through ConvNet

Input image

# Cropping Features: RoI Pool



Input Image
(e.g. 3 x 640 x 480)

Girshick, "Fast R-CNN", ICCV 2015.

# Cropping Features: RoI Pool



Input Image
(e.g. 3 x 640 x 480)

CNN

Image features
(e.g. 512 x 20 x 15)

Girshick, "Fast R-CNN", ICCV 2015.

# Cropping Features: RoI Pool

Project proposal
onto features

CNN

Input Image
(e.g. 3 x 640 x 480)

Image features
(e.g. 512 x 20 x 15)

Girshick, "Fast R-CNN", ICCV 2015.

# Cropping Features: RoI Pool



"Snap" to grid cells

Project proposal onto features

CNN

Input Image
(e.g. 3 x 640 x 480)

Image features
(e.g. 512 x 20 x 15)

Girshick, "Fast R-CNN", ICCV 2015.

# Cropping Features: RoI Pool

"Snap" to grid cells

Divide into 2x2 grid of (roughly) equal subregions

Project proposal onto features



CNN

Input Image
(e.g. 3 x 640 x 480)

Image features
(e.g. 512 x 20 x 15)

Girshick, "Fast R-CNN", ICCV 2015.

# Cropping Features: RoI Pool



Girshick, "Fast R-CNN", ICCV 2015.

"Snap" to grid cells

Project proposal onto features

Divide into 2x2 grid of (roughly) equal subregions

Max-pool within each subregion

CNN

Input Image
(e.g. 3 x 640 x 480)

Image features
(e.g. 512 x 20 x 15)

Region features
(here 512 x 2 x 2;
In practice e.g 512 x 7 x 7)

Region features always the same size even if input regions have different sizes!

# Cropping Features: RoI Pool

"Snap" to grid cells

Divide into 2x2 grid of (roughly) equal subregions

Project proposal onto features

Max-pool within each subregion



CNN

Input Image
(e.g. 3 x 640 x 480)

Image features
(e.g. 512 x 20 x 15)

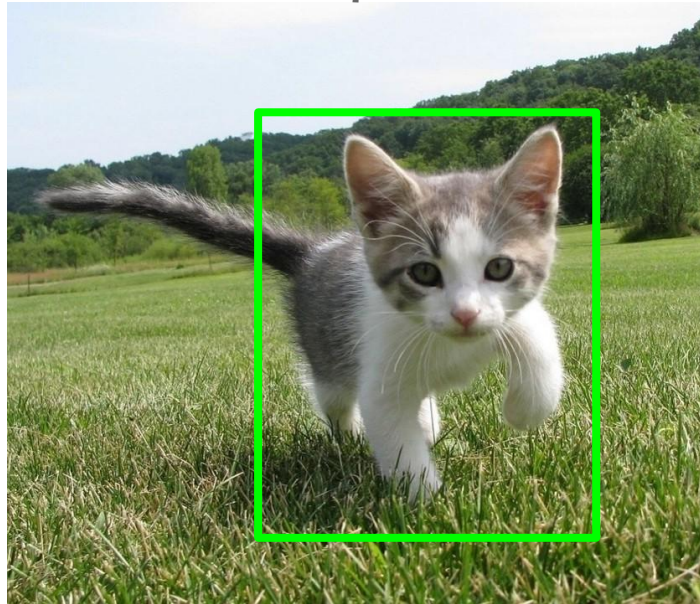Region features
(here 512 x 2 x 2;
In practice e.g 512 x 7 x 7)

Region features always the same size even if input regions have different sizes!

Girshick, "Fast R-CNN", ICCV 2015.

Problem: Slight misalignment due to snapping; different-sized subregions is weird

# Cropping Features: RoI <u>Align</u>

Divide into equal-sized subregions (may not be aligned to grid!)

Project proposal onto features

No "snapping"!

CNN

Input Image
(e.g. 3 x 640 x 480)

Image features
(e.g. 512 x 20 x 15)

He et al, "Mask R-CNN", ICCV 2017

# Cropping Features: RoI <u>Align</u>

Divide into equal-sized subregions
(may not be aligned to grid!)

Project proposal
onto features

No "snapping"!

Sample features at
regularly-spaced points
in each subregion using
**bilinear interpolation**

CNN

Input Image
(e.g. 3 x 640 x 480)

Image features
(e.g. 512 x 20 x 15)

He et al, "Mask R-CNN", ICCV 2017

# Cropping Features: RoI <u>Align</u>

Divide into equal-sized subregions
(may not be aligned to grid!)

No "snapping"!

Project proposal
onto features

Sample features at
regularly-spaced points
in each subregion using
**bilinear interpolation**

CNN

Input Image
(e.g. 3 x 640 x 480)

Image features
(e.g. 512 x 20 x 15)

Feature $f_{xy}$ for point (x, y) is a
linear combination of features
at its four neighboring grid cells:

# Cropping Features: RoI <u>Align</u>

Divide into equal-sized subregions (may not be aligned to grid!)

No "snapping"!

Project proposal onto features

Sample features at regularly-spaced points in each subregion using **bilinear interpolation**

CNN

After sampling, max-pool in each subregion

Input Image
(e.g. 3 x 640 x 480)

Image features
(e.g. 512 x 20 x 15)

Region features
(here 512 x 2 x 2;
In practice e.g 512 x 7 x 7)

He et al, "Mask R-CNN", ICCV 2017

# Fast R-CNN vs "Slow" R-CNN

**Fast R-CNN**: Apply differentiable cropping to shared image features
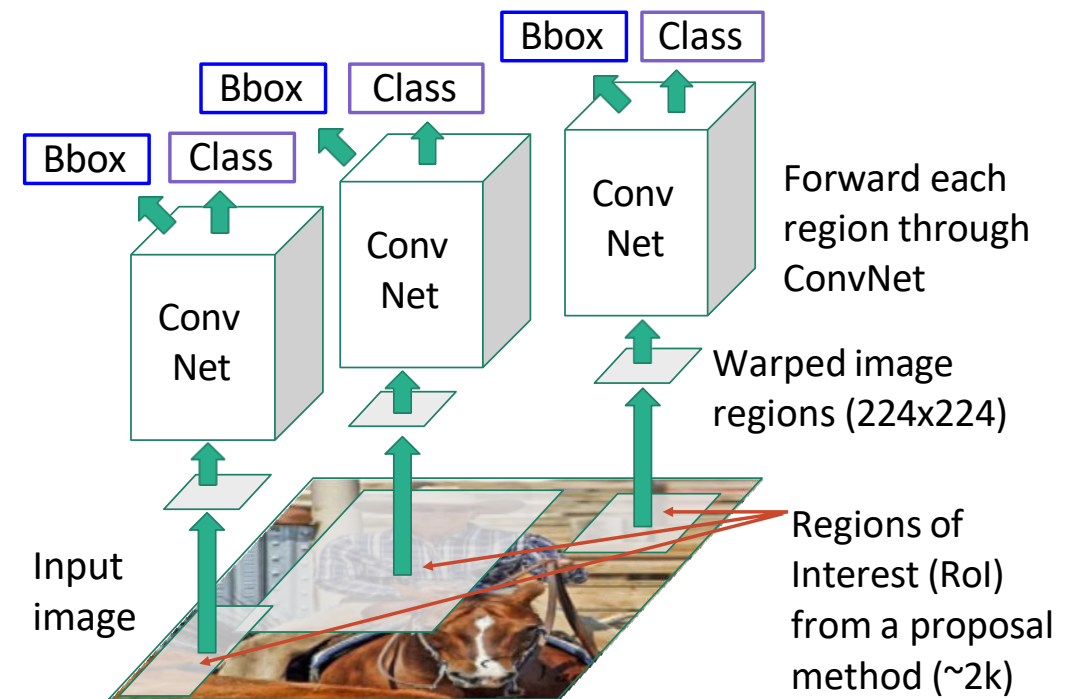
**"Slow" R-CNN**: Apply differentiable cropping to shared image features



Category and box transform per region

Regions of Interest (RoIs) from a proposal method

Per-Region Network

Crop + Resize features

Image features

"Backbone" network: AlexNet, VGG, ResNet, etc

Run whole image through ConvNet

Input image

Forward each region through ConvNet

Warped image regions (224x224)

Input image

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

# Fast R-CNN vs "Slow" R-CNN



Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.
He et al, "Spatial pyramid pooling in deep convolutional networks for visual recognition", ECCV 2014
Girshick, "Fast R-CNN", ICCV 2015

