# Week 15: Object Detection
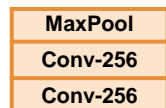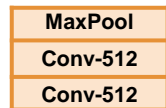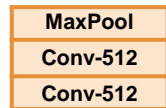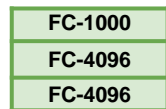
**Justin Johnson 2023**

# Transfer Learning with CNNs

Donahue et al, "DeCAF: A Deep Convolutional Activation Feature for Generic Visual Recognition", ICML 2014
Razavian et al, "CNN Features Off-the-Shelf: An Astounding Baseline for Recognition", CVPR Workshops 2014

1. Train on Imagenet

| FC-1000 |
| FC-4096 |
| FC-4096 |

| MaxPool |
| Conv-512 |
| Conv-512 |

| MaxPool |
| Conv-512 |
| Conv-512 |

| MaxPool |
| Conv-256 |
| Conv-256 |

| MaxPool |
| Conv-128 |
| Conv-128 |

| MaxPool |
| Conv-64 |
| Conv-64 |

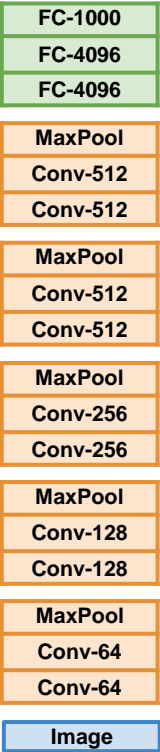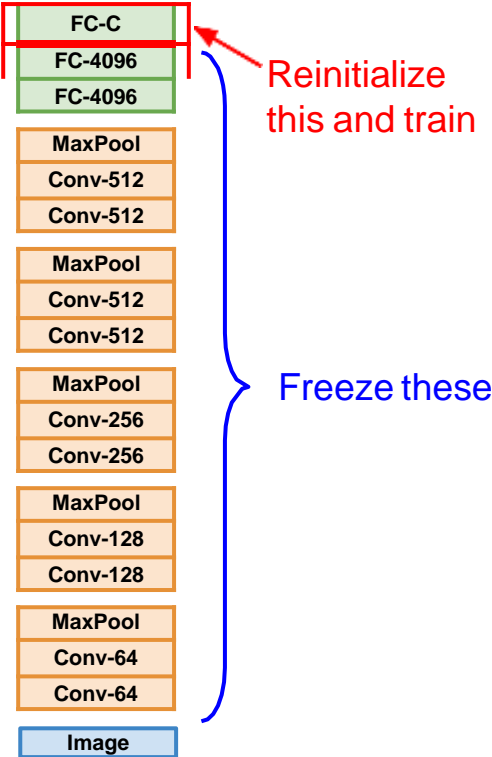| Image |

# Transfer Learning with CNNs

Donahue et al, "DeCAF: A Deep Convolutional Activation Feature for Generic Visual Recognition", ICML 2014
Razavian et al, "CNN Features Off-the-Shelf: An Astounding Baseline for Recognition", CVPR Workshops 2014



1. Train on Imagenet
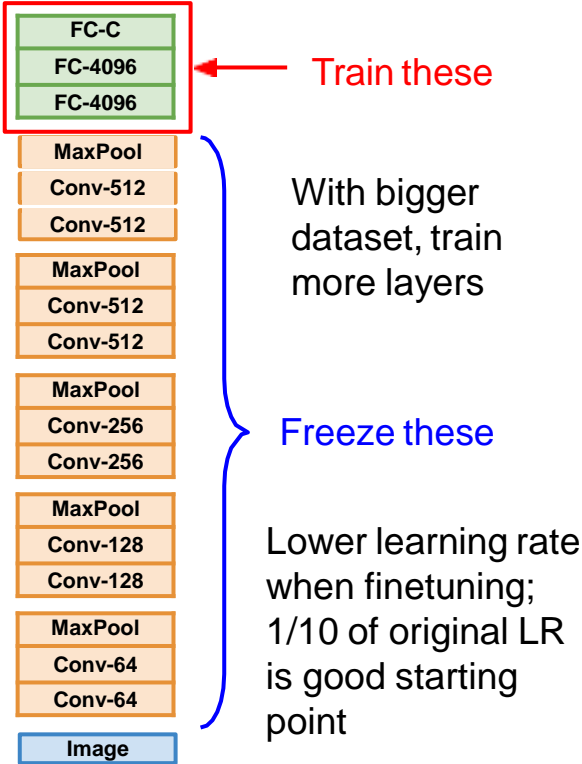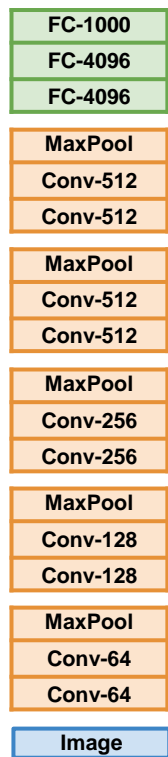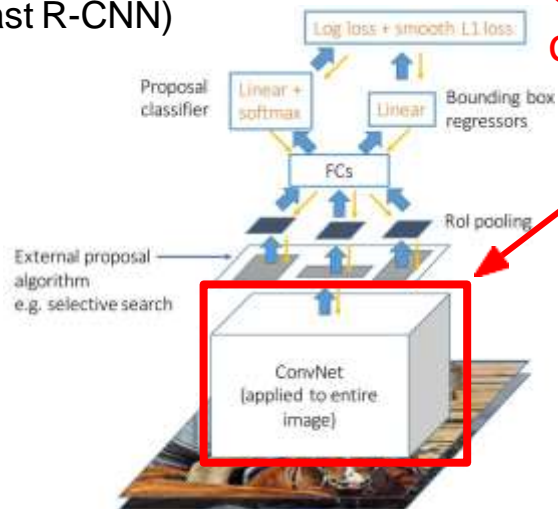
2. Small Dataset (C classes)

3. Bigger dataset

Reinitialize this and train

Train these

Freeze these

With bigger dataset, train more layers

Freeze these

Lower learning rate when finetuning; 1/10 of original LR is good starting point

| FC-1000 |
| FC-4096 |
| FC-4096 |

| MaxPool |
| Conv-512 |
| Conv-512 |

| MaxPool |
| Conv-512 |
| Conv-512 |

| MaxPool |
| Conv-256 |
| Conv-256 |

| MaxPool |
| Conv-128 |
| Conv-128 |

| MaxPool |
| Conv-64 |
| Conv-64 |

| Image |

More specific

More generic

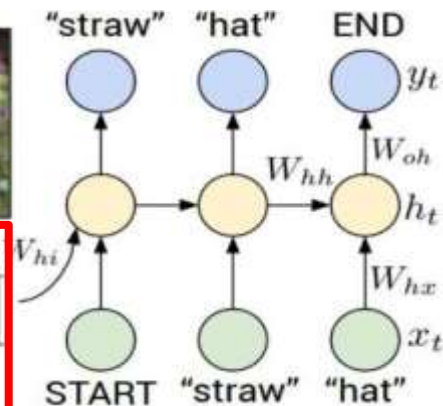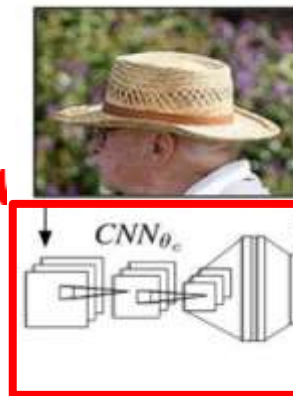|  | **very similar dataset** | **very different dataset** |
|---|---|---|
| **very little data** | Use Linear Classifier on top layer | You're in trouble… Try linear classifier from different stages |
| **quite a lot of data** | Finetune a few layers | Finetune a larger number of layers |

# Transfer learning with CNNs is pervasive…
## (it's the norm, not an exception)

Object Detection
(Fast R-CNN)

<span style="color:red">CNN pretrained
on ImageNet</span>

Image Captioning: CNN + RNN

# Transfer learning with CNNs - Architecture matters

Object detection on MSCOCO



Girshick, "The Generalized R-CNN Framework for Object Detection", ICCV 2017 Tutorial on Instance-Level Visual Recognition

# 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

# Bounding Boxes

Bounding boxes are
typically *axis-aligned*

# Bounding Boxes

Bounding boxes are typically *axis-aligned*

*Oriented* boxes are much less common

# Object Detection: Modal vs Amodal Boxes

Bounding boxes (usually) cover only the visible portion of the object

# Object Detection: Modal vs Amodal Boxes

Bounding boxes (usually) cover only the visible portion of the object

Amodal detection: box covers the entire extent of the object, even occluded parts

# 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"):

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

# 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"):

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

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",
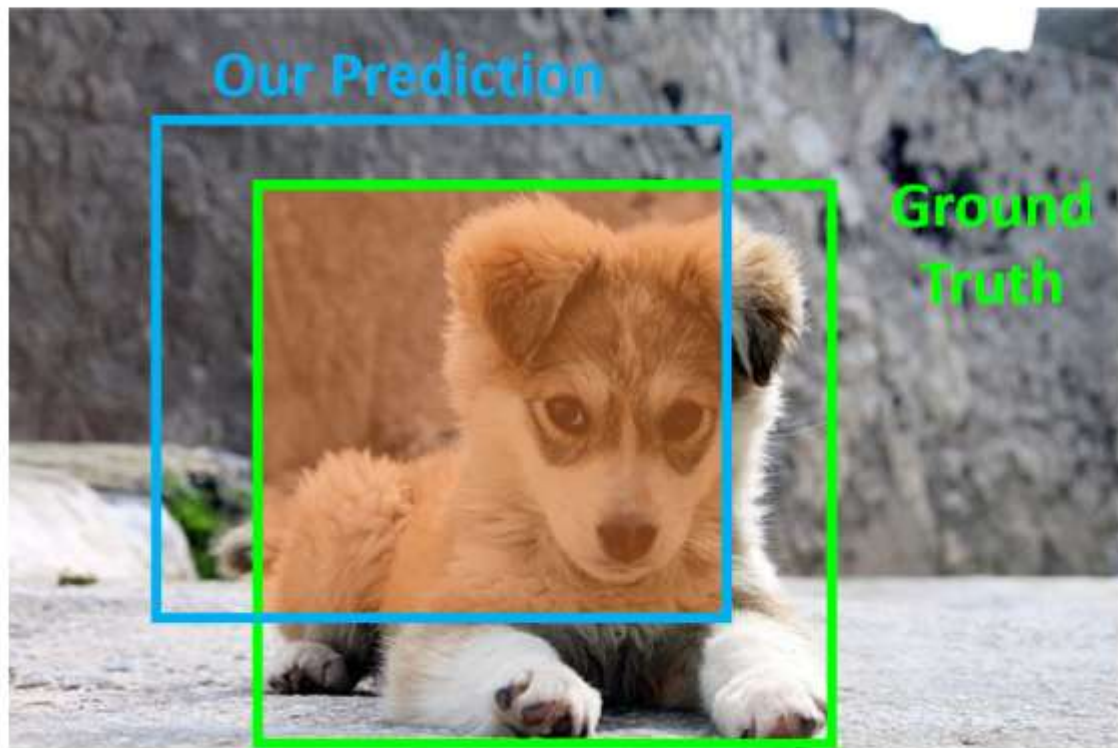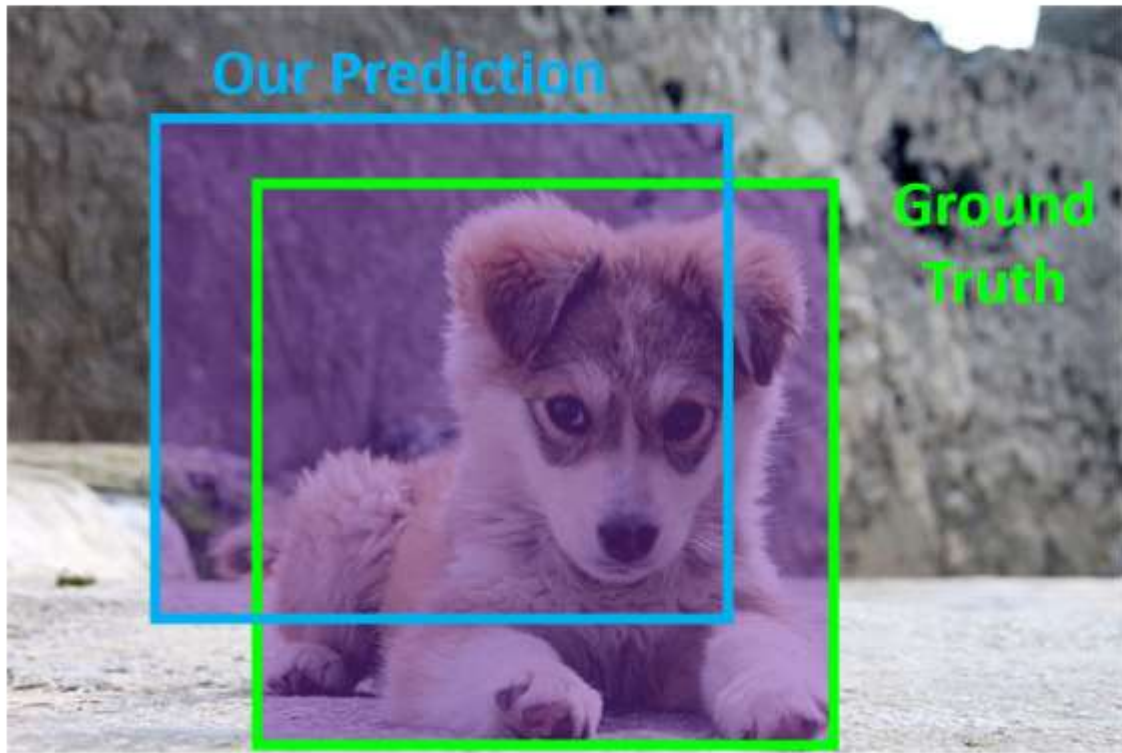IoU > 0.7 is "pretty good",

