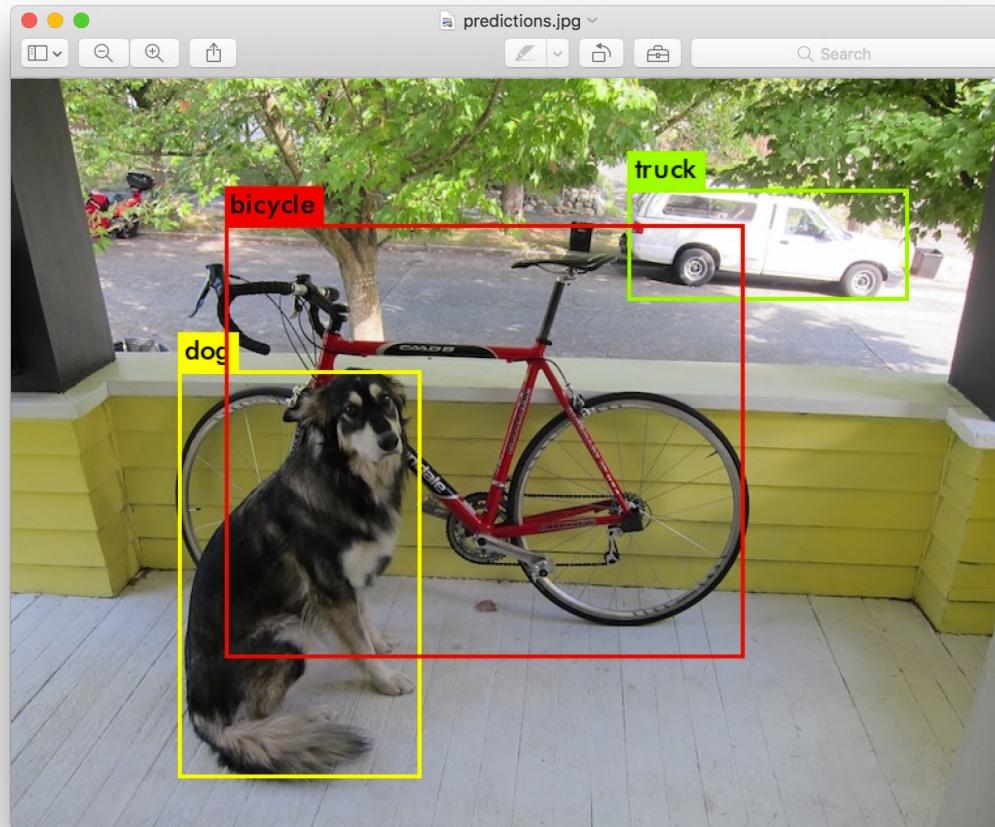


Detection and Localization



Many slides from: [Stanford CS231](#) , [Illinois cs498](#) and [Justin Johnson](#)

Computer Vision Tasks

2

Classification



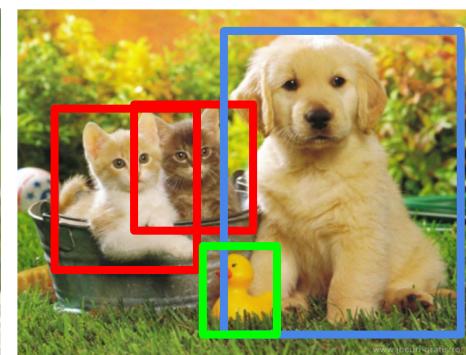
CAT

Classification + Localization



CAT

Object Detection



CAT, DOG, DUCK

Instance Segmentation



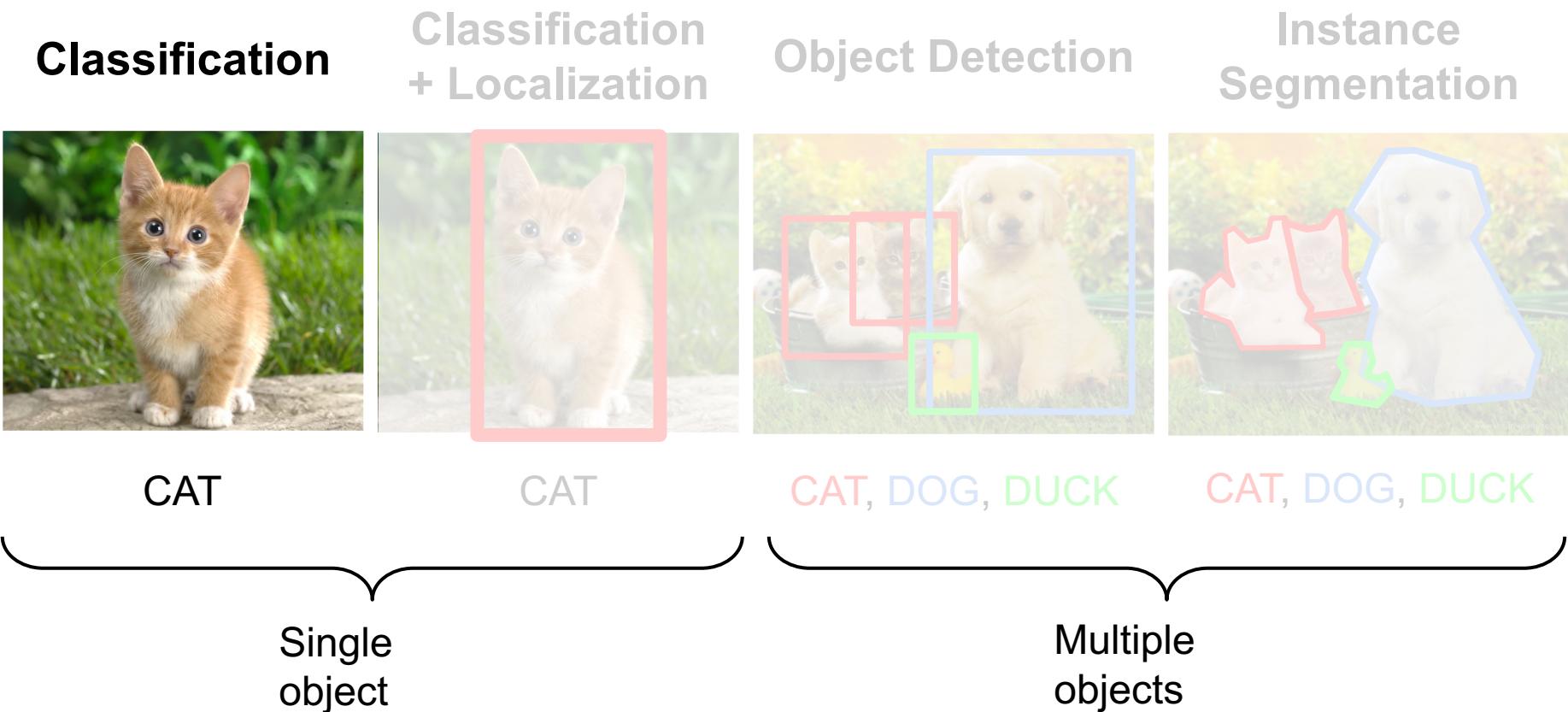
CAT, DOG, DUCK

Single
object

Multiple
objects

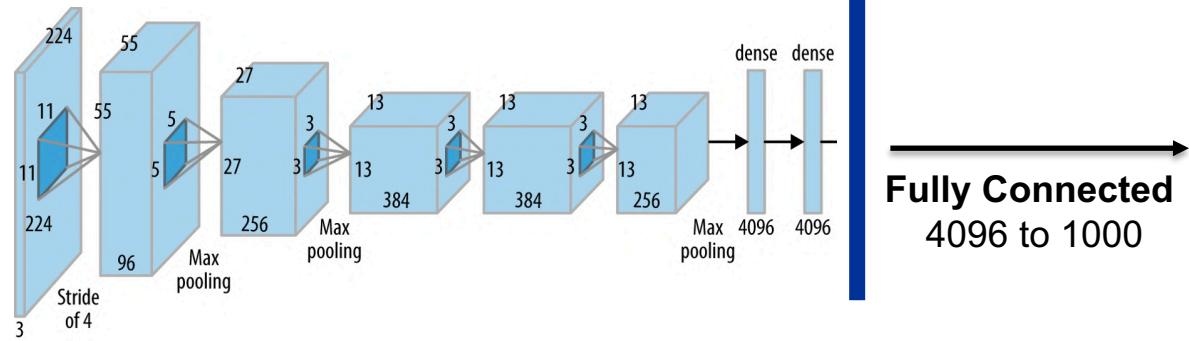
So Far: Image Classification

3



So Far: Image Classification

4



Class Scores:

Cat: 0.9
Dog: 0.05
Car: 0.01
...

Vector
4096

Computer Vision Tasks

5

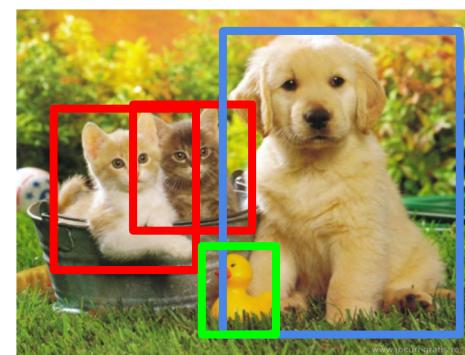
Classification



Classification + Localization



Object Detection



Instance Segmentation



CAT

CAT

CAT, DOG, DUCK

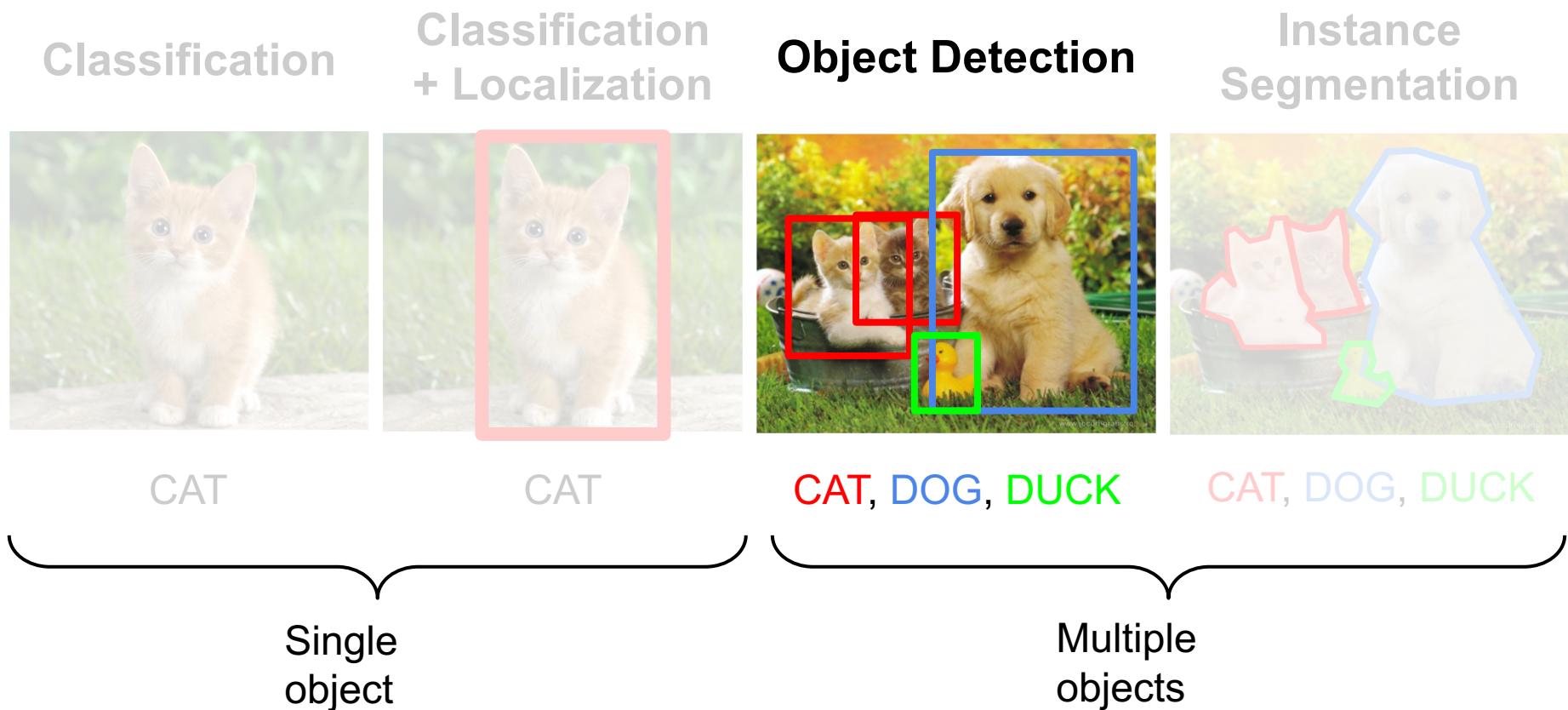
CAT, DOG, DUCK

Single
object

Multiple
objects

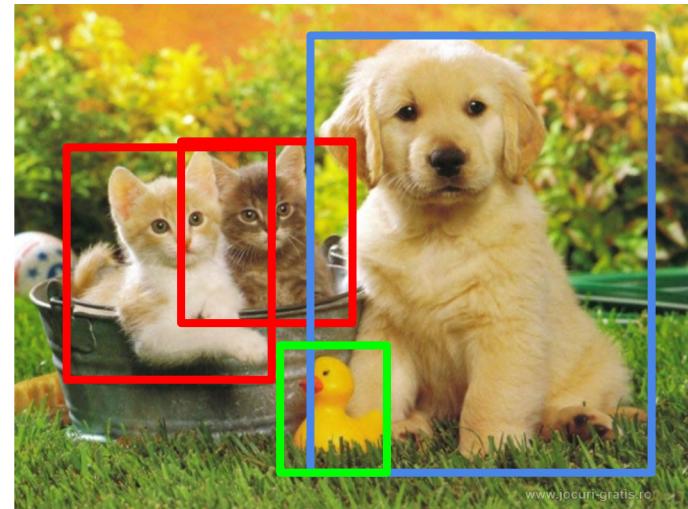
Today's Lesson: Object Detection

6



Object Detection: Task Definition

- **Input:** A single RGB image
- **Output:** A set of detected objects:
For each detected object:
 - Category label (from a fixed set)
 - Bounding box = $(x,y,\text{width},\text{height})$



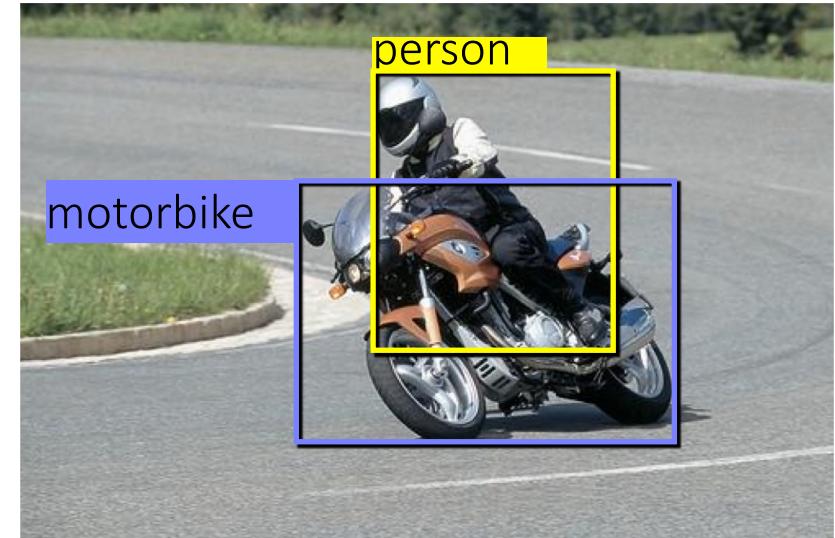
www.jocuri-gratis.ro

Object Detection: Task Definition

classes={airplane, bird, motorbike, person, sofa}



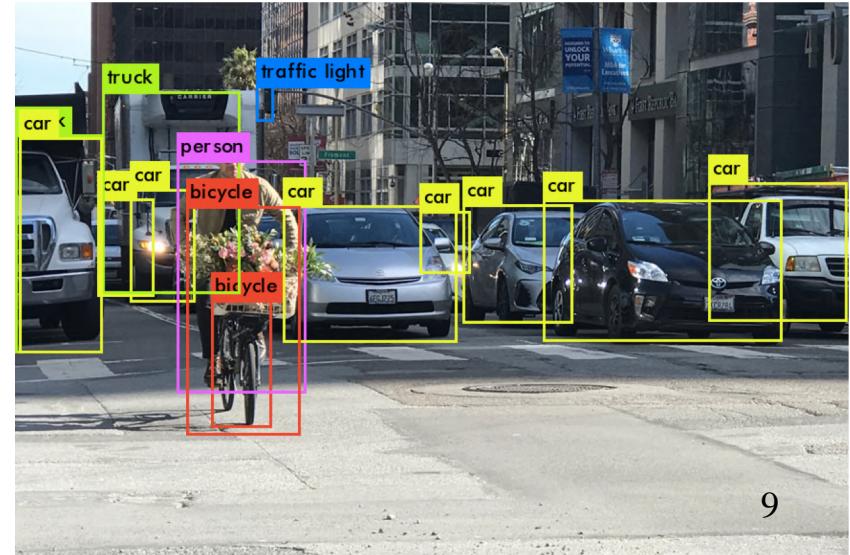
Input



Desired output

What are the challenges of an object detector?

- **Multiple Outputs:** Variable numbers of objects per image
- **Multiple types of Outputs:** Need to predict “what” (class) as well as “where” (bounding box)
- **Multiple Sizes:** multiple objects with variable sizes
- Multiple detections of the same object
- Multiple objects at the same location

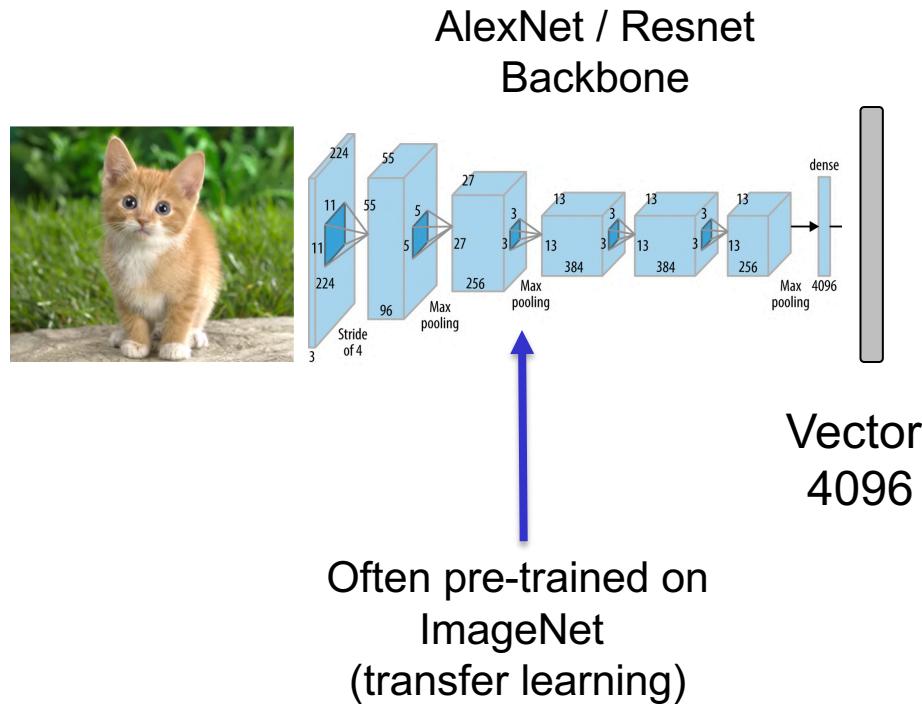


Outline

- Task definition and evaluation
- Conceptual approaches for detection
- Region-Proposal based Convolutional Networks
 - R-CNN
 - Fast R-CNN
 - Faster R-CNN
- Single Stage Object Detectors

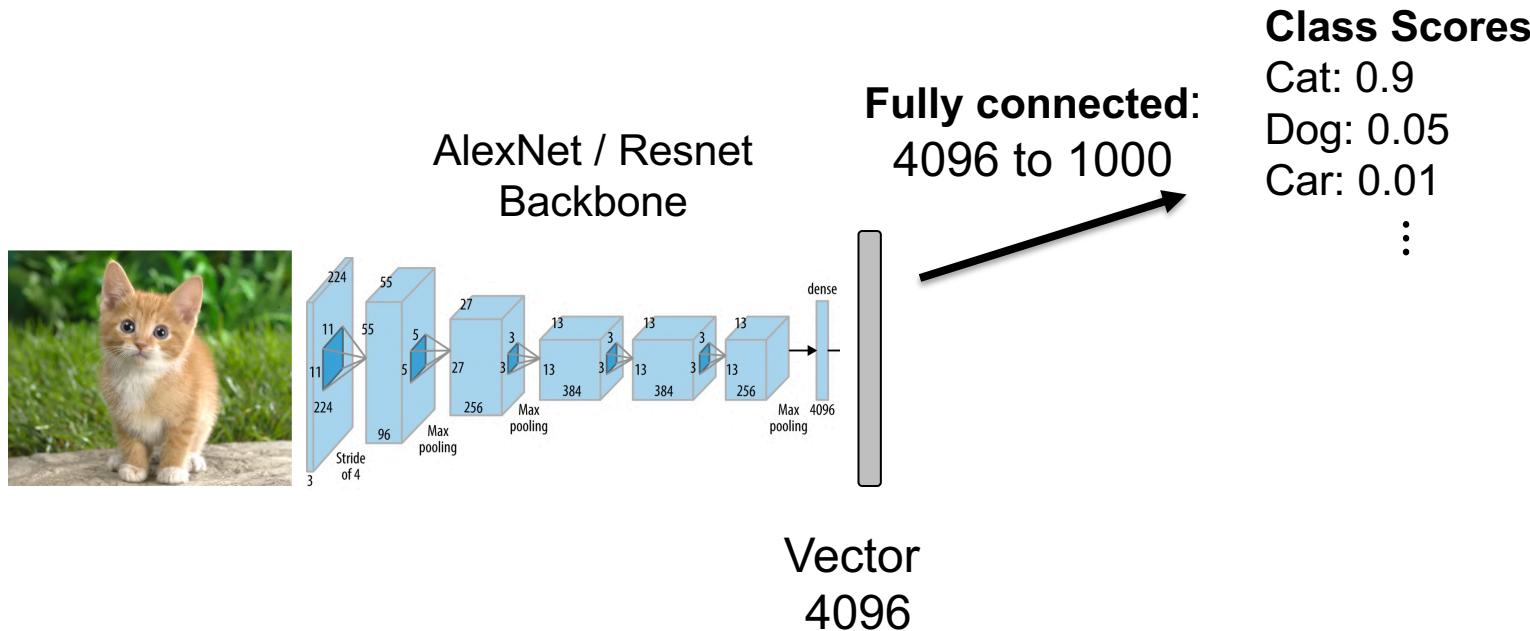
Object Detection: Single Object

- Can be treated as classification + localization



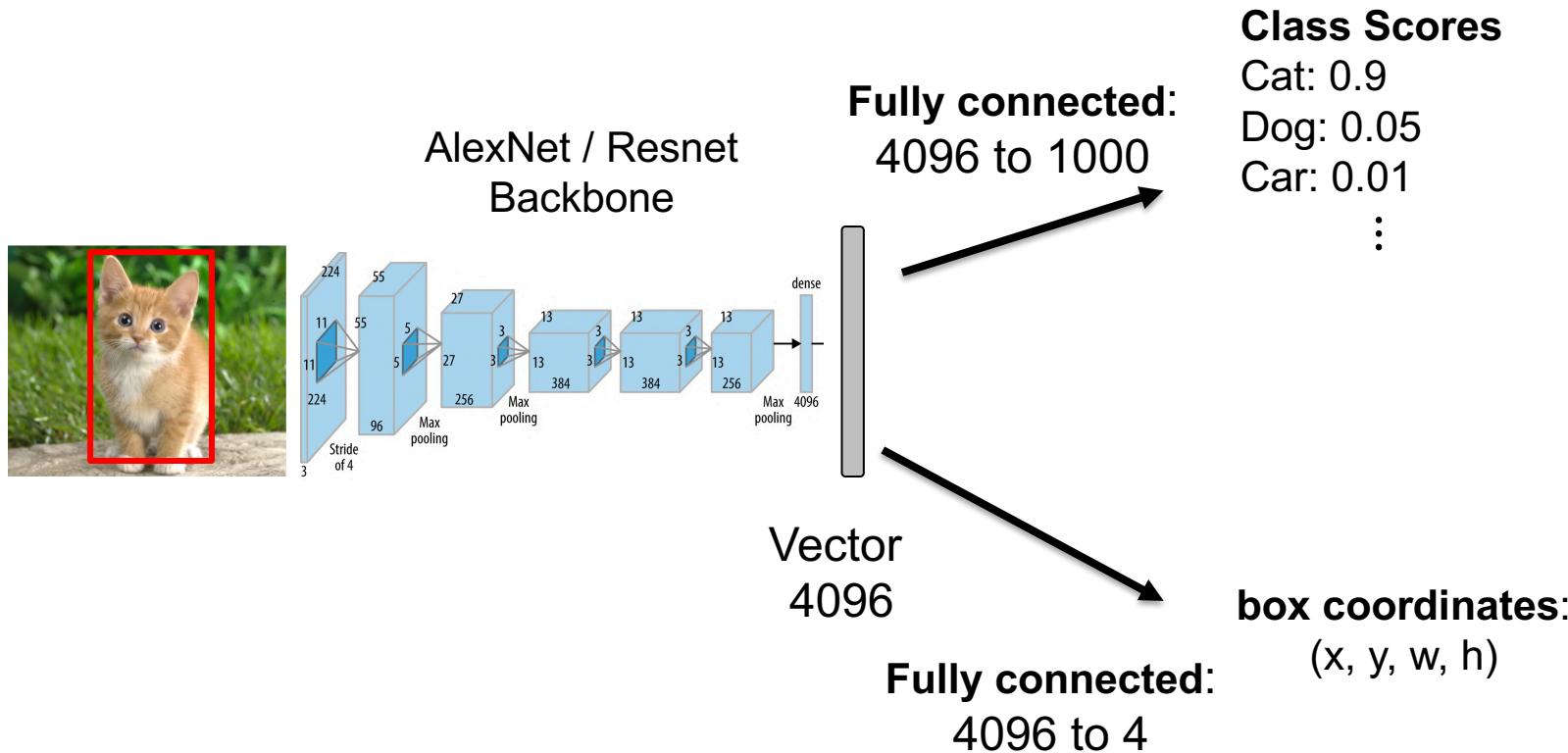
Object Detection: Single Object

- Can be treated as classification + localization



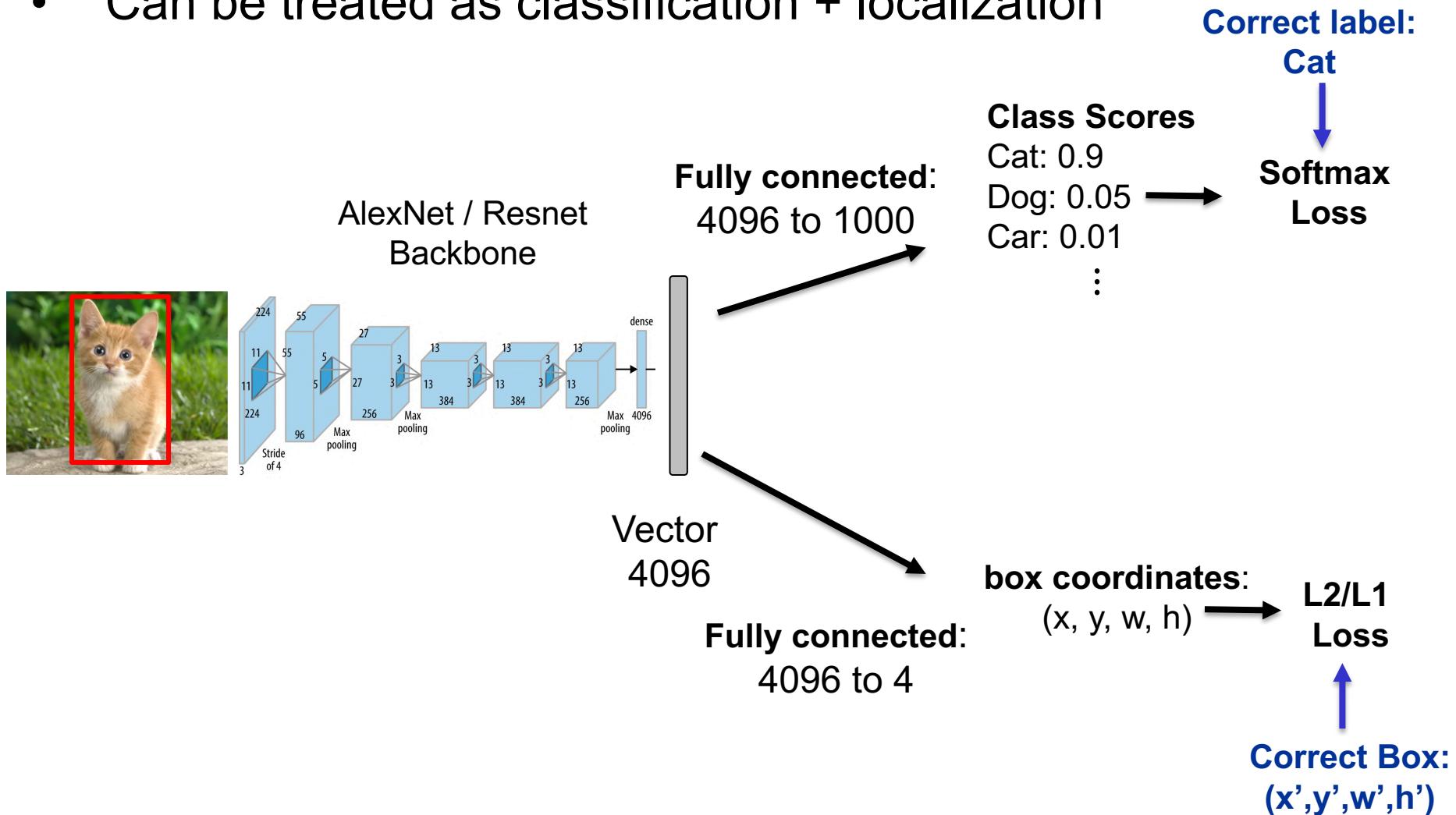
Object Detection: Single Object

- Can be treated as classification + localization



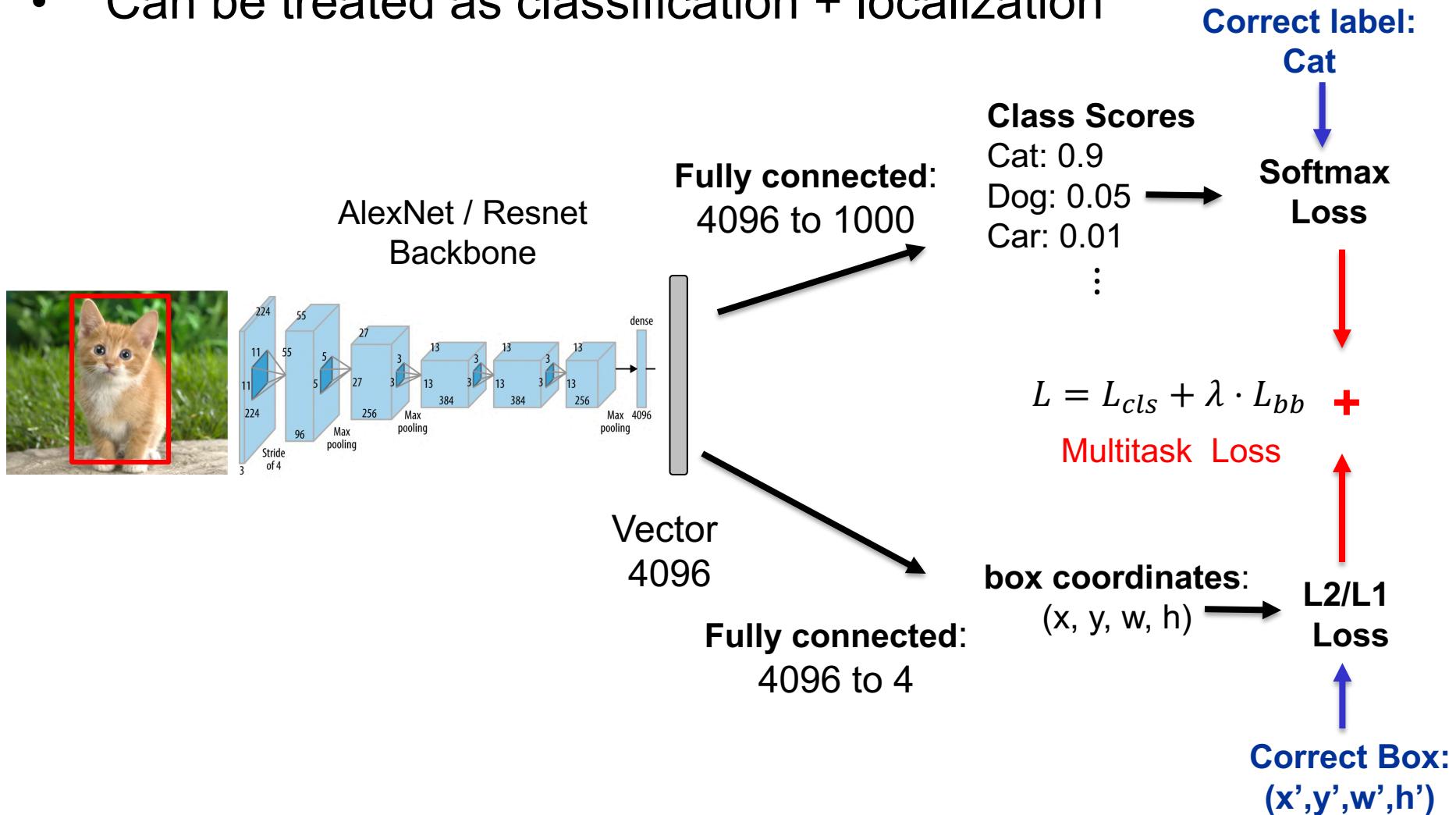
Object Detection: Single Object

- Can be treated as classification + localization



Object Detection: Single Object

- Can be treated as classification + localization



Multi-task loss

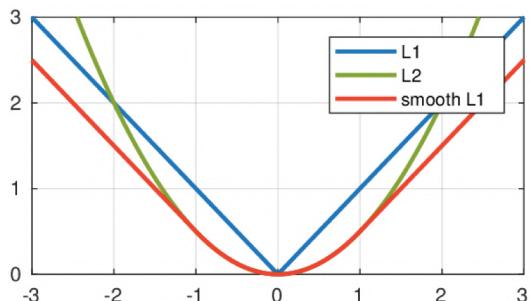
- Loss for ground truth class y , predicted class probabilities $P(y)$, ground truth box b , and predicted box b' :

$$L(y', y, P, b, b') = -\log P(y') + \lambda L_{\text{reg}}(b, b')$$

CE loss regression loss

- Regression loss:** *smooth L1* loss offsets relative to GT

$$L_{\text{reg}}(b, b') = \sum_{i=\{x,y,w,h\}} \text{smooth}_{L_1}(b_i - b'_i)$$



$$\text{smooth}_{L_1}(x) = \begin{cases} 0.5x^2 & \text{if } |x| < 1 \\ |x| - 0.5 & \text{otherwise} \end{cases}$$