# Fast R-CNN vs "Slow" R-CNN



**Training time (Hours)**

| Model | Hours |
|-------|-------|
| R-CNN | 84 |
| SPP-Net | 25.5 |
| Fast R-CNN | 8.75 |

**Test time (seconds)**

Including Region propos... | Excluding Region Propo...

| Model | Including | Excluding |
|-------|-----------|-----------|
| R-CNN | 49 | 47 |
| SPP-Net | 4.3 | 2.3 |
| Fast R-CNN | 2.3 | 0.32 |

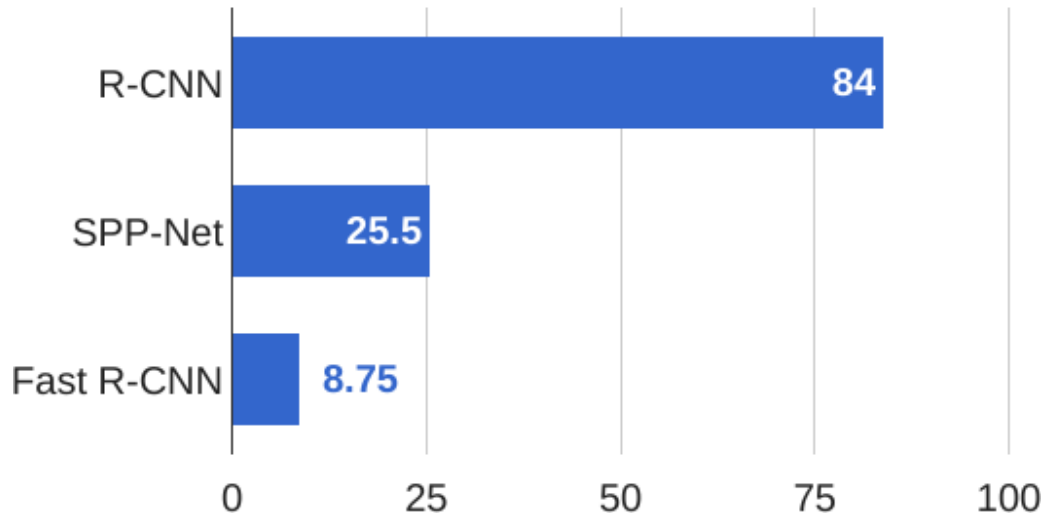**Problem**: Runtime dominated by region proposals

Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.
He et al, "Spatial pyramid pooling in deep convolutional networks for visual recognition", ECCV 2014
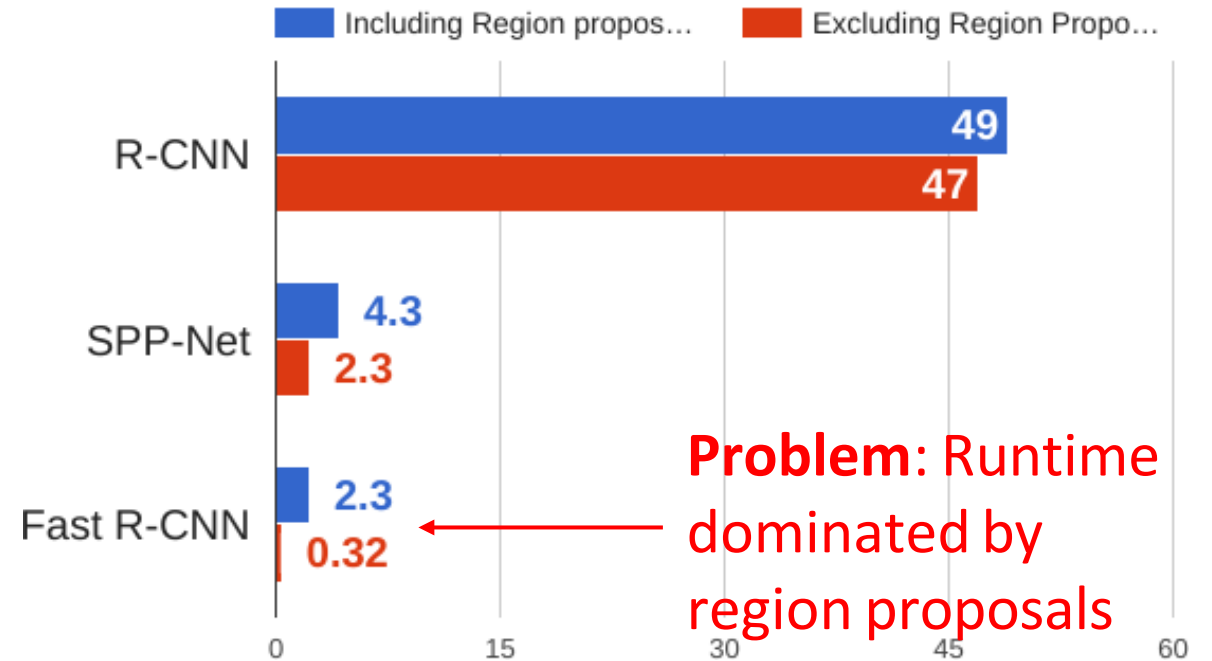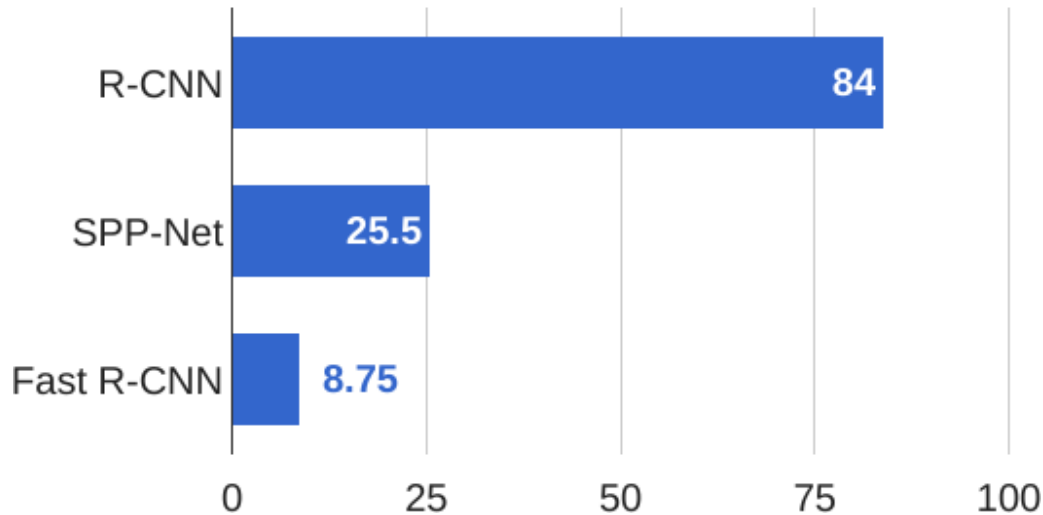Girshick, "Fast R-CNN", ICCV 2015

# Fast R-CNN vs "Slow" R-CNN



**Training time (Hours)**

R-CNN — 84
SPP-Net — 25.5
Fast R-CNN — 8.75

**Test time (seconds)**

Including Region propos... / Excluding Region Propo...