# 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"):

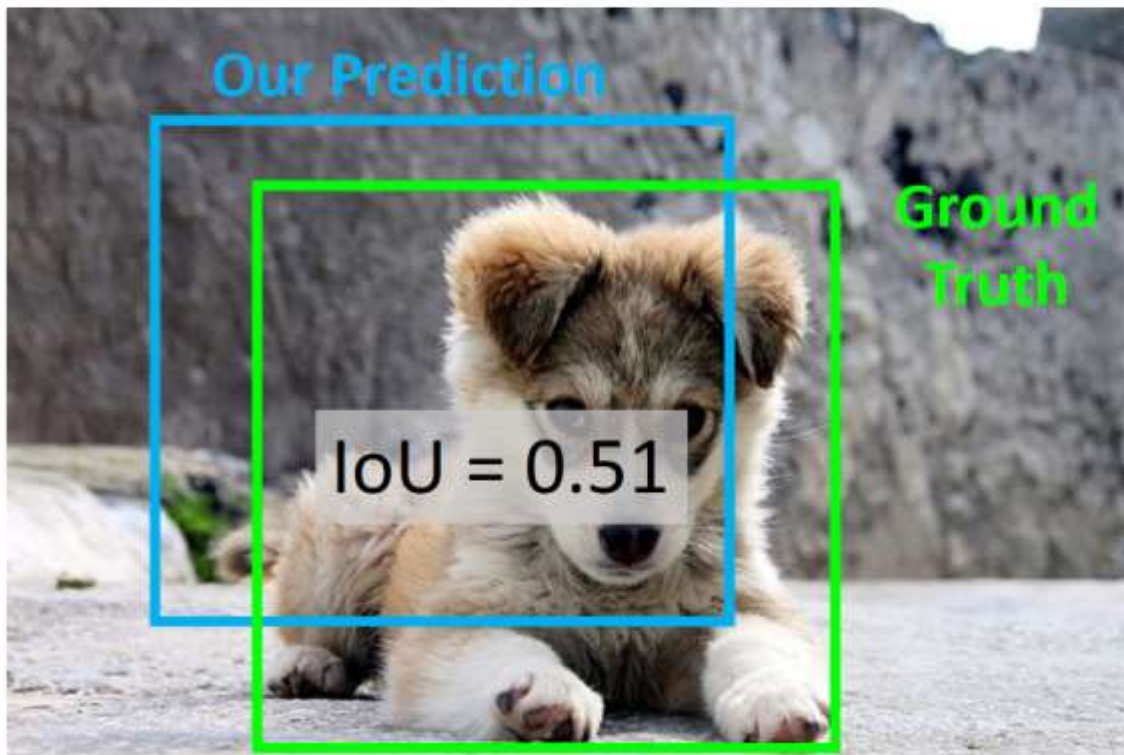$$\frac{\textit{Area of Intersection}}{\textit{Area of Union}}$$

IoU > 0.5 is "decent",
IoU > 0.7 is "pretty good",
IoU > 0.9 is "almost perfect"

# Detecting a single object



This image is CC0 public domain

Treat localization as a regression problem!

**Vector:**
4096

# Detecting a single object "What"



This image is CCO public domain

Treat localization as a regression problem!

**Vector:**
4096

**Fully Connected:**
4096 to 1000

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

**Correct label:**
Cat

**Softmax Loss**

# Detecting a single object  "What"

**Correct label:** Cat

**Fully Connected:** 4096 to 1000

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

**Softmax Loss**

This image is CC0 public domain

**Vector: 4096**

Treat localization as a regression problem!

**Fully Connected:** 4096 to 4

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

**L2 Loss**

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

"Where"

# Detecting a single object

"What"

**Correct label:**
Cat

This image is CCO public domain

Treat localization as a regression problem!

"Where"

**Vector: 4096**

**Fully Connected:** 4096 to 1000

**Fully Connected:** 4096 to 4

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

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

**Softmax Loss**

**L2 Loss**

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

Multitask Loss

**Weighted Sum**

**Loss**

$$L = L_{cls} + \lambda L_{reg}$$

# Detecting Multiple Objects

Need different numbers of outputs per image

CAT: (x, y, w, h)    4 numbers

DOG: (x, y, w, h)
DOG: (x, y, w, h)    12 numbers
CAT: (x, y, w, h)

DUCK: (x, y, w, h)    Many
DUCK: (x, y, w, h)    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

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

# Detecting Multiple Objects: Sliding Window



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

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

Consider a box of size h x w:
Possible x positions: W − w + 1
Possible y positions: H − h + 1
Possible positions:
(W − w + 1) * (H − h + 1)

# Detecting Multiple Objects: Sliding Window

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

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

Consider a box of size h x w:
Possible x positions: W − w + 1
Possible y positions: H − h + 1
Possible positions:
(W − w + 1) * (H − h + 1)

Total possible boxes:

$$= \frac{H(H+1)}{2} \frac{W(W+1)}{2}$$
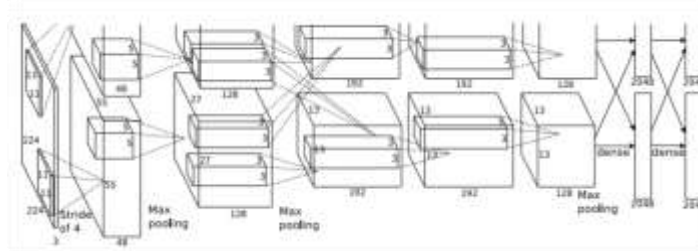
# Detecting Multiple Objects: Sliding Window

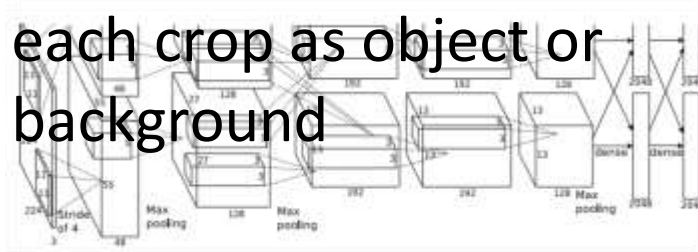Apply a CNN to many different crops of the image, CNN classifies each crop as object or background

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

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

Consider a box of size h x w:
Possible x positions: W − w + 1
Possible y positions: H − h + 1
Possible positions:
(W − w + 1) * (H − h + 1)

Total possible boxes:

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

# 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.
Figure copyright Ross Girshick, 2015; source. Reproduced with permission.

# R-CNN: Region-Based CNN



Input image

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

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

# R-CNN: Region-Based CNN



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.
Figure copyright Ross Girshick, 2015; source. Reproduced with permission.

# R-CNN: Region-Based CNN



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.
Figure copyright Ross Girshick, 2015; source. Reproduced with permission.

# R-CNN: Region-Based CNN

Class

Class

Class

ConvNet

ConvNet

ConvNet

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.
Figure copyright Ross Girshick, 2015; source. Reproduced with permission.

# R-CNN: Region-Based CNN

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



Bbox   Class

Bbox   Class

Bbox   Class

ConvNet

ConvNet

ConvNet

Forward each region through ConvNet

Warped image regions (224x224)

Input image

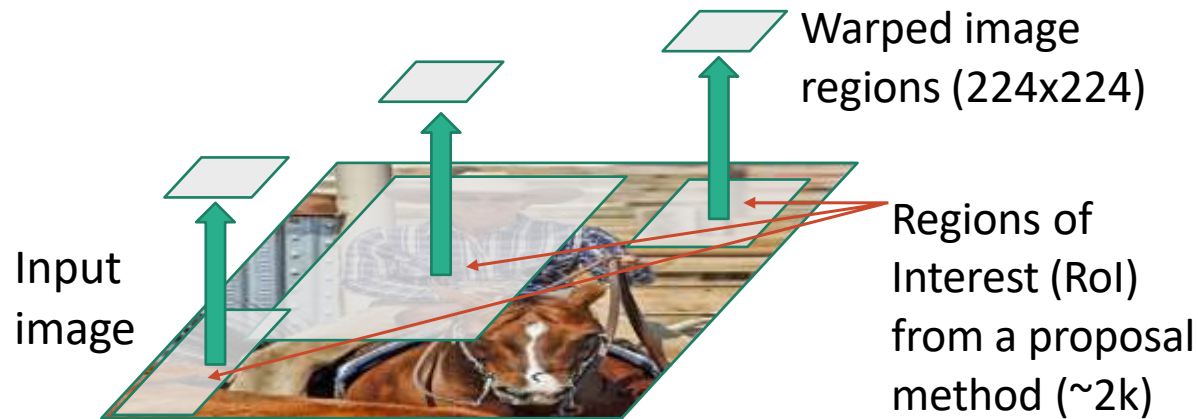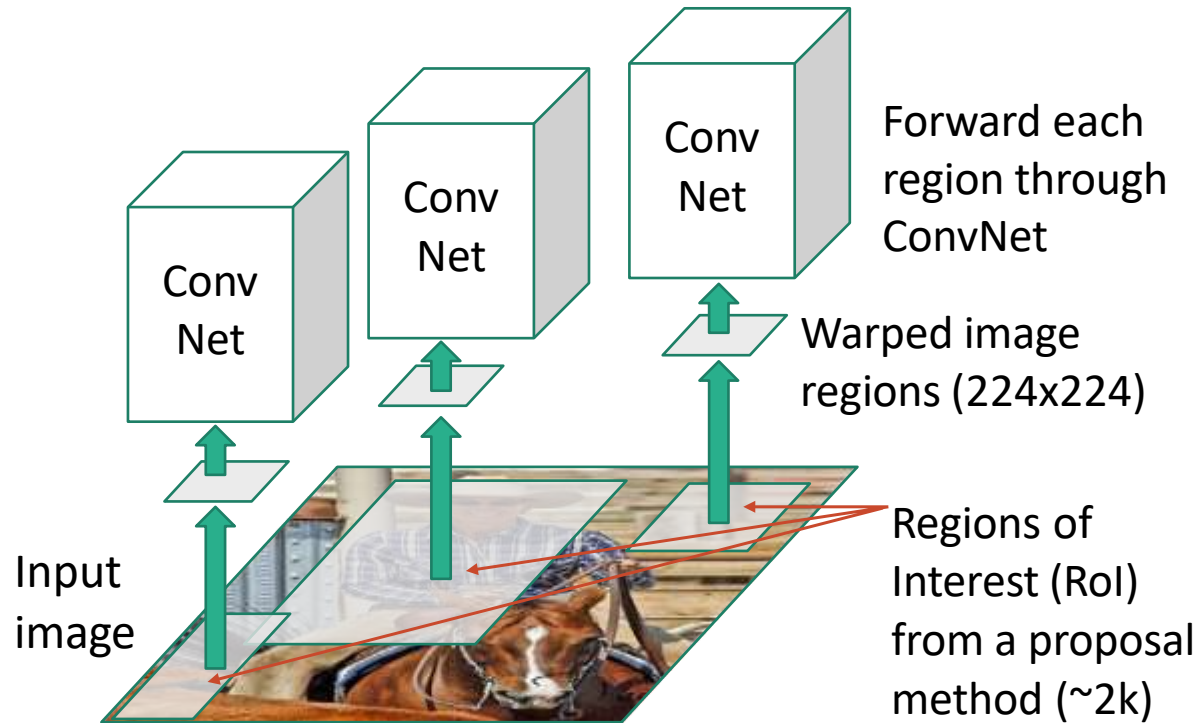Regions of Interest (RoI) from a proposal method (~2k)