Single object Evaluation

classification + localization: Regress both

Classification: C classes

Input: Image

Output: Class label

Evaluation metric: Accuracy



CAT

Localization:

Input: Image

Output: Box in the image (x, y, w, h)

Evaluation metric: Intersection over Union

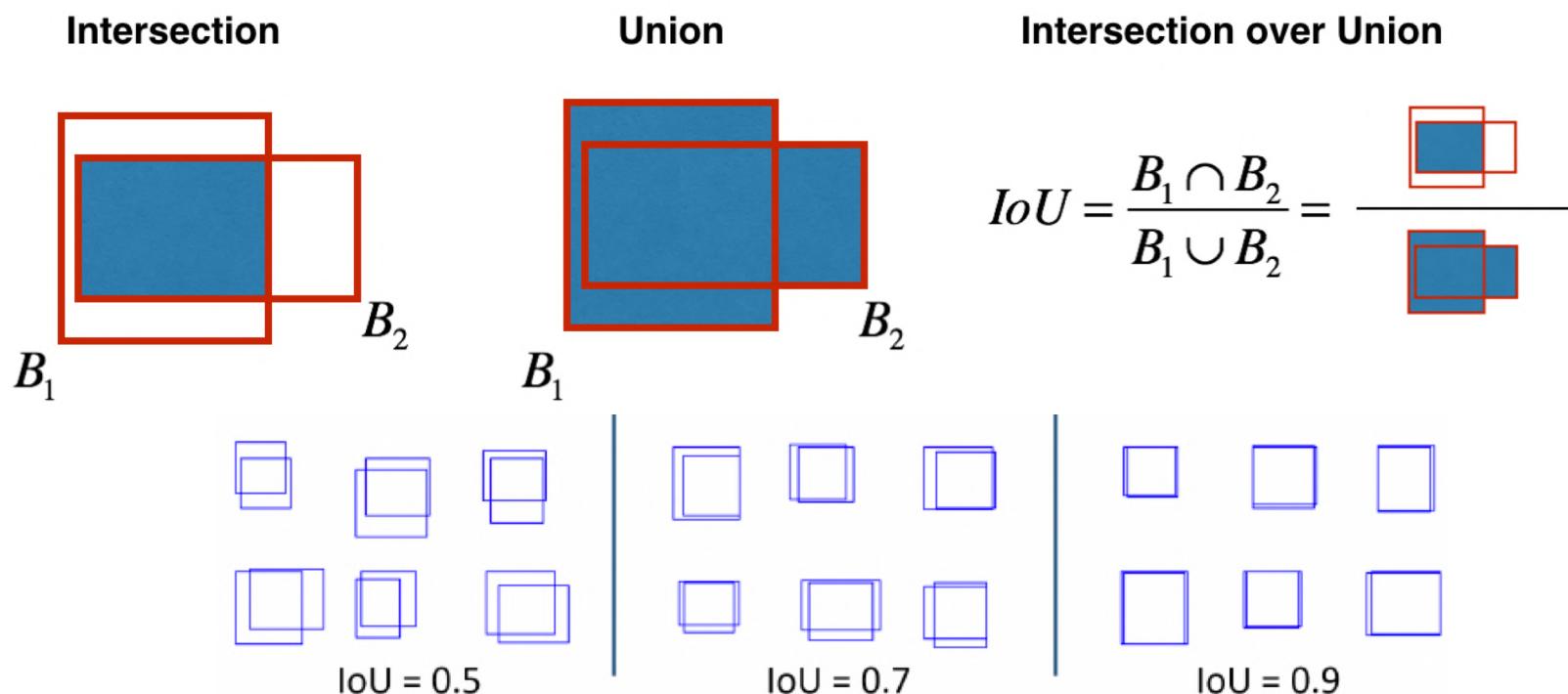


(x, y, w, h)

Comparing Boxes

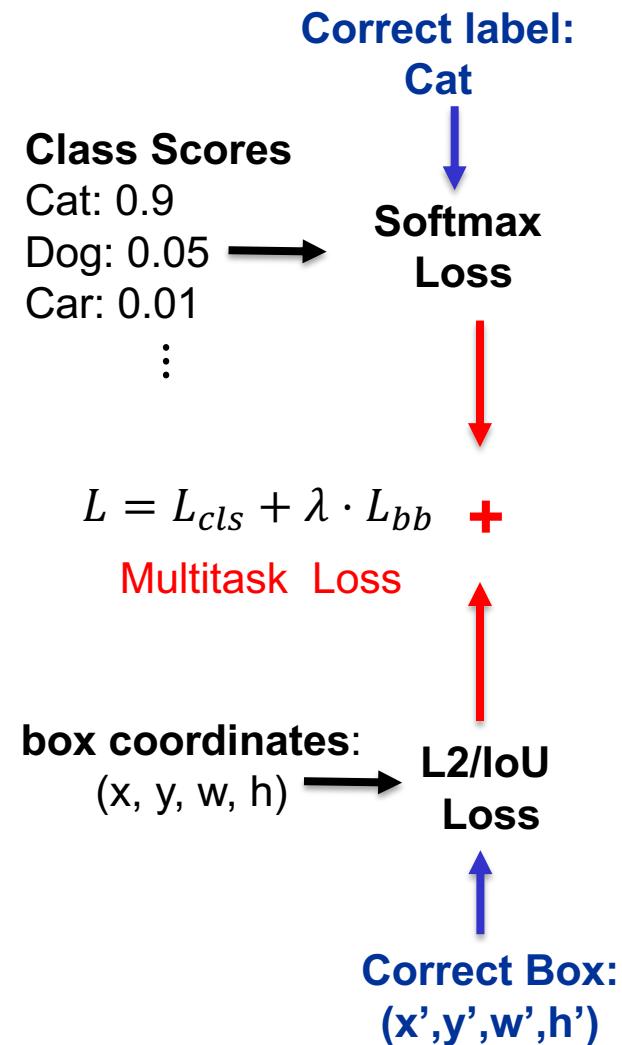
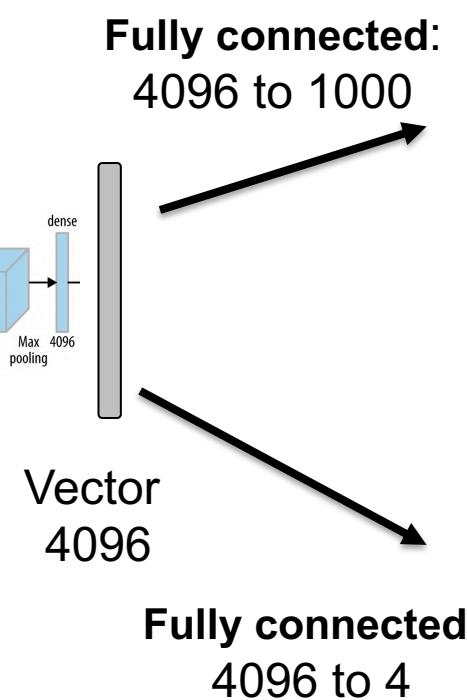
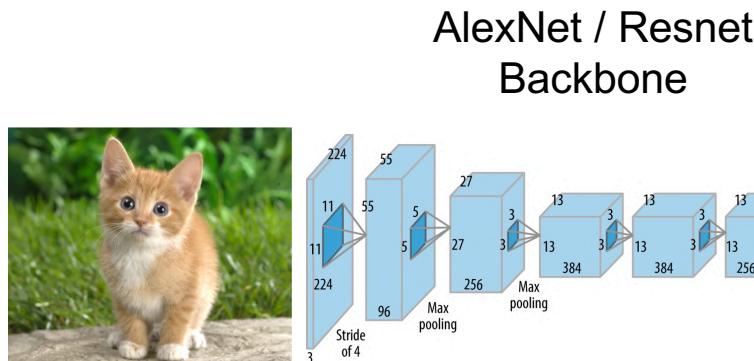
Intersection over Union (IoU).

- The IoU score is between [0..1] and is scale invariant.
- IoU is calculated between predicted box & GT box
- IoU > 0.5 is “decent”, IoU > 0.7 is “very good”, and IoU > 0.9 is “almost perfect”.



Object Detection: Single object

- A single object detection can be treated as classification + localization



Problem: Images can have more than a single object

Object Detection: Multiple objects



→ CAT, (x, y, w, h) = 4 numbers



→ DOG, (x, y, w, h)
CAT, (x, y, w, h) = 8 numbers

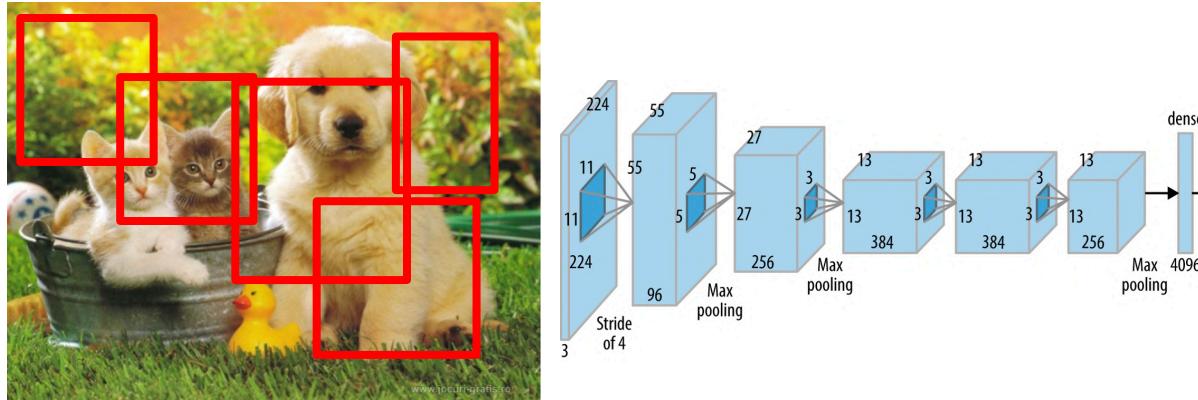


→ DOG, (x, y, w, h)
CAT, (x, y, w, h)
CAT, (x, y, w, h)
DUCK (x, y, w, h) = 16 numbers

Object Detection: Sliding Window

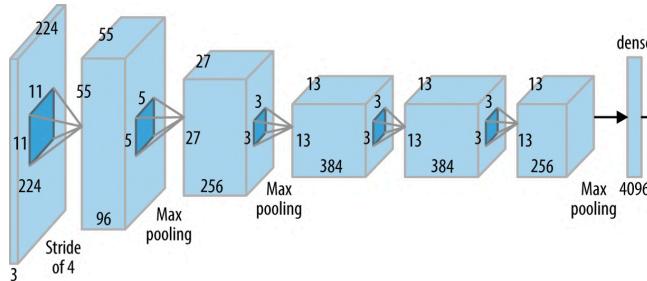
21

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



Object Detection: Sliding Window

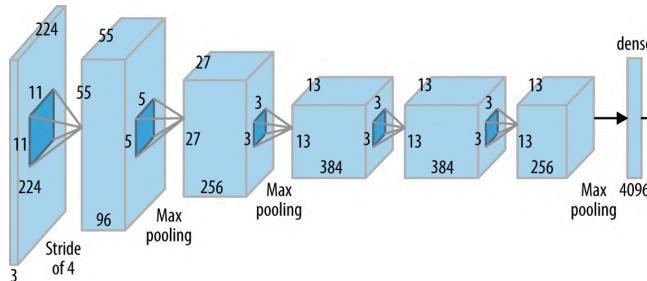
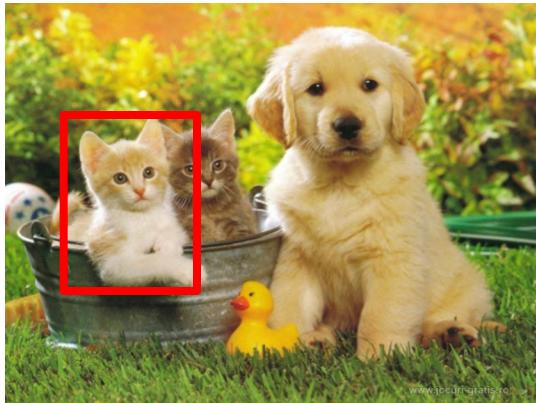
- Apply CNN to many different crops of the image.
- CNN classifies each crop as object or background



Dog: No
Cat: No
Duck: No
Background: Yes

Object Detection: Sliding Window

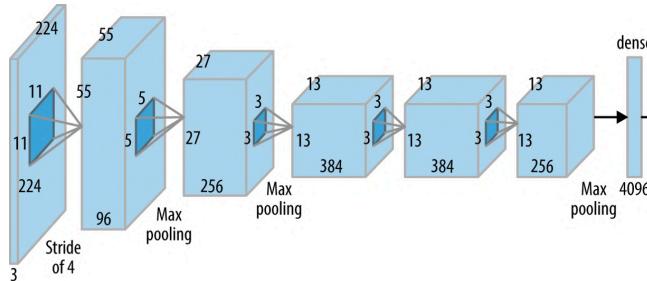
- Apply CNN to many different crops of the image.
- CNN classifies each crop as object or background



Dog: No
Cat: Yes
Duck: No
Background: No

Object Detection: Sliding Window

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



Dog: Yes

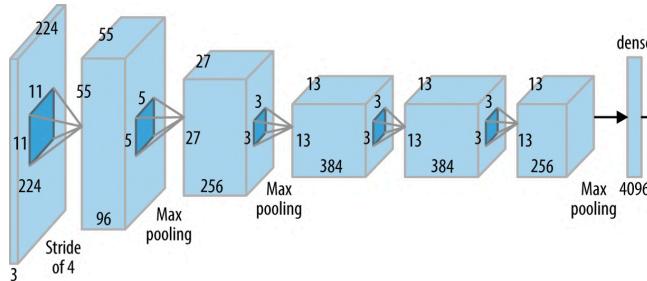
Cat: No

Duck: No

Background: No

Object Detection: Sliding Window

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



Dog: Yes
 Cat: No
 Duck: No
 Background: No

Answer:

Consider an image of size $W \times H$

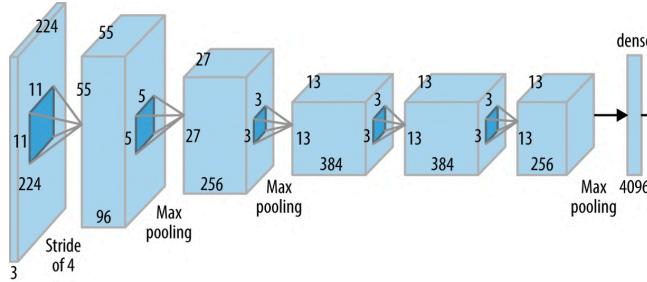
Total Possible boxes: $= O(W^2H^2)$

Question:

How many possible crops are there in an image of size $H \times W$?

Object Detection: Sliding Window

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

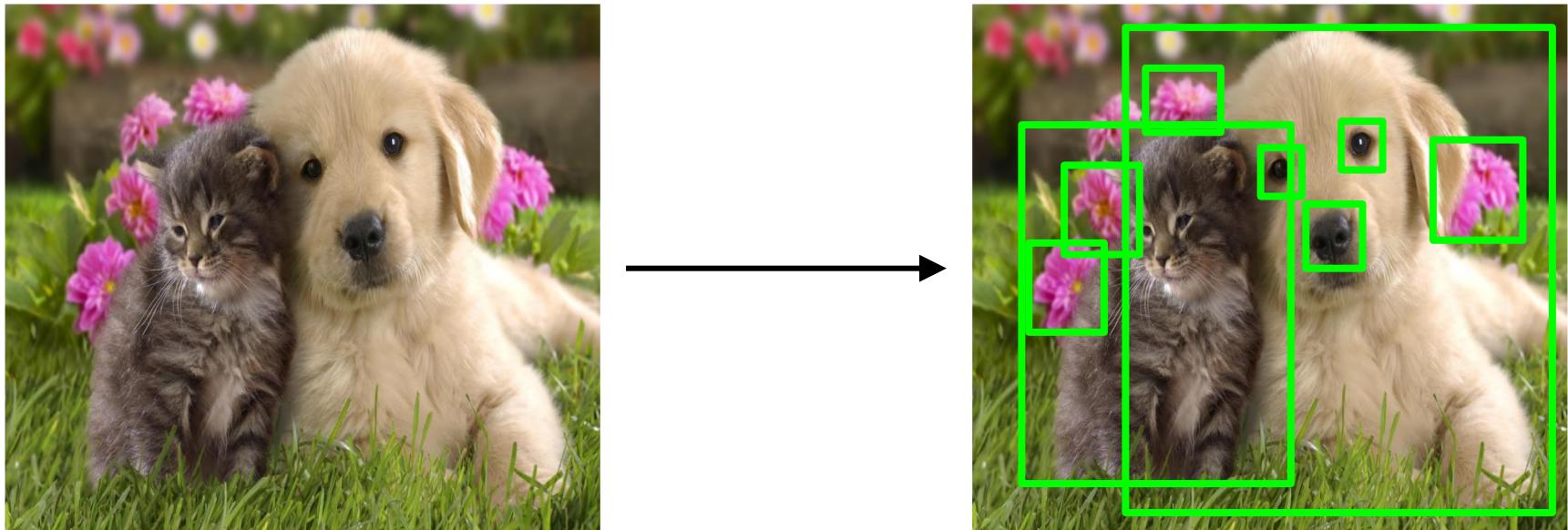


Dog: Yes
Cat: No
Duck: No
Background: No

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

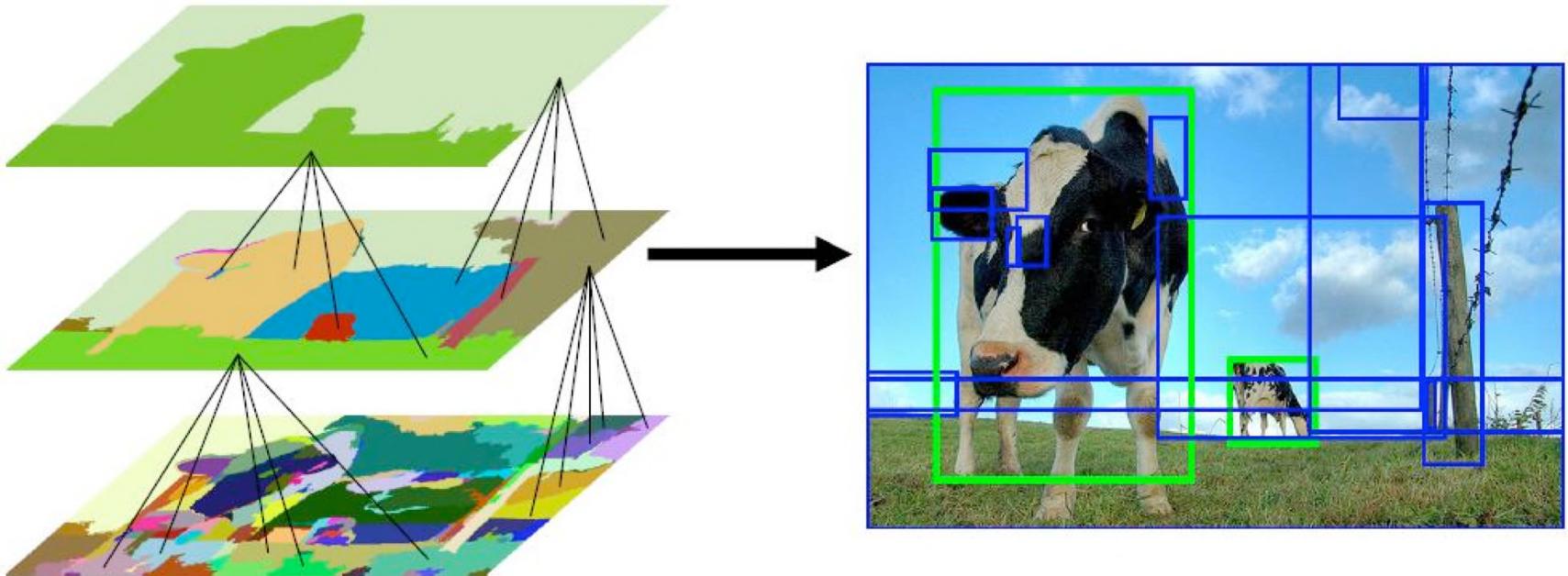
Solution: Region proposal

- Find “blobby” image regions that are likely to contain objects
- Relatively fast to run; e.g. **Selective Search** gives 2000 regions
- Region proposals in a few seconds on CPU



Region proposal: Selective Search

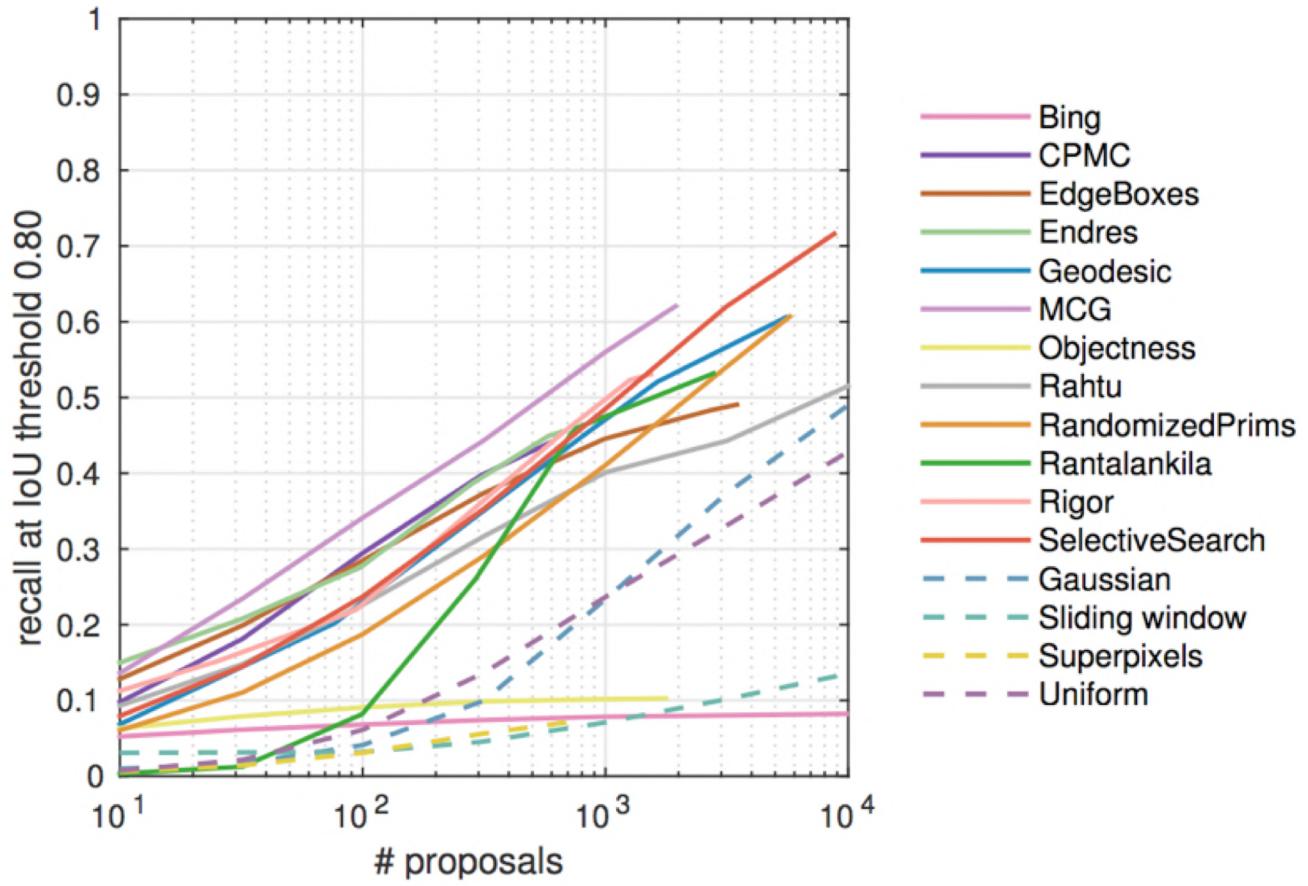
- **Selective Search:** Use hierarchical segmentation; Start with small *superpixels* and group segments based on homogeneity cues



Region proposal: Many other choices

- Tens of ways of generating candidates (“proposals”)

What fraction of GT
objects have
proposals near them?



Object Detection: Region proposal

We will cover three Region Proposal approaches:

1. R-CNN (Region CNN):

- Region proposal using classical image processing approaches
- Independent CNN for each ROI

2. Fast R-CNN

- Region proposal using classical image processing approaches
- Shared CNN and feature maps