R-CNN — 49 / 47
SPP-Net — 4.3 / 2.3
Fast R-CNN — 2.3 / 0.32

**Problem**: Runtime dominated by region proposals

**Recall**: Region proposals computed by heuristic "Selective Search" algorithm on CPU -- let's learn them with a CNN instead

Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.
He et al, "Spatial pyramid pooling in deep convolutional networks for visual recognition", ECCV 2014
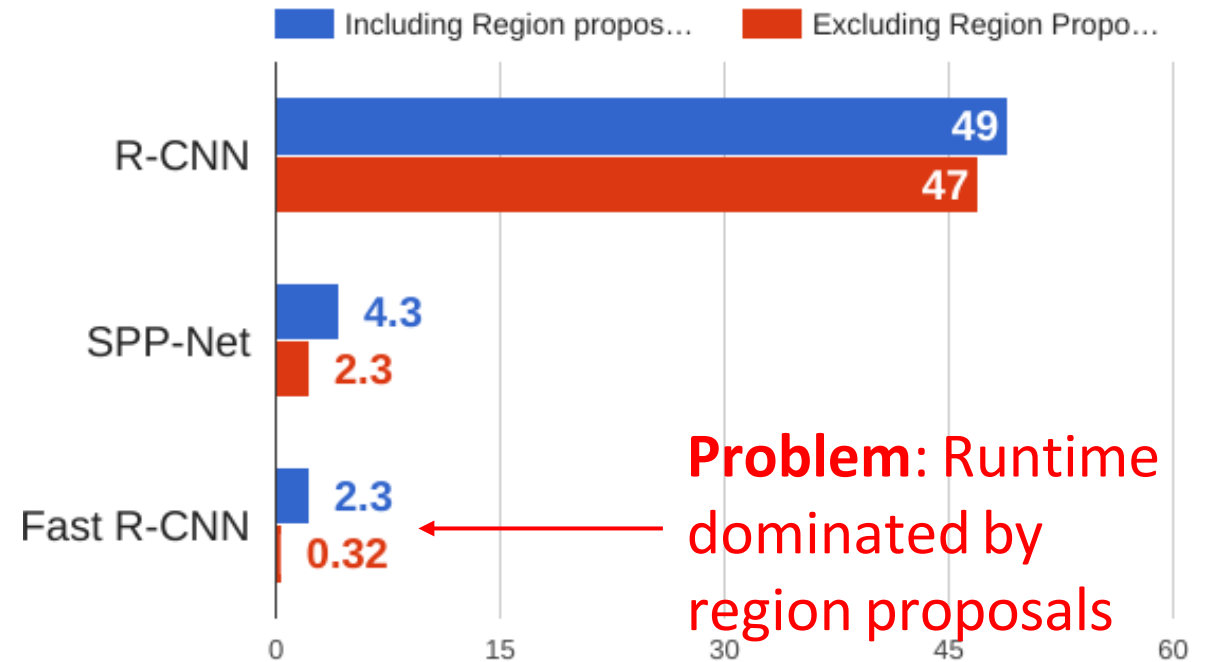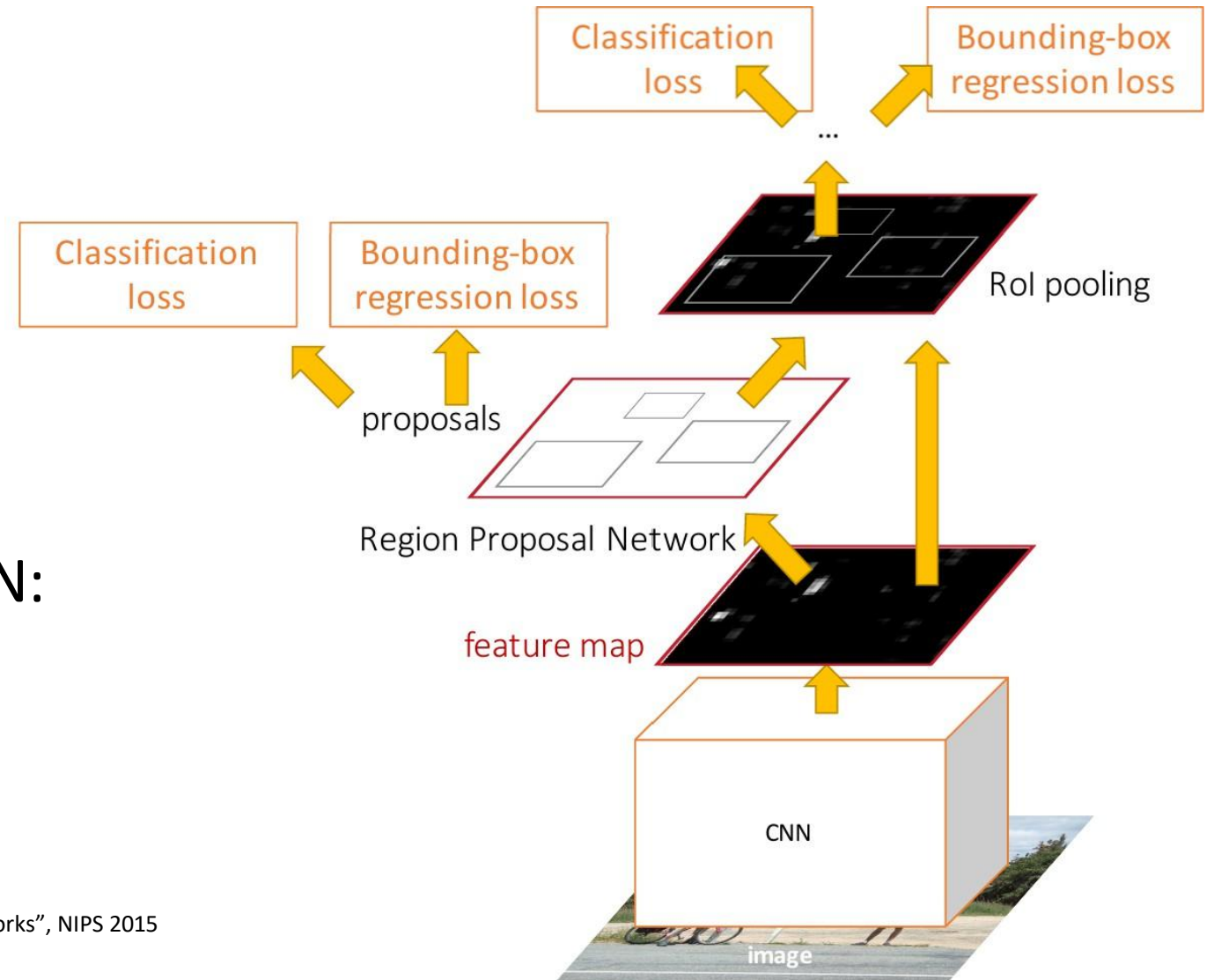Girshick, "Fast R-CNN", ICCV 2015

# Faster R-CNN: Learnable Region Proposals

Insert **Region Proposal Network (RPN)** to predict proposals from features

Otherwise same as Fast R-CNN:
Crop features for each proposal, classify each one



Classification loss

Bounding-box regression loss

...

RoI pooling

Classification loss

Bounding-box regression loss

proposals

Region Proposal Network

feature map

CNN

image

# Region Proposal Network (RPN)

Run backbone CNN to get
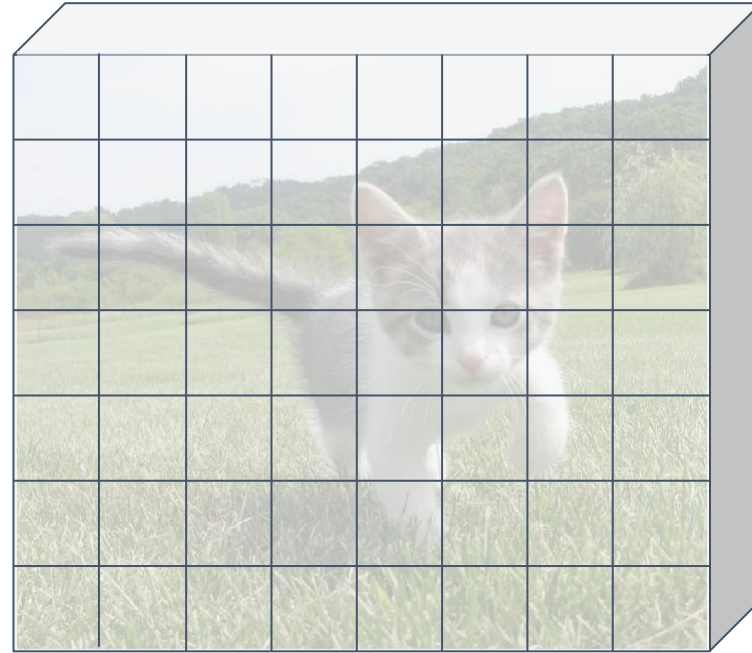features aligned to input image