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

# R-CNN Test-Time
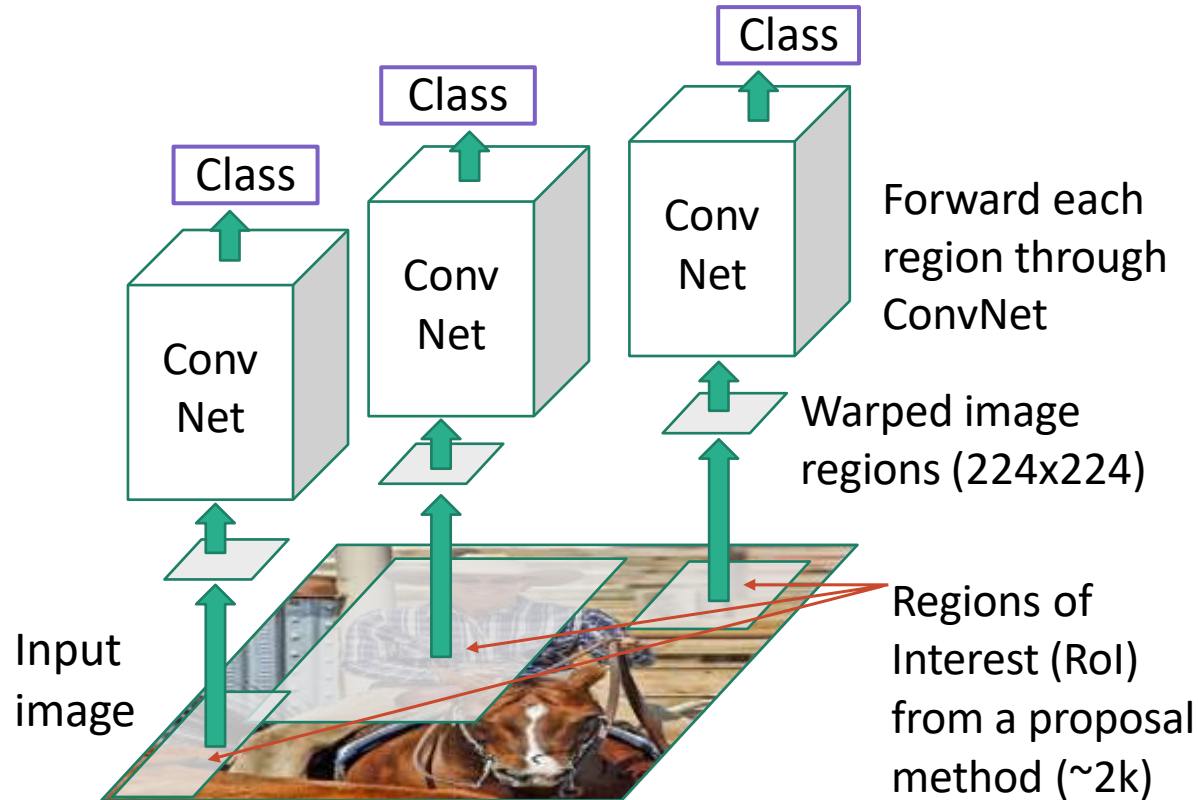
Input Image



Region Proposals

1. Run proposal method
2. Run CNN on each proposal to get class scores, transforms
3. Threshold class scores to get a set of detections

2 problems:
- CNN often outputs overlapping boxes
- How to set thresholds?

# Overlapping Boxes

**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**



P(dog) = 0.9

P(dog) = 0.75

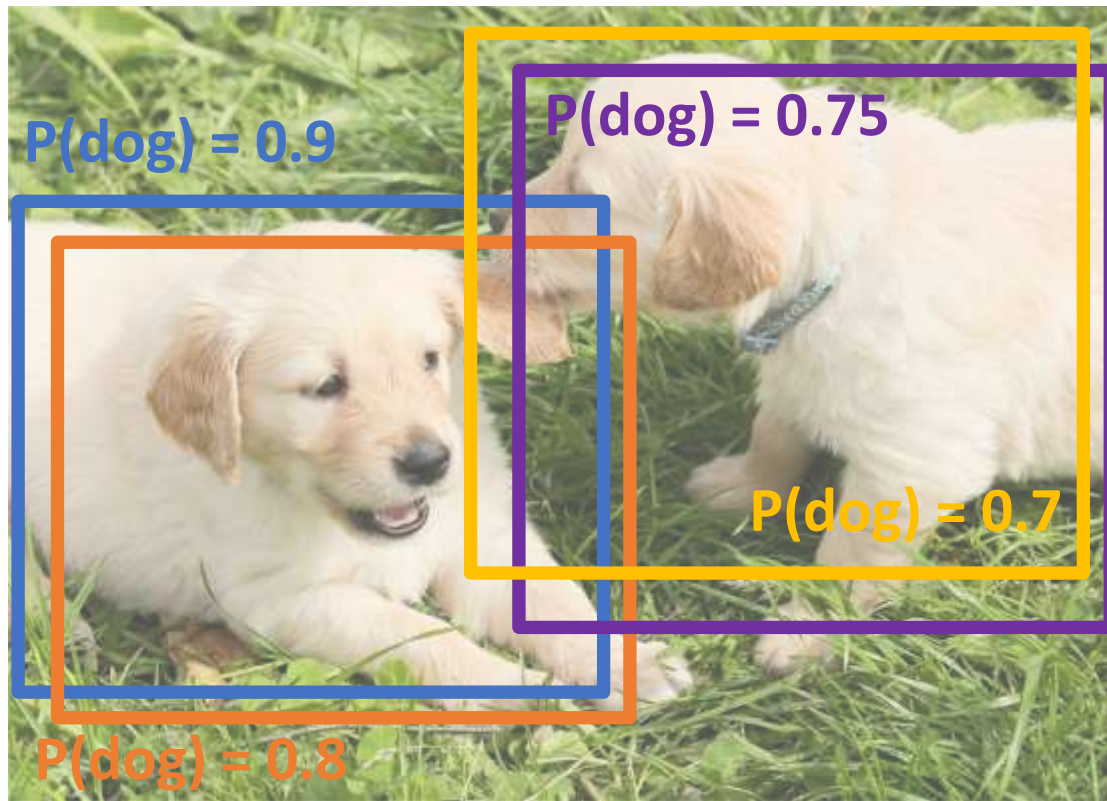P(dog) = 0.7

# Overlapping Boxes: Non-Max Suppression (NMS)

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

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

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



P(dog) = 0.9

P(dog) = 0.75

# Overlapping Boxes: Non-Max Suppression (NMS)

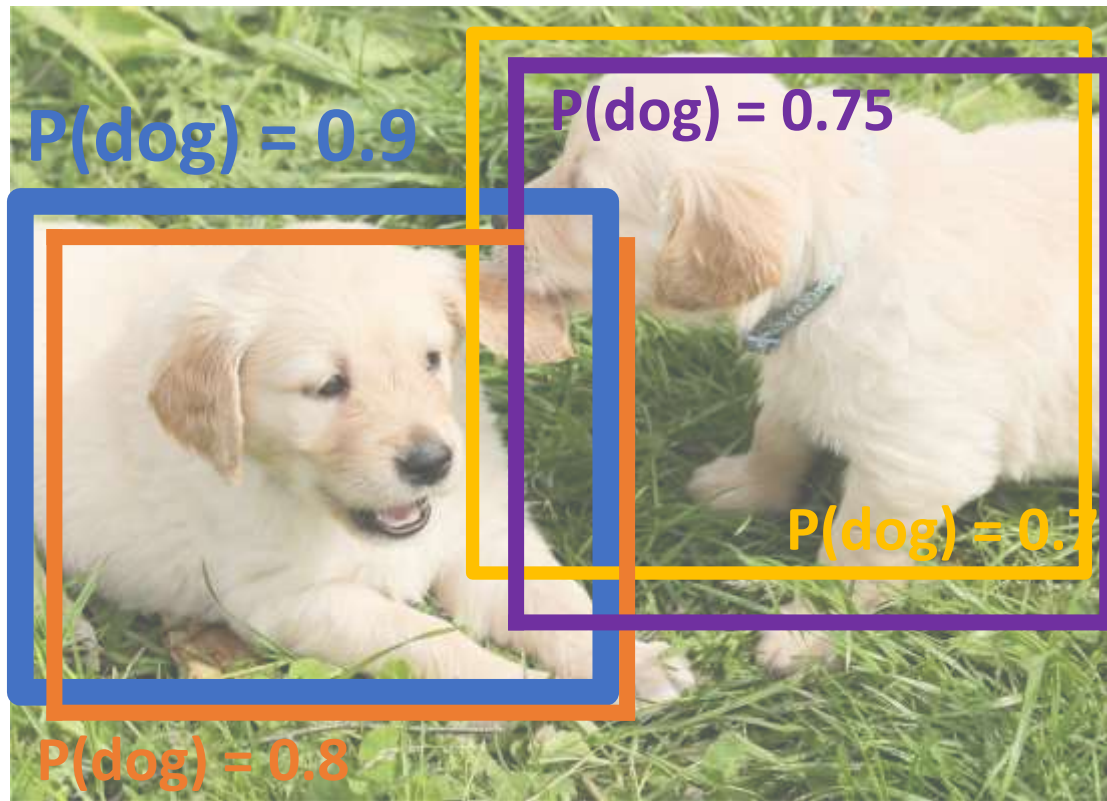**Solution**: Post-process raw detections using **Non-Max Suppression (NMS)**

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... no good solution =(

# Summary

**Transfer learning** allows us to re-use a trained network for new tasks

**Object detection** is the task of localizing objects with bounding boxes

**Intersection over Union (IoU)** quantifies differences between bounding boxes

The **R-CNN** object detector processes **region proposals** with a CNN

At test-time, eliminate overlapping detections using **non-max suppression (NMS)**

Evaluate object detectors using **mean average precision (mAP)**

# Last Time: R-CNN



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

Classify each region

Bbox   Class

Bbox   Class

Bbox   Class

ConvNet

ConvNet

ConvNet

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.
Figure copyright Ross Girshick, 2015; source. Reproduced with permission.