3. Faster R-CNN

- Region proposal using CNN

R-CNN: Region-Based CNN

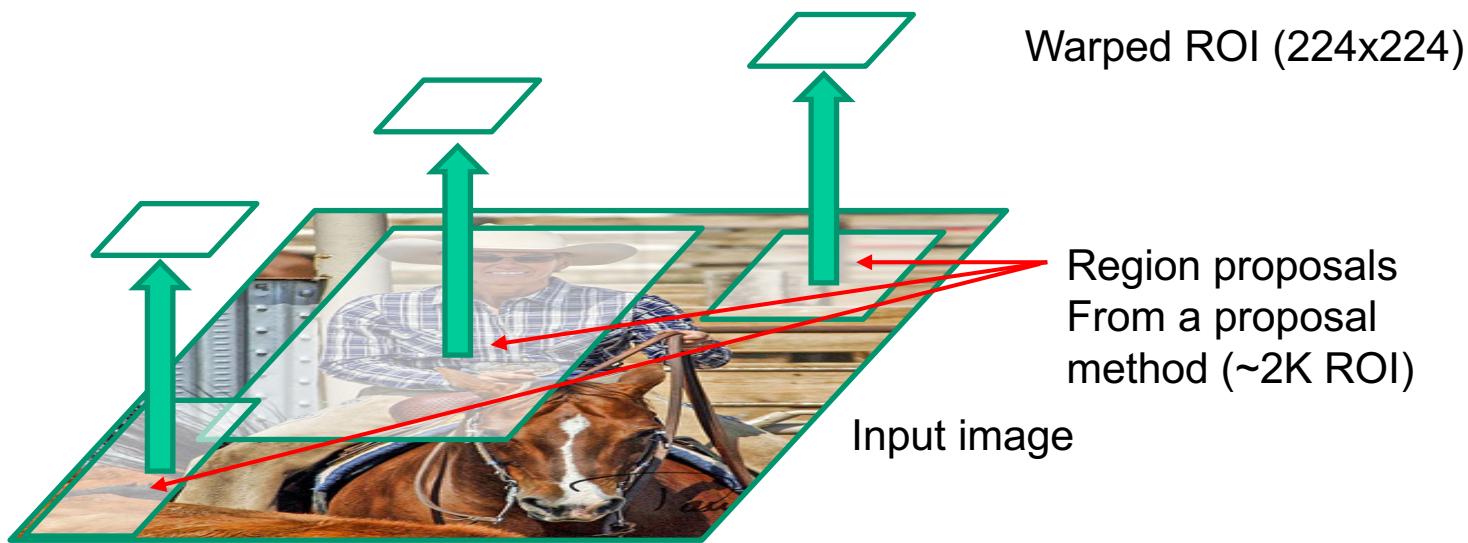


Input image

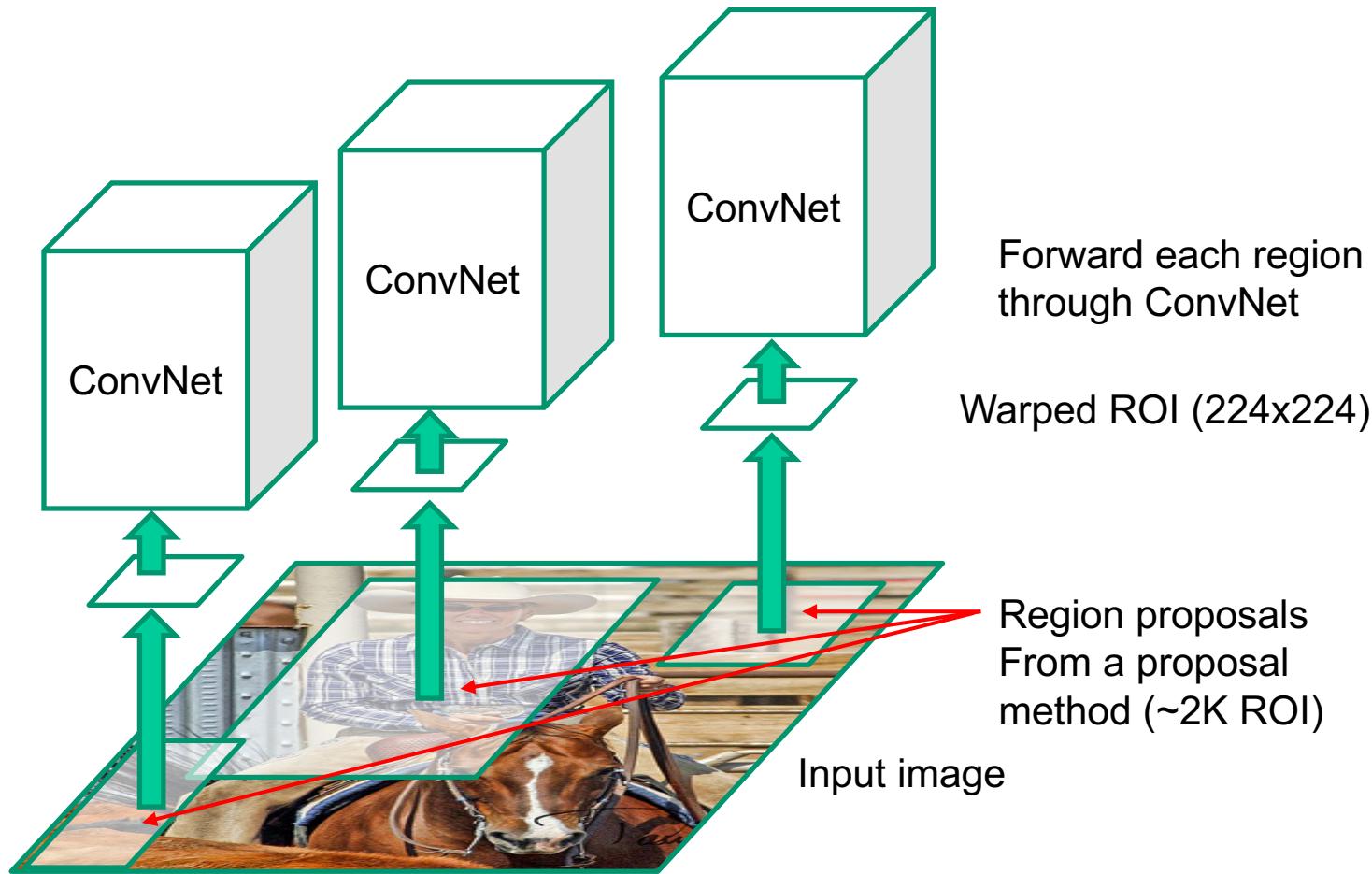
R-CNN: Region-Based CNN



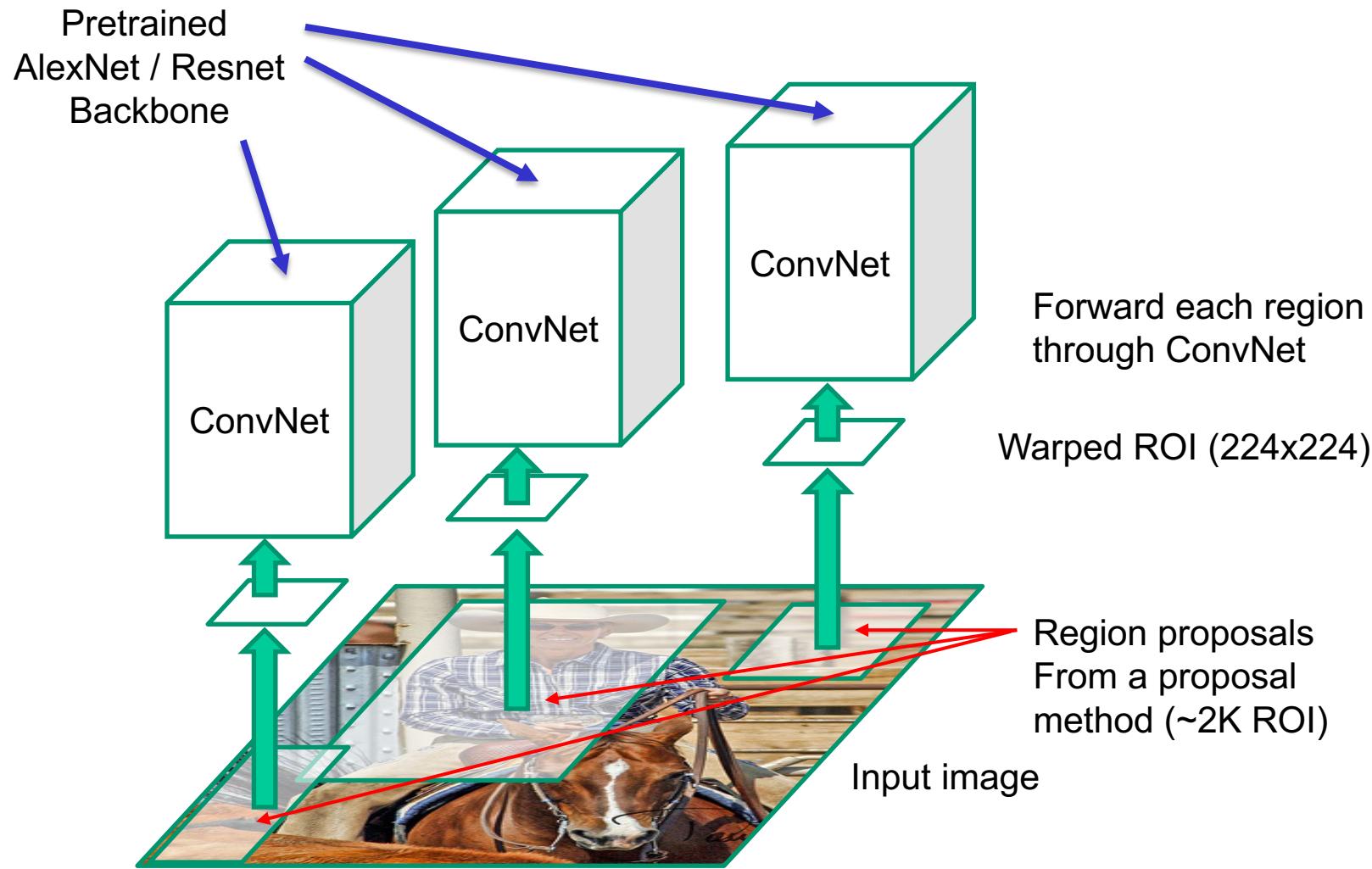
R-CNN: Region-Based CNN



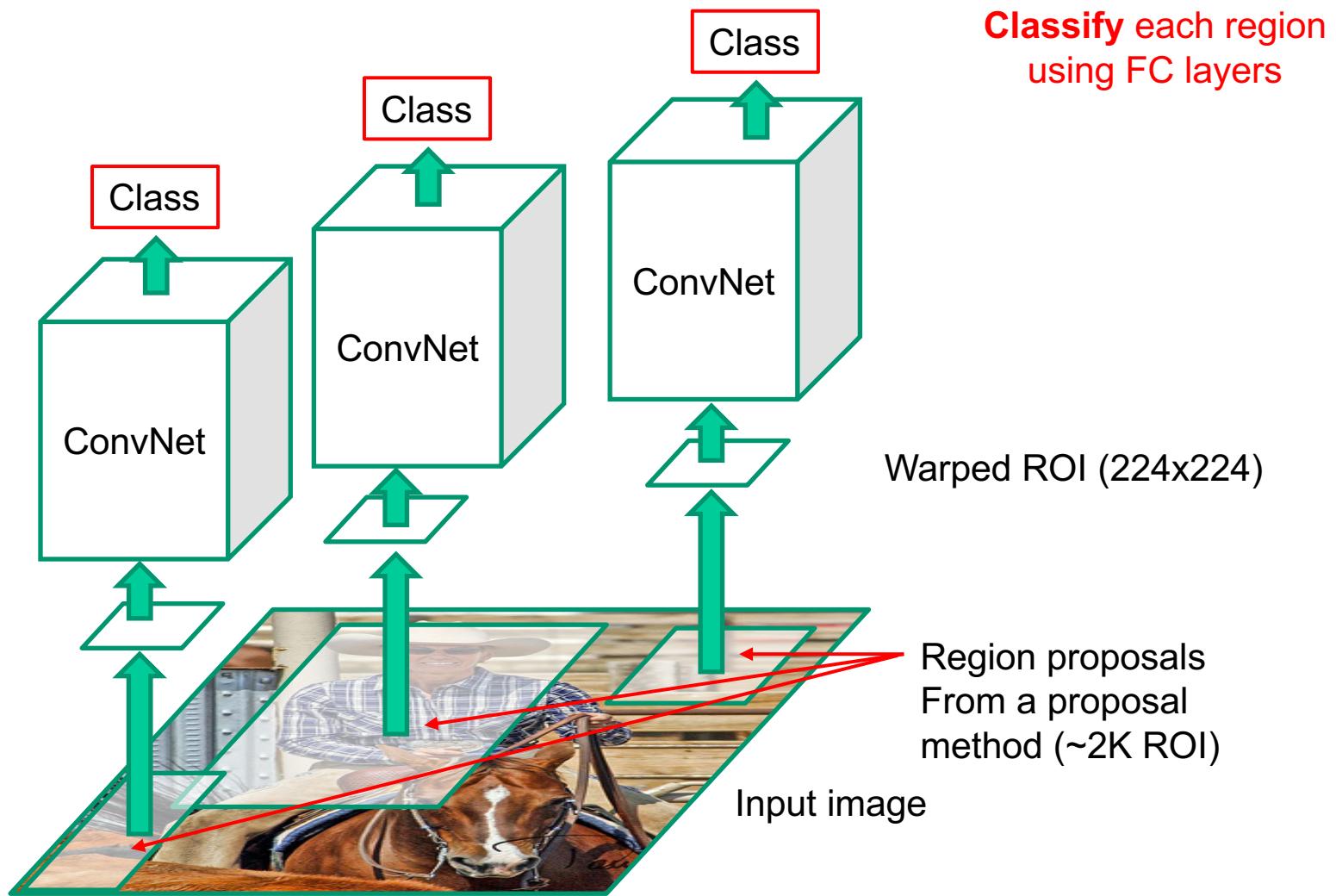
R-CNN: Region-Based CNN



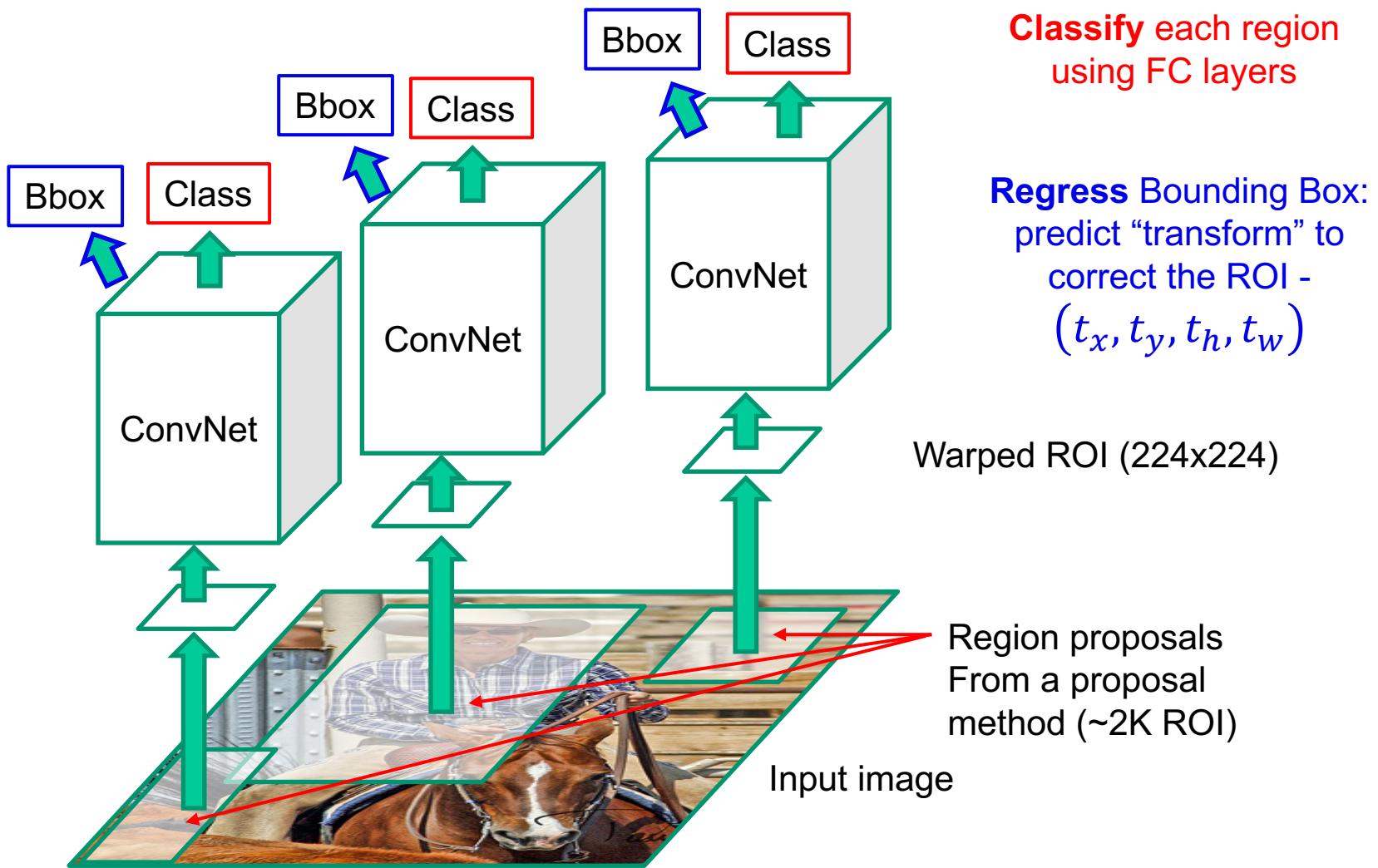
R-CNN: Region-Based CNN



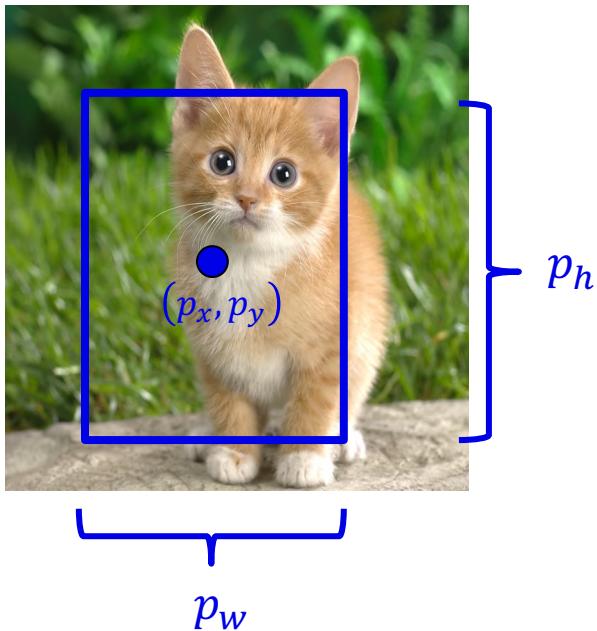
R-CNN: Region-Based CNN



R-CNN: Region-Based CNN

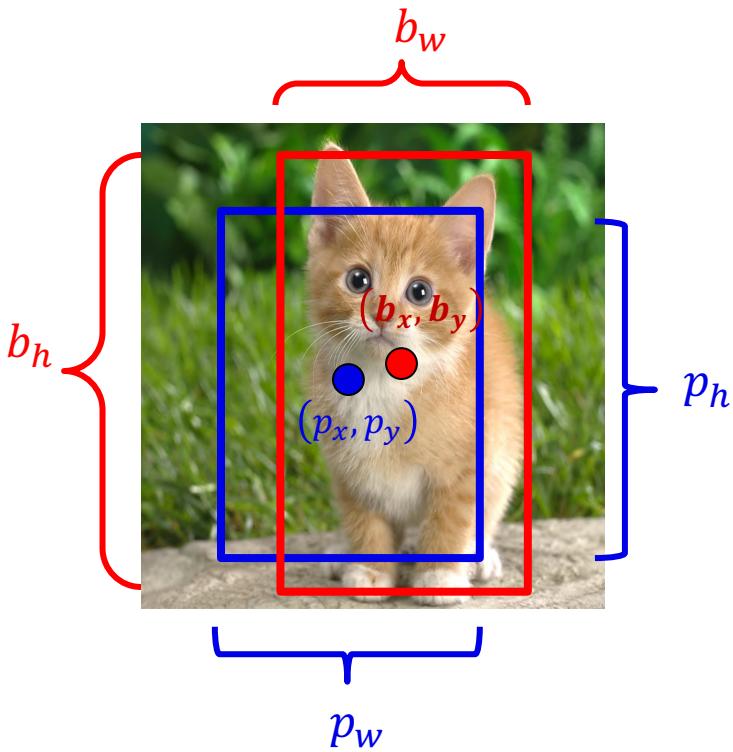


R-CNN: Box Regression



- Consider a Region Proposal with center (p_x, p_y) and size (p_h, p_w)
- The model predicts a transform to correct the region proposal: (t_x, t_y, t_w, t_h)

R-CNN: Box Regression



- Consider a **Region Proposal** with center (p_x, p_y) and size (p_h, p_w)
- The model predicts **a transform** to correct the region proposal: (t_x, t_y, t_w, t_h)
- The **output box** is defined by:

$$b_x = p_x + p_w * t_x$$

Shift center by amount relative to the RP size

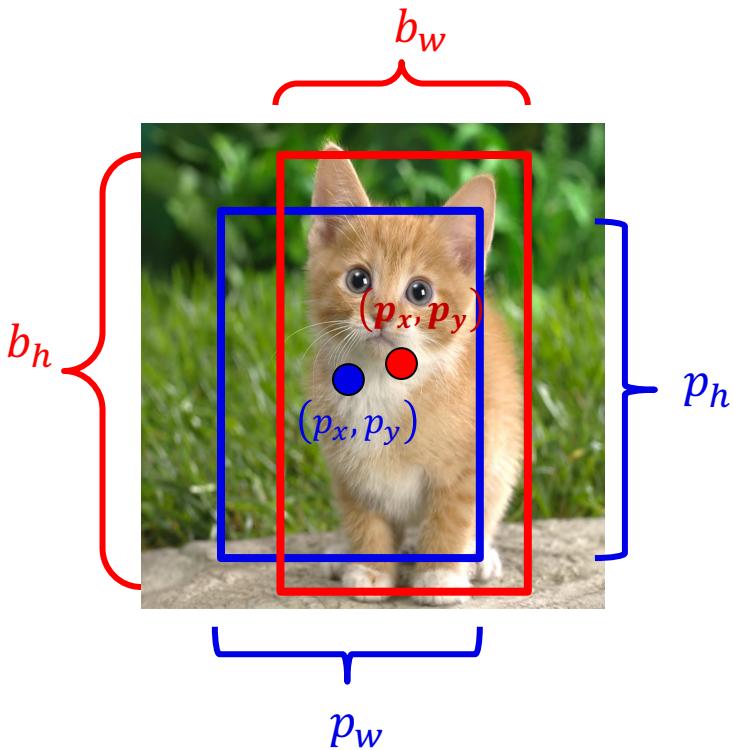
$$b_y = p_y + p_h * t_y$$

Scale RP ; exp ensures that scaling factor > 0

$$b_w = p_w * \exp(t_w)$$

$$b_h = p_h * \exp(t_h)$$

R-CNN: Box Regression



- Consider a **Region Proposal** with center (p_x, p_y) and size (p_h, p_w)
- The model predicts **a transform** to correct the region proposal: (t_x, t_y, t_w, t_h)
- The **output box** is defined by:

$$b_x = p_x + p_w * t_x$$

When transform=0
output=proposal

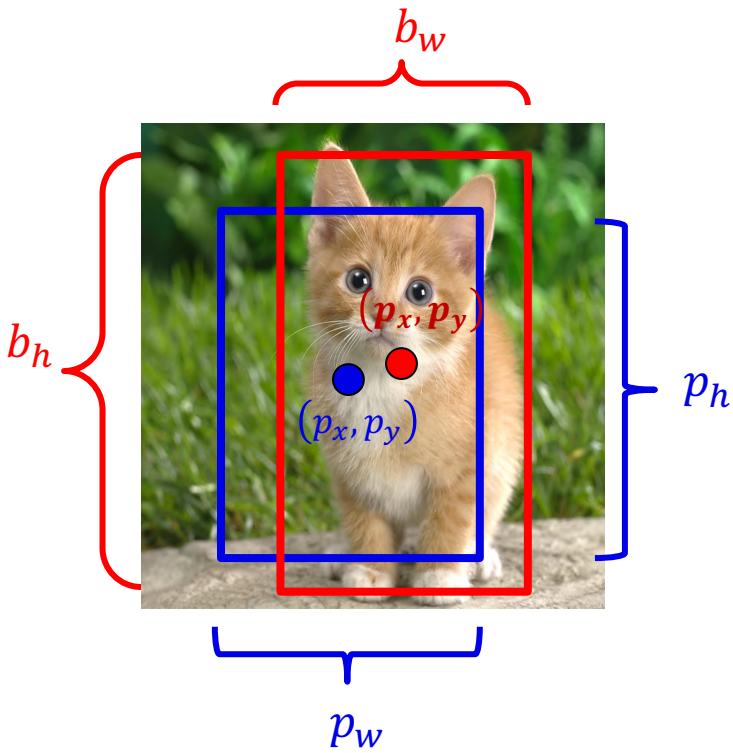
$$b_y = p_y + p_h * t_y$$

L2 regularization
encourages leaving the
proposal unchanged

$$b_w = p_w * \exp(t_w)$$

$$b_h = p_h * \exp(t_h)$$

R-CNN: Box Regression



- Consider a Region Proposal with center (p_x, p_y) and size (p_h, p_w)
- The model predicts a transform to correct the region proposal: (t_x, t_y, t_w, t_h)
- The output box is defined by:

$$b_x = p_x + p_w * t_x$$

$$b_y = p_y + p_h * t_y$$

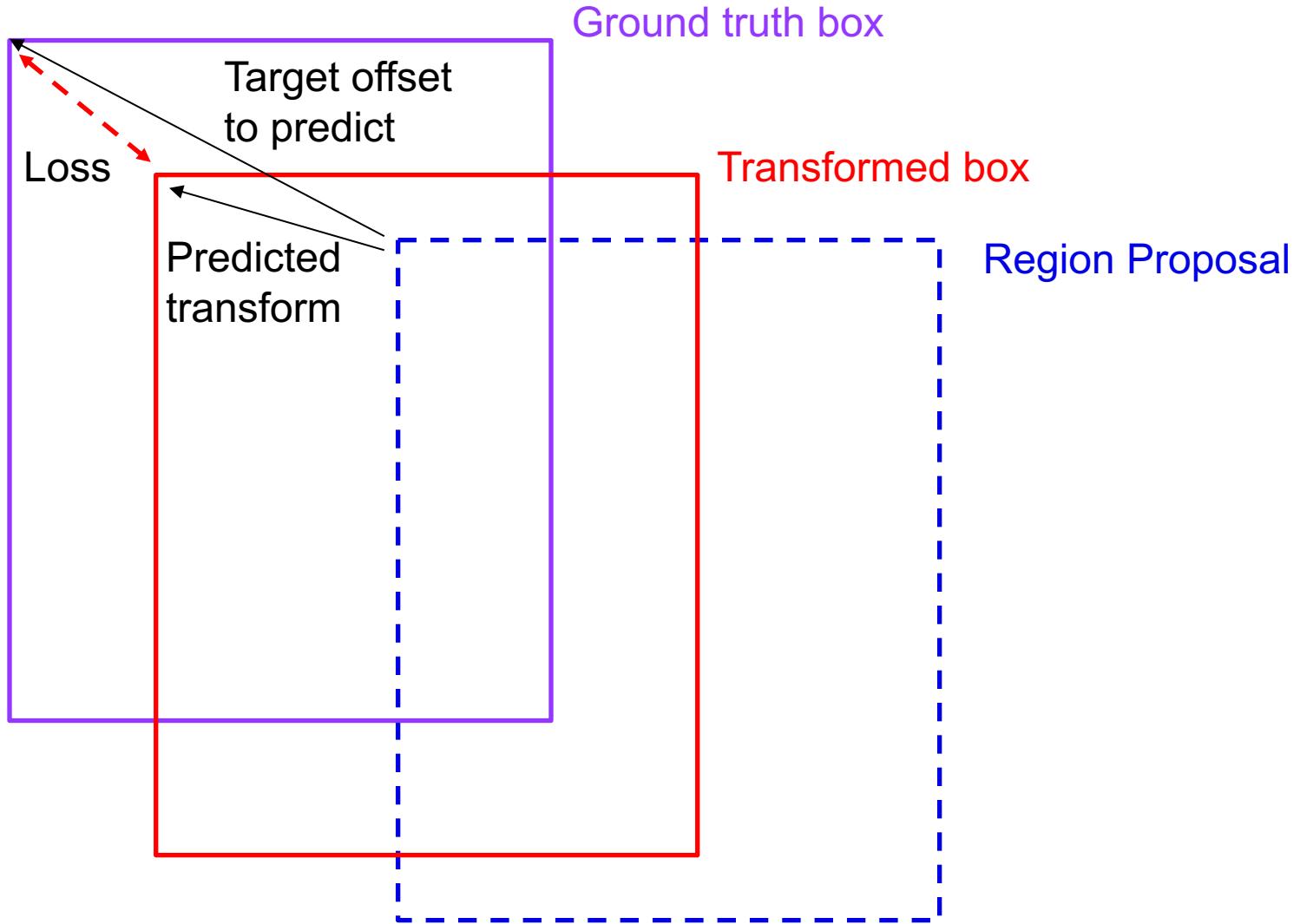
$$b_w = p_w * \exp(t_w)$$

$$b_h = p_h * \exp(t_h)$$

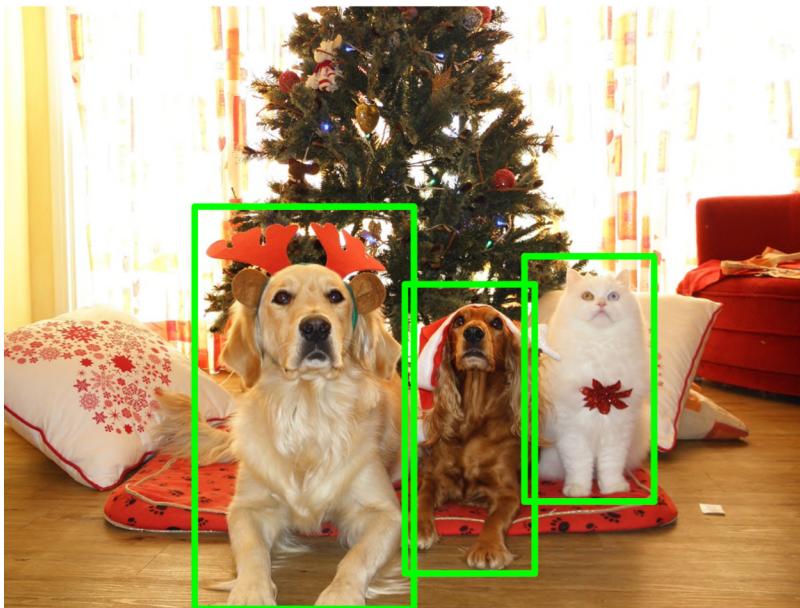
The transform encodes relative difference between proposal and the output;

Note, CNN doesn't see absolute size or position after cropping and warping

Bounding box regression Loss

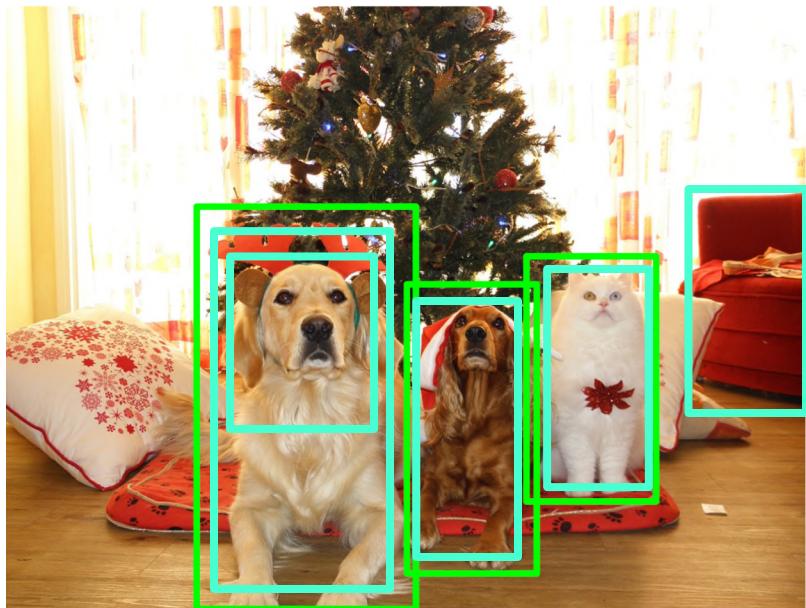


R-CNN Training



GT Boxes

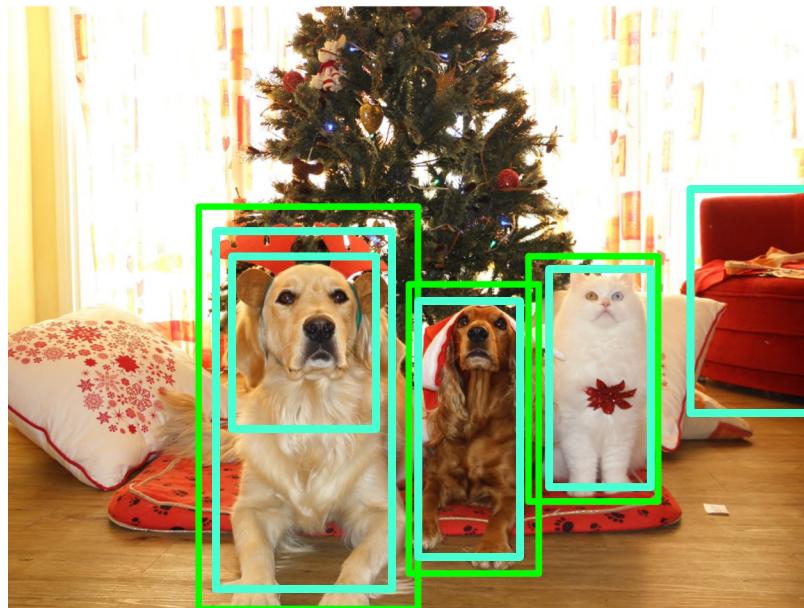
R-CNN Training



GT Boxes

Region Proposals

R-CNN Training



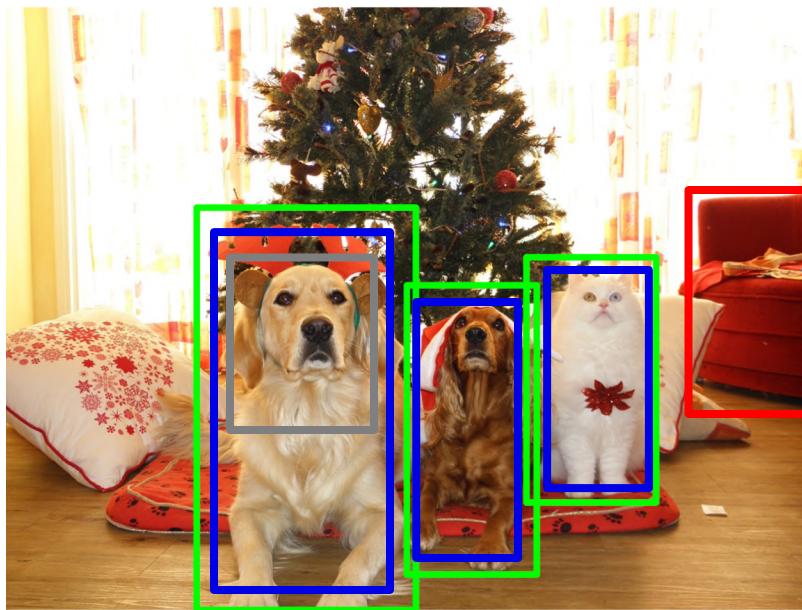
GT Boxes

Region Proposals

Categorize each RP as:

- Positive: $\text{IOU} > 0.5$ with GT
- Negative: $\text{IOU} < 0.3$ with GT
- Neutral: Otherwise

R-CNN Training



Categorize each RP as:

- Positive: $\text{IOU} > 0.5$ with GT
- Negative: $\text{IOU} < 0.3$ with GT
- Neutral: Otherwise

GT Boxes

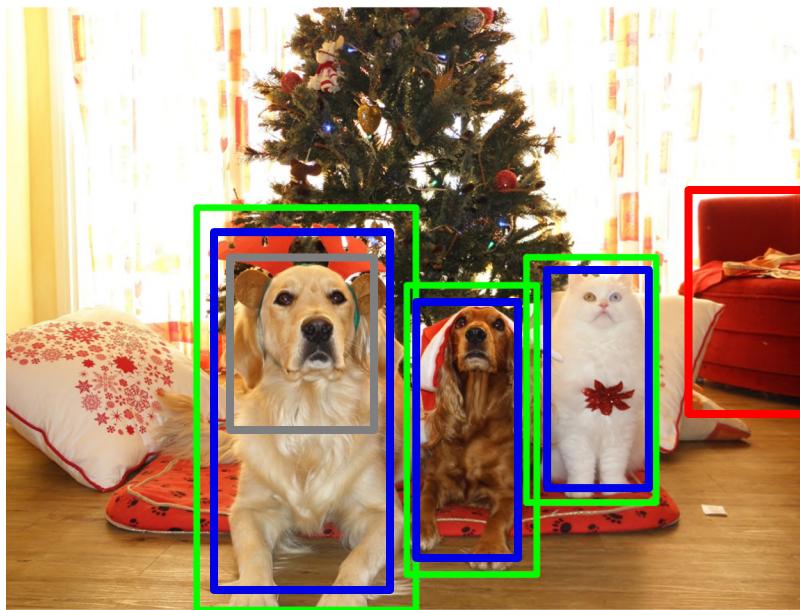
Positive

Negative

Neutral

R-CNN Training

Crop pixels from each **positive** and **negative** proposals and resize to 224x224:



GT Boxes

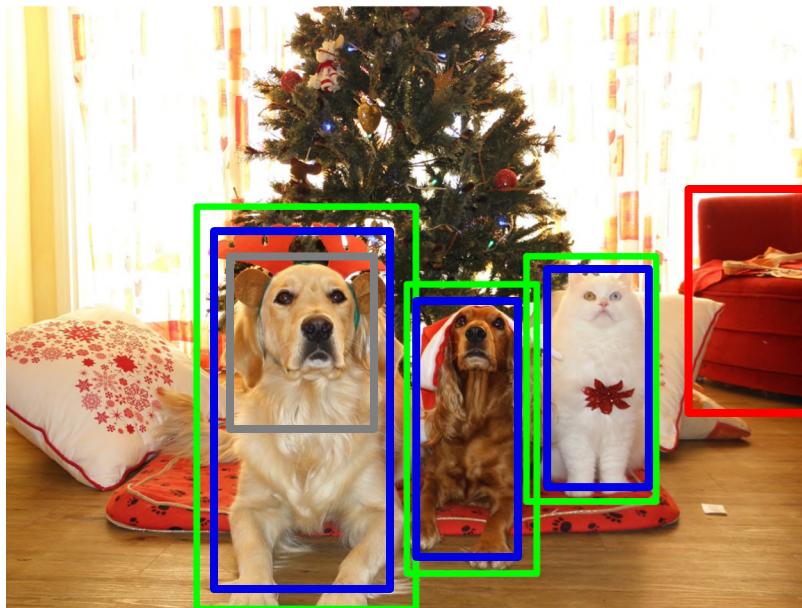
Positive

Negative

Neutral

R-CNN Training

Run each crop through CNN



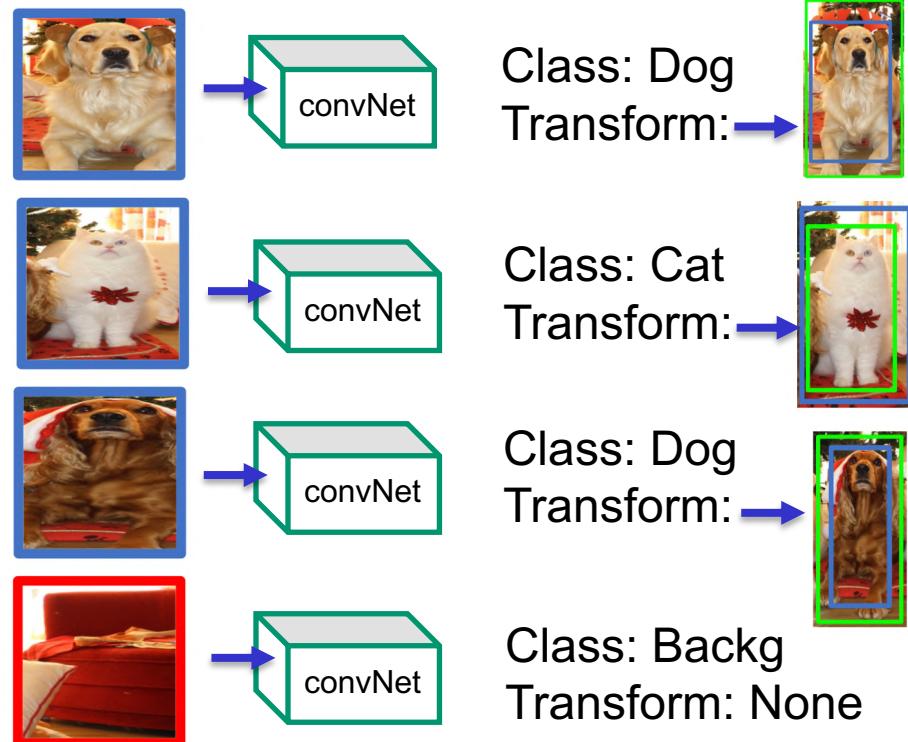
GT Boxes

Positive

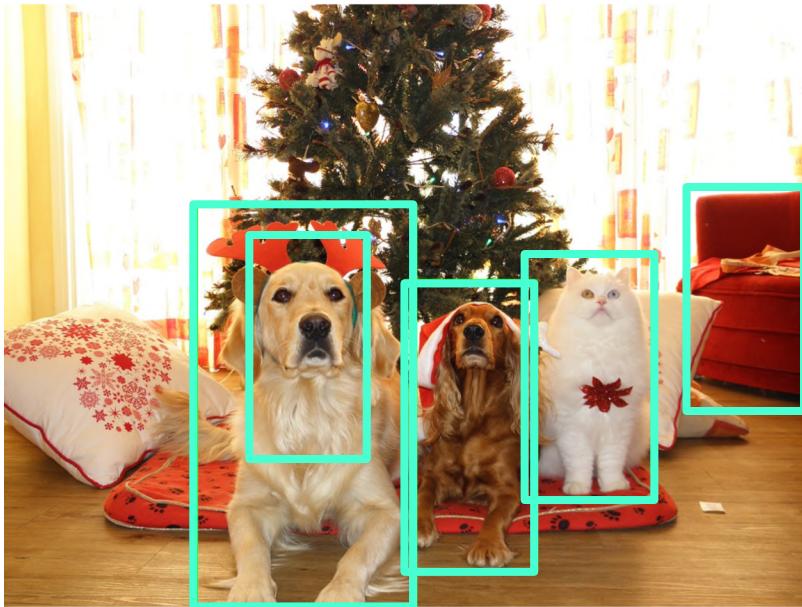
Negative

Neutral

Positive regions: predicts class + transform
Negative regions: predict only class (backg)



R-CNN Test-Time



Input: an image

1. Run proposal method \Rightarrow 2000 RPs.
2. Run CNN on each RP to get Class scores and Transform.
3. Threshold class scores to get a set of detections.