CNN

Input Image
(e.g. 3 x 640 x 480)
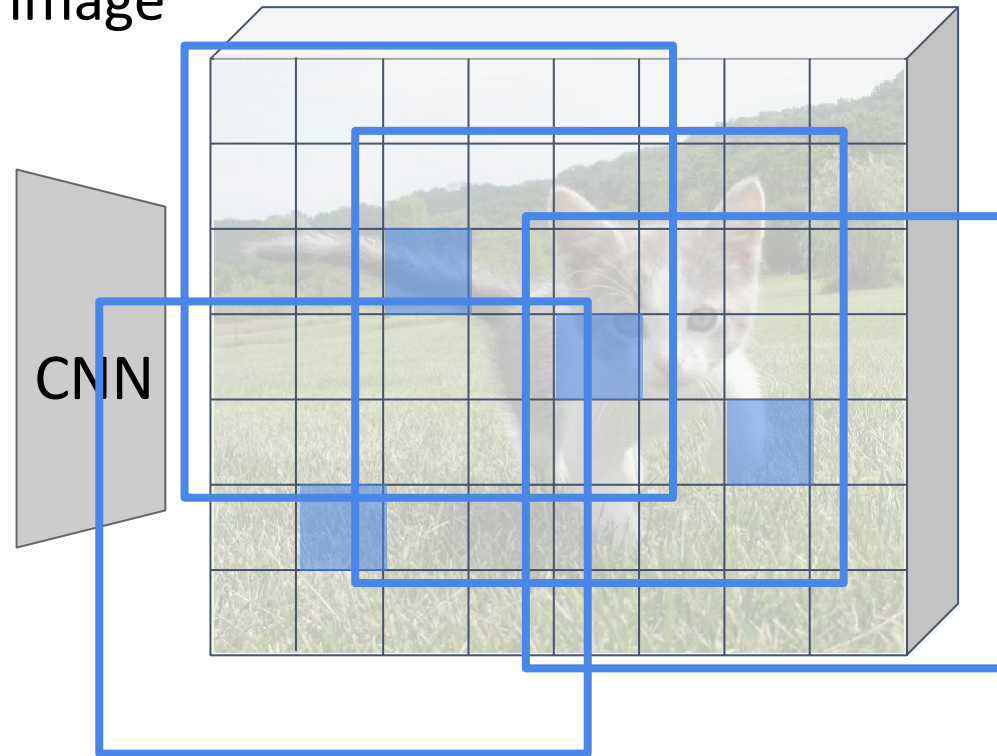
Image features
(e.g. 512 x 20 x 15)

# Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image

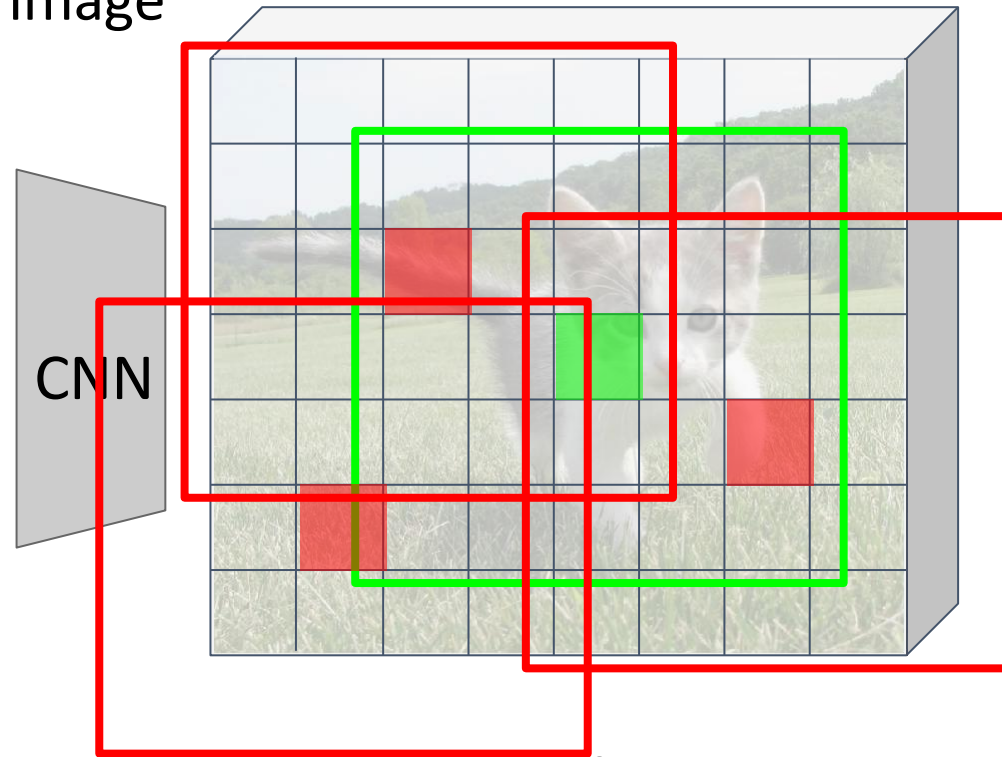Imagine an **anchor box** of fixed size at each point in the feature map



CNN

Input Image
(e.g. 3 x 640 x 480)

Image features
(e.g. 512 x 20 x 15)

# Region Proposal Network (RPN)

Imagine an anchor box of fixed size at each point in the feature map

Run backbone CNN to get features aligned to input image



CNN

Input Image
(e.g. 3 x 640 x 480)

Image features
(e.g. 512 x 20 x 15)

Conv

Anchor is an object?
1 x 20 x 15

At each point, predict whether the corresponding anchor contains an object (per-cell logistic regression, predict scores with conv layer)

# Region Proposal Network (RPN)

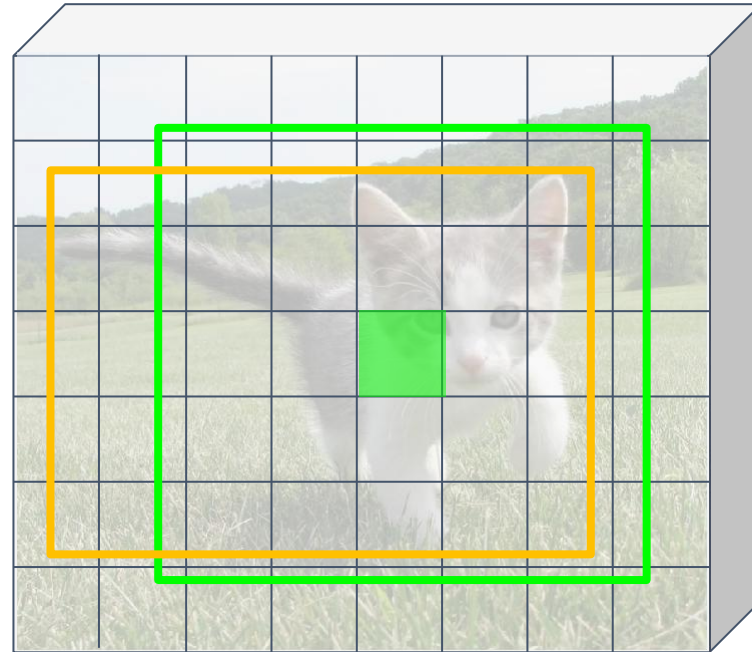Run backbone CNN to get features aligned to input image

Imagine an anchor box of fixed size at each point in the feature map



CNN

Input Image
(e.g. 3 x 640 x 480)

Image features
(e.g. 512 x 20 x 15)

Conv

Anchor is an object?
1 x 20 x 15

Box transforms
4 x 20 x 15

For positive boxes, also predict a box transform to regress from **anchor box** to **object box**