# Last Time: R-CNN



Input image

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

Warped image regions (224x224)

Forward each region through ConvNet

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 2000 forward passes through CNN per image

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

# Last Time: R-CNN



Bbox Class

Bbox Class

Bbox Class

Forward each region through ConvNet

ConvNet ConvNet ConvNet

Warped image regions (224x224)
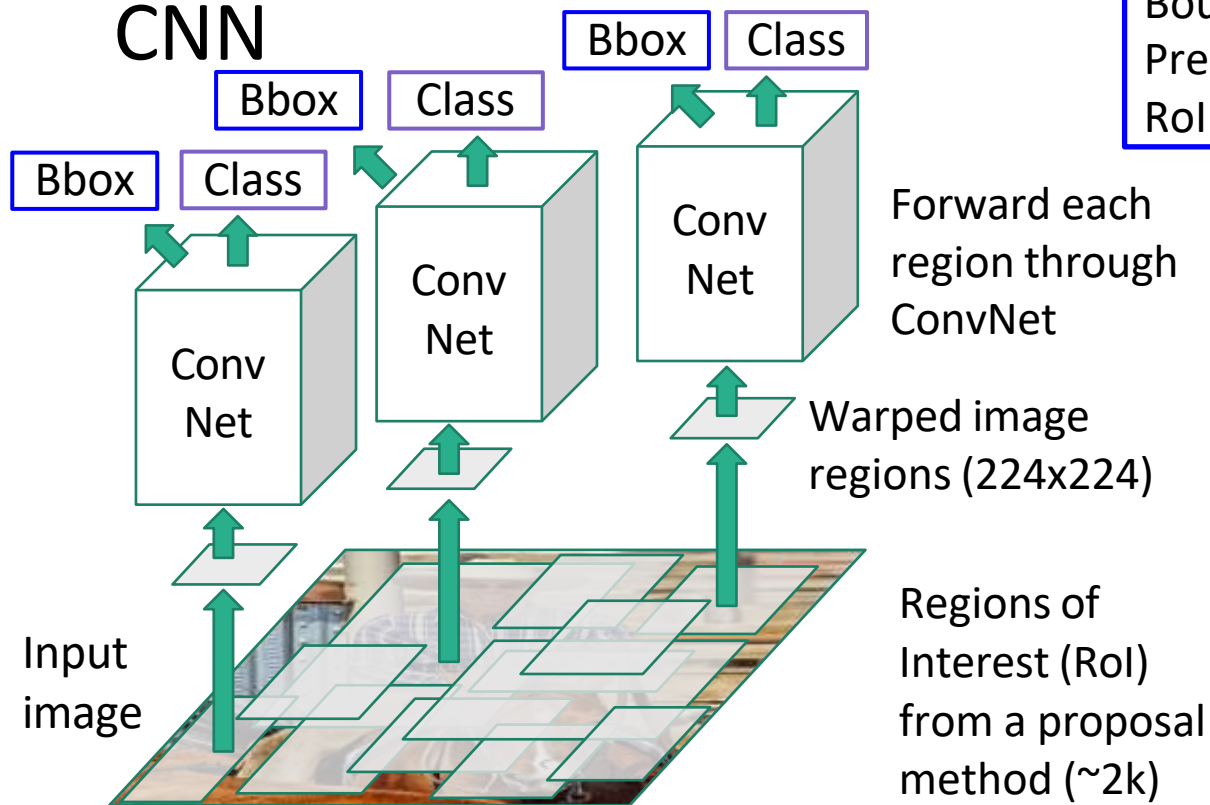
Input image

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

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 2000 forward passes through CNN per image

**Idea**: Overlapping proposals cause a lot of repeated work: same pixels processed many times. Can we avoid this?

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

"Slow" R-CNN
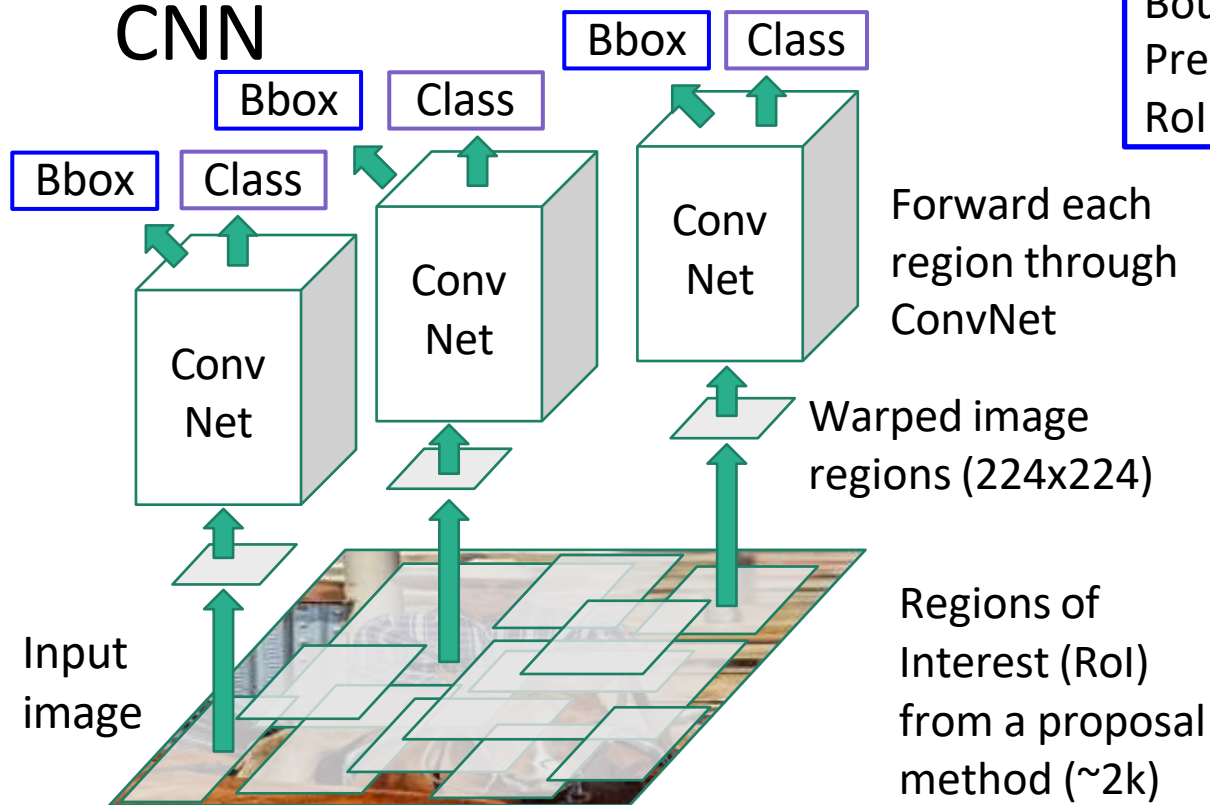Process each region
independently

# Fast R-CNN



"Slow" R-CNN
Process each region independently

Input image

Girshick, "Fast R-CNN", ICCV 2015. Figure copyright Ross Girshick, 2015; source. Reproduced with permission.

# Fast R-CNN



"Slow" R-CNN
Process each region independently

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

Image features

Run whole image through ConvNet

ConvNet

Input image

Bbox | Class
Bbox | Class
Bbox | Class

Conv Net
Conv Net
Conv Net

Input image

Girshick, "Fast R-CNN", ICCV 2015. Figure copyright Ross Girshick, 2015; source. Reproduced with permission.

# 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

Run whole image through ConvNet

ConvNet

Input image

Bbox    Class

Bbox    Class

Bbox    Class

ConvNet

ConvNet

ConvNet

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

Girshick, "Fast R-CNN", ICCV 2015. Figure copyright Ross Girshick, 2015; source. Reproduced with permission.

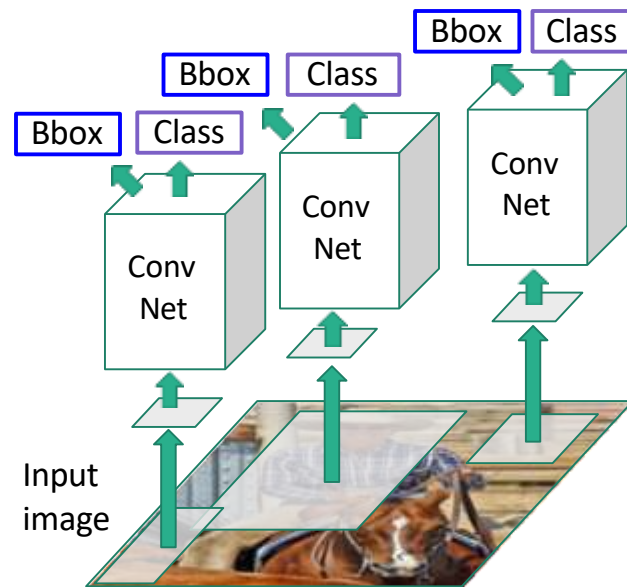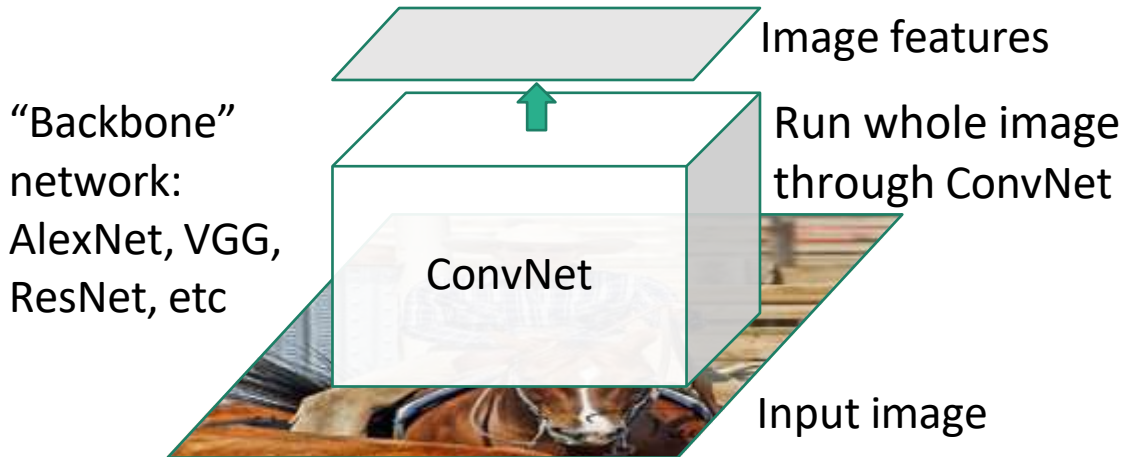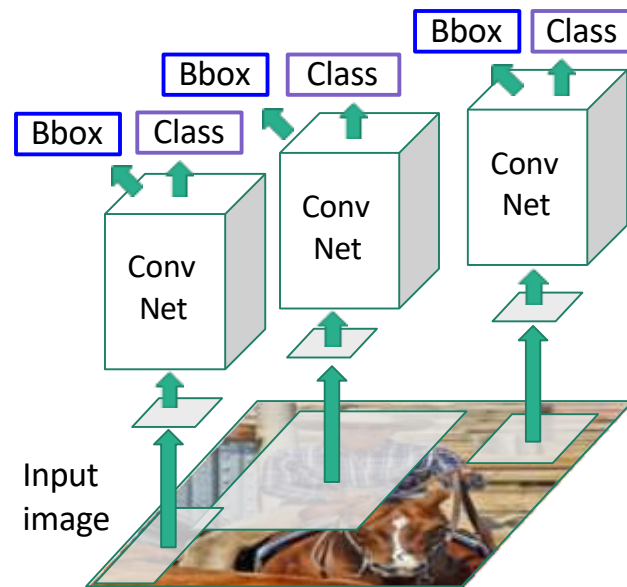# Fast R-CNN