Problems:

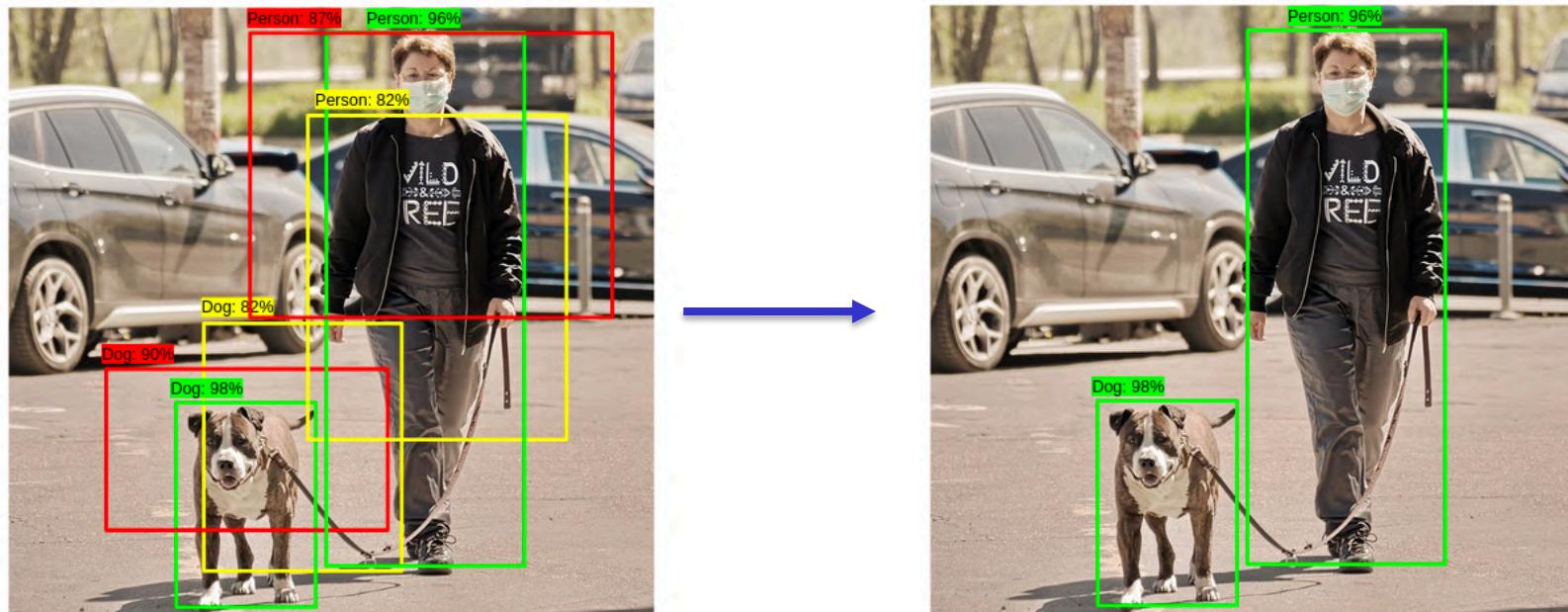
- CNN often outputs overlapping boxes
- How to set the Thresholds?

Post Processing: Non-Max Supression

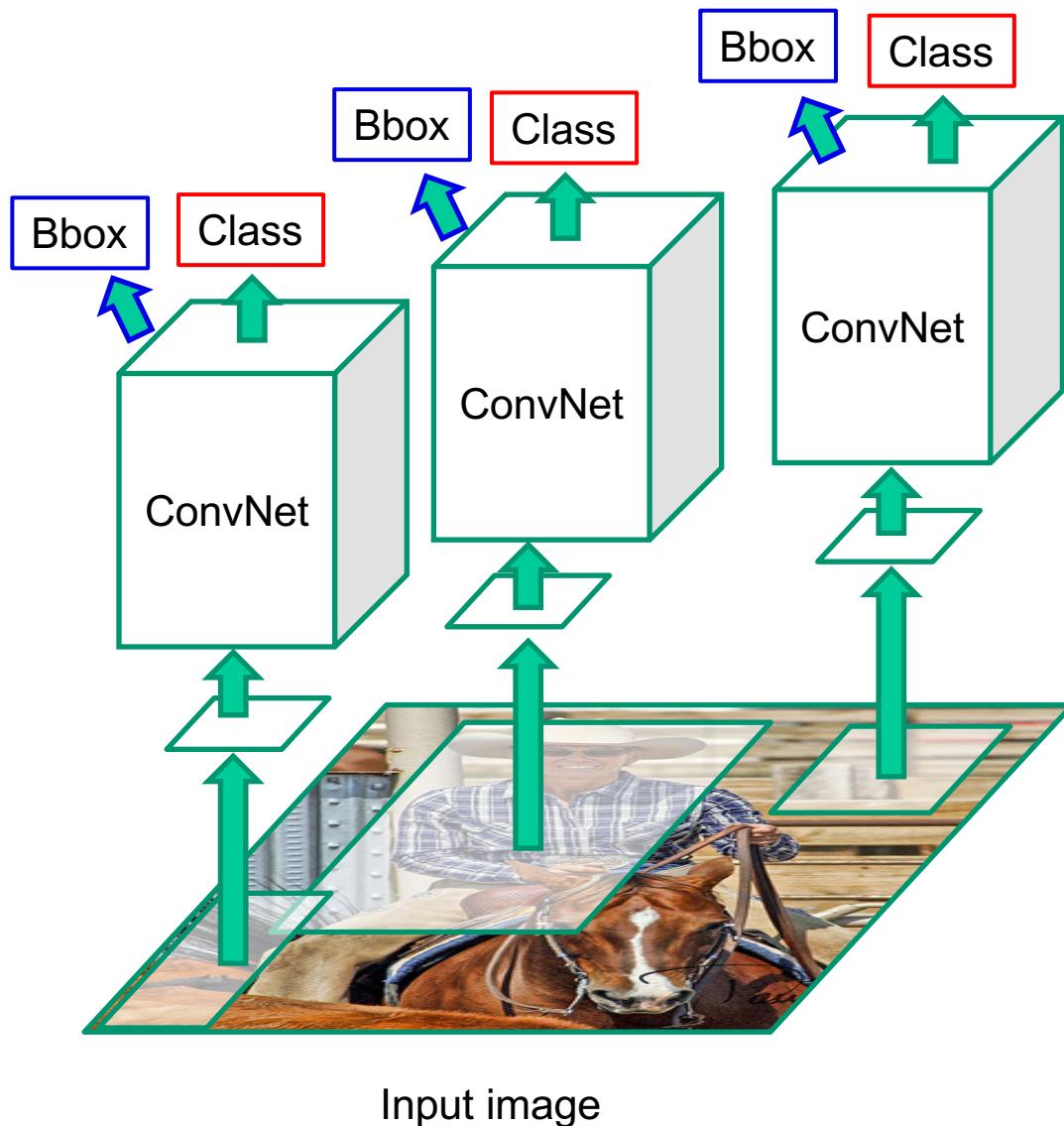
Non-max suppression:

Sort the predictions by the confidence scores.

1. Select the next highest-score box,
2. Eliminate any current prediction if it has the same class and $\text{IoU} > 0.5$ with the current prediction.
3. Repeat step 1 until all predictions are checked.



R-CNN: Problems

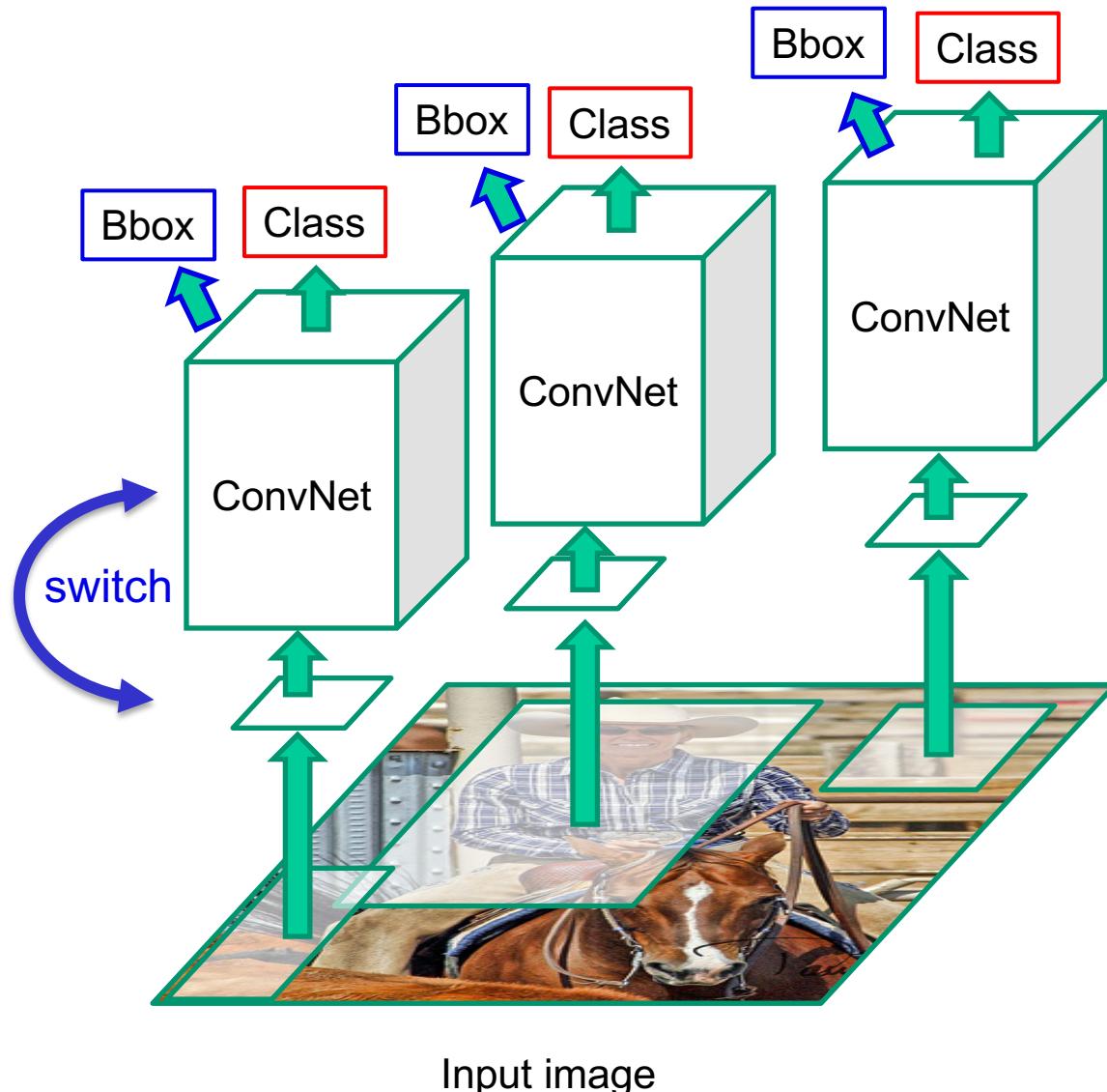


Problem: Very slow!

Need to run ~2000
forward passes for each
image.

Waste calculations for
overlapping RPs

R-CNN: Problems



Problem: Very slow!

Need to run ~2000
forward passes for each
image.

Waste calculations for
overlapping RPs

Solution: Run CNN
before cropping.

Fast R-CNN

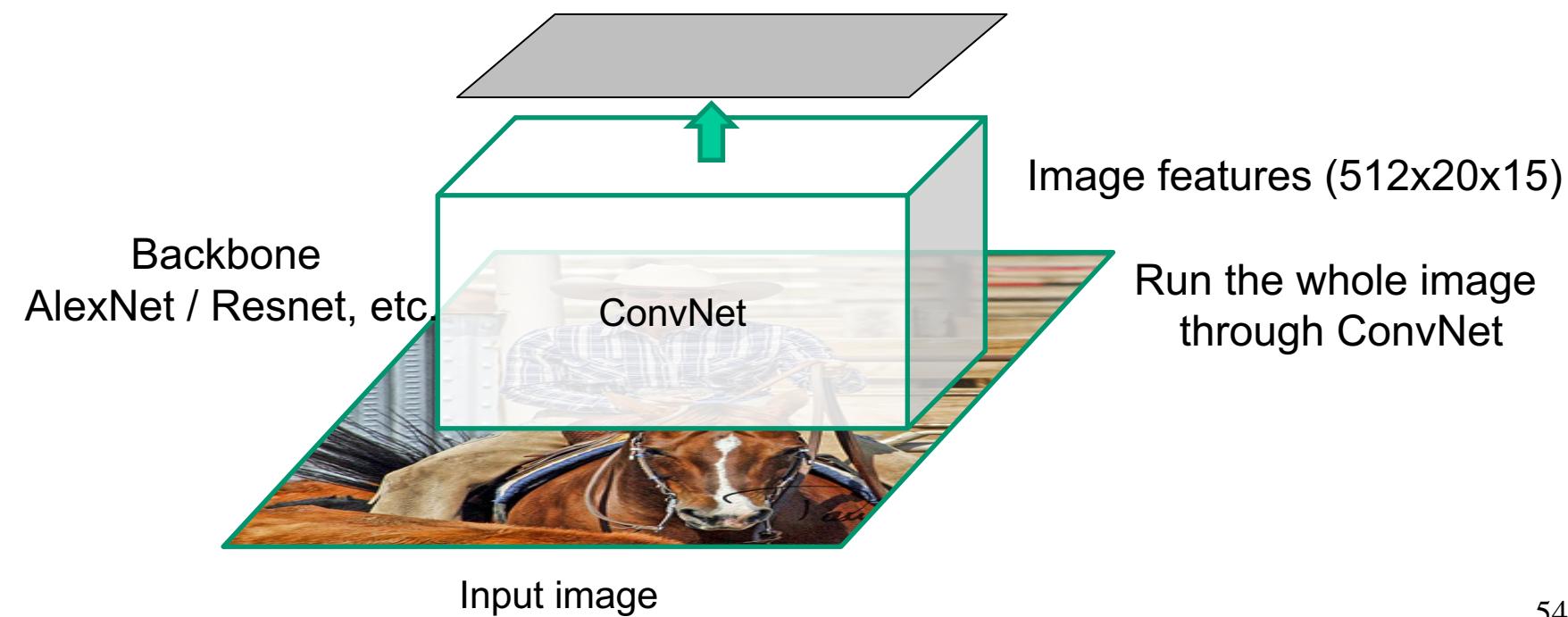


Input image

53

R. Girshick, [Fast R-CNN](#), ICCV 2015

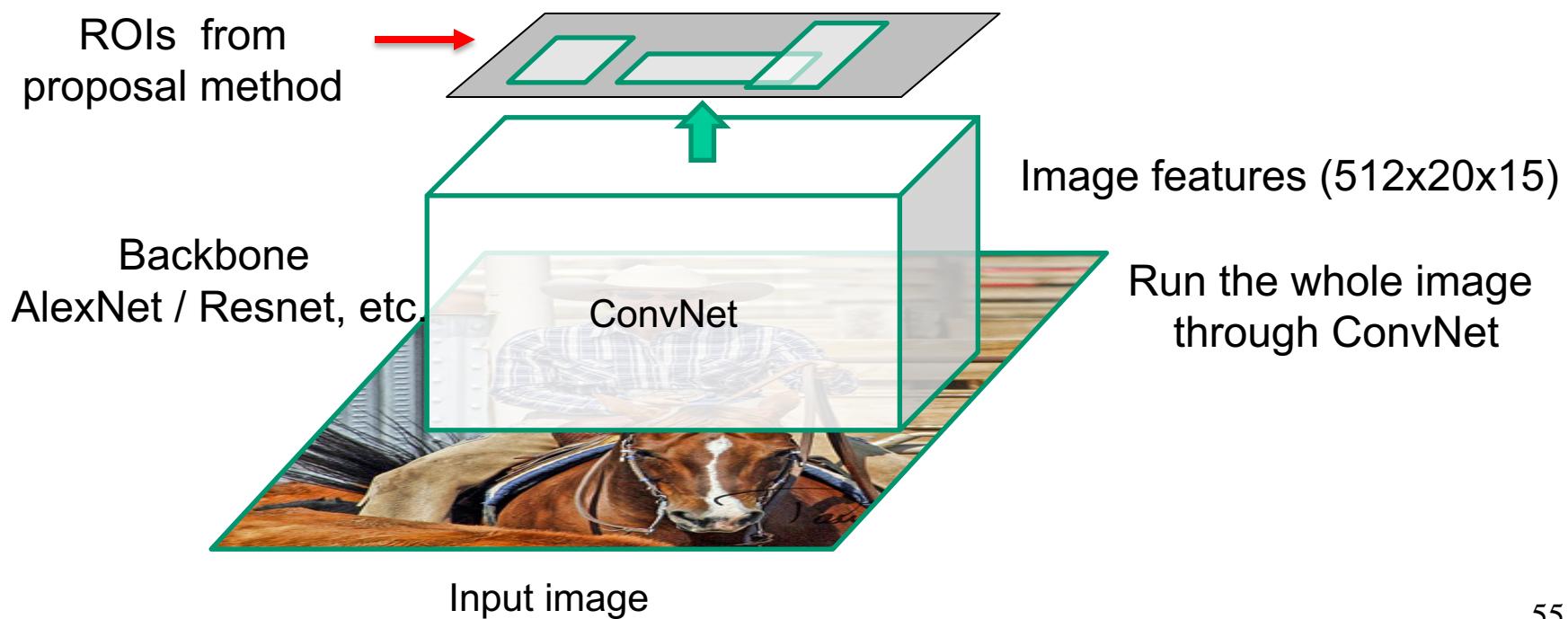
Fast R-CNN



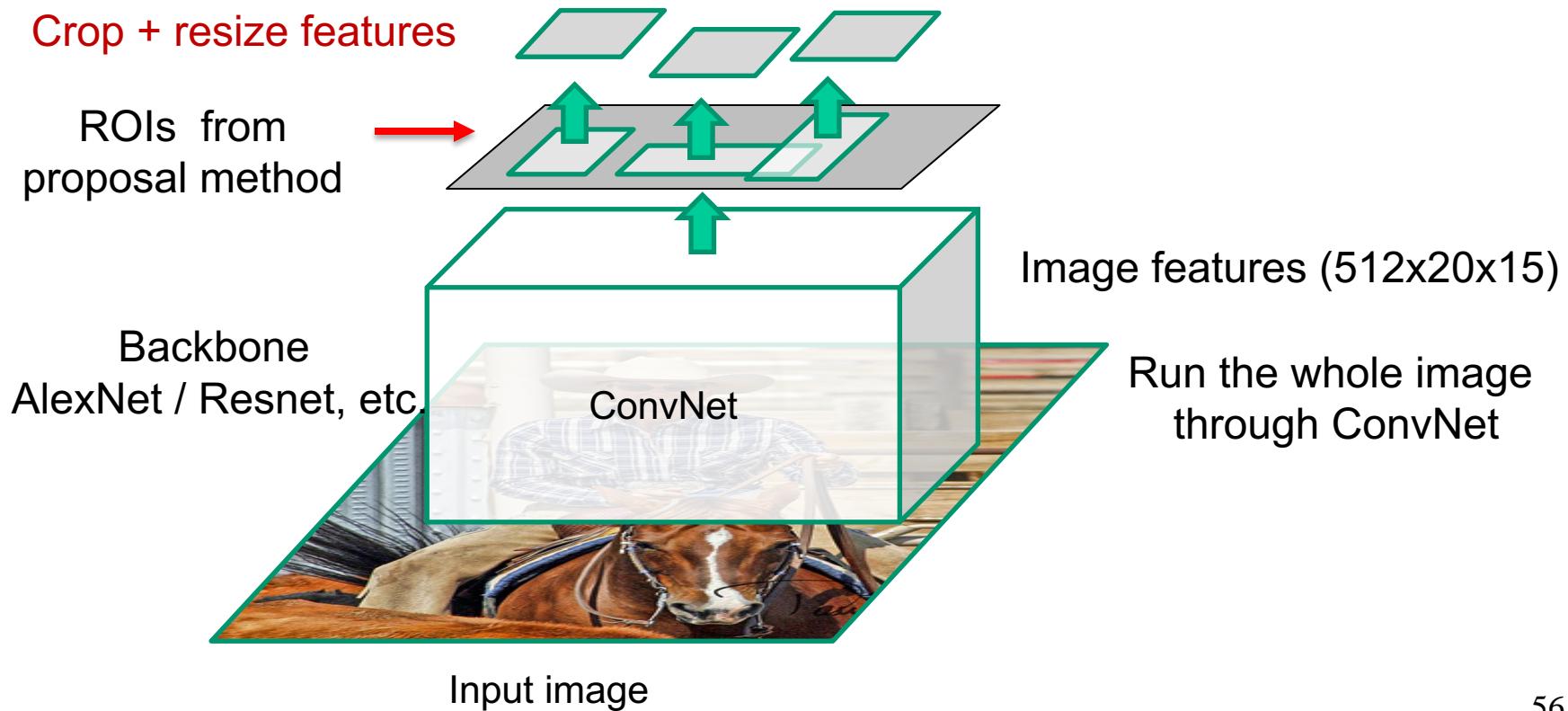
54

R. Girshick, [Fast R-CNN](#), ICCV 2015

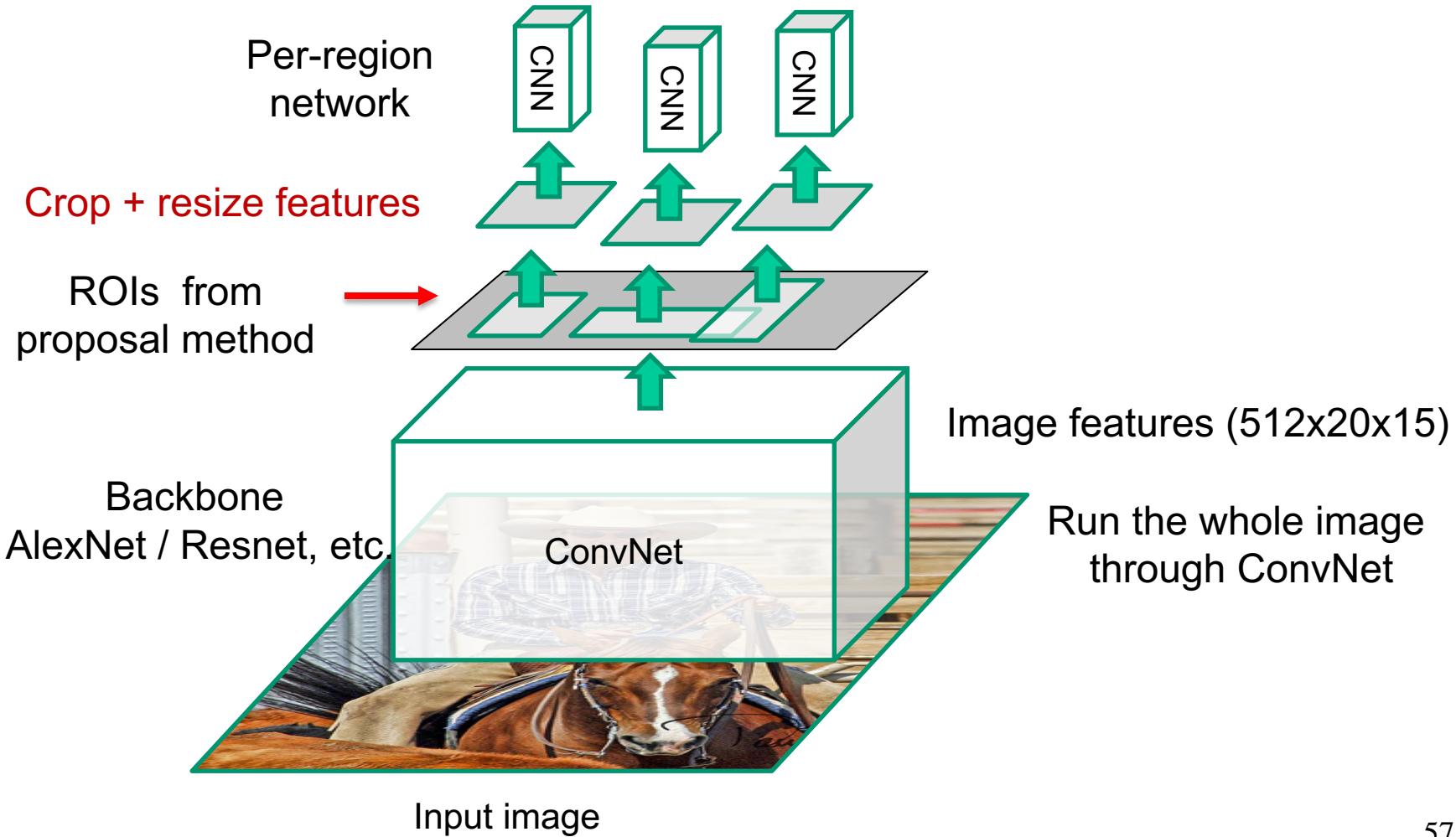
Fast R-CNN



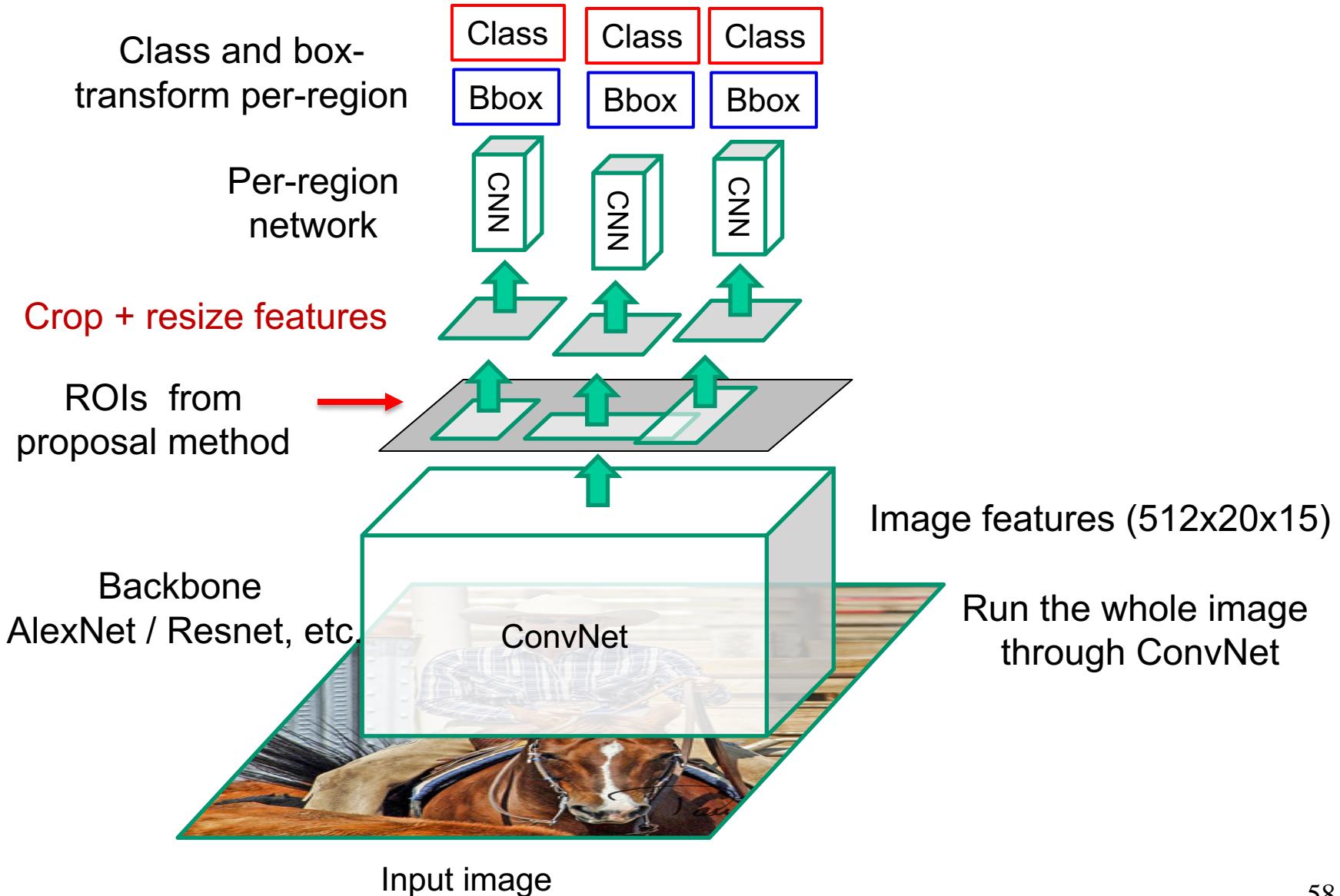
Fast R-CNN



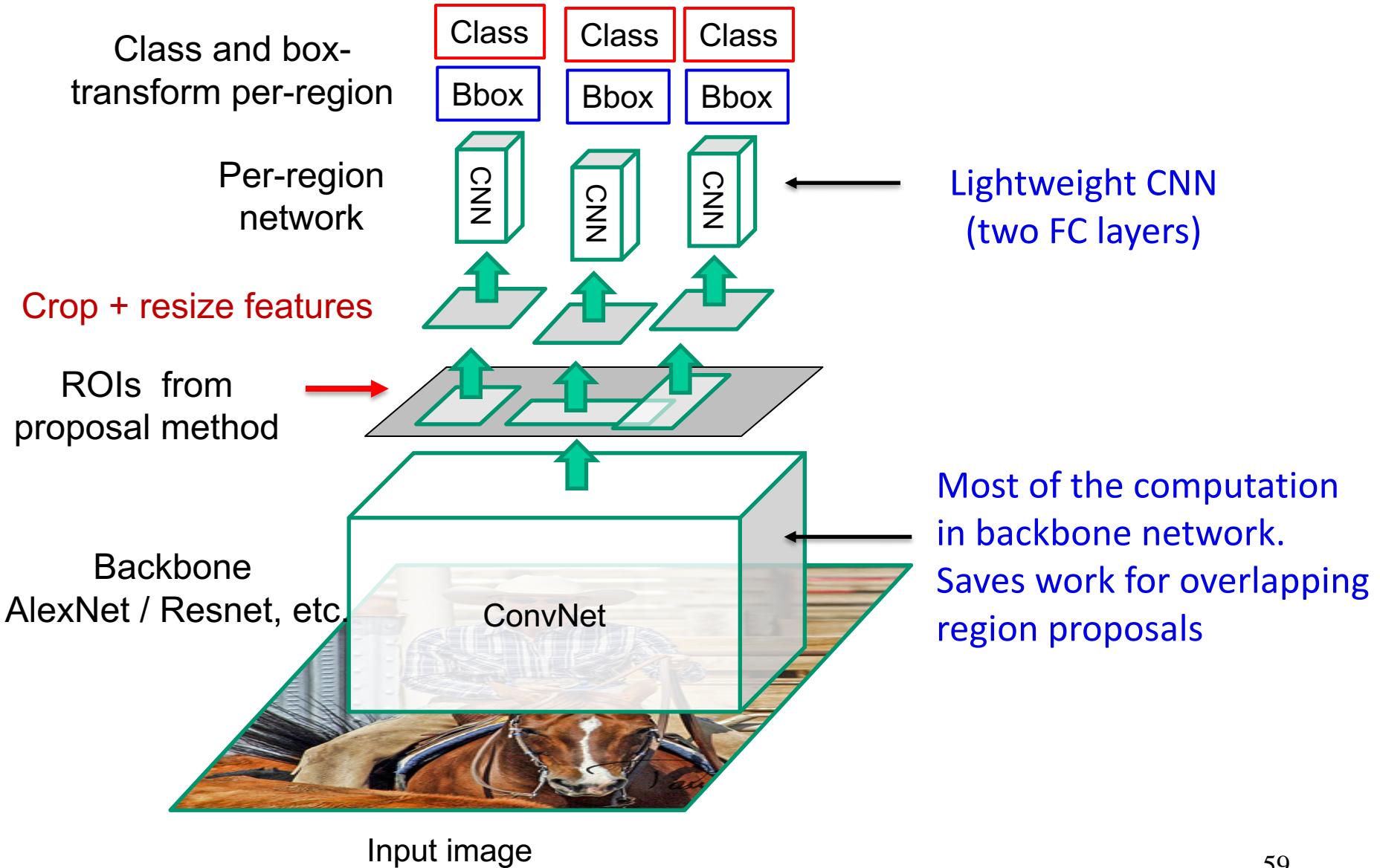
Fast R-CNN



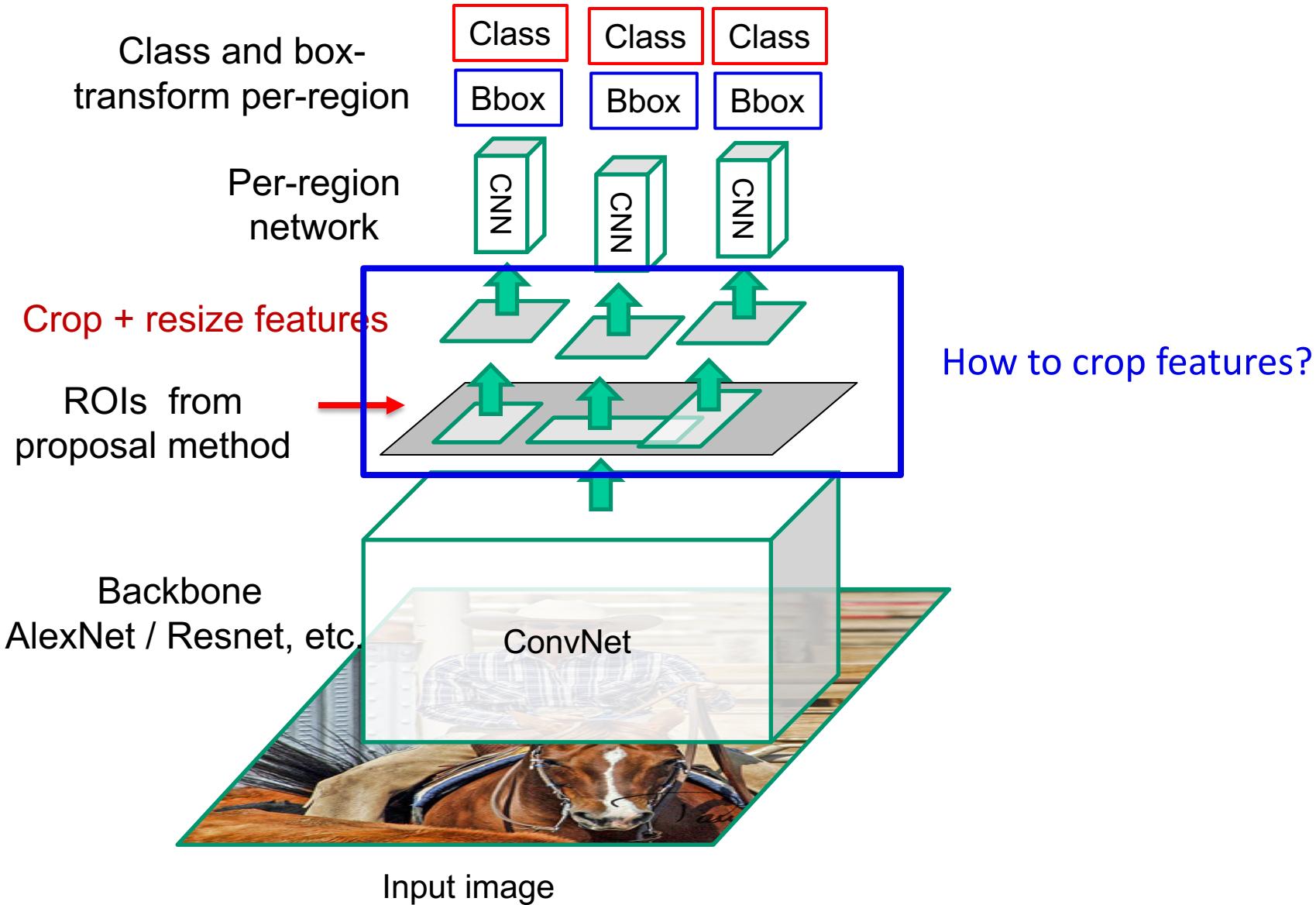
Fast R-CNN



Fast R-CNN



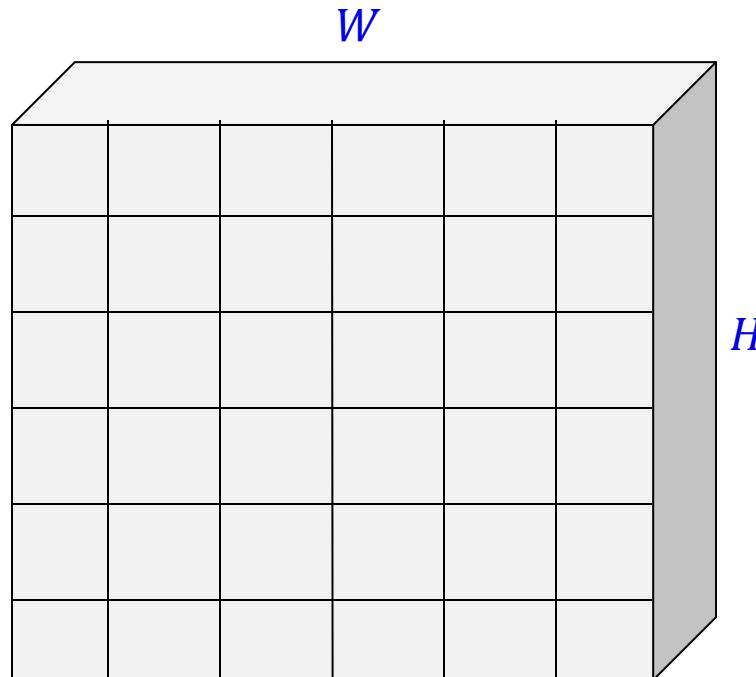
Fast R-CNN



Fast R-CNN: Cropping using ROI Align



Input image
(e.g. 3x640x480)