# Region Proposal Network (RPN)

Run backbone CNN to get
features aligned to input image

**Problem**: Anchor box may
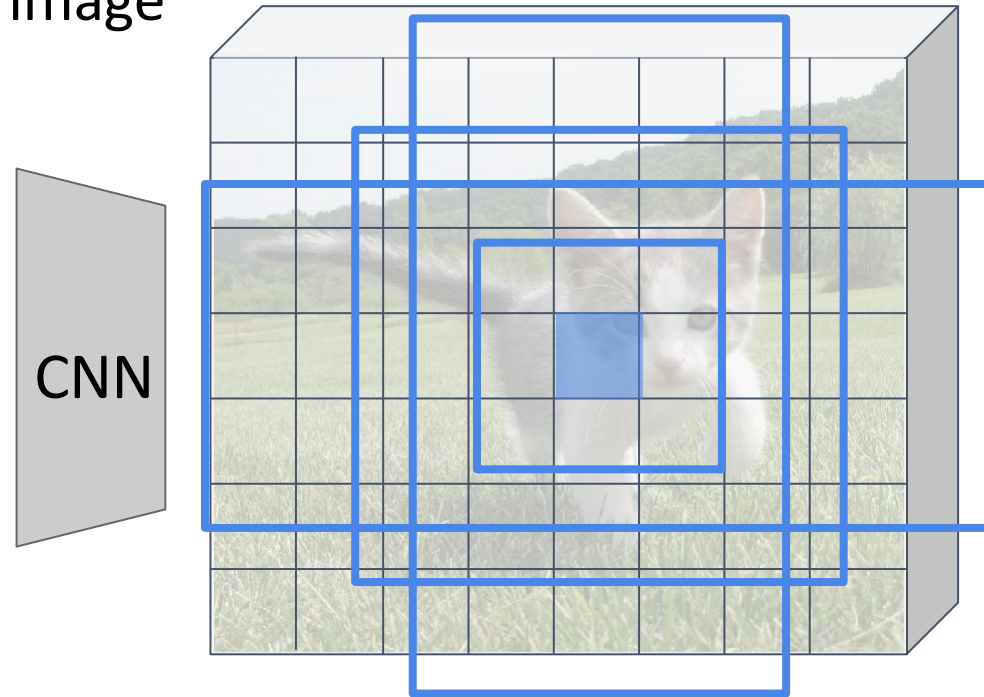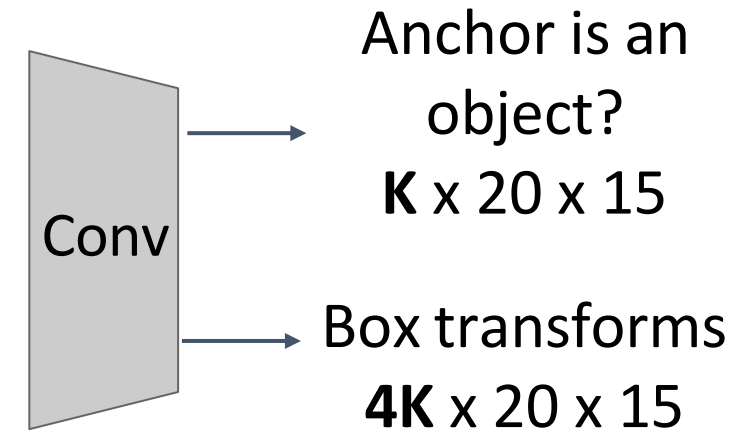have the wrong size / shape
**Solution**: Use **K** different
anchor boxes at each point!



Input Image
(e.g. 3 x 640 x 480)

CNN

Image features
(e.g. 512 x 20 x 15)

Conv

Anchor is an
object?
**K** x 20 x 15

Box transforms
**4K** x 20 x 15

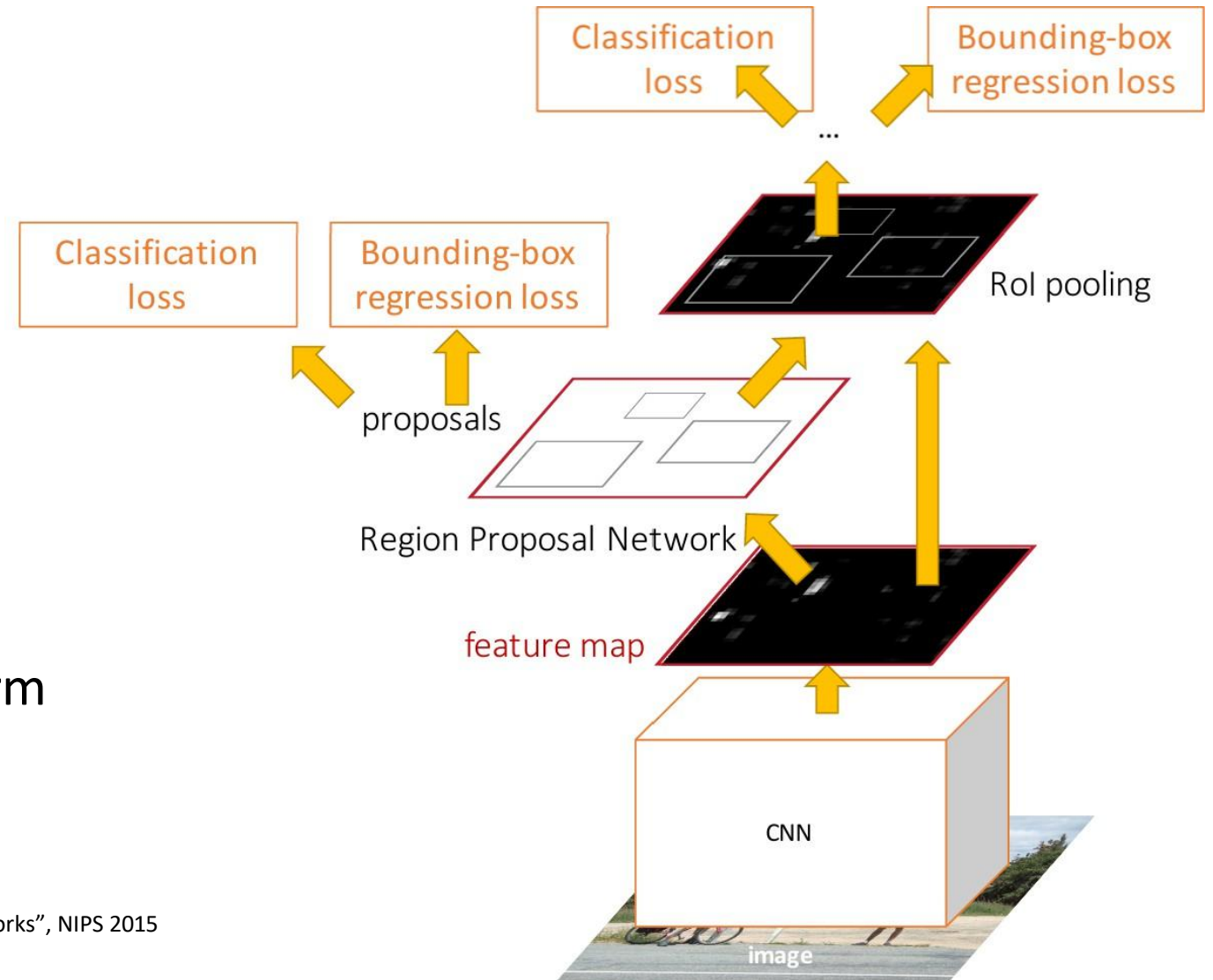At test time: sort all
K*20*15 boxes by their
score, and take the top ~300
as our region proposals
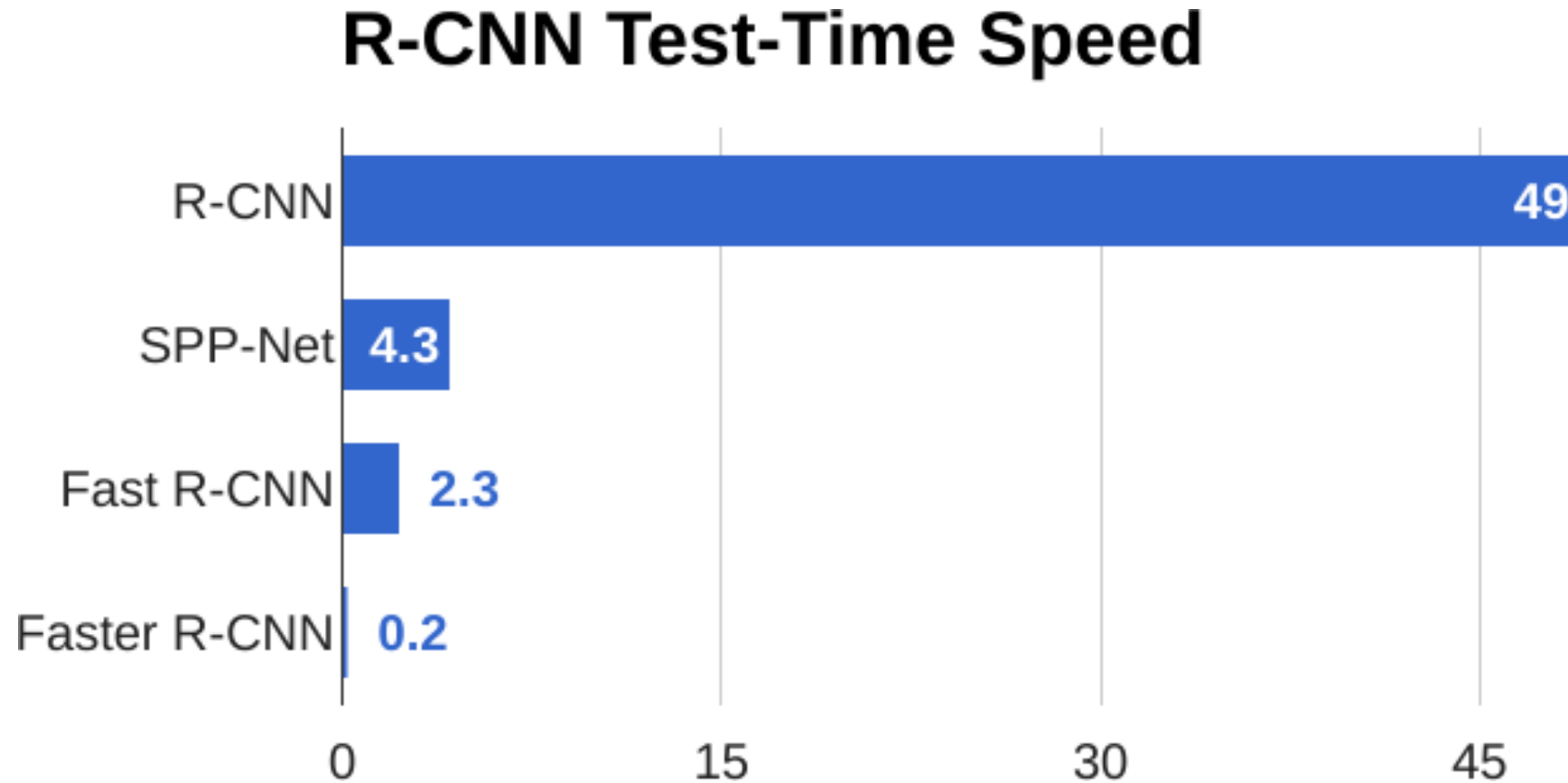
# Faster R-CNN: Learnable Region Proposals