"Slow" R-CNN
Process each region
independently

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



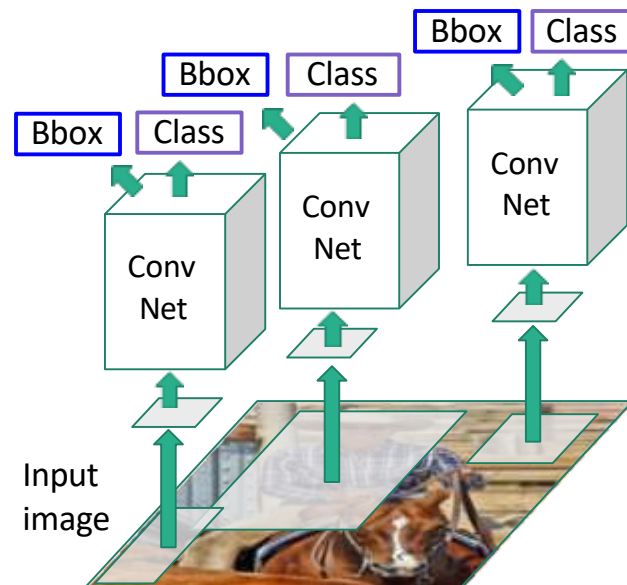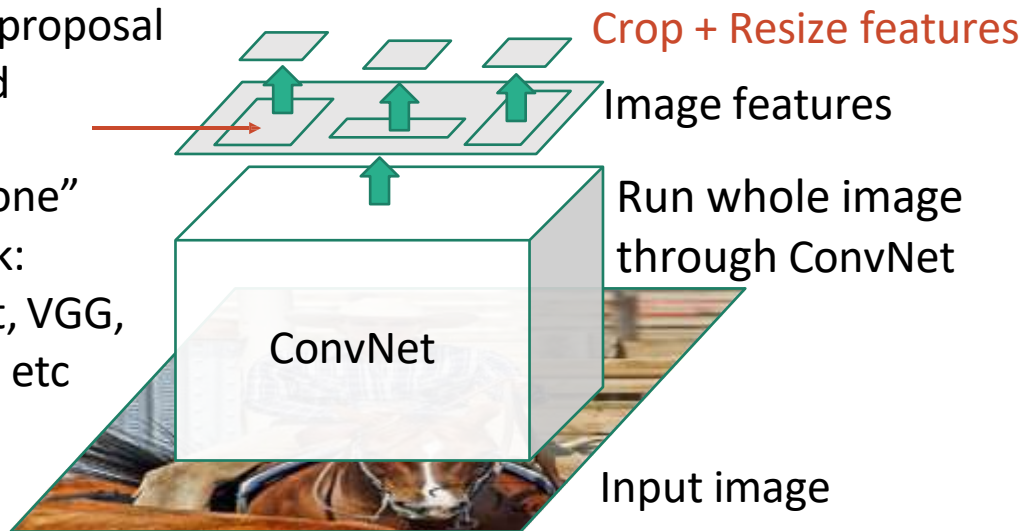Girshick, "Fast R-CNN", ICCV 2015. Figure copyright Ross Girshick, 2015; source. Reproduced with permission.

# Fast R-CNN



Bbox   Bbox   Bbox
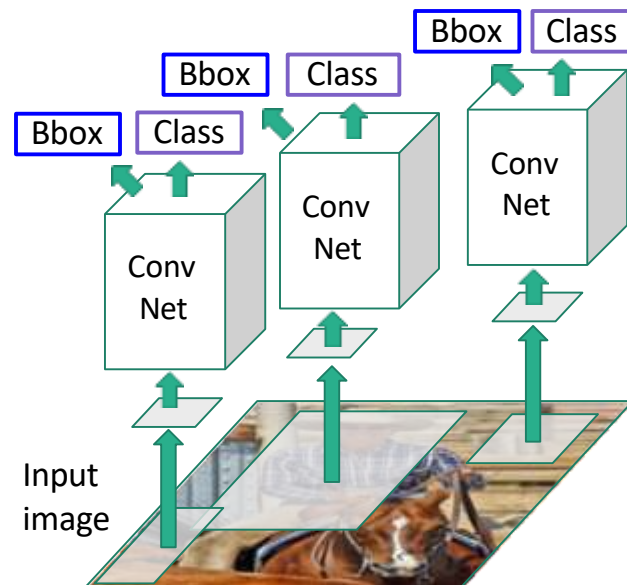
Class   Class   Class

**Category and box transform per region**

Regions of Interest (RoIs) from a proposal method

Per-Region Network

**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

Conv Net

Input image

# Fast R-CNN



Bbox  Bbox  Bbox

Class  Class  Class

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

ConvNet

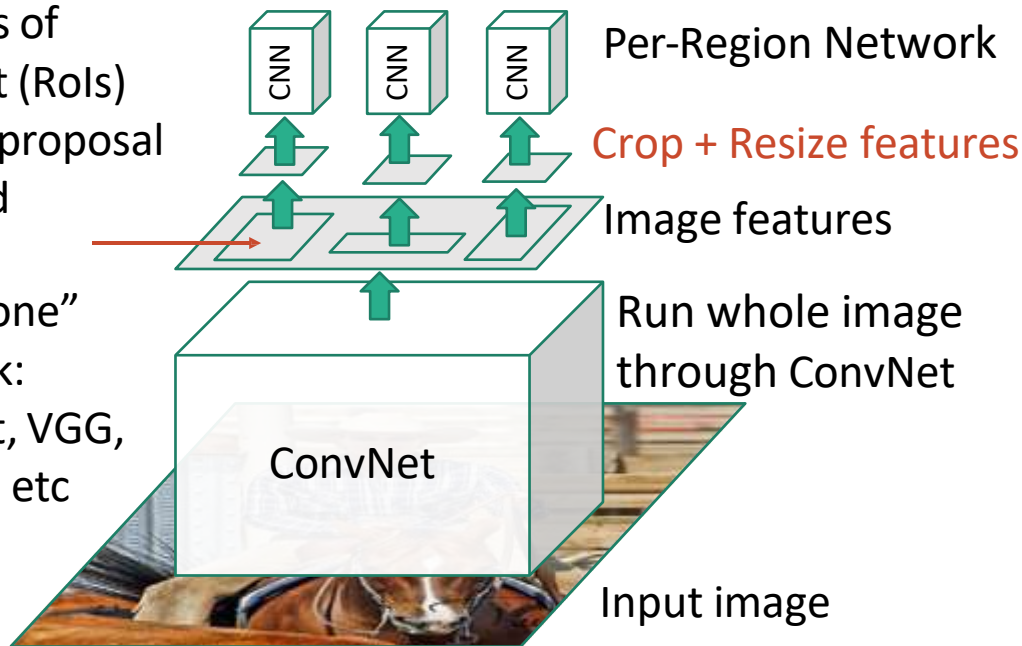Most of the computation happens in backbone network; this saves work for overlapping region proposals

Input image

Girshick, "Fast R-CNN", ICCV 2015. Figure copyright Ross Girshick, 2015; source. Reproduced with permission.

# Fast R-CNN

Bbox   Bbox   Bbox

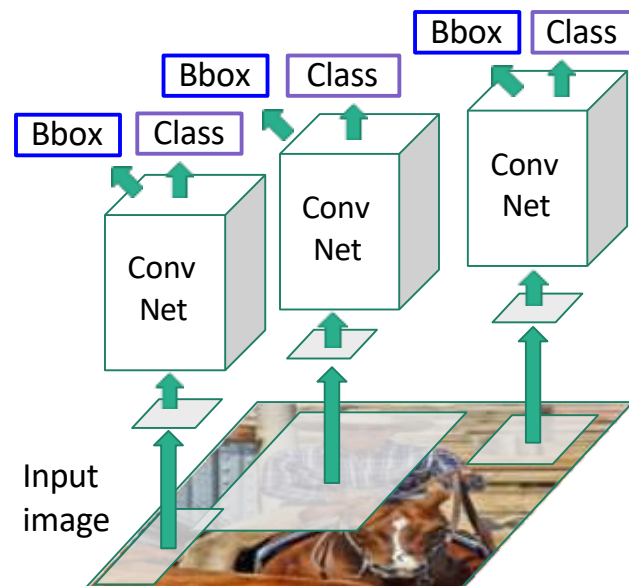Class   Class   Class

**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

ConvNet

Run whole image through ConvNet

Input image

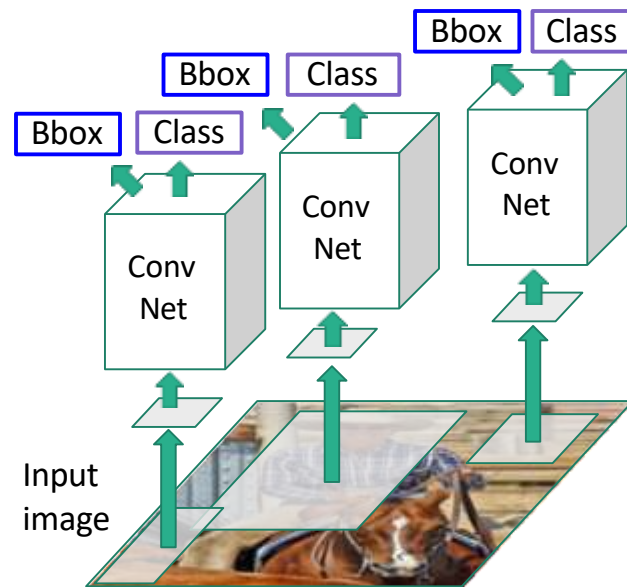Example: When using AlexNet for detection, five conv layers are used for backbone and two FC layers are used for per-region network

Girshick, "Fast R-CNN", ICCV 2015. Figure copyright Ross Girshick, 2015; source. Reproduced with permission.
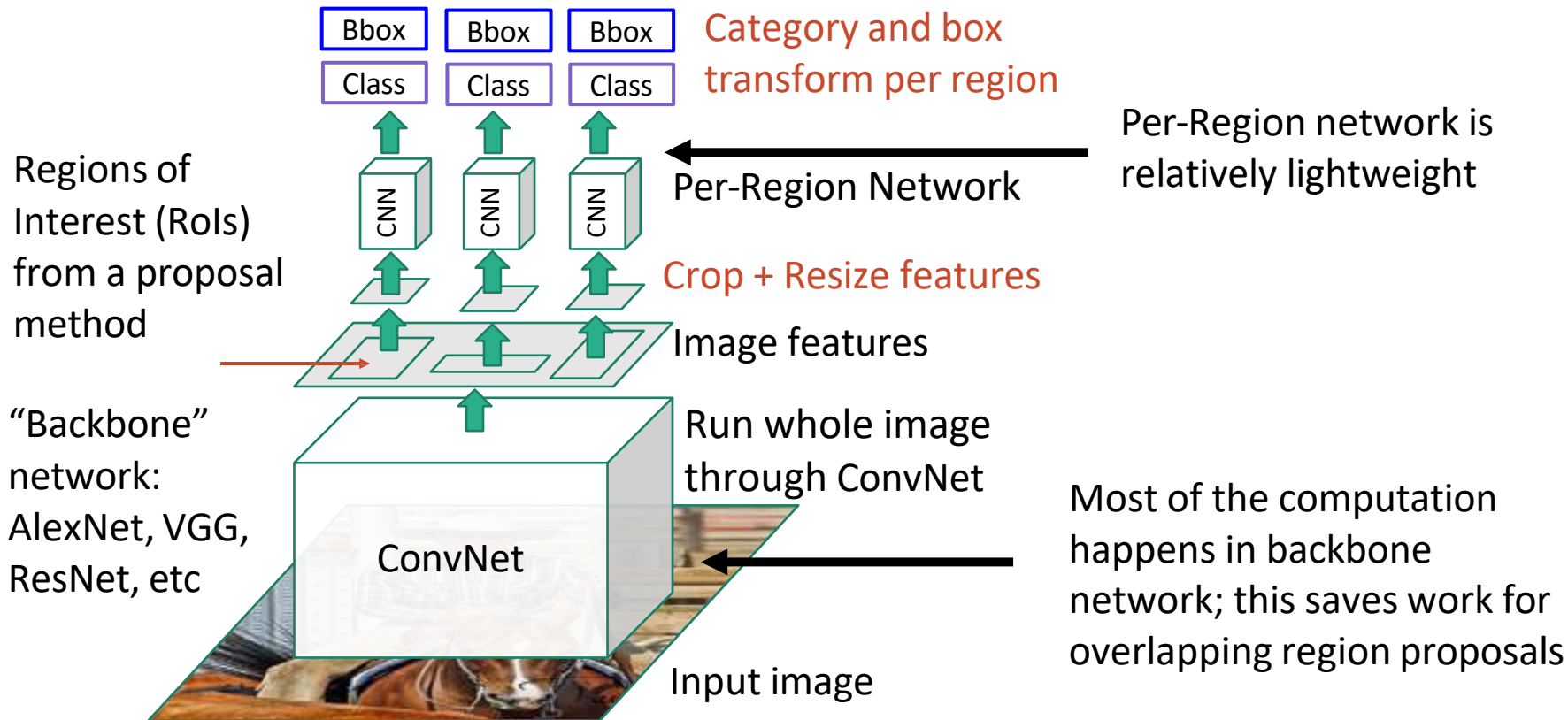
# Fast R-CNN