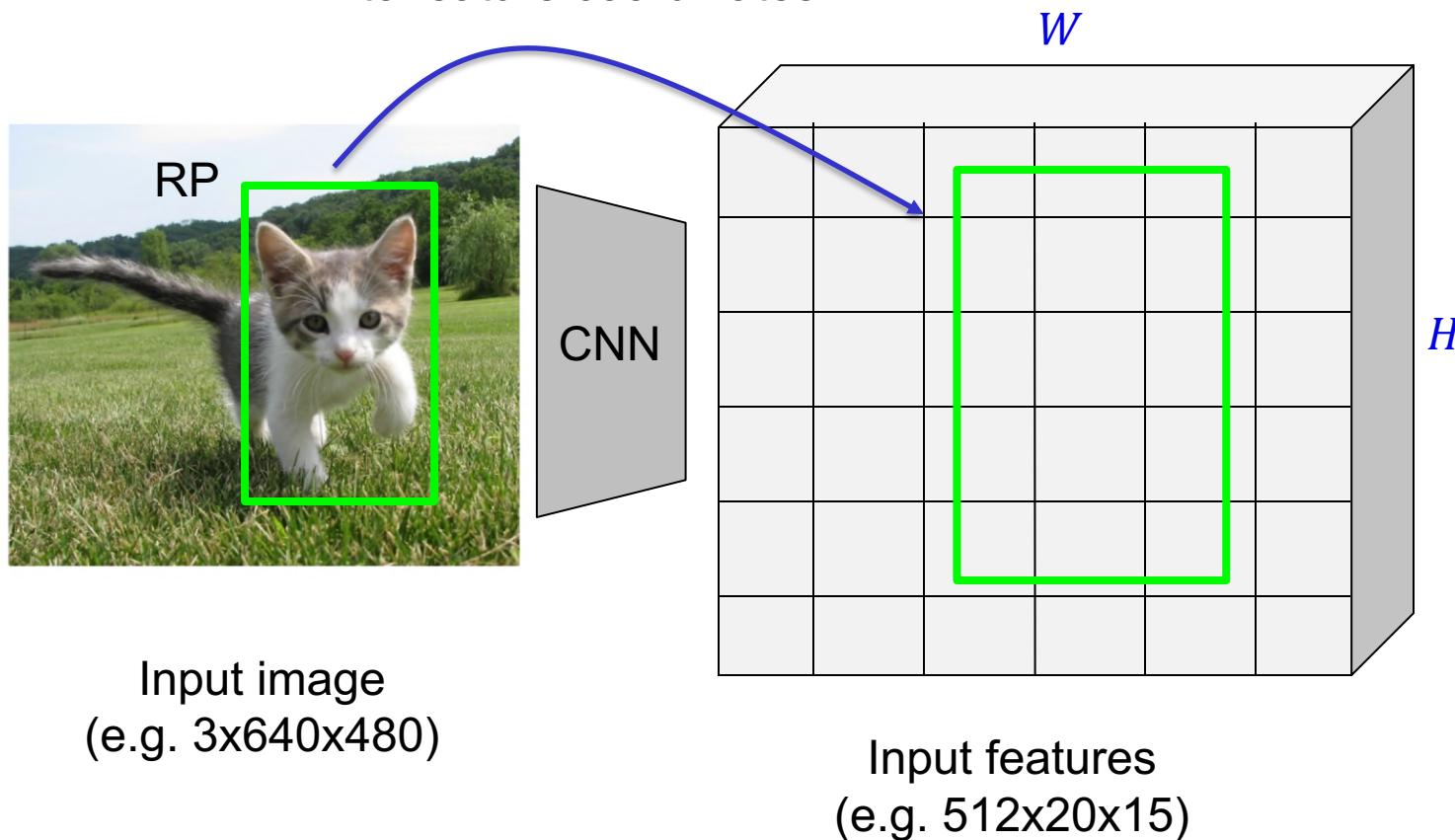


Input features
(e.g. 512x20x15)

Fast R-CNN: Cropping using ROI Align

62

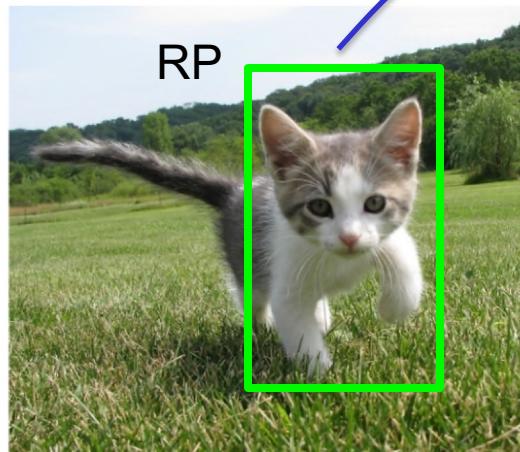
Map RP from image coordinates
to feature coordinates



Fast R-CNN: Cropping using ROI Align

63

Map RP from image coordinates
to feature coordinates



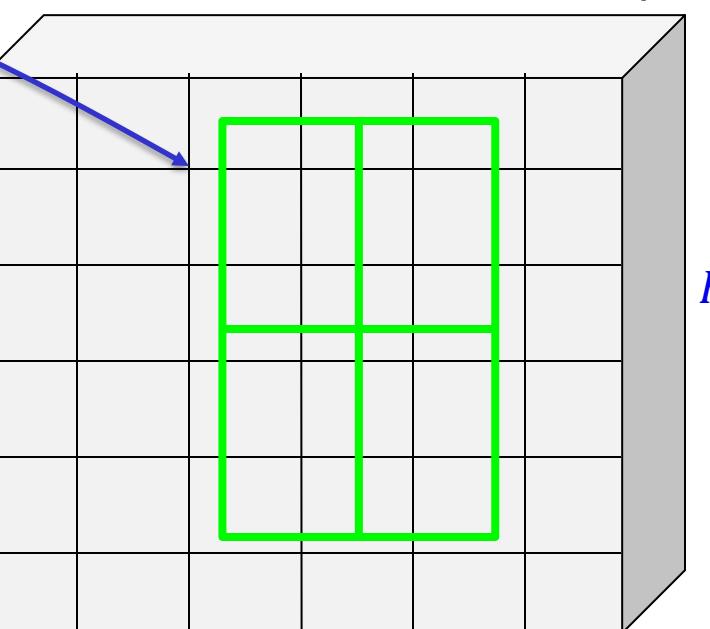
Input image
(e.g. 3x640x480)

CNN

W

H

Divide RP into $k \times k$
subregions
(here 512x2x2, in
practice 512x7x7)

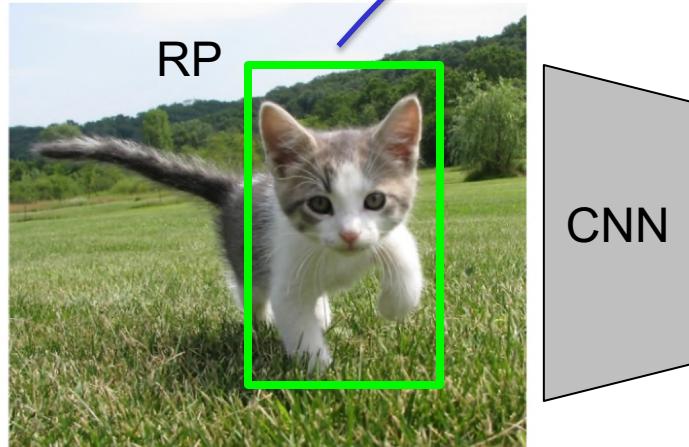


Input features
(e.g. 512x20x15)

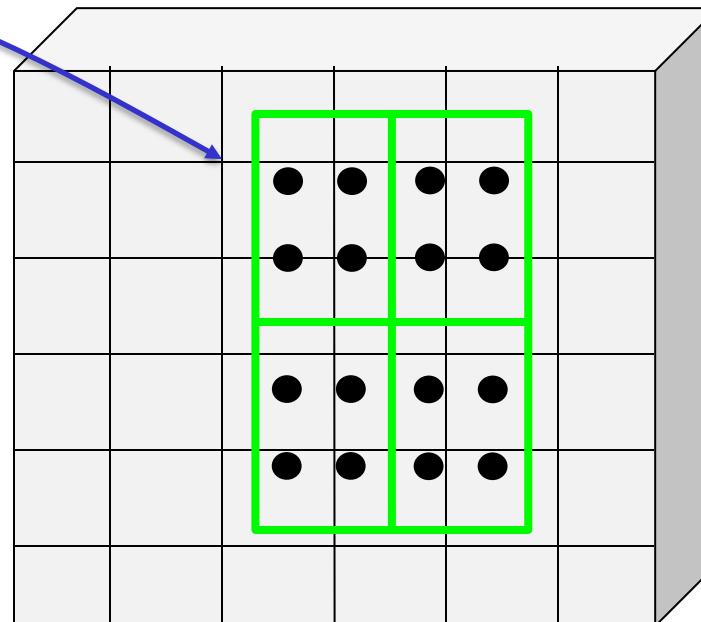
Fast R-CNN: Cropping using ROI Align

64

Map RP from image coordinates
to feature coordinates



Input image
(e.g. 3x640x480)



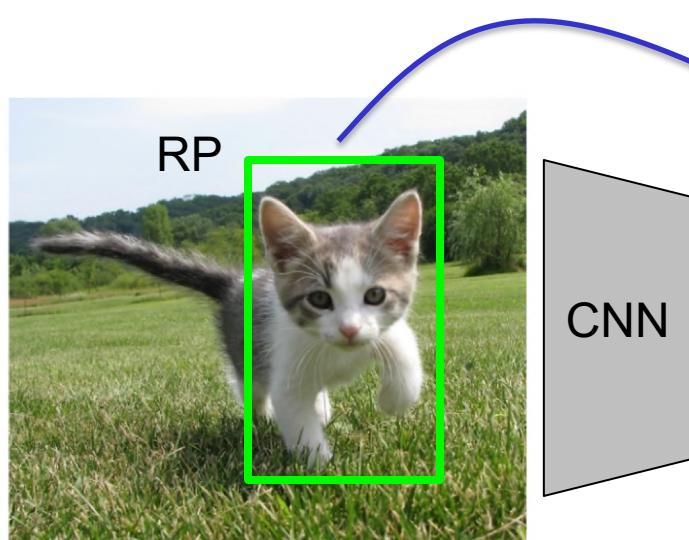
Input features
(e.g. 512x20x15)

Divide RP into $k \times k$
subregions
(here 512x2x2, in
practice 512x7x7)

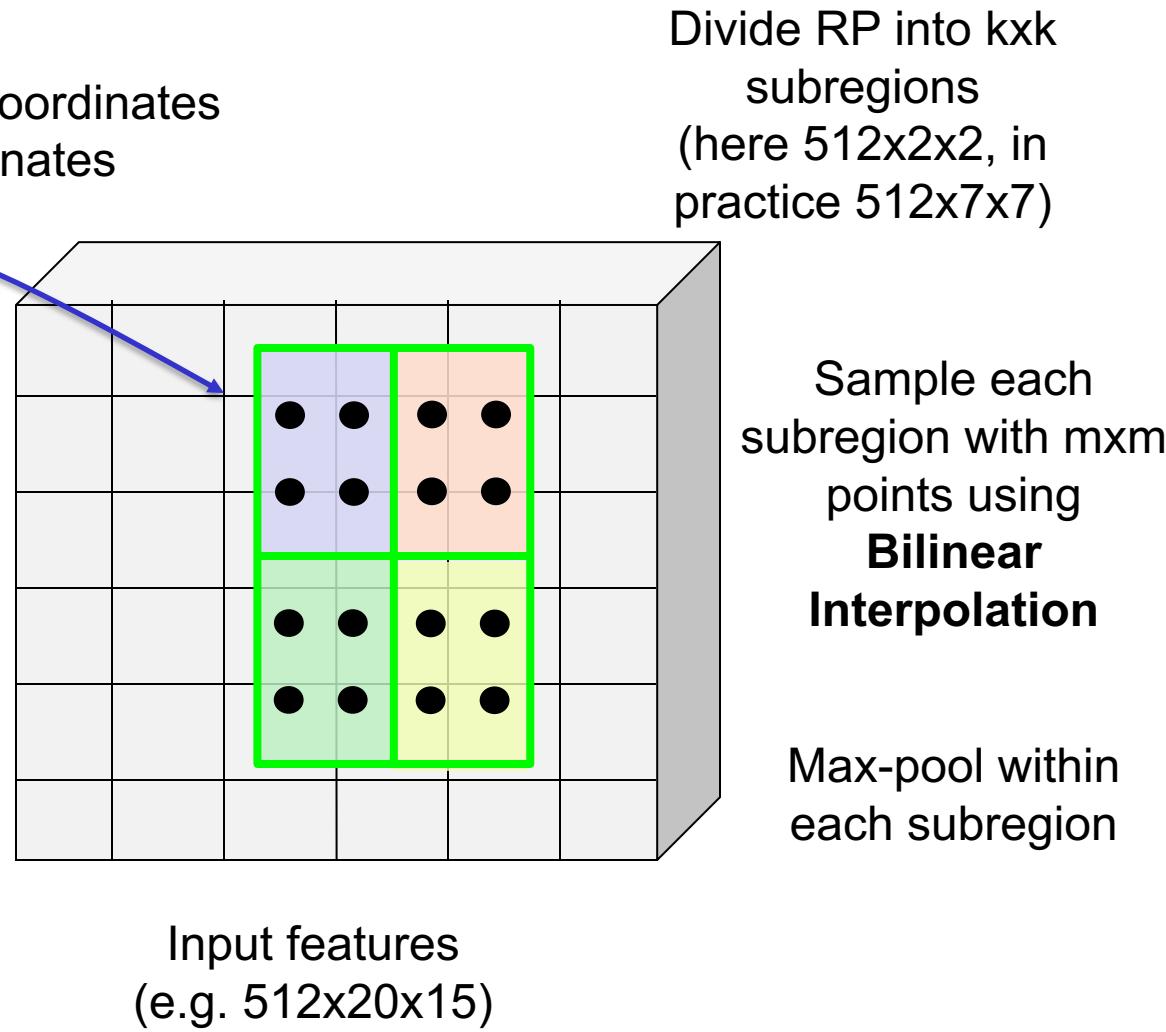
Sample each
subregion with $m \times m$
points using
**Bilinear
Interpolation**

Fast R-CNN: Cropping using ROI Align

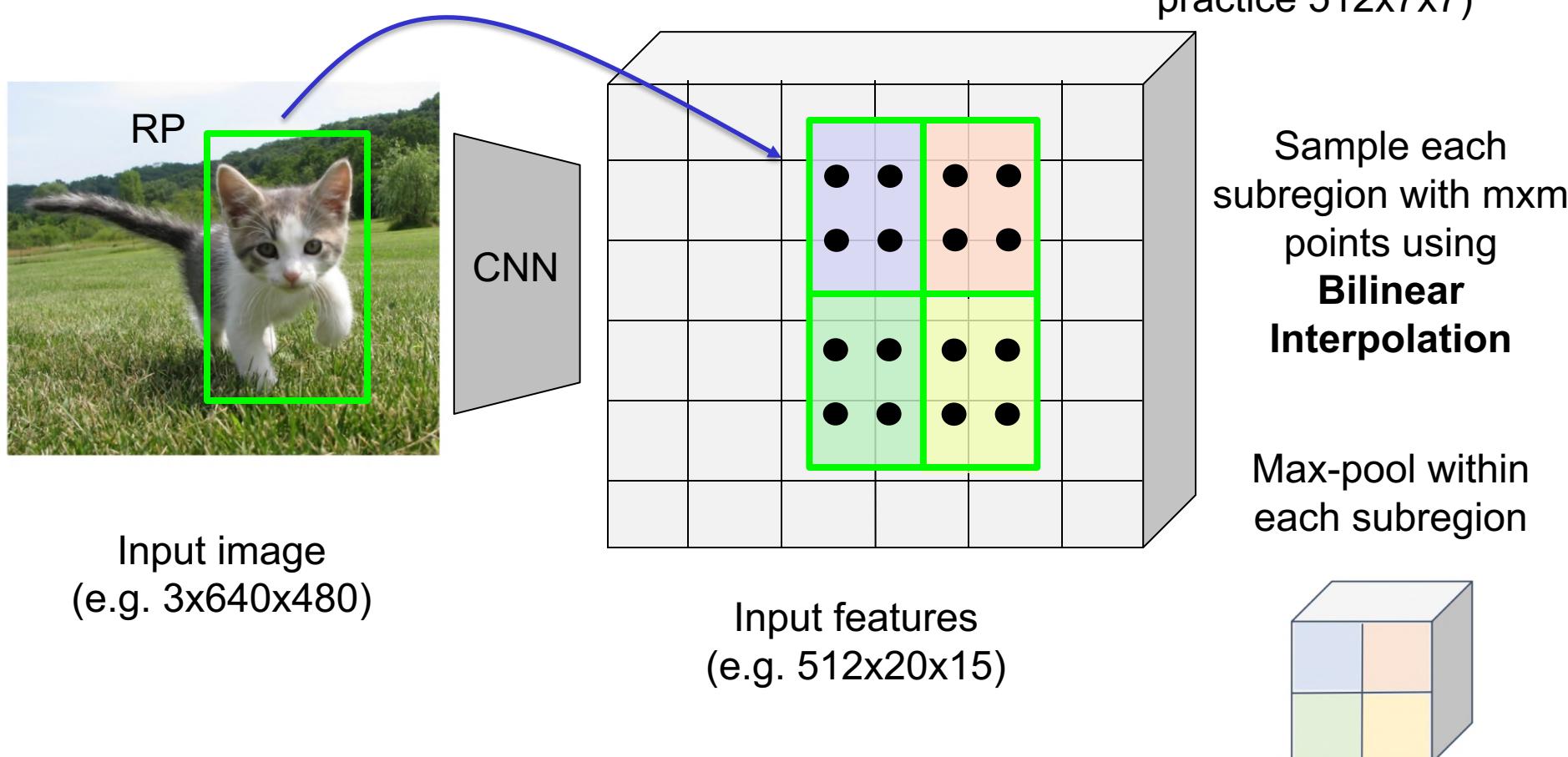
Map RP from image coordinates
to feature coordinates



Input image
(e.g. 3x640x480)

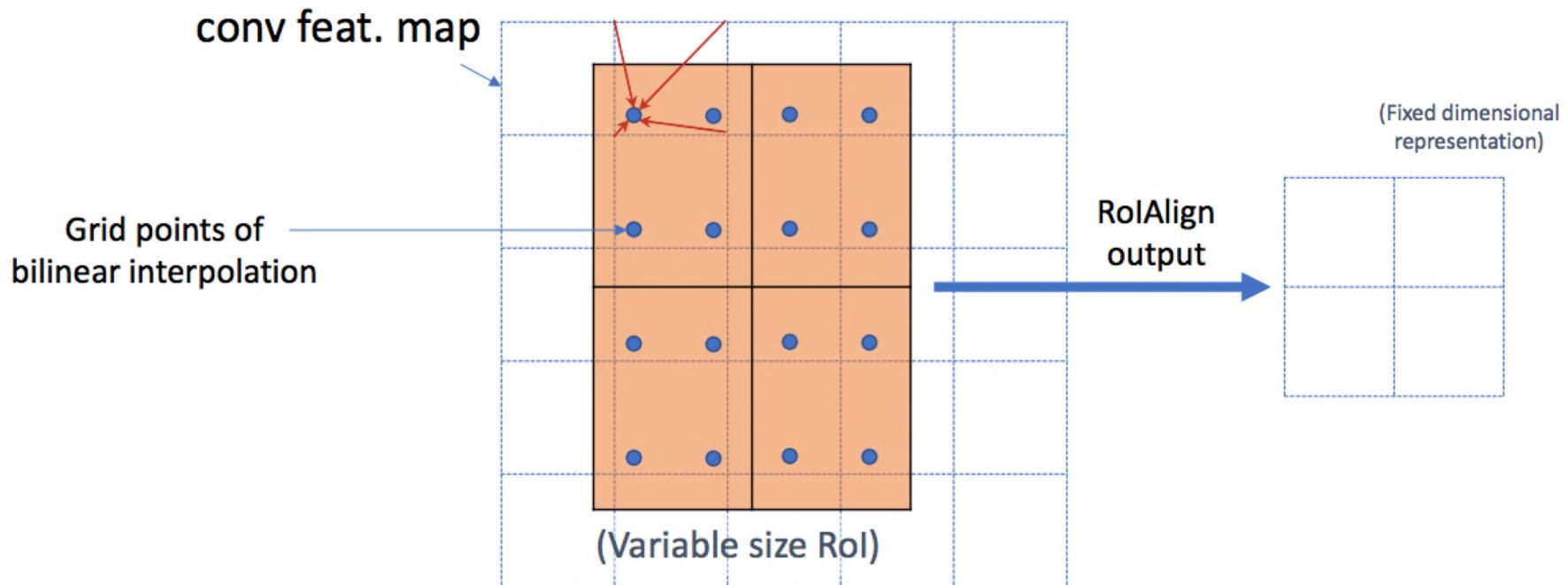


Fast R-CNN: Cropping using ROI Align



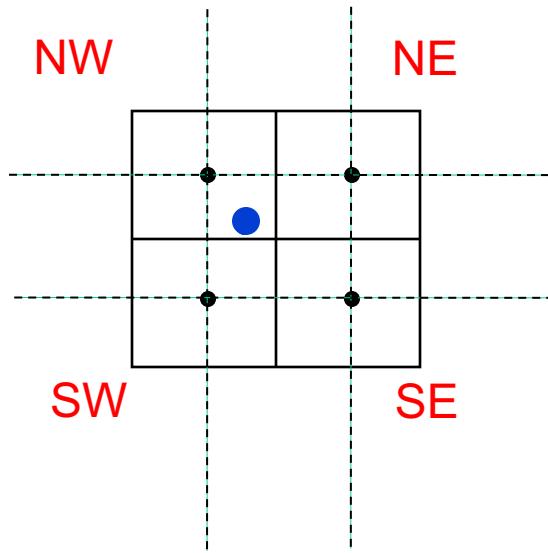
ROI Align

- Bilinear interpolation then Max-pooling

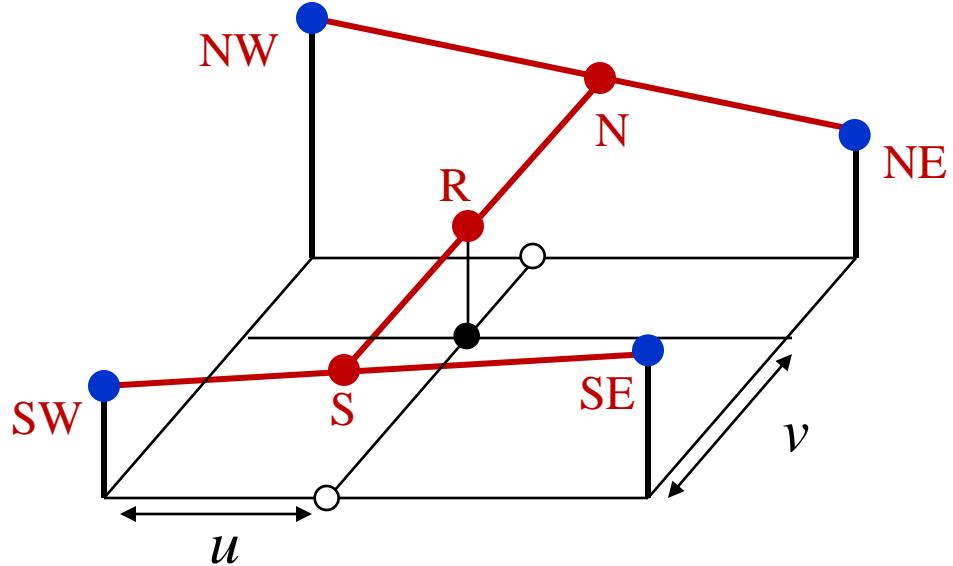


Bilinear interpolation

- Each sample has 4 nearest neighbor values.



Bilinear interpolation



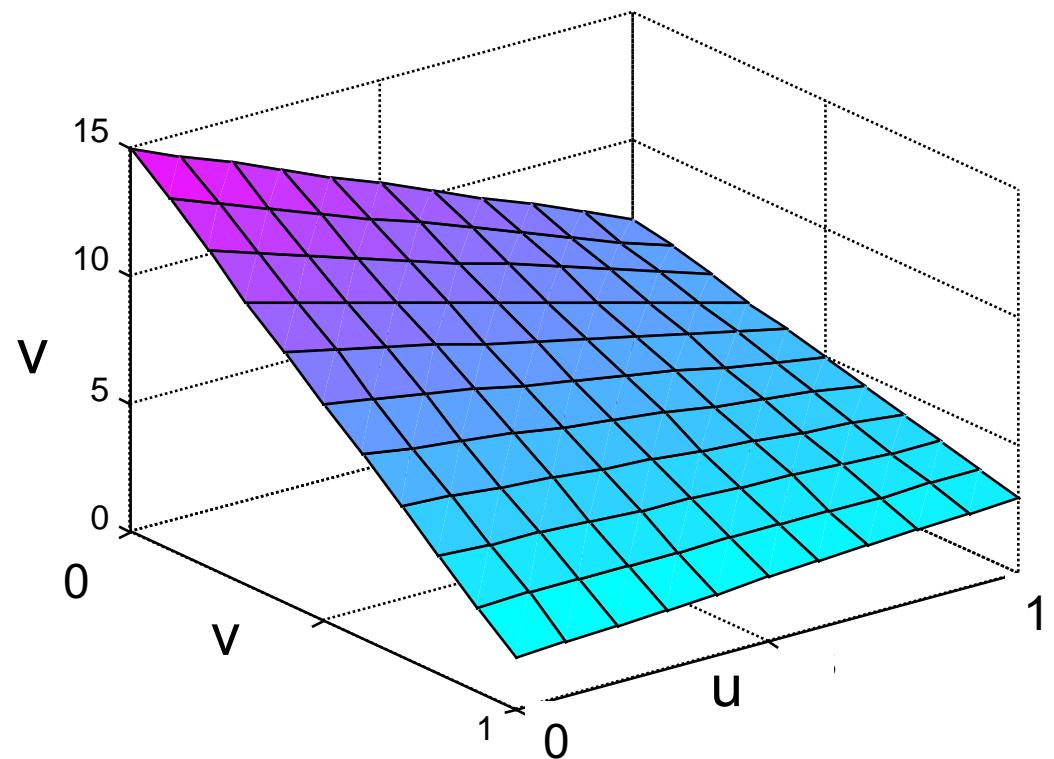
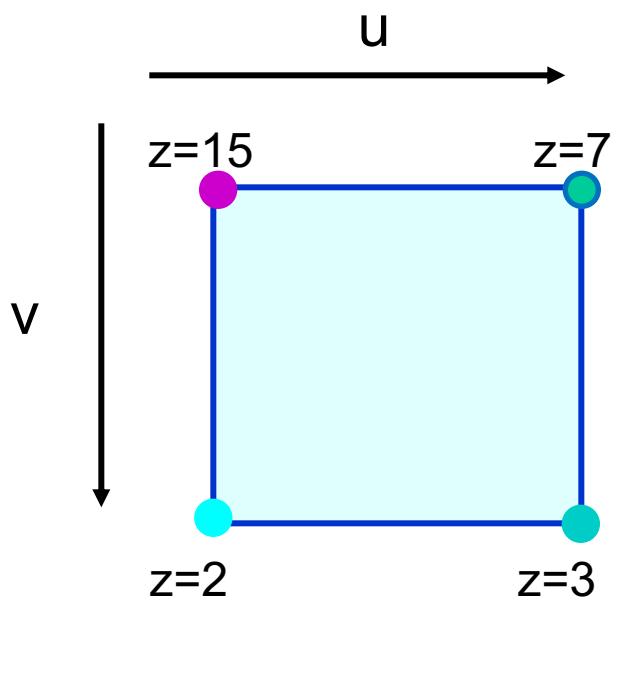
$$u, v \in [0,1]$$

$$S(u) = u SE + (1 - u) SW$$

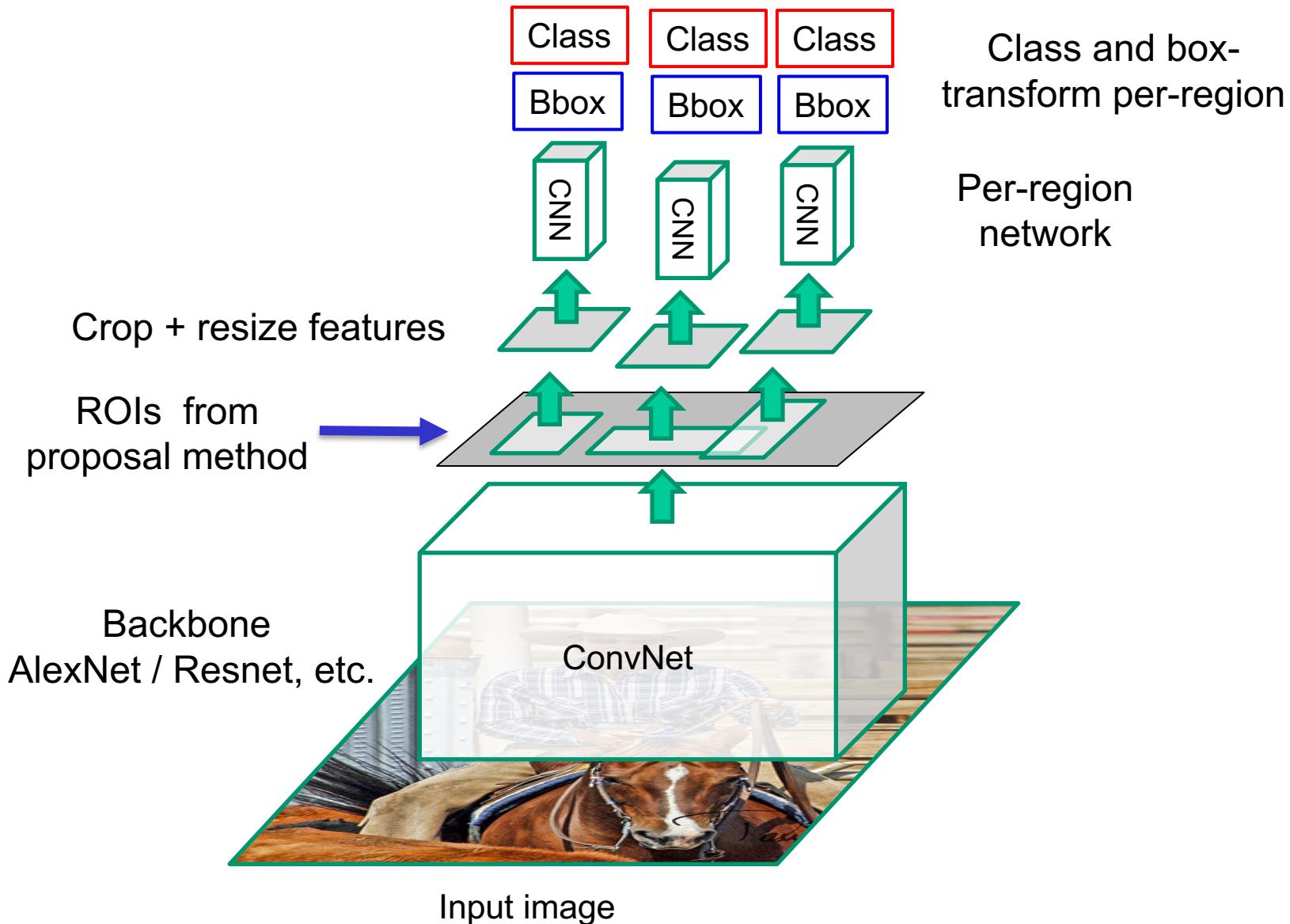
$$N(u) = u NE + (1 - u) NW$$

$$R(u, v) = v N(u) + (1 - v) S(u)$$

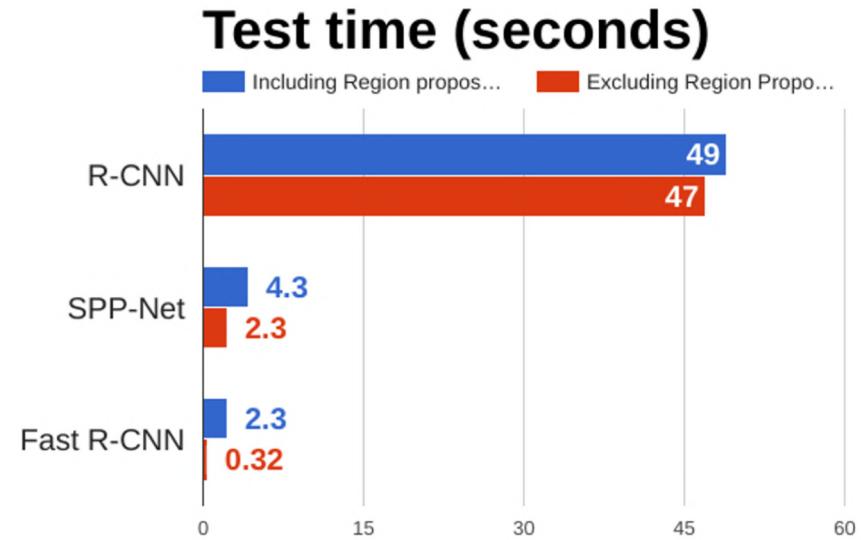
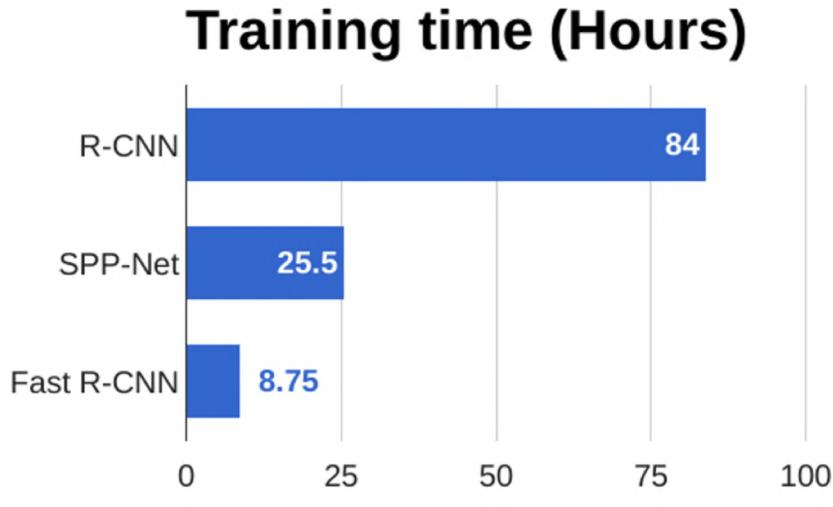
Bilinear interpolation: Example



Fast R-CNN



R-CNN vs. Fast R-CNN



What about the performance?