Jointly train with 4 losses:

1. **RPN classification**: anchor box is object / not an object
2. **RPN regression**: predict transform from anchor box to proposal box
3. **Object classification**: classify proposals as background / object class
4. **Object regression**: predict transform from proposal box to object box



Ren et al, "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks", NIPS 2015
Figure copyright 2015, Ross Girshick; reproduced with permission

# Faster R-CNN: Learnable Region Proposals



**R-CNN Test-Time Speed**

| | |
|---|---|
| R-CNN | 49 |
| SPP-Net | 4.3 |
| Fast R-CNN | 2.3 |
| Faster R-CNN | 0.2 |

0    15    30    45

# Faster R-CNN: Learnable Region Proposals

Faster R-CNN is a
**Two-stage object detector**

First stage: Run once per image
- Backbone network
- Region proposal network

Second stage: Run once per region
- Crop features: RoI pool / align
- Predict object class
- Prediction bbox offset

# Faster R-CNN: Learnable Region Proposals

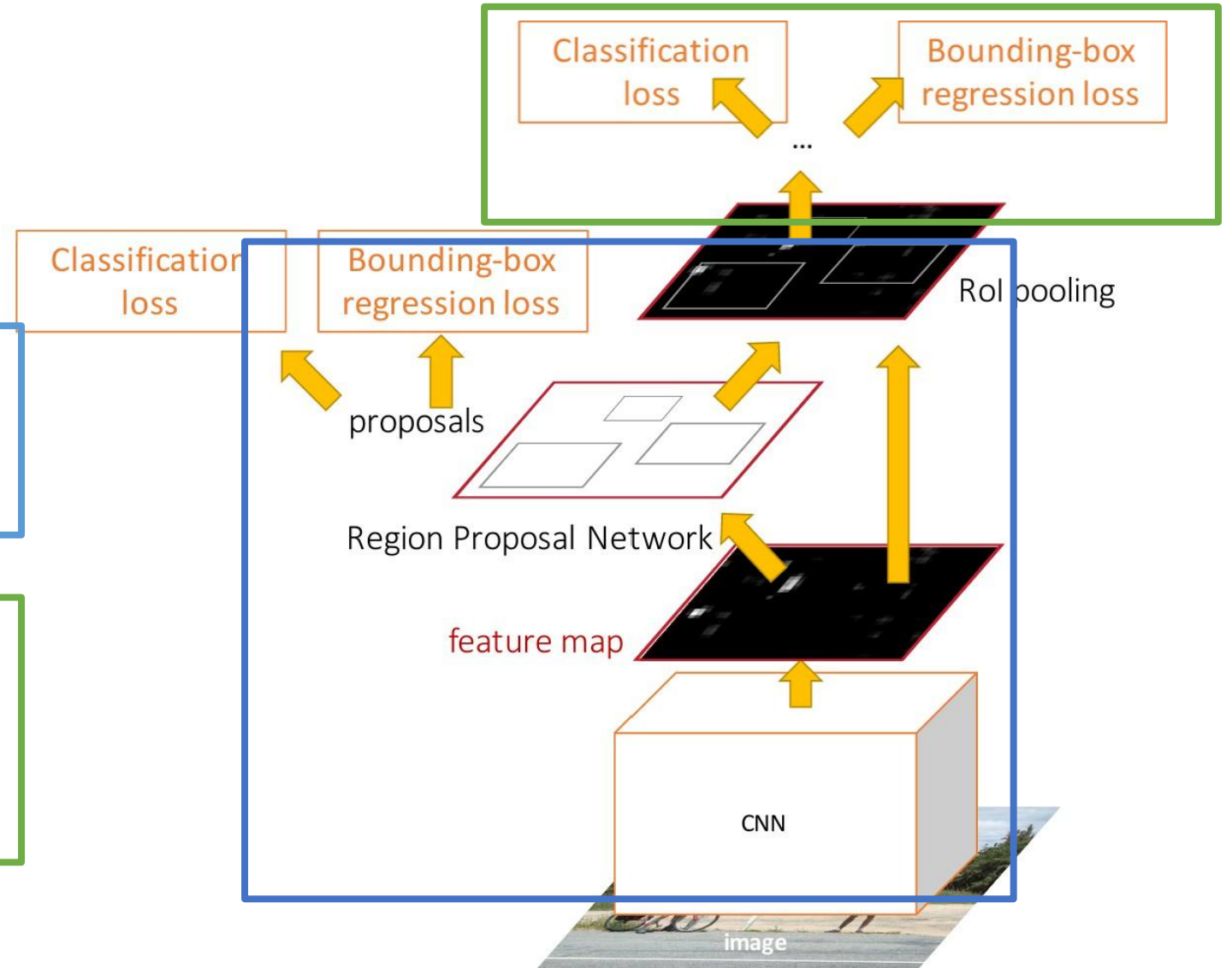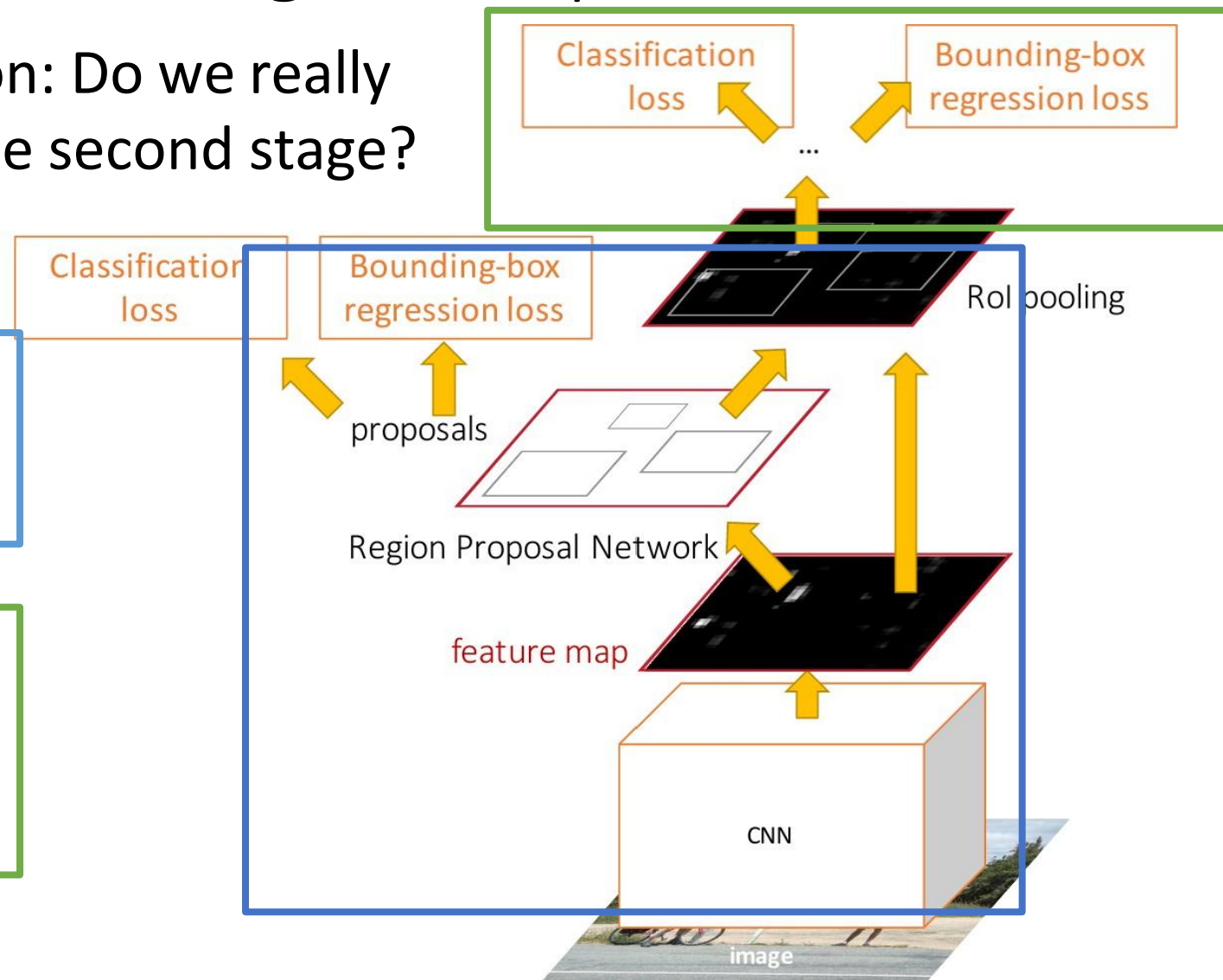Question: Do we really need the second stage?