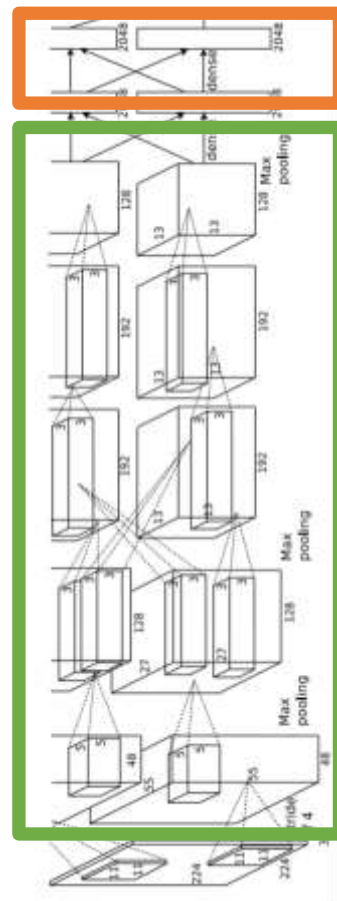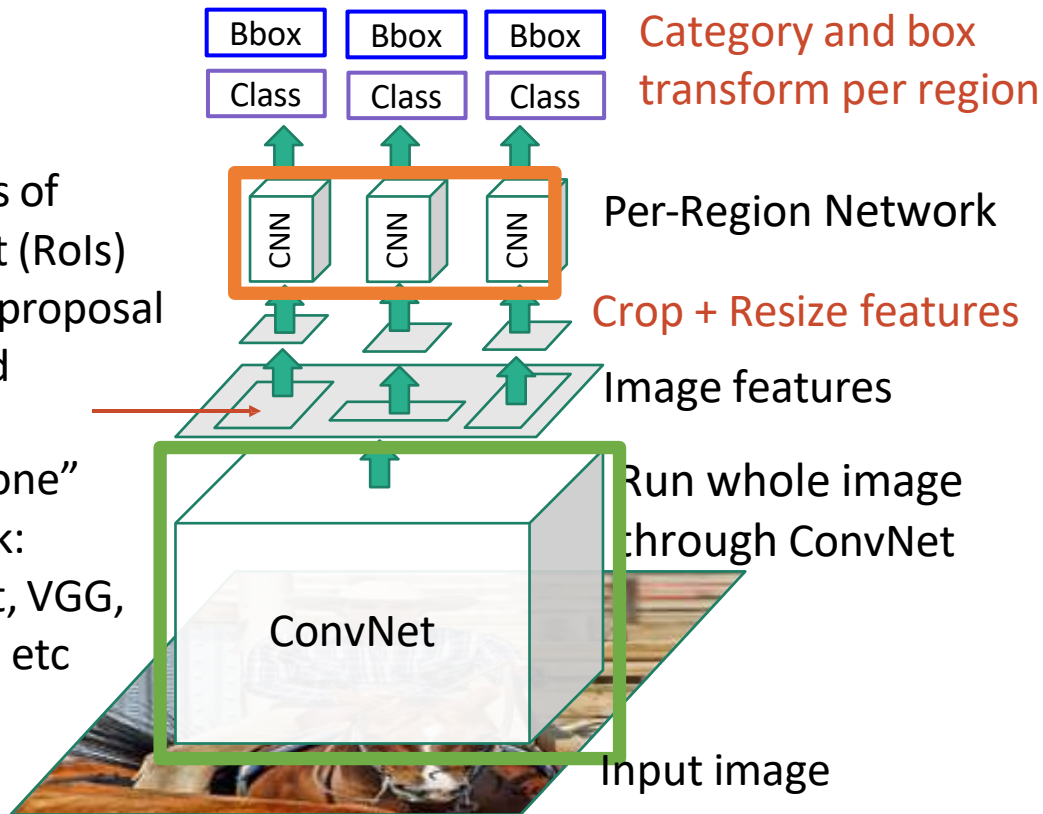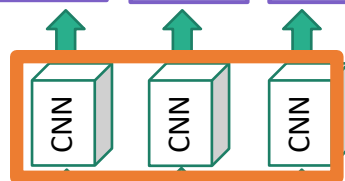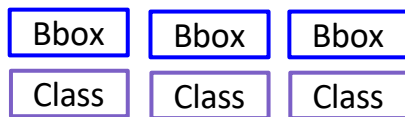| Bbox | Bbox | Bbox |
|---|---|---|
| Class | Class | Class |

Category and box transform per region

Regions of Interest (RoIs) from a proposal method

Per-Region Network

Crop + Resize features

Image features

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

Run whole image through ConvNet

ConvNet

Input image

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

Softmax
FC 1000
Pool
3x3 conv, 512
3x3 conv, 512
3x3 conv, 512
3x3 conv, 512
3x3 conv, 512
3x3 conv, 512, / 2

3x3 conv, 128
3x3 conv, 128
3x3 conv, 128
3x3 conv, 128
3x3 conv, 128
3x3 conv, 128, / 2
3x3 conv, 64
3x3 conv, 64
3x3 conv, 64
3x3 conv, 64
3x3 conv, 64
3x3 conv, 64
Pool
7x7 conv, 64, / 2
Input

# 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

# 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)**

- 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!

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)**

**Test time (seconds)**

**Problem**: Runtime dominated by region proposals!

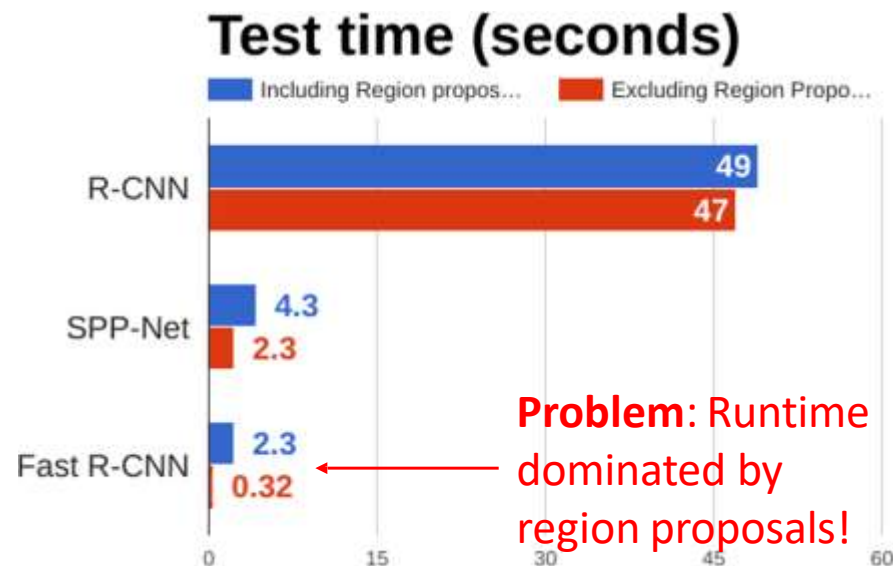**Recall**: Region proposals computed by heuristic "Selective Search" algorithm on CPU -- let's learn them with a CNN instead!
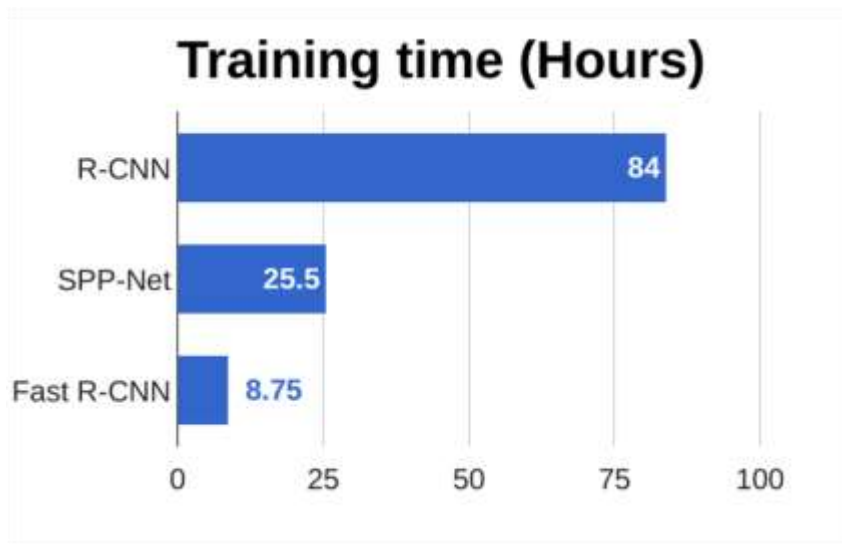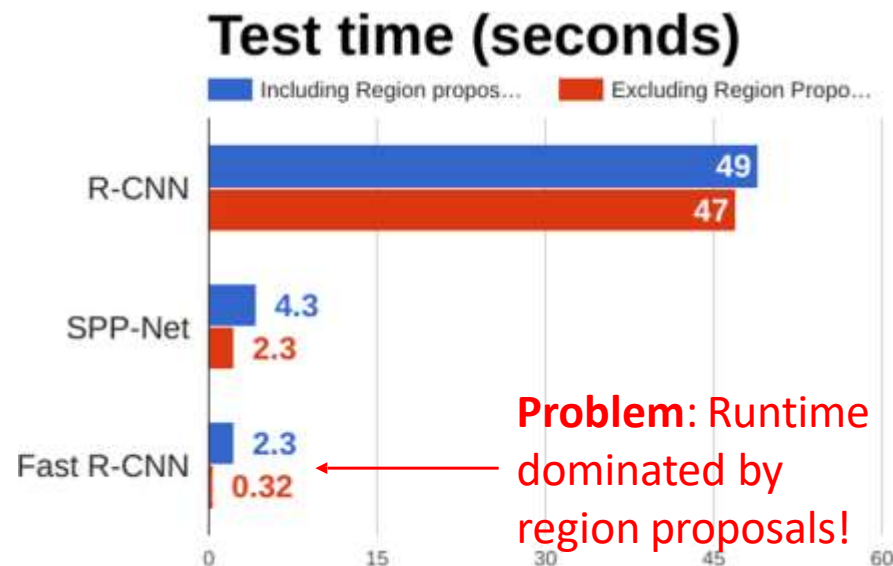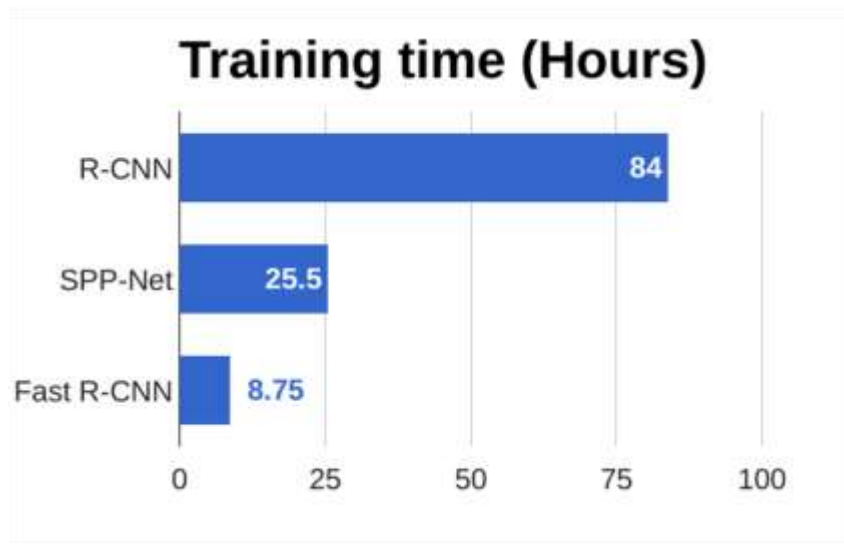
Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.
He et al, "Spatial pyramid pooling in deep convolutional networks for visual recognition", ECCV 2014
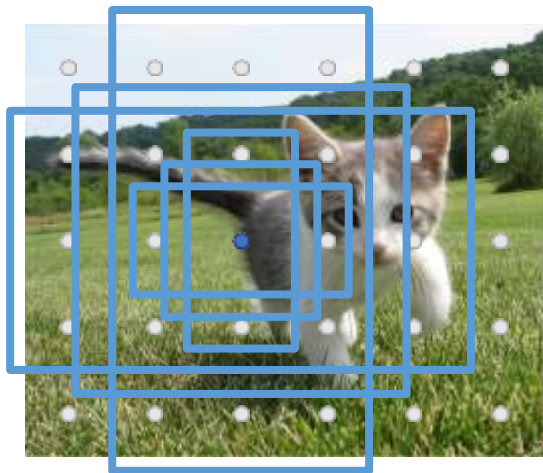Girshick, "Fast R-CNN", ICCV 2015

# Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image

Each feature corresponds to a point in the input

In practice: Rather than using one anchor per point, instead consider K different anchors with different size and scale (here K = 6)



CNN

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

Image features
(e.g. 512 x 5 x 6)

Conv

Anchor is object?
2K x 5 x 6

Anchor transforms
4K x 5 x 6

At test-time, sort all K*5*6 boxes by their positive score, take top 300 as our region proposals

Ren et al, "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks", NIPS 2015

# Faster R-CNN: Learnable Region Proposals

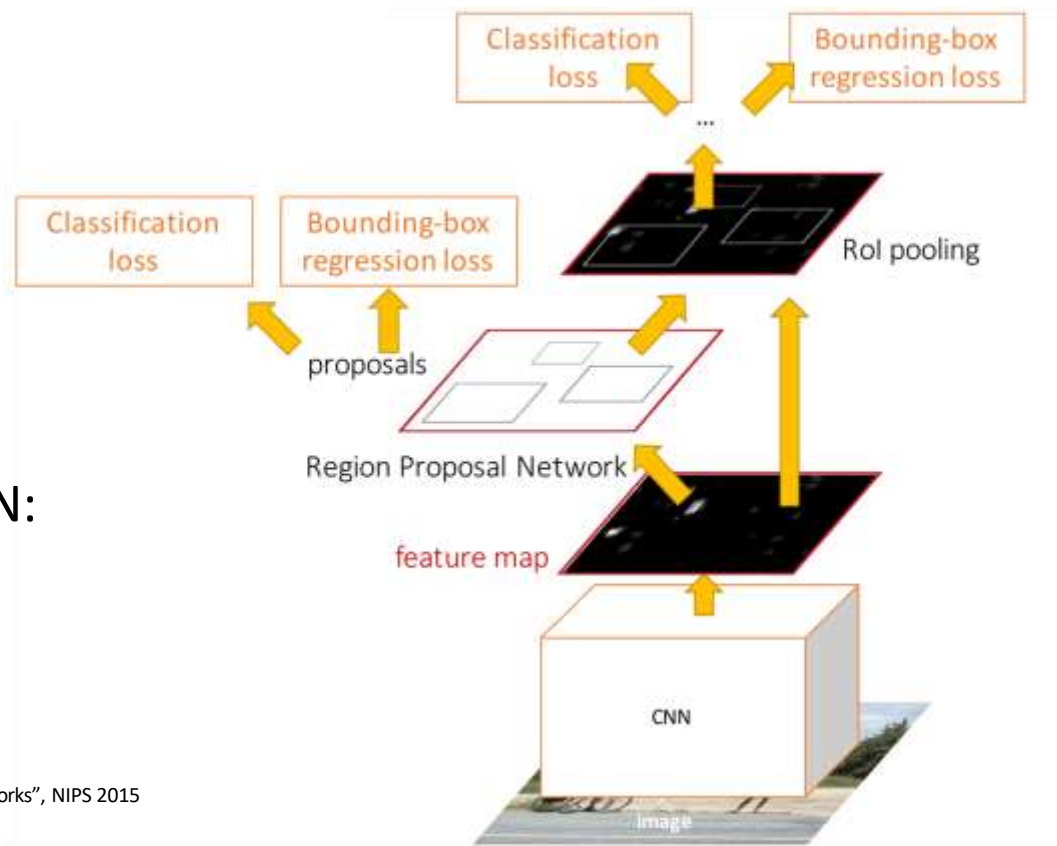Insert **Region Proposal Network (RPN)** to predict proposals from features

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

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

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



Classification loss

Bounding-box regression loss

RoI pooling

Classification loss

Bounding-box regression loss

proposals

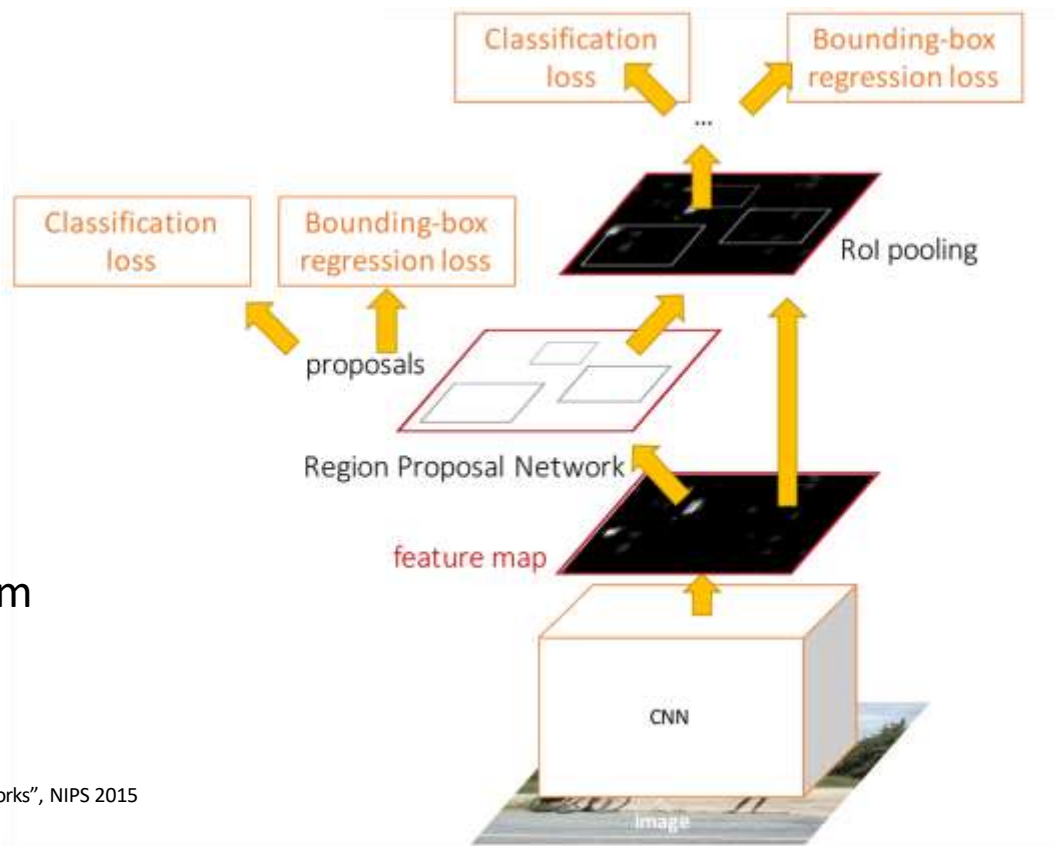Region Proposal Network

feature map

CNN

Image

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
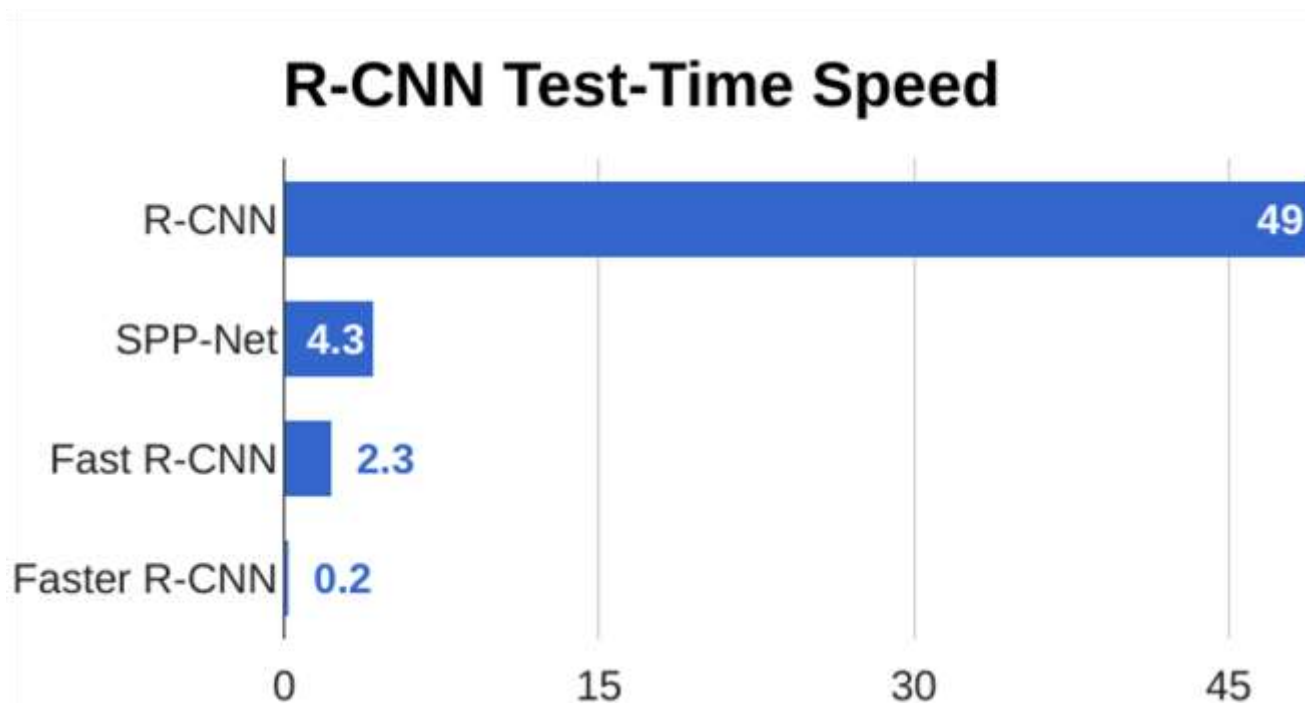
# Faster R-CNN: Learnable Region Proposals