Evaluating a detector



Test image (previously unseen)

Evaluating a detector

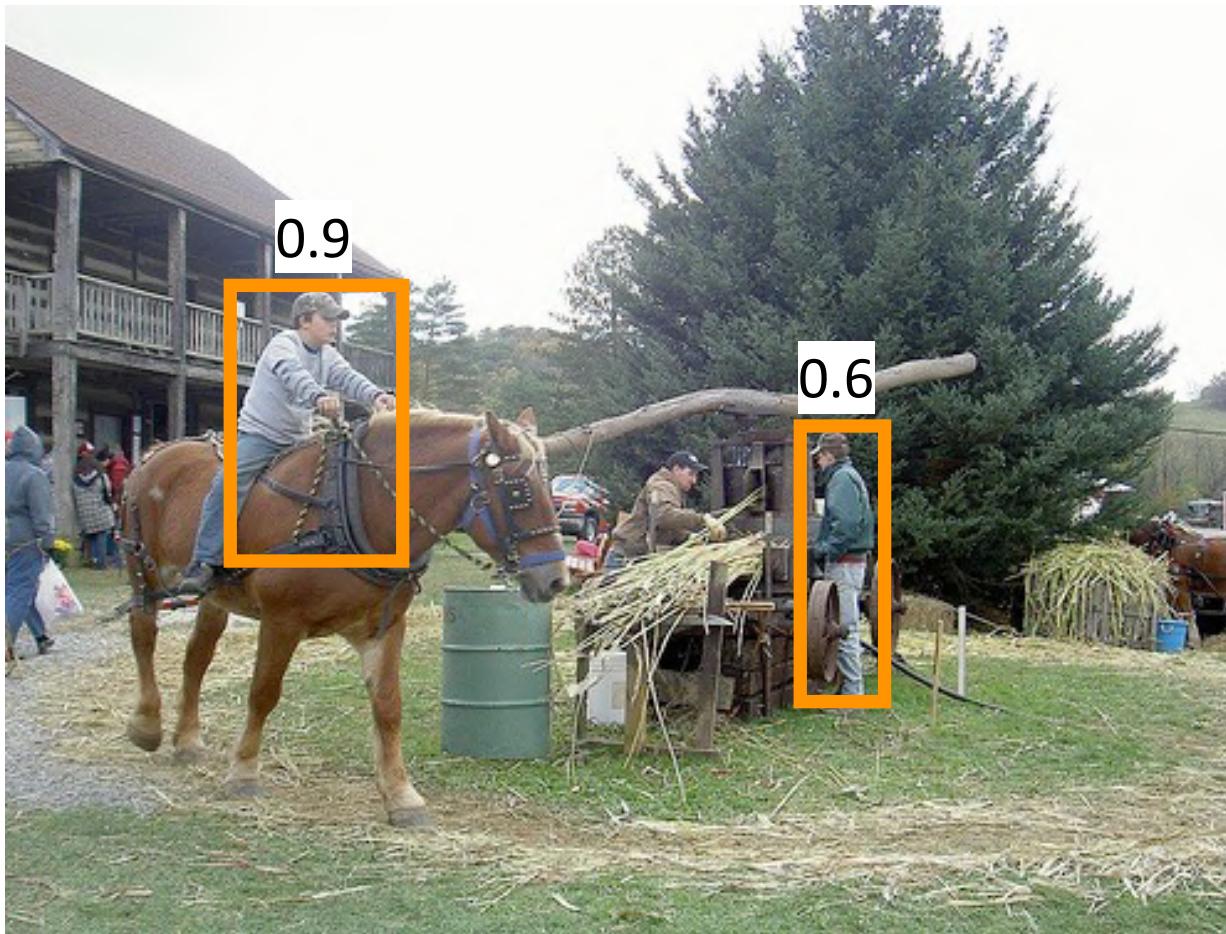
First
Detection



- ‘person’ detector predictions

Evaluating a detector

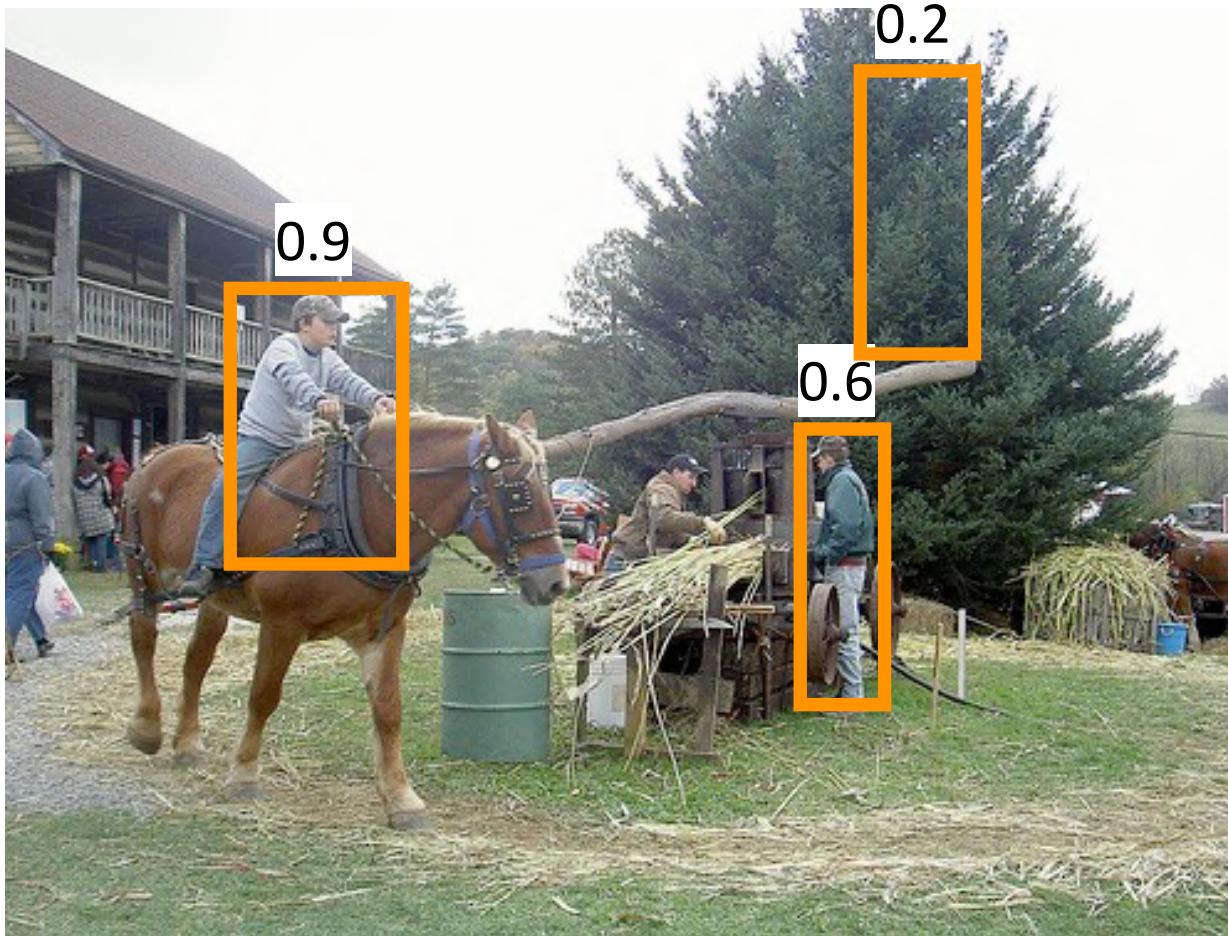
Second
Detection



‘person’ detector predictions

Evaluating a detector

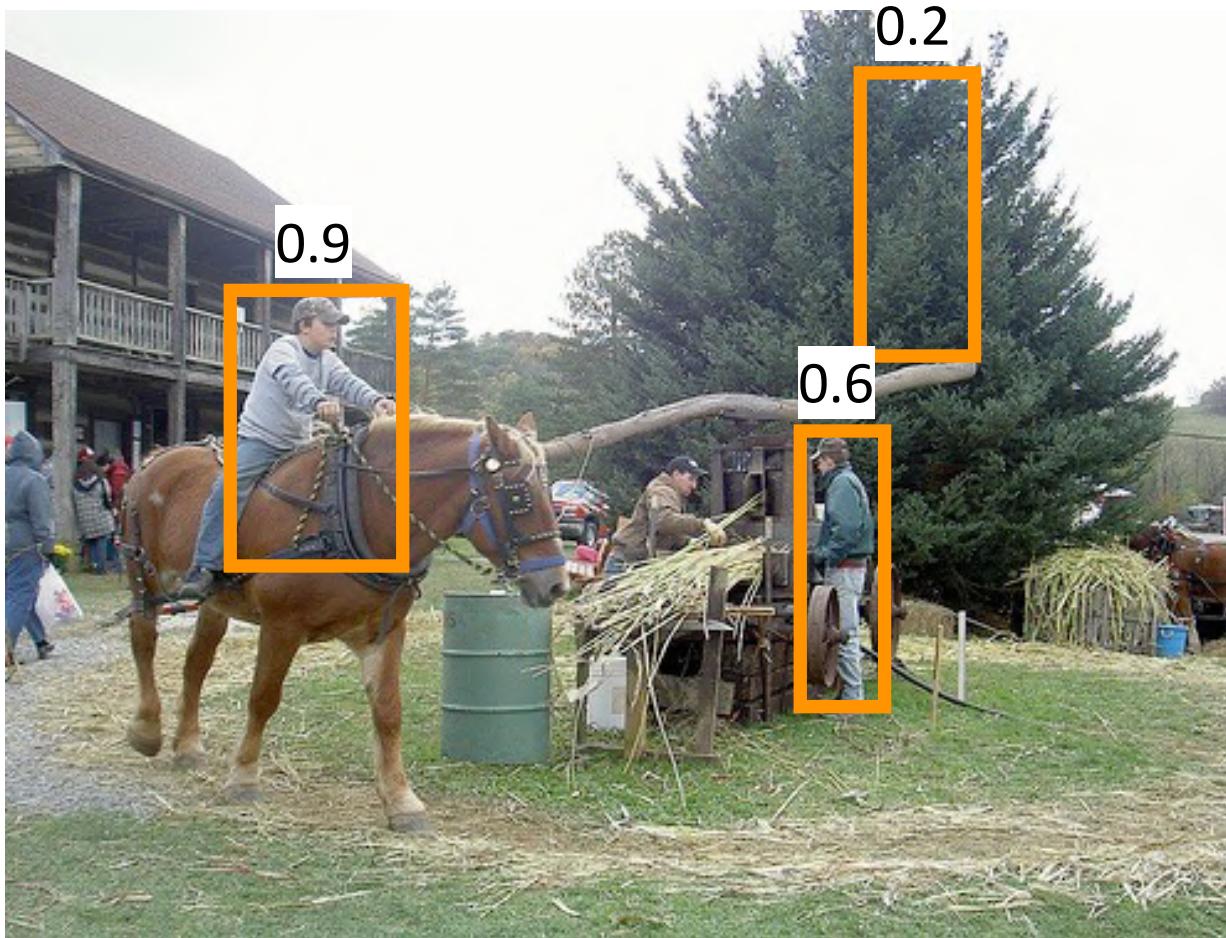
Third
Detection



‘person’ detector predictions

Evaluating a detector

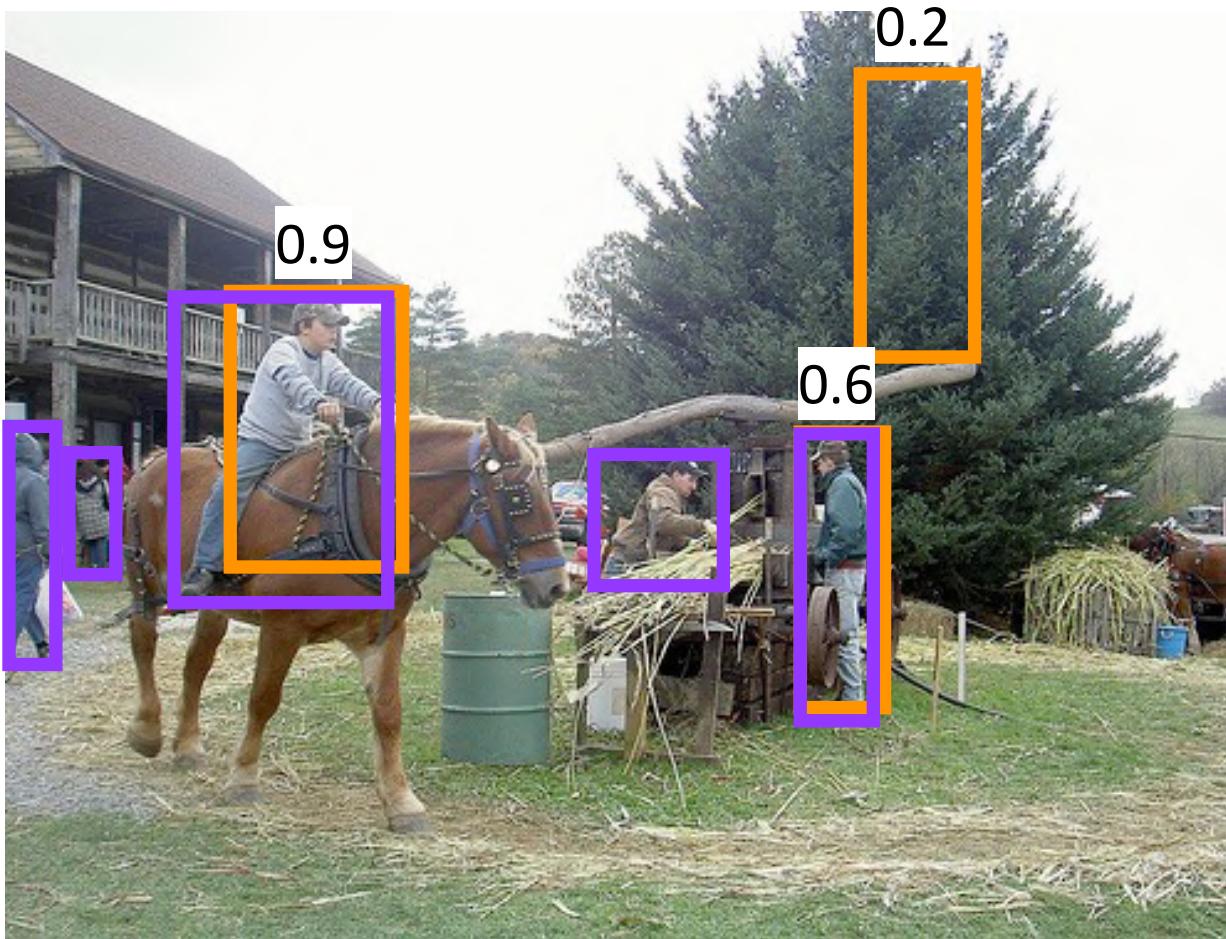
Compared
to GT



- ‘person’ detector predictions
- ground truth ‘person’ boxes

Evaluating a detector

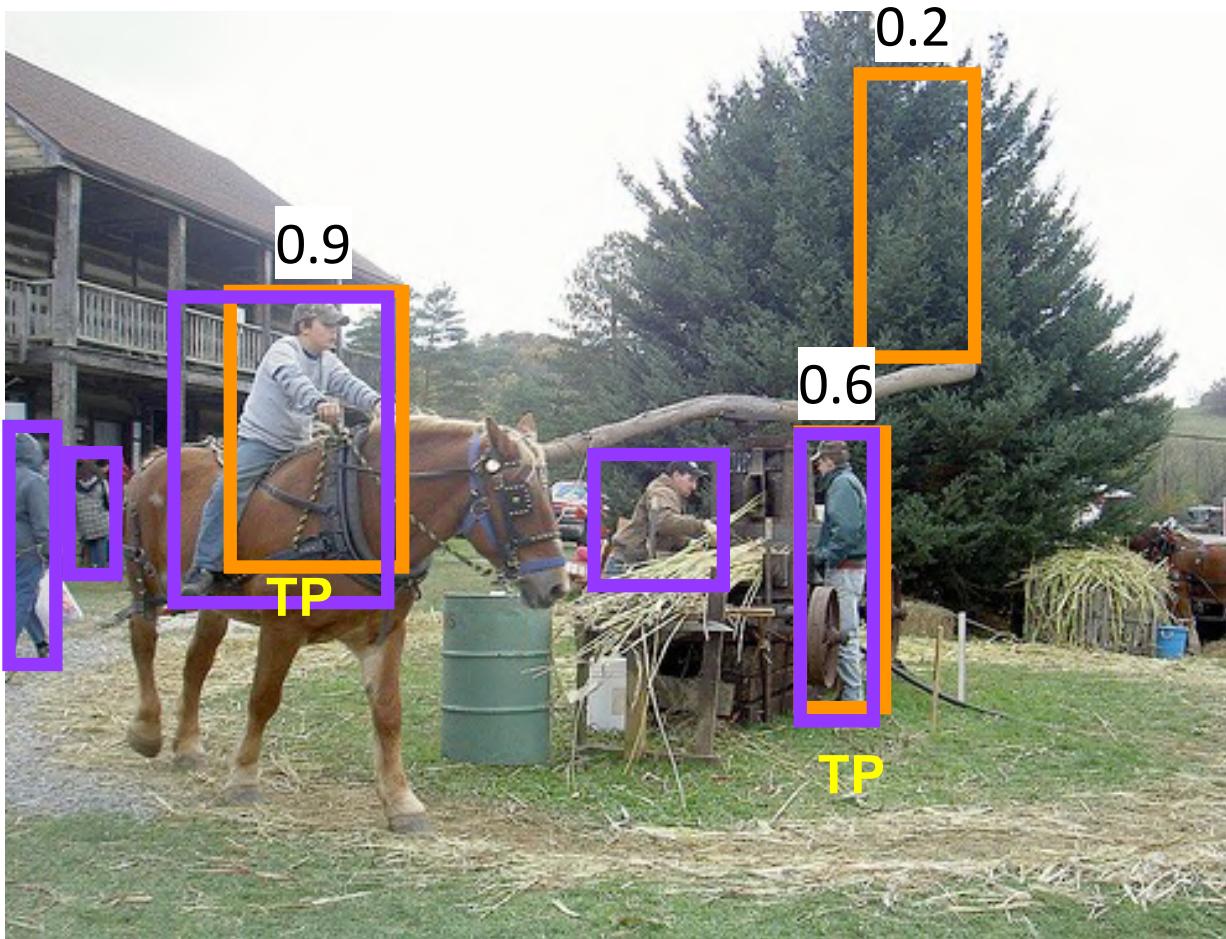
Compared
to GT



- ‘person’ detector predictions
- ground truth ‘person’ boxes

Evaluating a detector

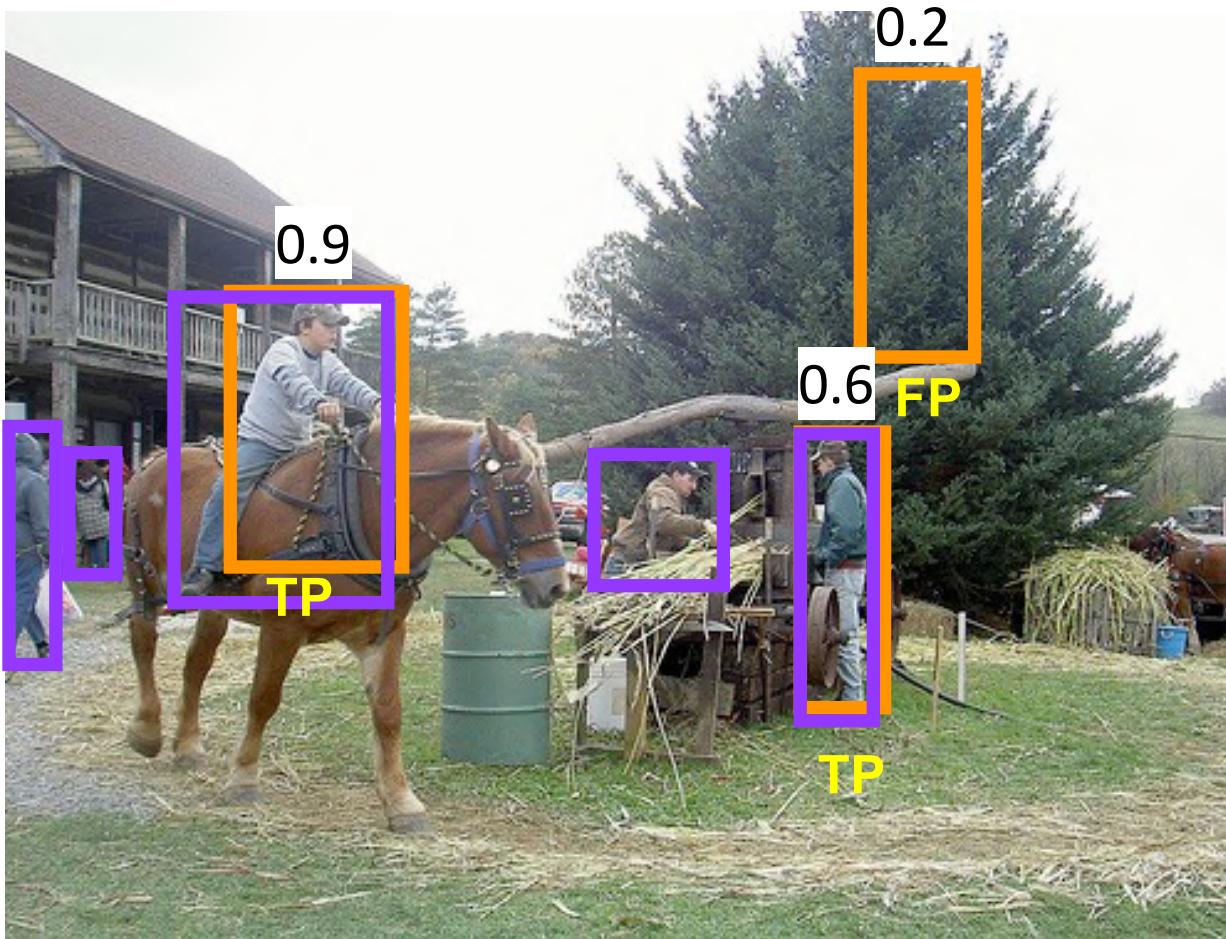
Compared
to GT



- ‘person’ detector predictions
- ground truth ‘person’ boxes

Evaluating a detector

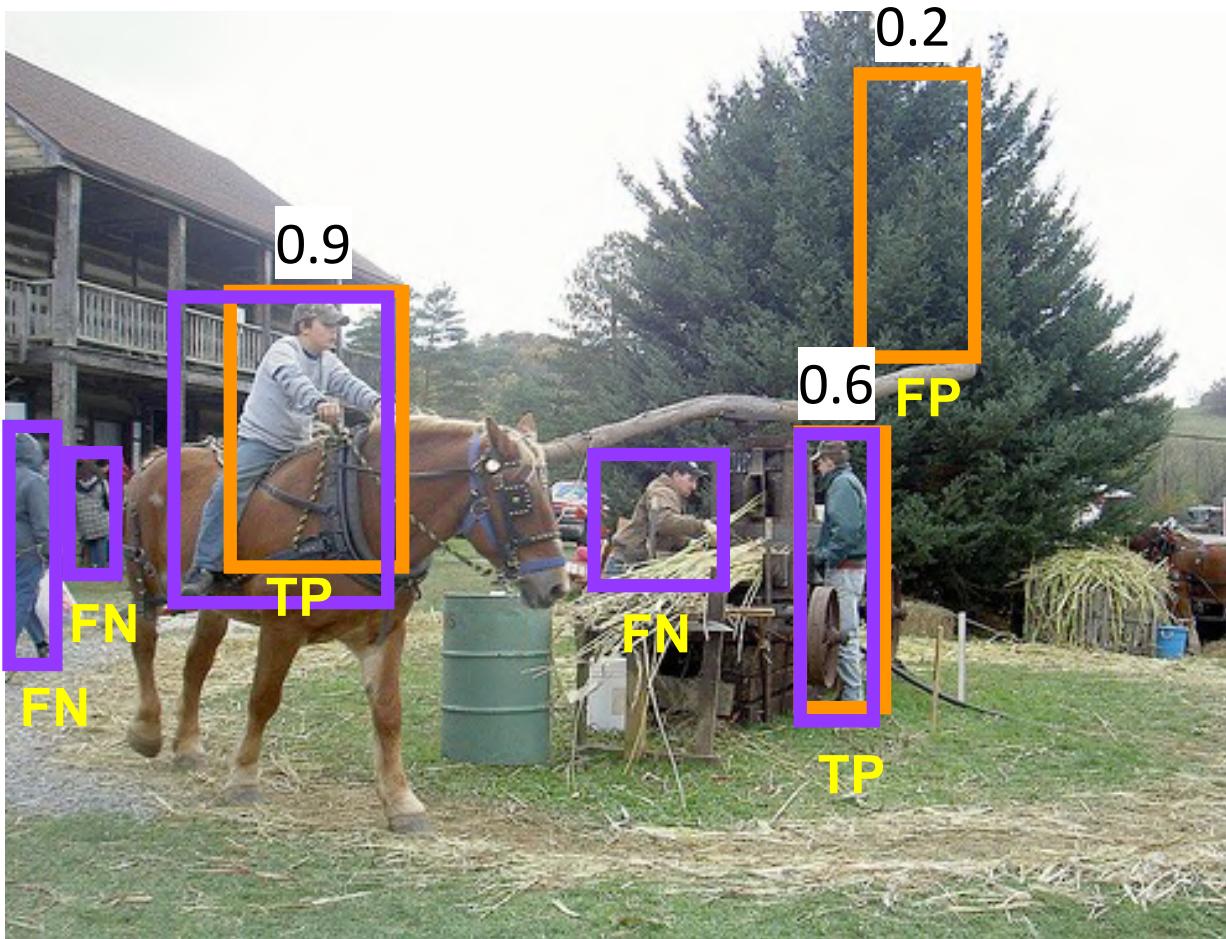
Compared
to GT



- ‘person’ detector predictions
- ground truth ‘person’ boxes

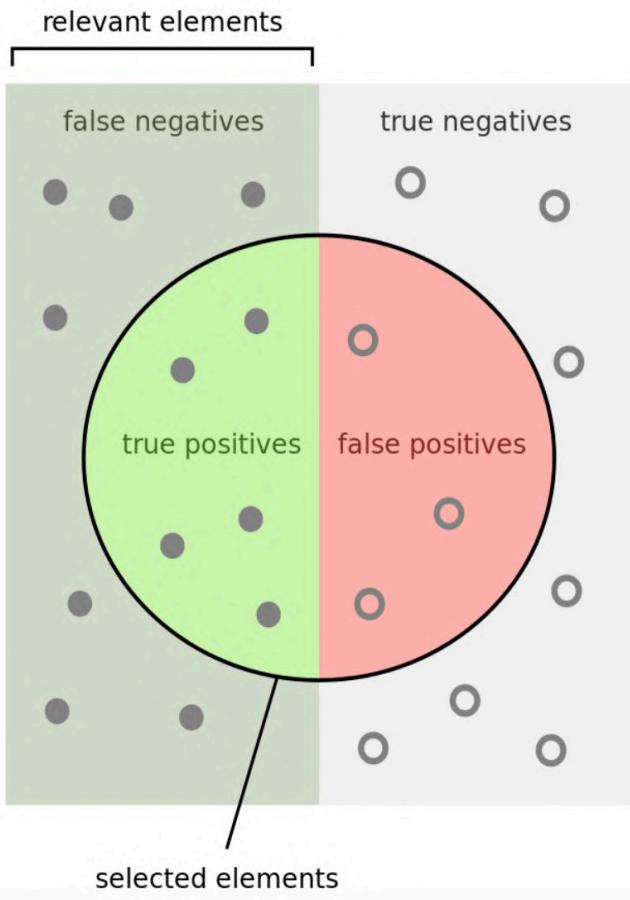
Evaluating a detector

Compared
to GT



- ‘person’ detector predictions
- ground truth ‘person’ boxes

Evaluating a detector



Precision vs. Recall

How many selected items are relevant?

$$\text{Precision} = \frac{\text{true positives}}{\text{true positives} + \text{false positives}}$$



How many relevant items are selected?

$$\text{Recall} = \frac{\text{true positives}}{\text{true positives} + \text{false negatives}}$$



$$precision = \frac{TP}{TP + FP}$$

$$recall = \frac{TP}{TP + FN}$$

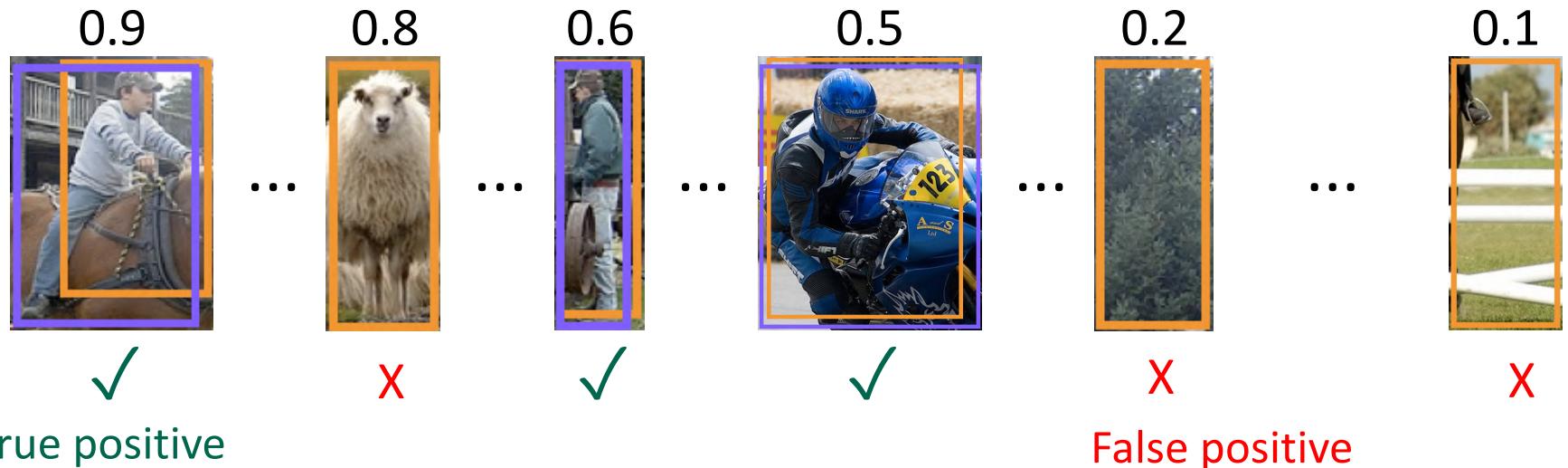
Evaluating a detector

- Sort predictions by confidence
- $\text{TP} = \text{IoU}(\text{prediction}, \text{some GT}) \geq 0.5$
- $\text{FP} = \text{IoU}(\text{prediction}, \text{any GT}) < 0.5$