Faster R-CNN is a
**Two-stage object detector**

First stage: Run once per image
- Backbone network
- Region proposal network

Second stage: Run once per region
- Crop features: RoI pool / align
- Predict object class
- Prediction bbox offset



Classification loss

Bounding-box regression loss

...

RoI pooling

Classification loss

Bounding-box regression loss

proposals

Region Proposal Network

feature map

CNN

image

# Single-Stage Object Detection

Run backbone CNN to get features aligned to input image
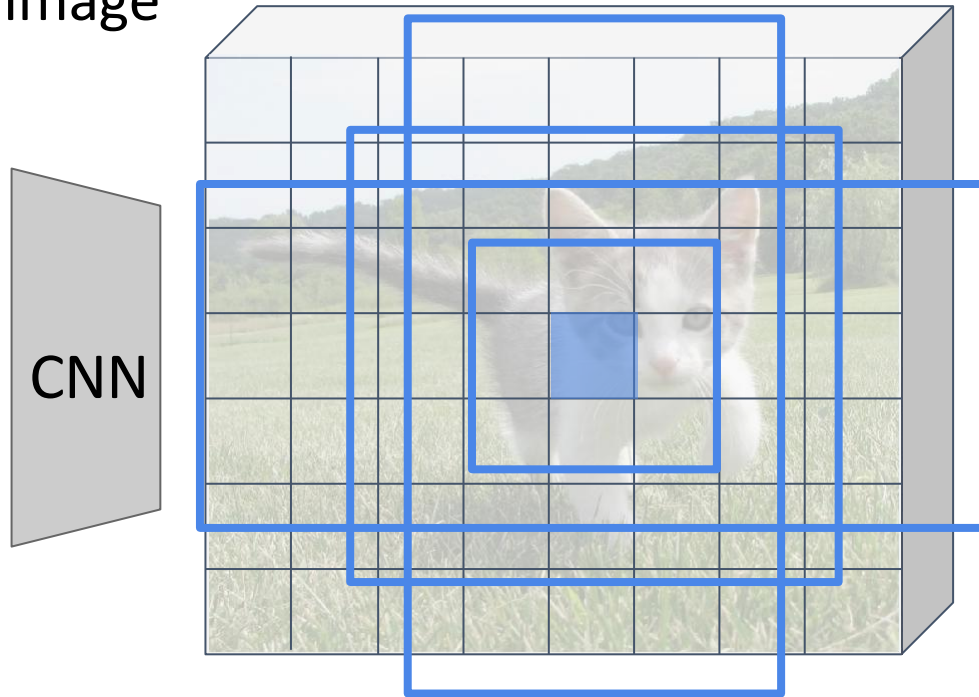


Input Image
(e.g. 3 x 640 x 480)

CNN

Image features
(e.g. 512 x 20 x 15)

Conv

**RPN**: Classify each anchor as object / not object

**Single-Stage Detector**: Classify each object as one of C categories (or background)

Anchor category
→ (C+1) x K x 20 x 15

Box transforms
4K x 20 x 15

Remember: K anchors at each position in image feature map

# Single-Stage Object Detection

Run backbone CNN to get features aligned to input image
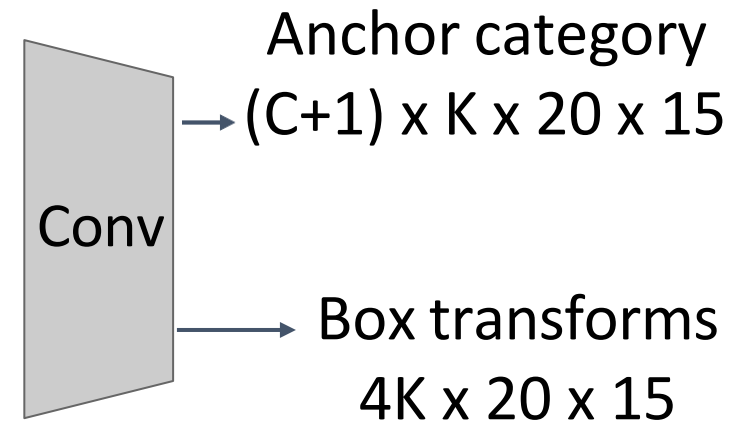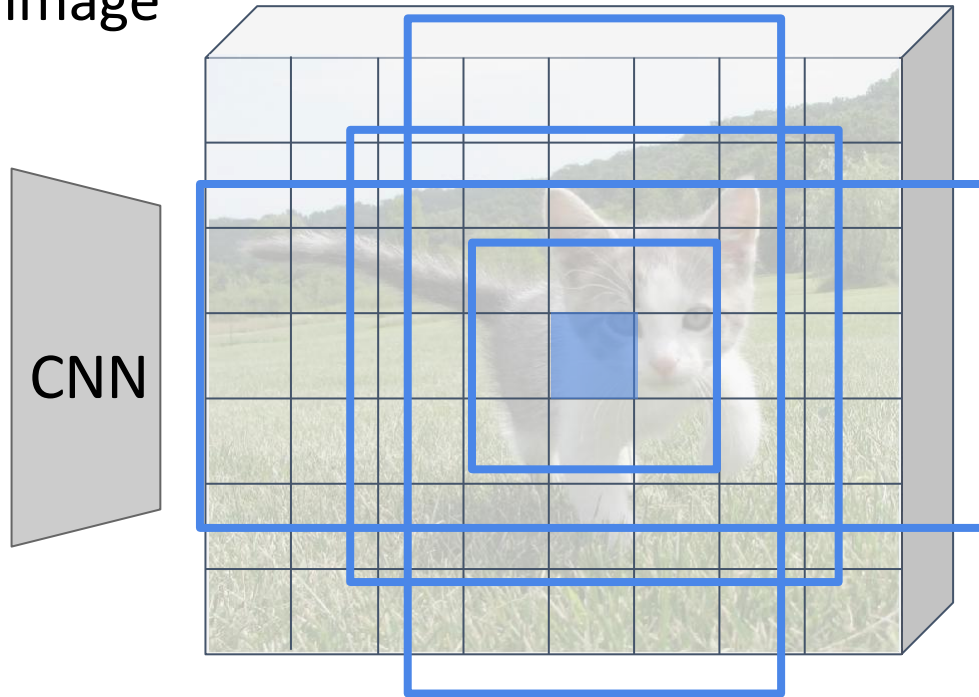