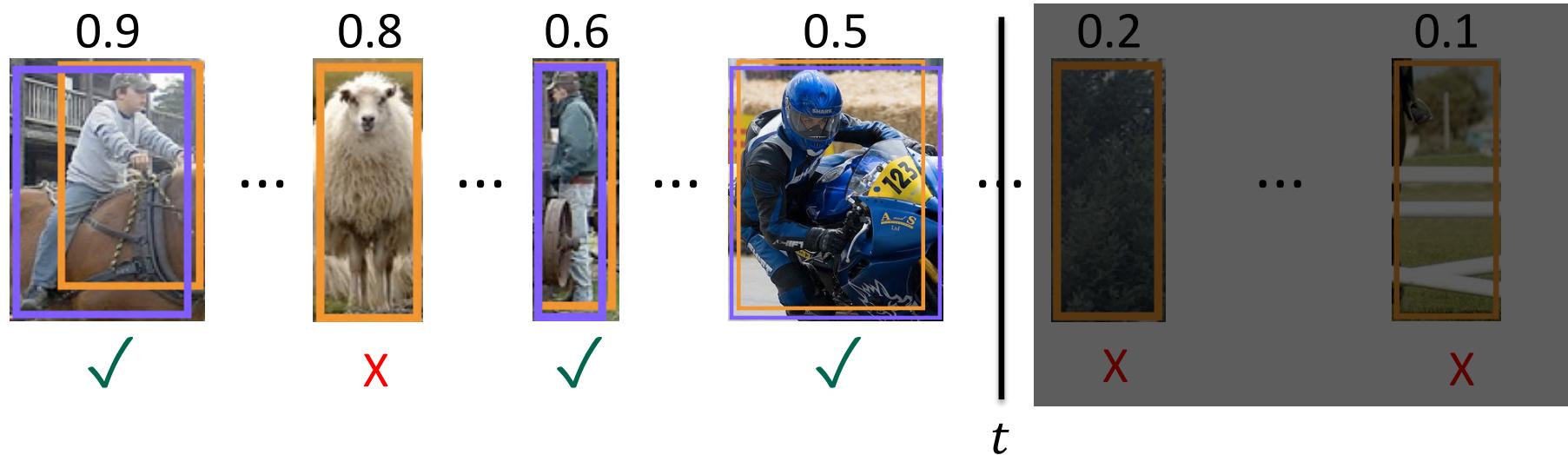
Evaluating a detector

- Sort predictions by confidence



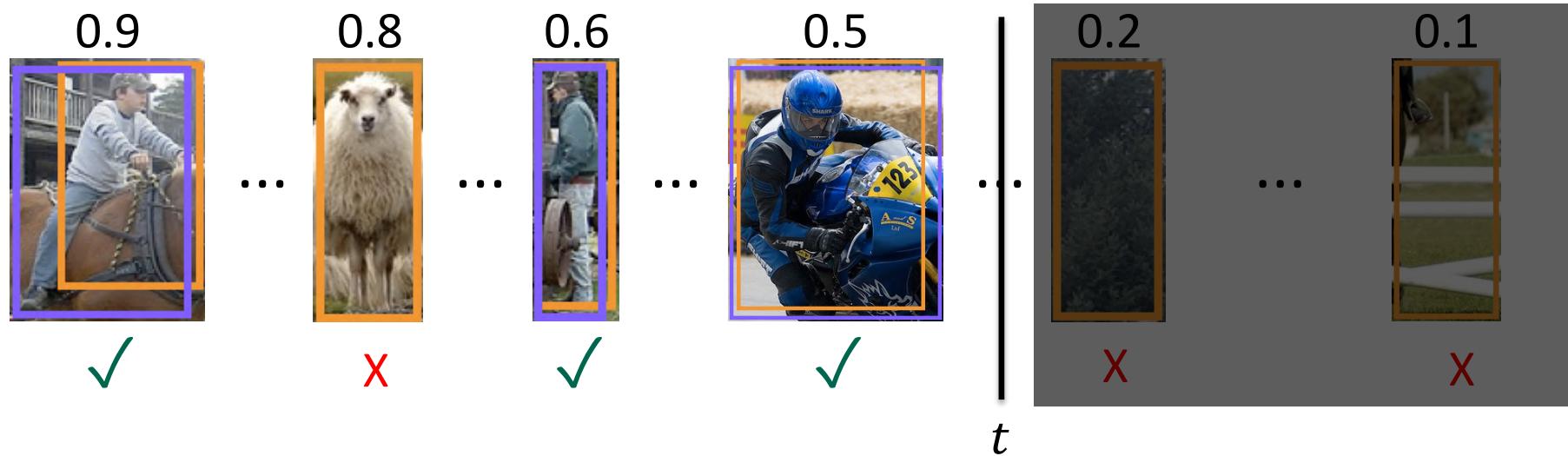
Evaluating a detector

- Sort predictions by confidence
- Choose a confidence threshold



Evaluating a detector

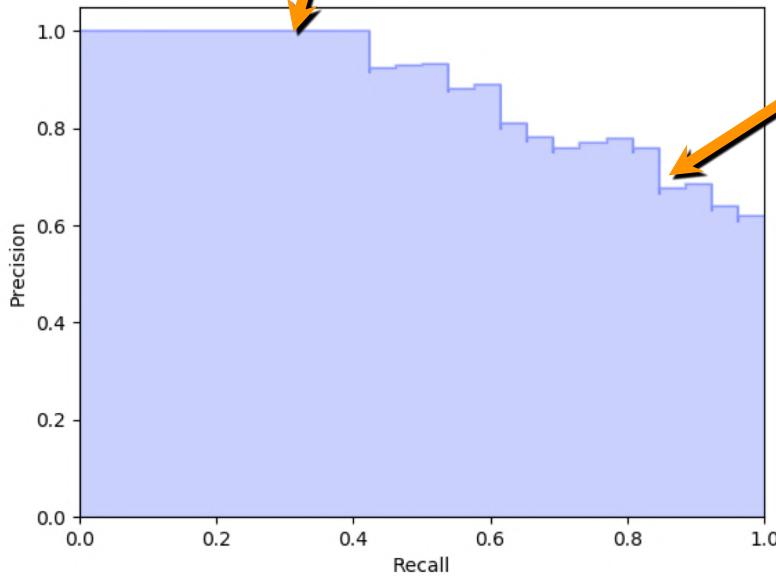
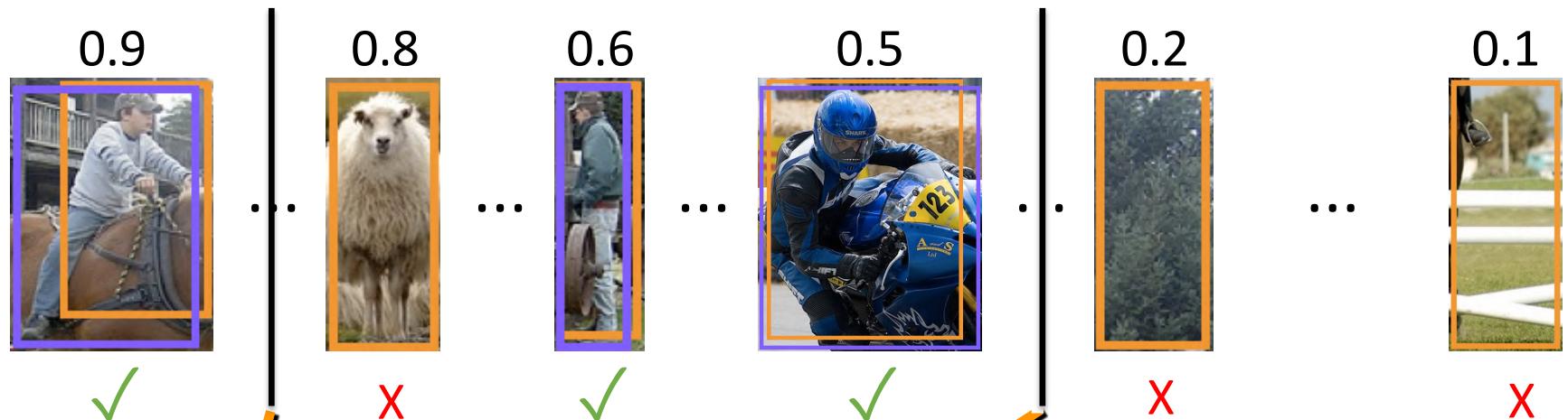
- Sort predictions by confidence
- Choose a confidence threshold
- Evaluate metric



$$\text{precision}(t) = \frac{TP(t)}{TP(t) + FP(t)} = \frac{\text{✓}}{\text{✓} + \text{✗}}$$

$$\text{recall}(t) = \frac{TP(t)}{GT} = \frac{TP(t)}{TP(t) + FN(t)} = \frac{\text{✓}}{\square}$$

Evaluating a detector



Average Precision (AP):
The area under the PR curve

0% is worst
100% is best

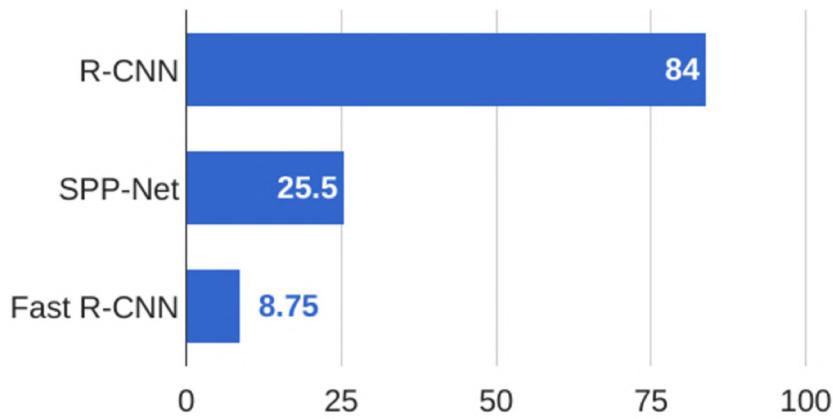
mean AP (mAP):
The mean of AP over all classes

mAP: Fast R-CNN vs. R_CNN

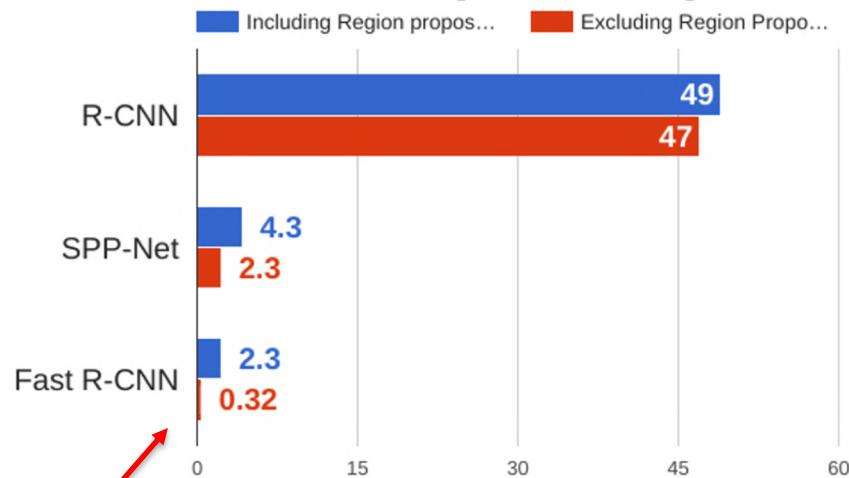
	R-CNN	Fast R-CNN
Test time per image (with proposals)	50 seconds	2 seconds
(Speedup)	1x	25x
mAP (VOC 2007)	66.0	66.9

R-CNN vs. Fast R-CNN

Training time (Hours)



Test time (seconds)



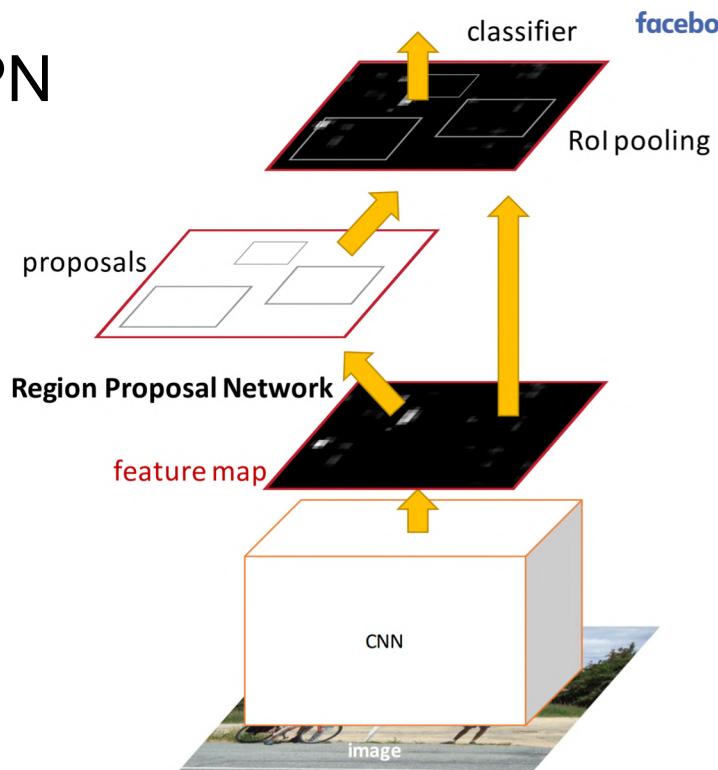
Problem:

Runtime is dominated by
the Region Proposals

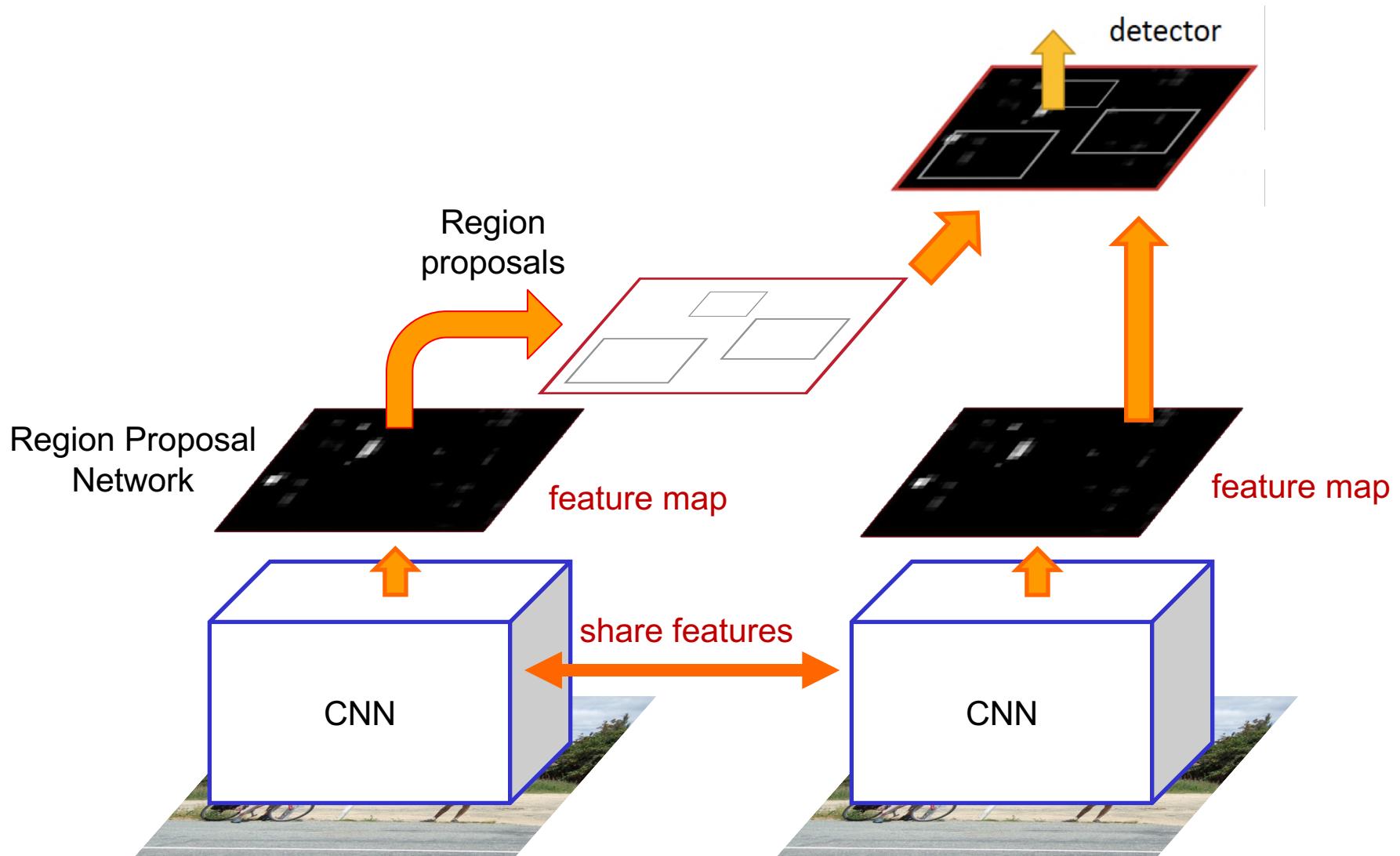
Faster R-CNN

Main idea:

- Use a **region proposal network (RPN)** to predict RP from CNN features.
- Use the same features for the RPN and for class/box regressions.



Faster R-CNN

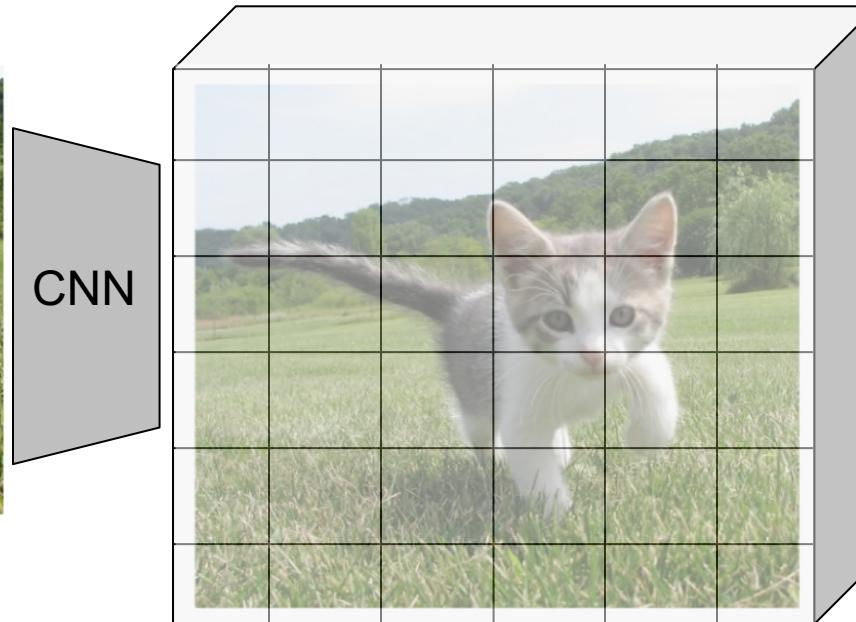


Faster R-CNN: RPN

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



Input image
(e.g. 3x640x480)



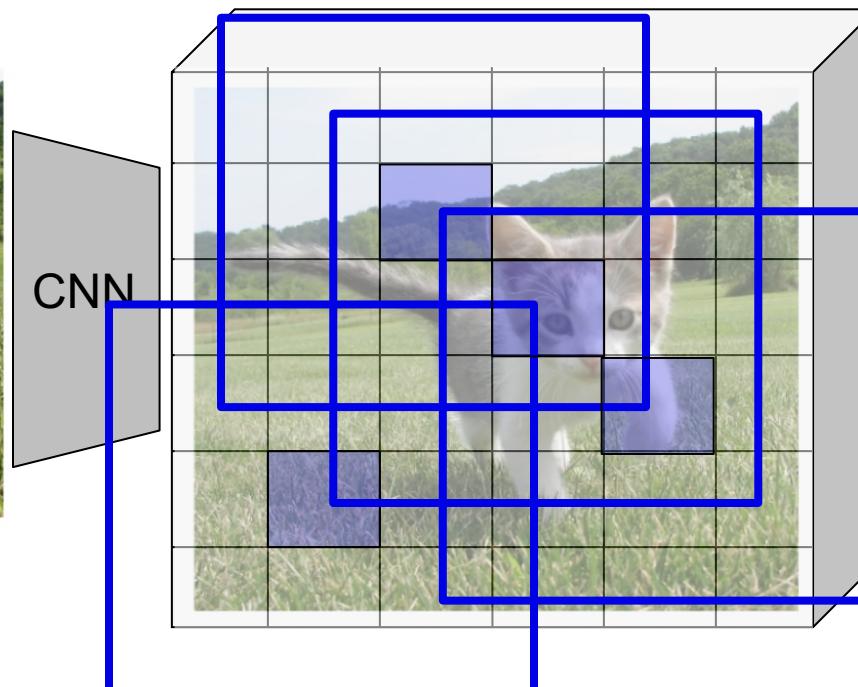
Input features
(e.g. 512x20x15)

Faster R-CNN: RPN

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



Input image
(e.g. 3x640x480)



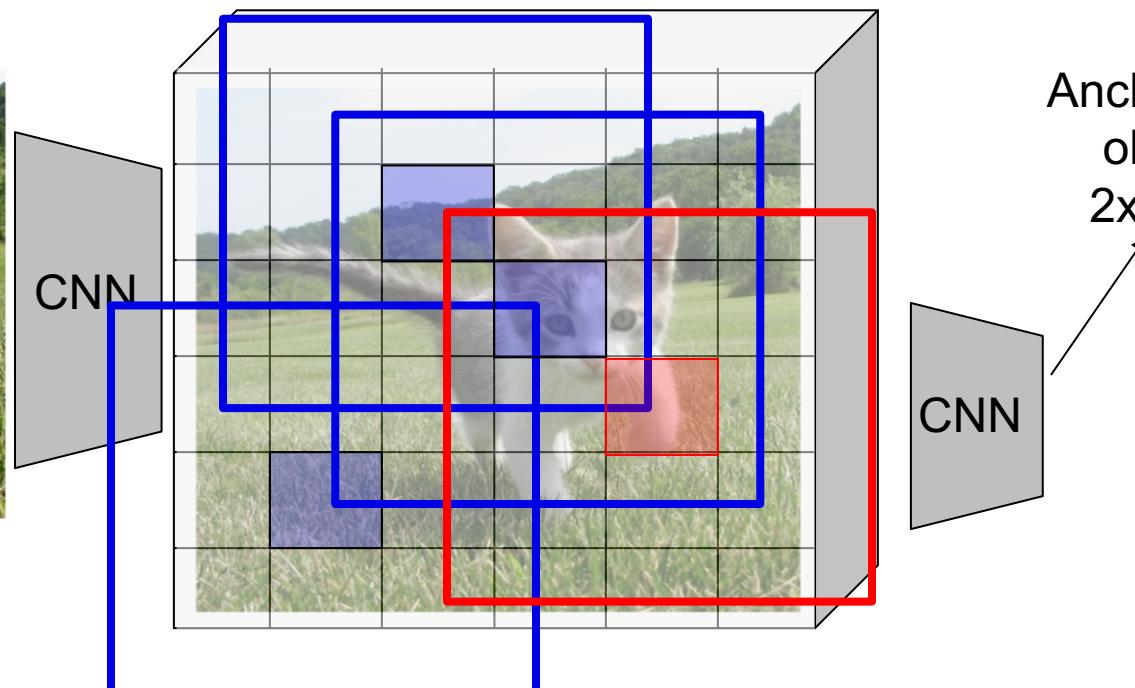
Input features
(e.g. 512x20x15)

Faster R-CNN: RPN

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



Input image
(e.g. 3x640x480)



Input features
(e.g. 512x20x15)

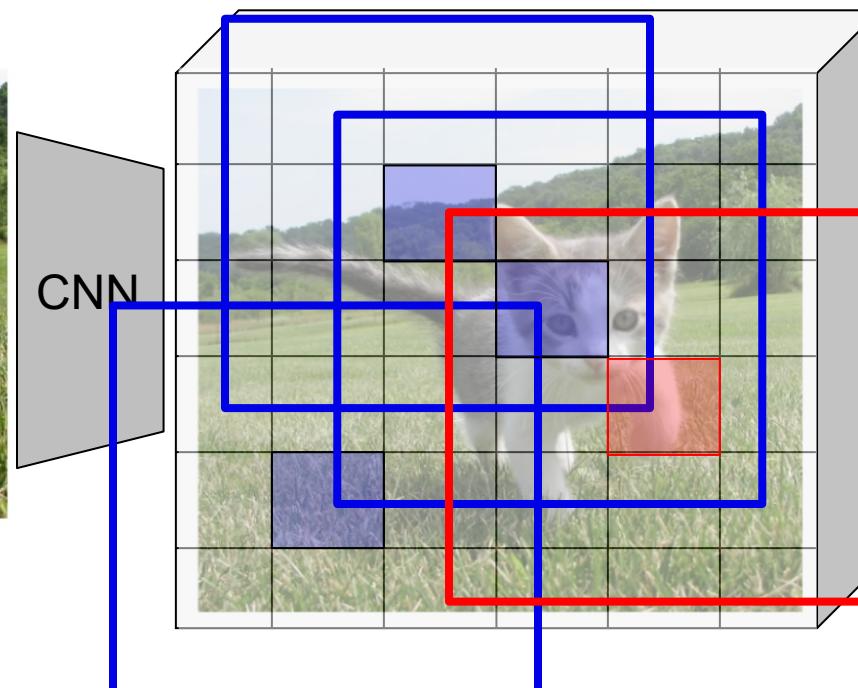
Anchor is an
object?
 $2 \times 20 \times 15$

Faster R-CNN: RPN

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



Input image
(e.g. 3x640x480)



(e.g. 512x20x15)

Anchor is an object?
 $2 \times 20 \times 15$

Box transform
 $4 \times 20 \times 15$

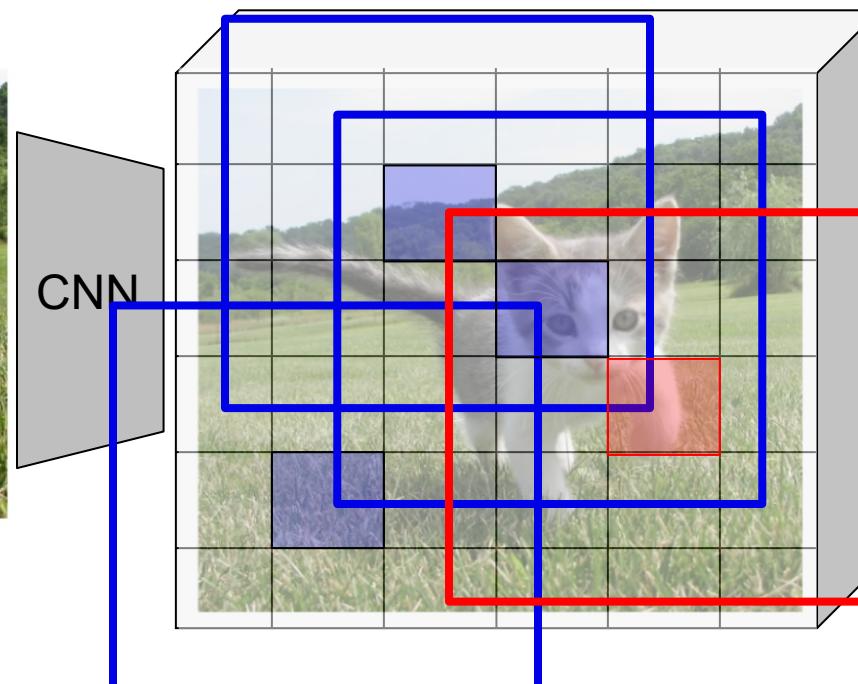
Faster R-CNN: RPN

96

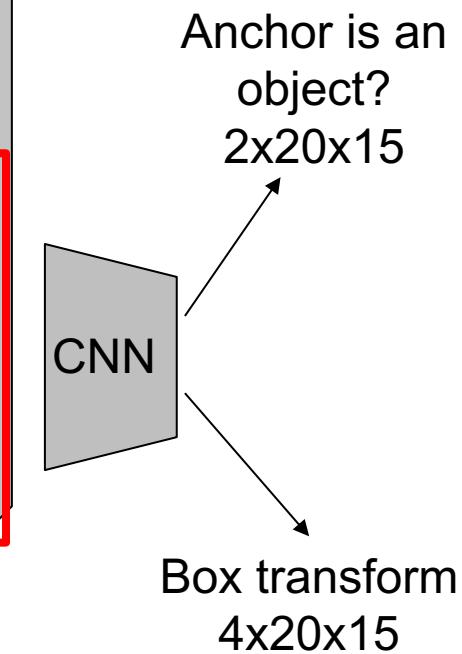
- **Problem:** Anchor box may have the wrong size / shape
- **Solution:** Use k different anchor boxes at each point.



Input image
(e.g. 3x640x480)



Input features
(e.g. 512x20x15)



Anchor is an
object?
2x20x15

Box transform
4x20x15

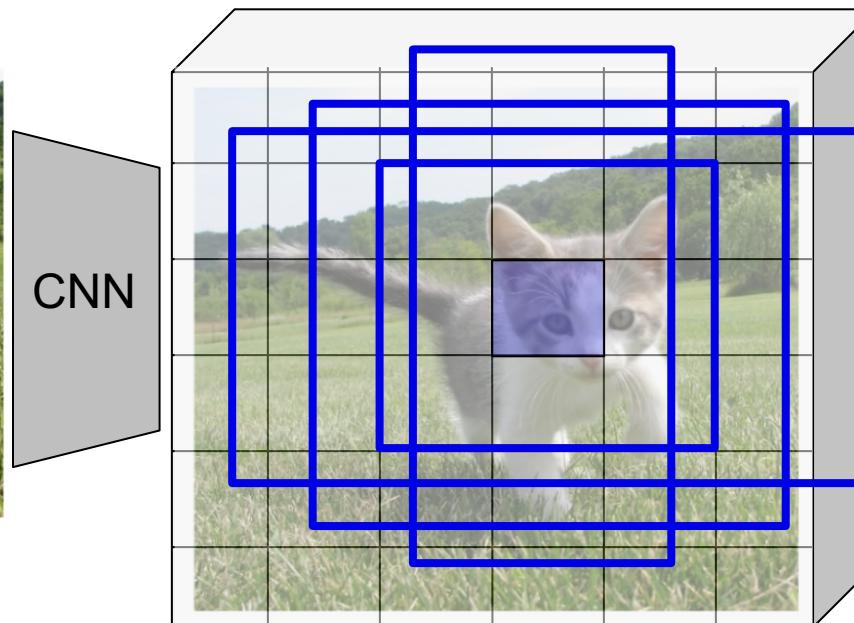
Faster R-CNN: RPN

97

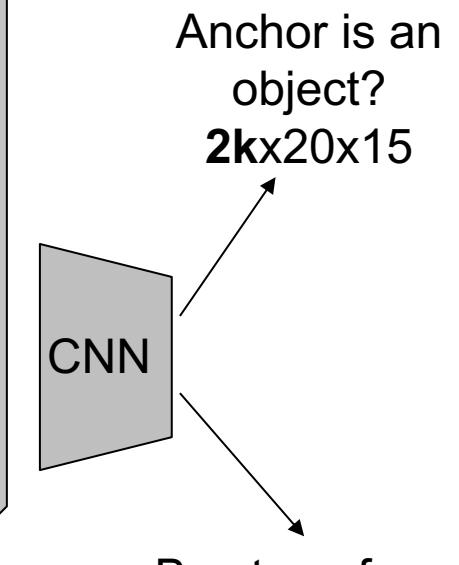
- **Problem:** Anchor box may have the wrong size / shape
- **Solution:** Use k different anchor boxes at each point.



Input image
(e.g. 3x640x480)



Input features
(e.g. 512x20x15)



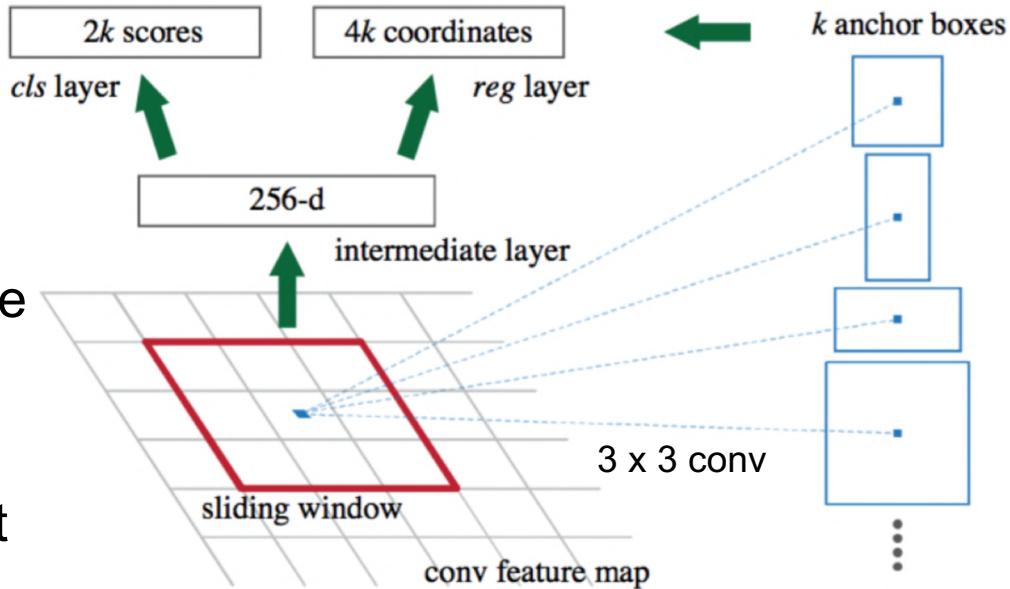
Anchor is an
object?
 $2k \times 20 \times 15$

Box transform
 $4k \times 20 \times 15$

Faster R-CNN: RPN

98

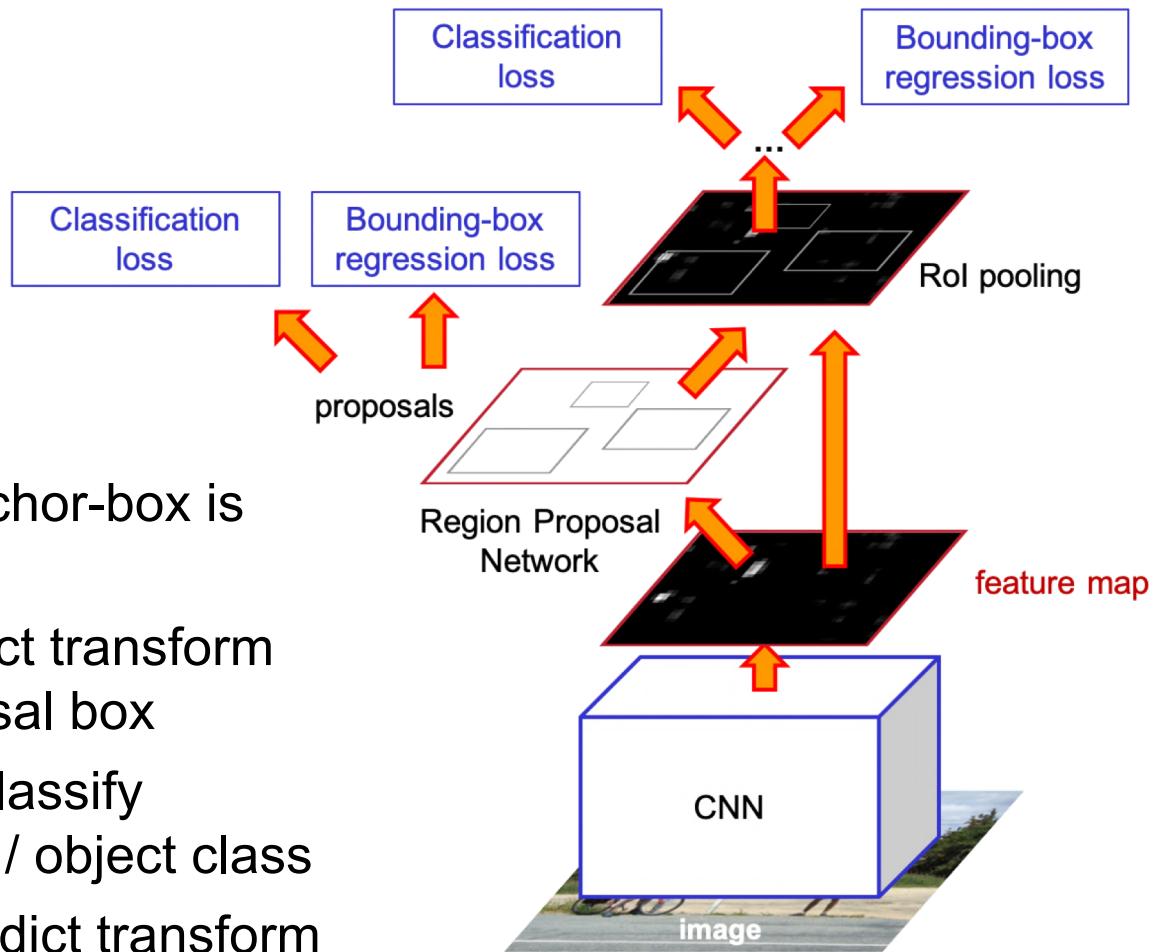
- At each location, hypothesize the k anchor boxes
- Build a small network for:
 - classifying object/non-object
 - regressing bbox locations
- Position of the sliding window provides localization information with reference to the image
- Box regression provides finer localization information with reference to this sliding window



Faster R-CNN: Training

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 proposal as background / object class
4. **Object Regression:** predict transform from proposal box to object box



Faster R-CNN: Testing

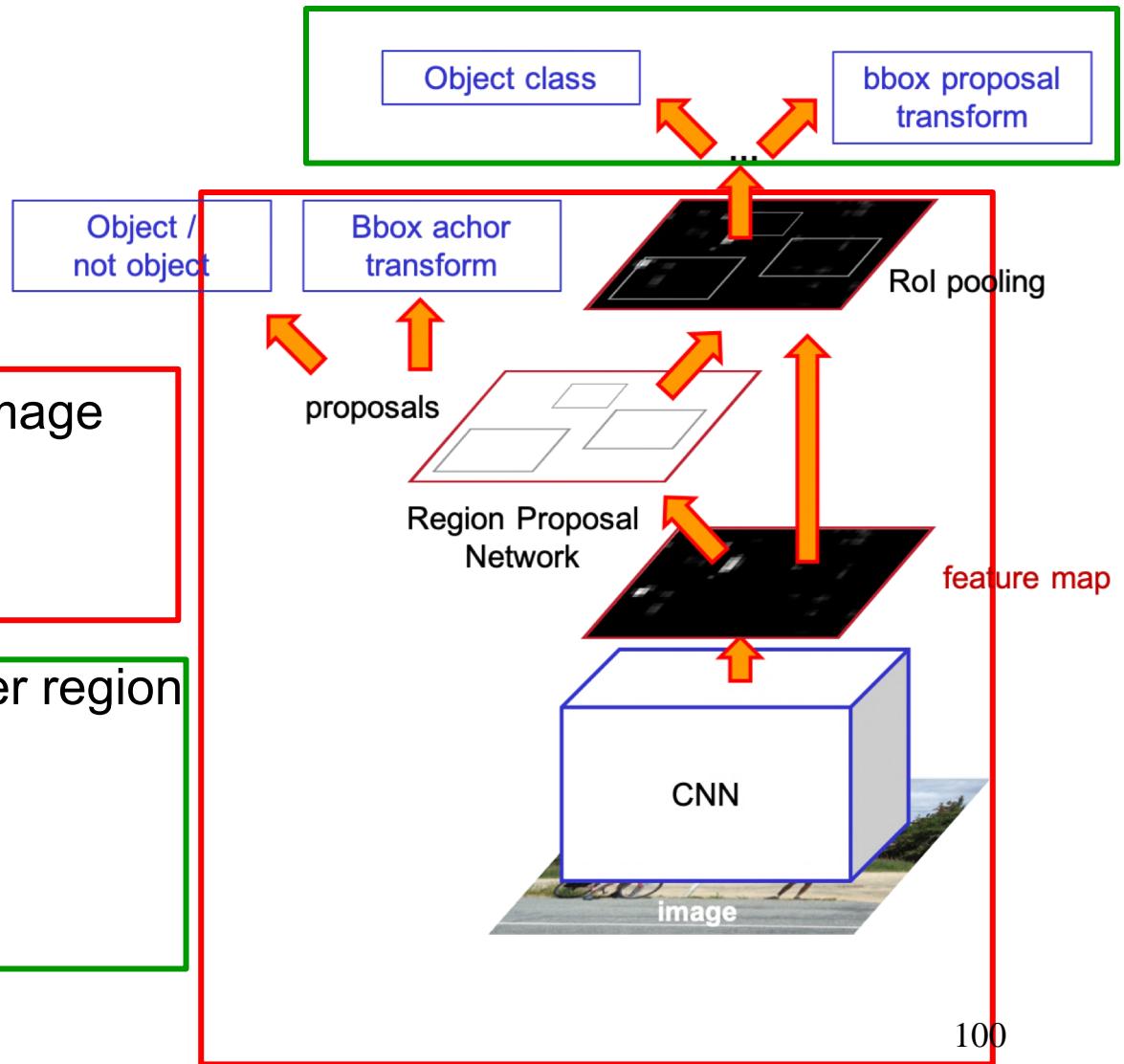
Two stage object detector:

First stage: Run once per image

- Backbone network
- Region proposal network

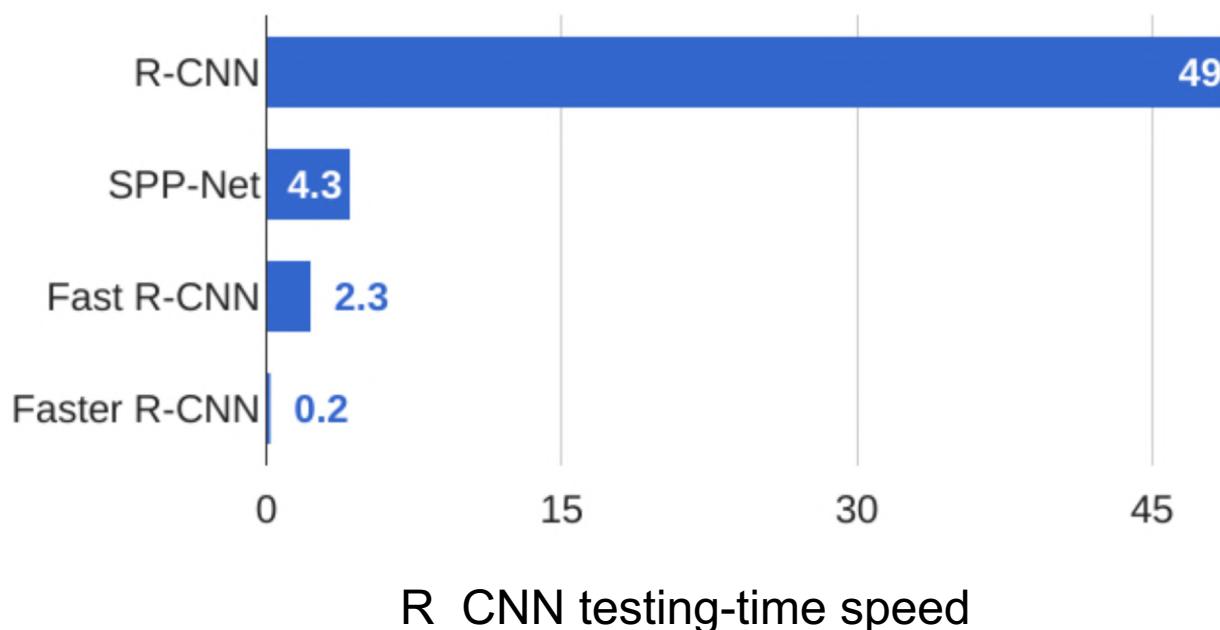
Second stage: Run once per region

- Crop feature / ROI align
- Predict object class
- Predict bbox offset



Faster R-CNN: Results

	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



Single Stage Object Detection

Do we really need two stages?

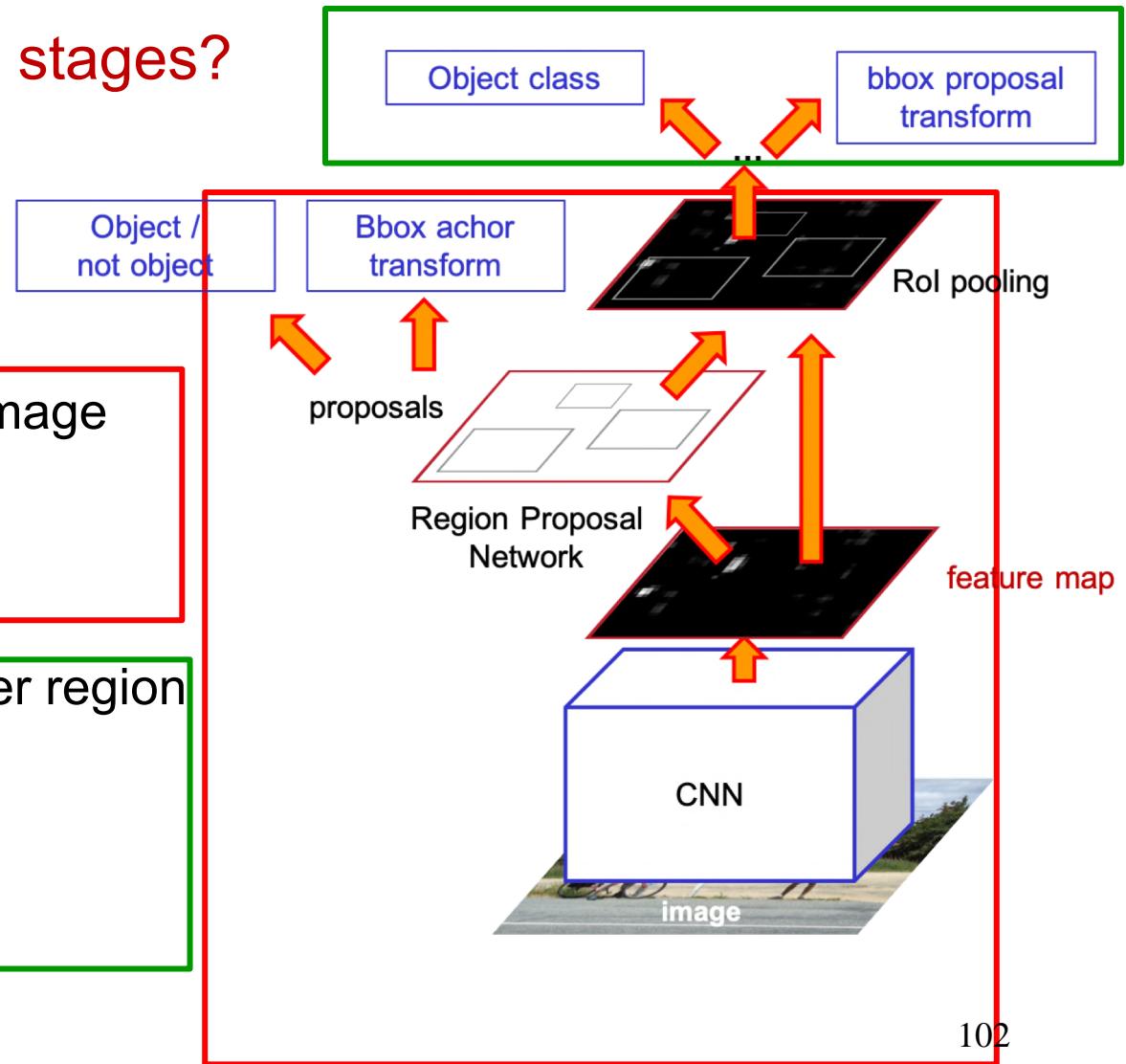
Two stage object detector:

First stage: Run once per image

- Backbone network
- Region proposal network

Second stage: Run once per region

- Crop feature / ROI align
- Predict object class
- Predict bbox offset



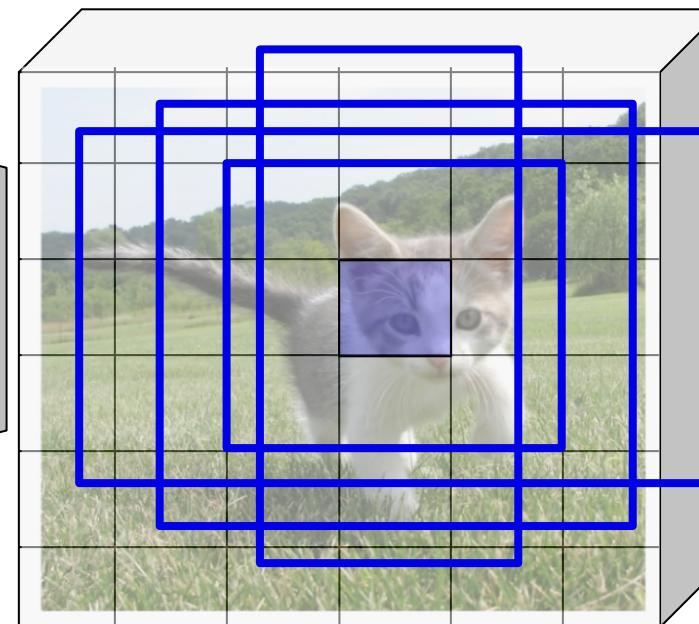
Single Stage Object Detection

103

- **RPN:** Classify each anchor box as object / non-object
- **Single stage detector:** Classify each anchor box as one of C categories (or background)



Input image
(e.g. 3x640x480)



Input features
(e.g. 512x20x15)

RPN

Anchor is an object?
 $2k \times 20 \times 15$



Box transform
 $4k \times 20 \times 15$

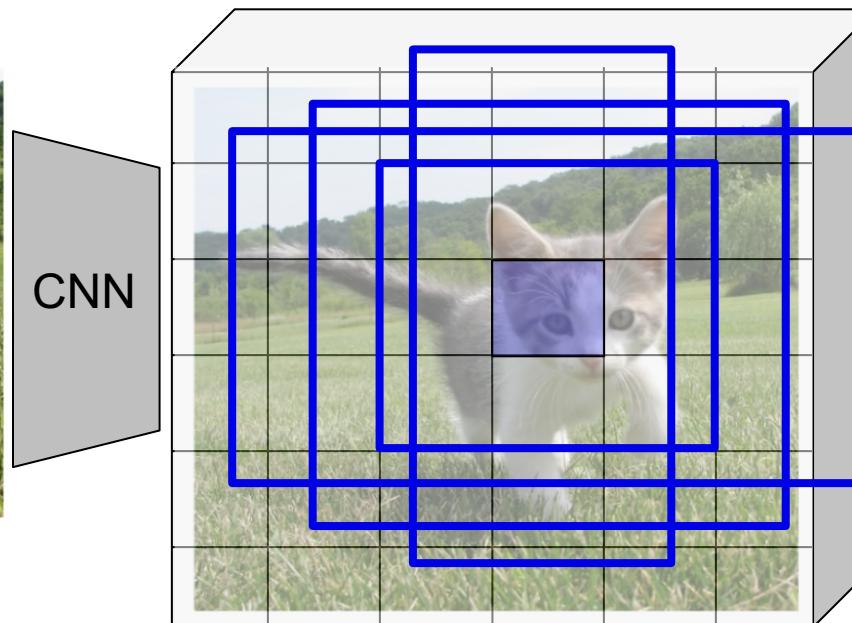
Single Stage Object Detection

104

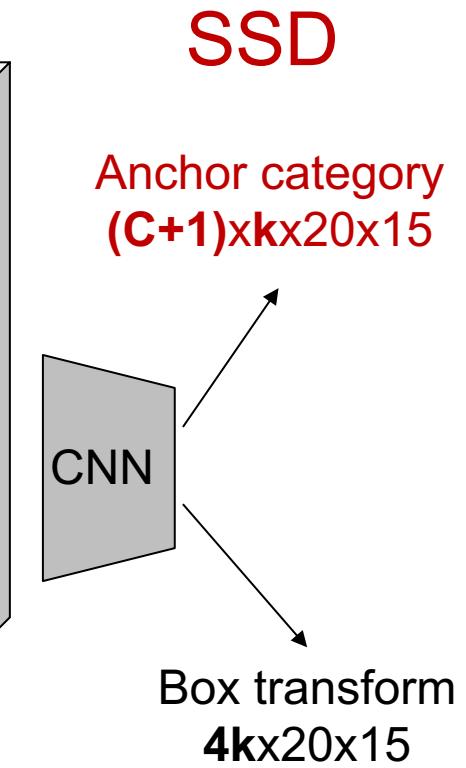
- **RPN:** Classify each anchor box as object / non-object
- **Single stage detector:** Classify each anchor box as one of C categories (or background)



Input image
(e.g. 3x640x480)



Input features
(e.g. 512x20x15)



Detectors without proposal

1. YOLO (You Only Look once)

- Object detection system targeted for real-time processing
- Divides the input image into an $S \times S$ grid and look at each grid cell a single object.

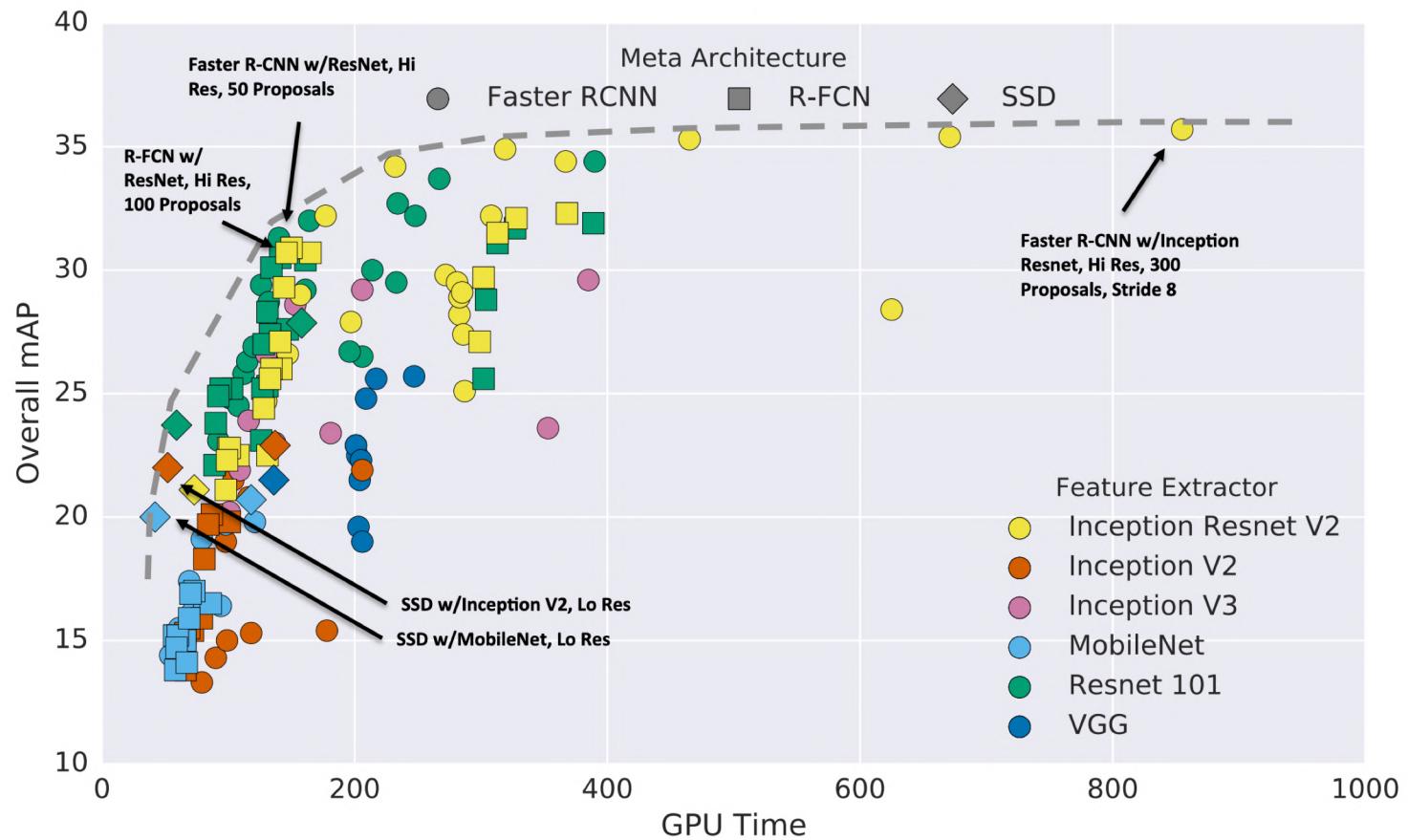
2. SSD (Single Shot Detector)

- Similar to YOLO
- Hypothesize k anchor boxes at each cell

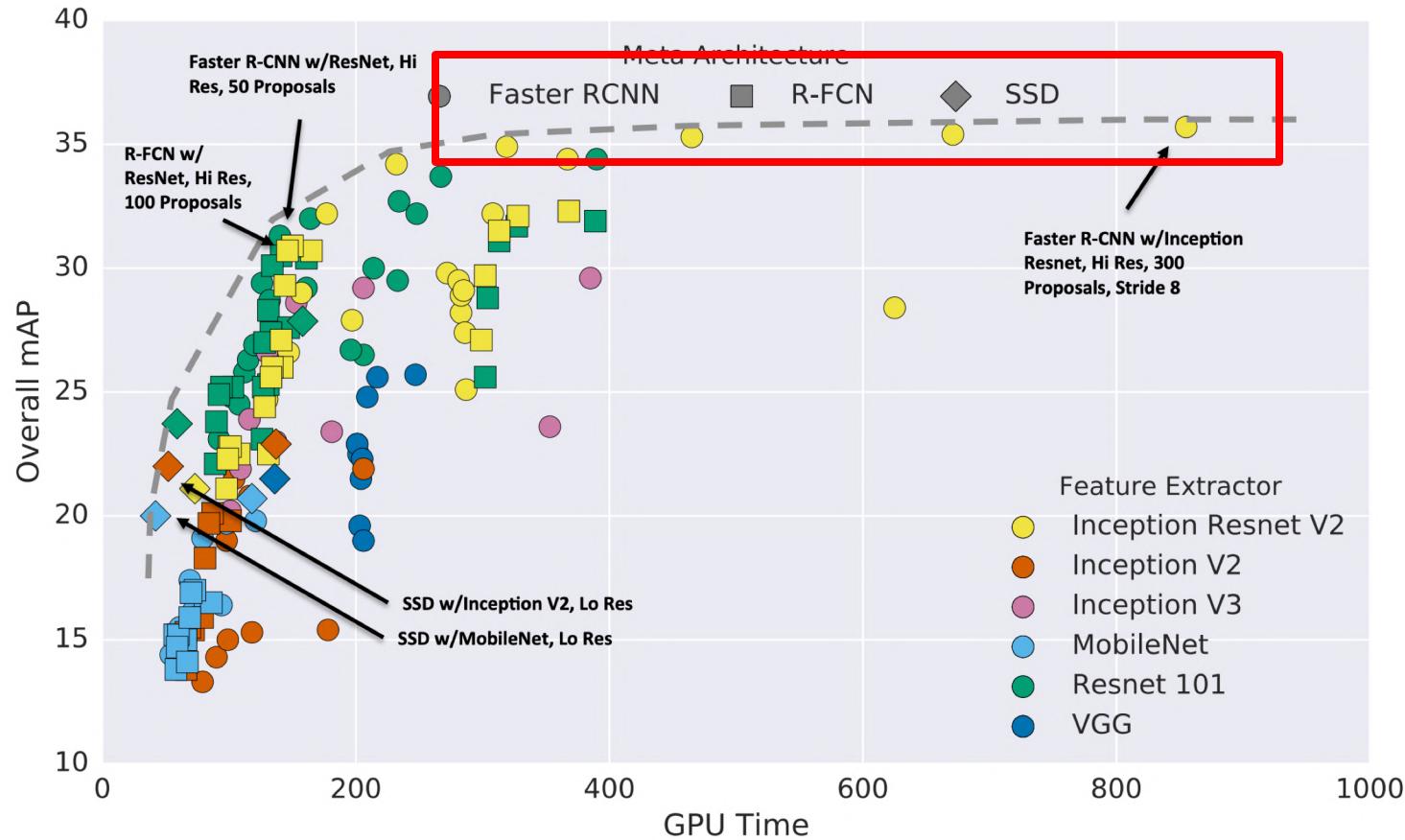
W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. Berg, [SSD: Single Shot MultiBox Detector](#), ECCV 2016.

J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, [You Only Look Once: Unified, Real-Time Object Detection](#), CVPR 2016

Object Detectors

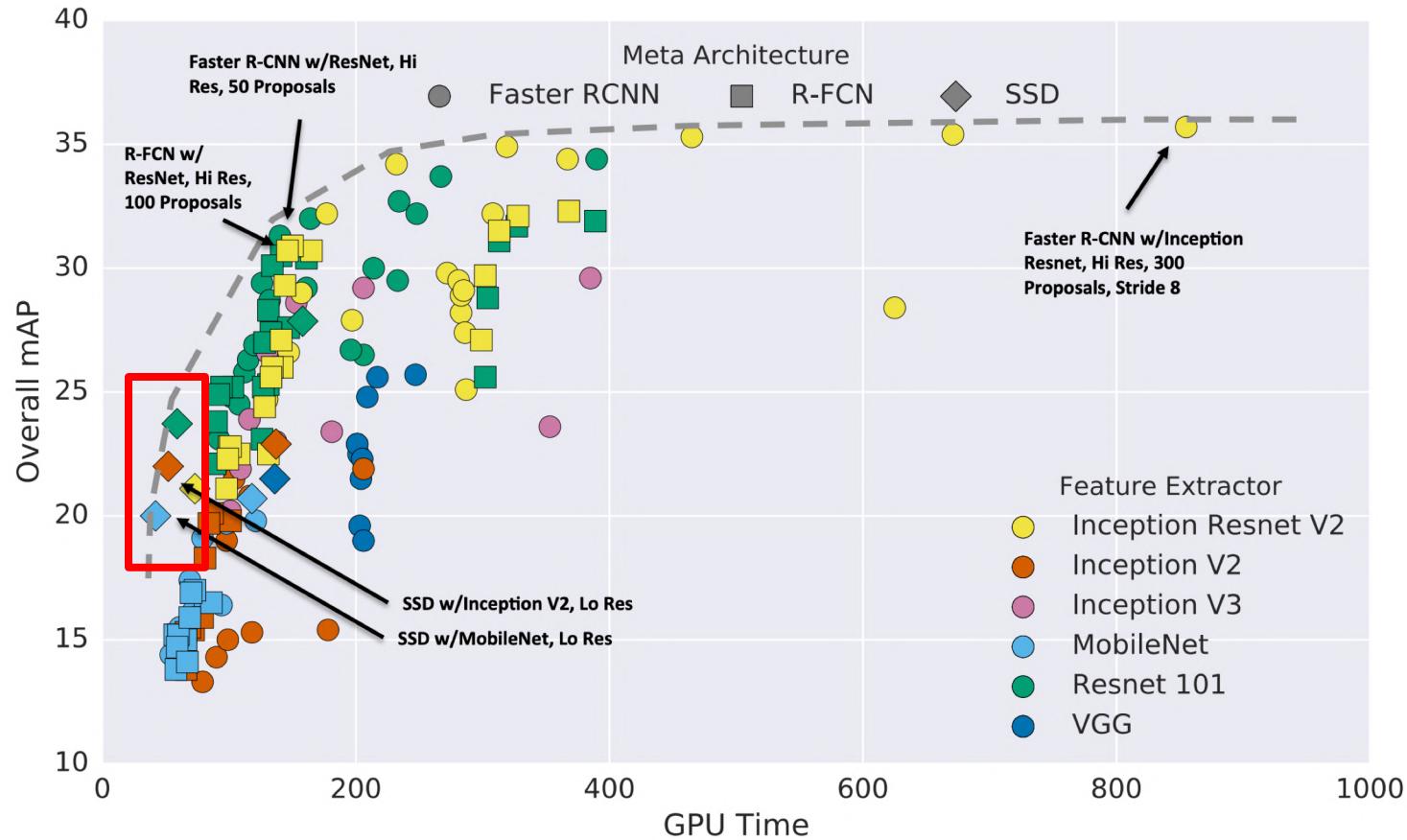


Object Detectors



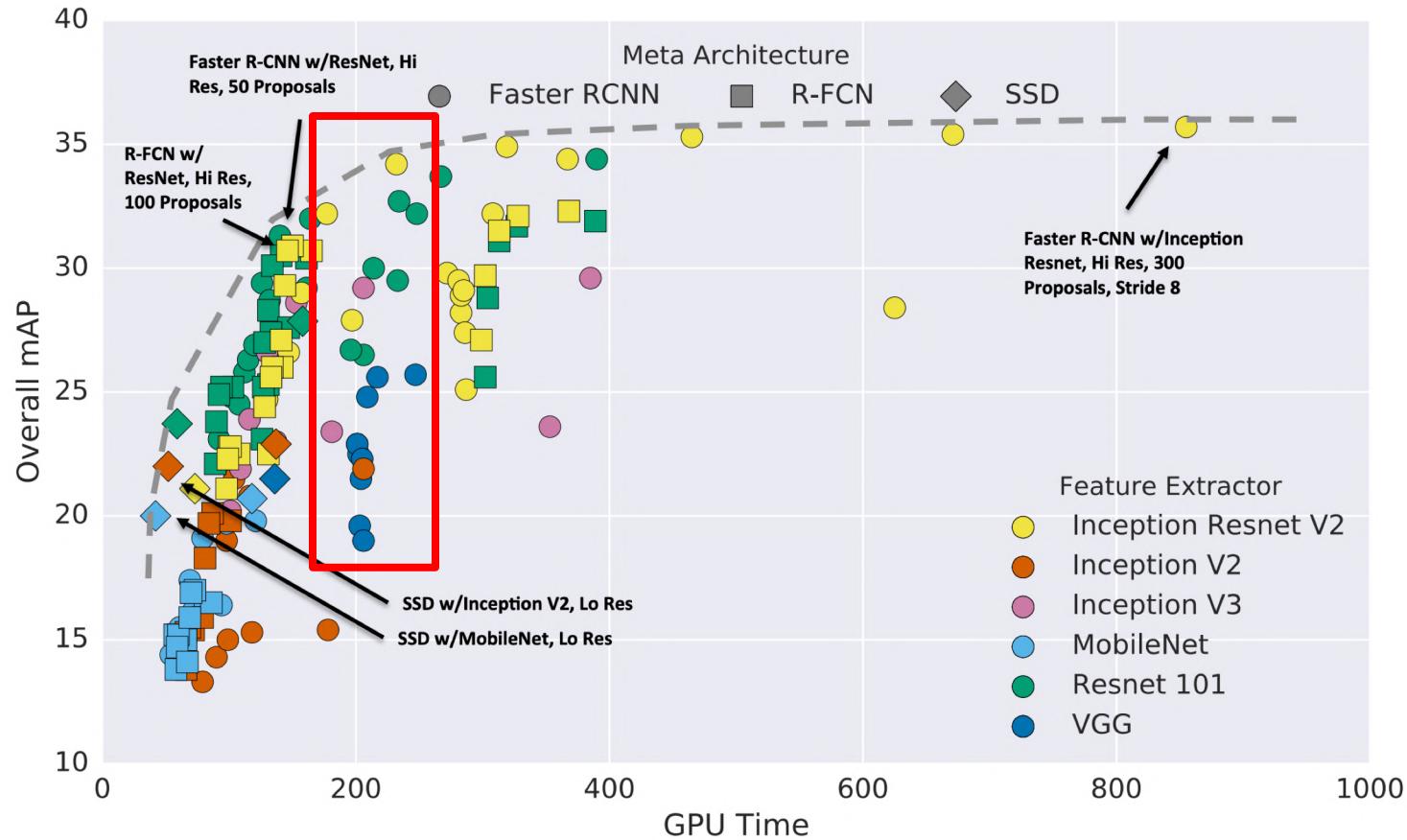
- Two stages detectors get the best accuracy, but slower

Object Detectors



- One stage detectors are faster, but don't perform well

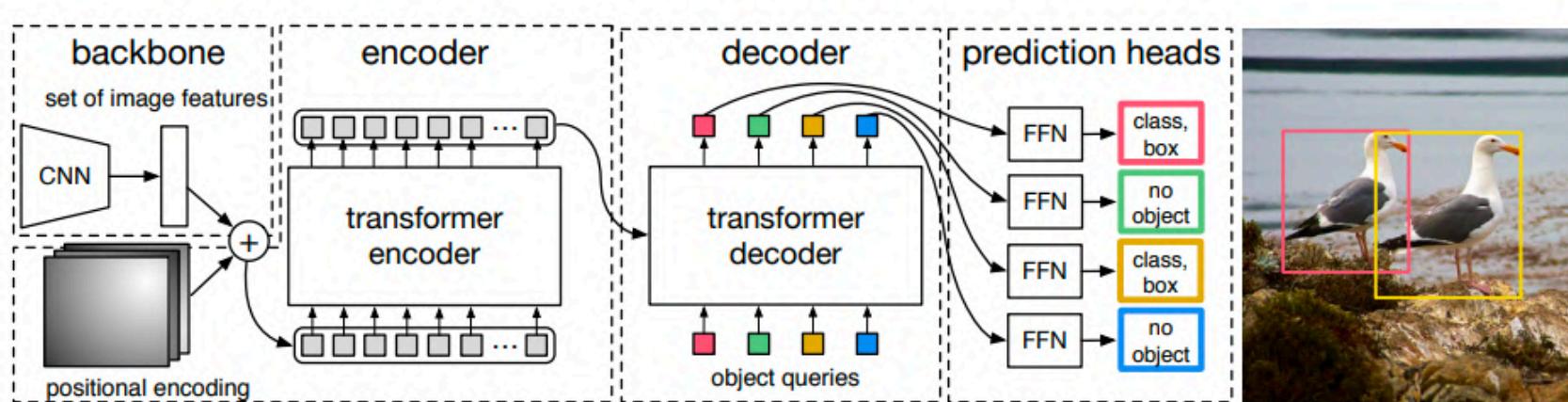
Object Detectors



- Bigger backbone improve performance, but slower

Object Detectors

- These results are few years old...
- Nowadays the best object detectors use **Transformer based backbone**



THE END