Input Image
(e.g. 3 x 640 x 480)

CNN

Image features
(e.g. 512 x 20 x 15)

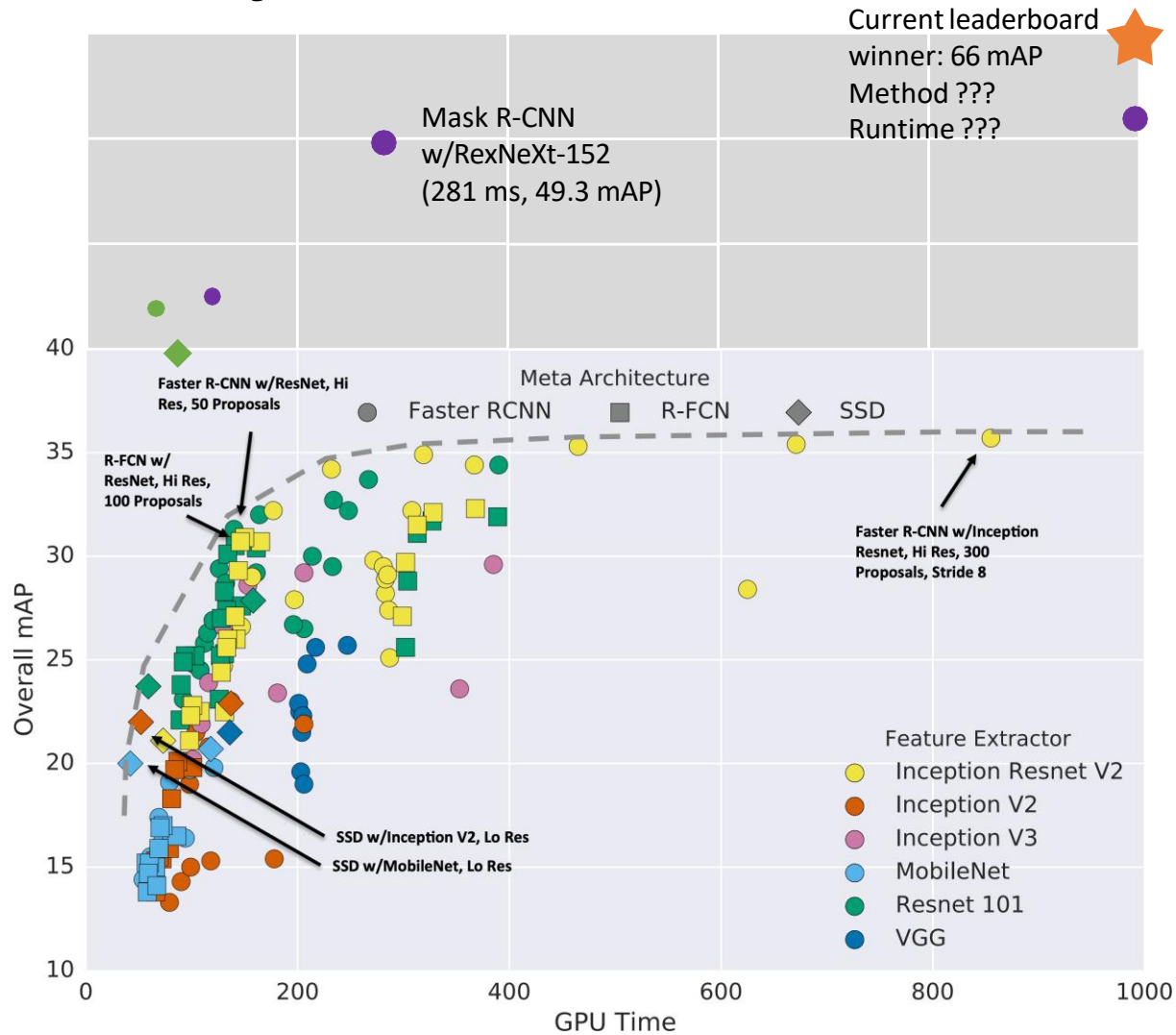**RPN**: Classify each anchor as object / not object
**Single-Stage Detector**: Classify each object as one of C categories (or background)

Conv

Anchor category
→ (C+1) x K x 20 x 15

Box transforms
**C** x 4K x 20 x 15

Sometimes use **category-specific regression**: Predict different box transforms for each category

Redmon et al, "You Only Look Once: Unified, Real-Time Object Detection", CVPR 2016
Liu et al, "SSD: Single-Shot MultiBox Detector", ECCV 2016
Lin et al, "Focal Loss for Dense Object Detection", ICCV 2017

# Object Detection: Lots of variables!



Current leaderboard winner: 66 mAP
Method ???
Runtime ???

Mask R-CNN w/RexNeXt-152 (281 ms, 49.3 mAP)

Faster R-CNN w/ResNet, Hi Res, 50 Proposals

R-FCN w/ ResNet, Hi Res, 100 Proposals

Meta Architecture
Faster RCNN   R-FCN   SSD

Faster R-CNN w/Inception Resnet, Hi Res, 300 Proposals, Stride 8

Feature Extractor
- Inception Resnet V2
- Inception V2
- Inception V3
- MobileNet
- Resnet 101
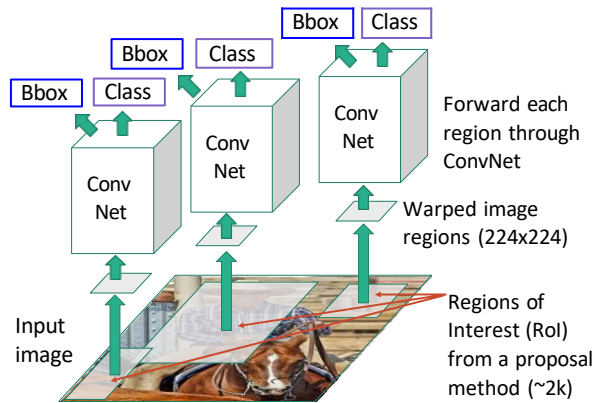- VGG

SSD w/Inception V2, Lo Res
SSD w/MobileNet, Lo Res

These results are a few years old … since then GPUs have gotten faster, and we've improved performance with many tricks:
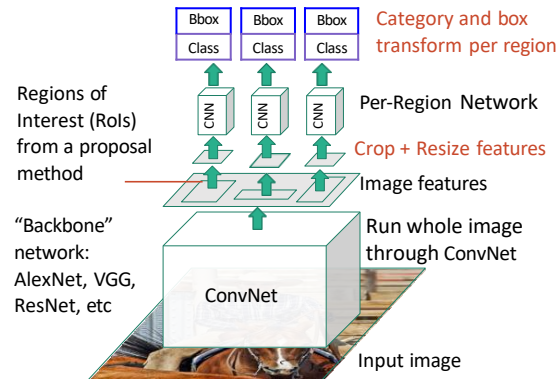- Train longer!
- Multiscale backbone: Feature Pyramid Networks
- Better backbone: ResNeXt
- Single-Stage methods have improved
- Very big models work better
- Test-time augmentation pushes numbers up
- Big ensembles, more data, etc

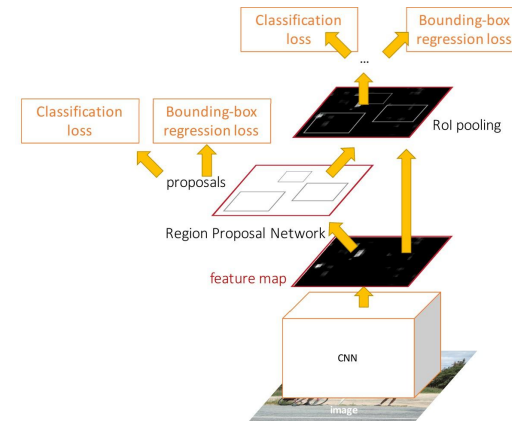Huang et al, "Speed/accuracy trade-offs for modern convolutional object detectors", CVPR 2017

https://codalab.lisn.upsaclay.fr/competitions/7384#results

# Summary

**"Slow" R-CNN**: Run CNN independently for each region

**Fast R-CNN**: Apply differentiable cropping to shared image features

**Faster R-CNN**: Compute proposals with CNN

**Single-Stage**: Fully convolutional detector