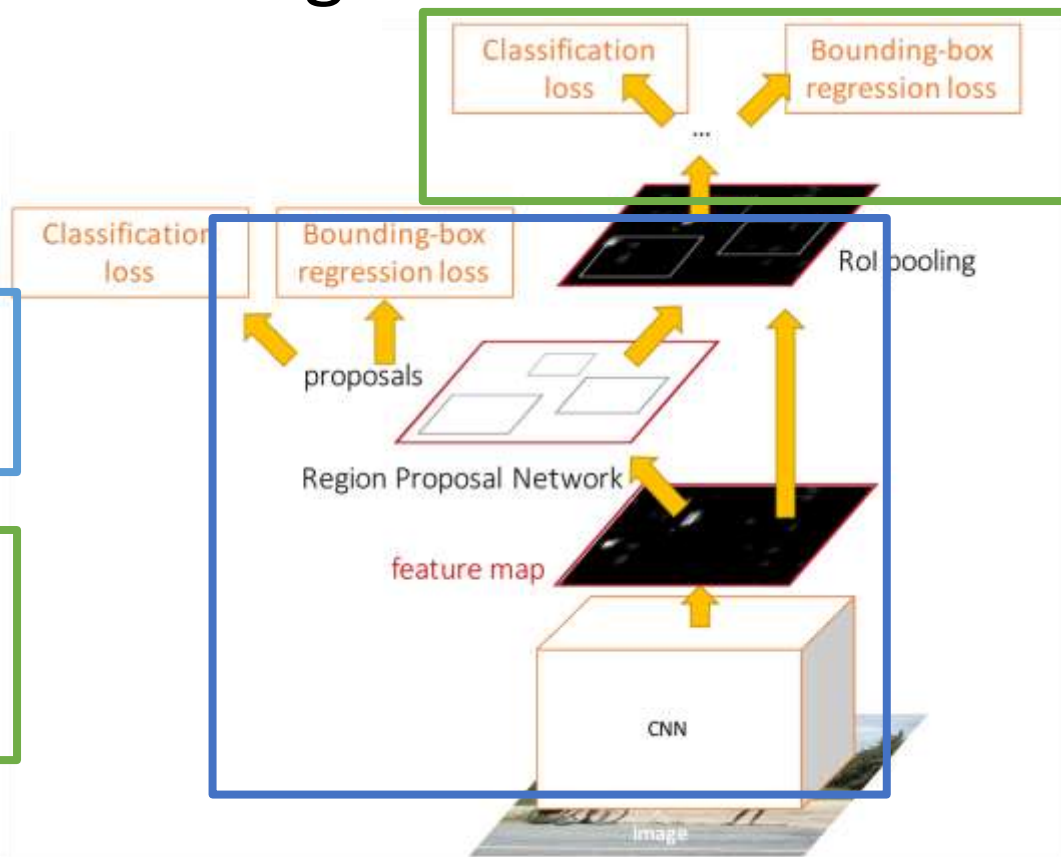# Faster R-CNN: Learnable Region Proposals

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

# Dealing with Scale

We need to detect objects of many different scales.
How to improve *scale invariance* of the detector?

Dealing with Scale: Image Pyramid

Classic idea: build an *image pyramid* by resizing the image to different scales, then process each image scale independently.



Object Detector

Object Detector

Object Detector

Dealing with Scale: Image Pyramid

Classic idea: build an *image pyramid* by resizing the image to different scales, then process each image scale independently.

Problem: Expensive! Don't share any computation between scales

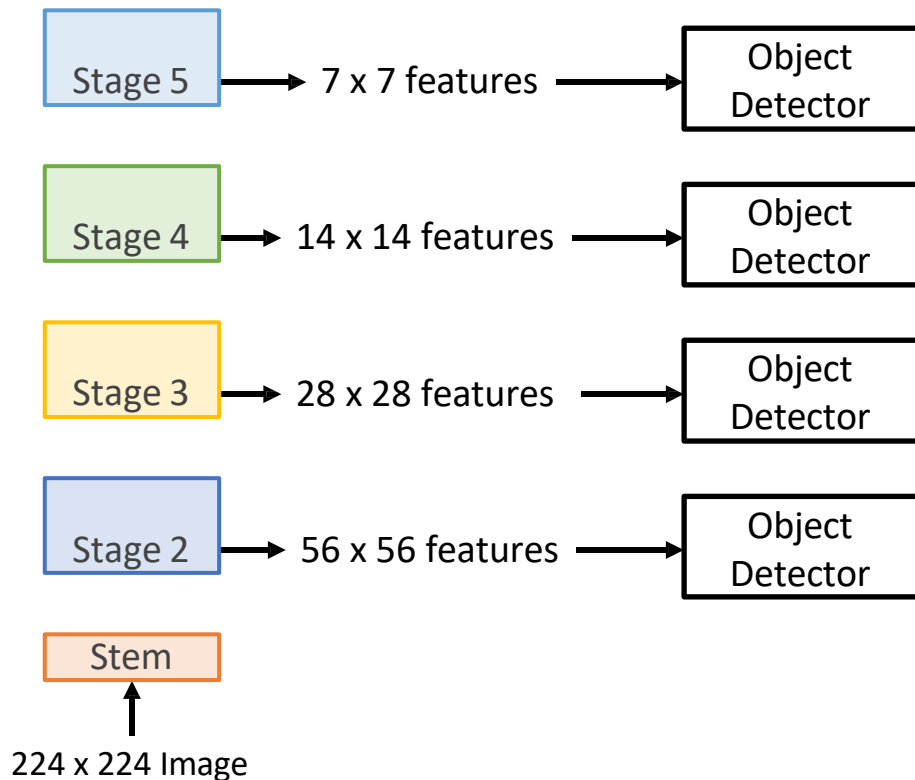# Dealing with Scale: Multiscale Features

CNNs have multiple *stages* that operate at different resolutions. Attach an independent detector to the features at each level

Lin et al, "Feature Pyramid Networks for Object Detection", ICCV 2017
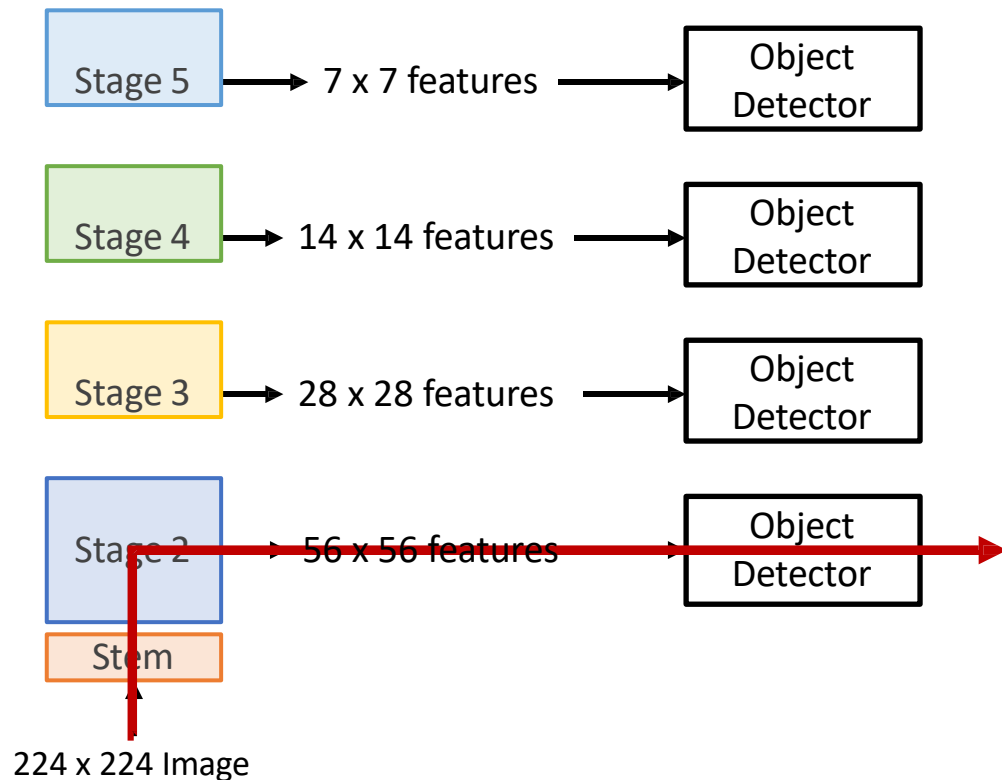
# Dealing with Scale: Multiscale Features

CNNs have multiple *stages* that operate at different resolutions. Attach an independent detector to the features at each level

Problem: detector on early features doesn't make use of the entire backbone; doesn't get access to high-level features

| Stage 5 | → 7 x 7 features → | Object Detector |
|---------|--------------------|-----------------|
| Stage 4 | → 14 x 14 features → | Object Detector |
| Stage 3 | → 28 x 28 features → | Object Detector |
| Stage 2 | → 56 x 56 features → | Object Detector |

Stem

224 x 224 Image

# Dealing with Scale: Feature Pyramid Network

Add *top down connections* that feed information from high level features back down to lower level features



Stage 5 → 7 x 7 feats → Object Detector

Stage 4 → 14 x 14 feats

Stage 3 → 28 x 28 feats

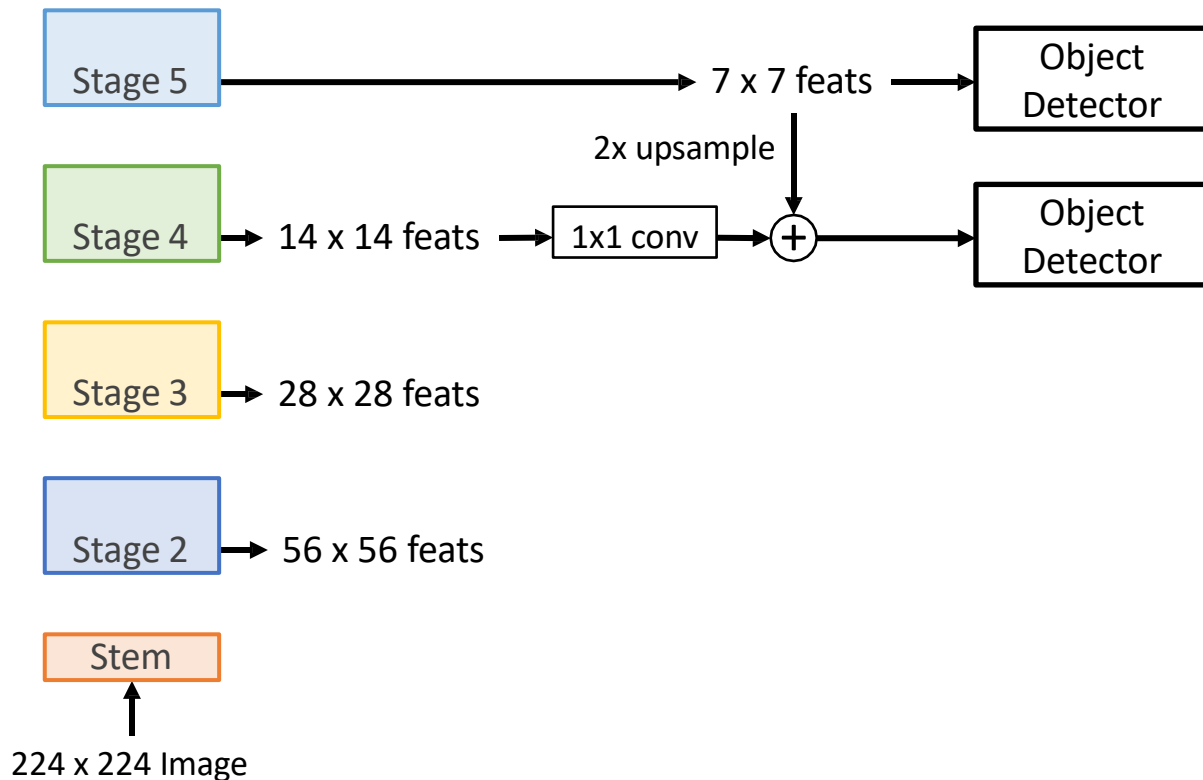Stage 2 → 56 x 56 feats

Stem

224 x 224 Image

# Dealing with Scale: Feature Pyramid Network

Add *top down connections* that feed information from high level features back down to lower level features

# Dealing with Scale: Feature Pyramid Network

Add *top down connections* that feed information from high level features back down to lower level features
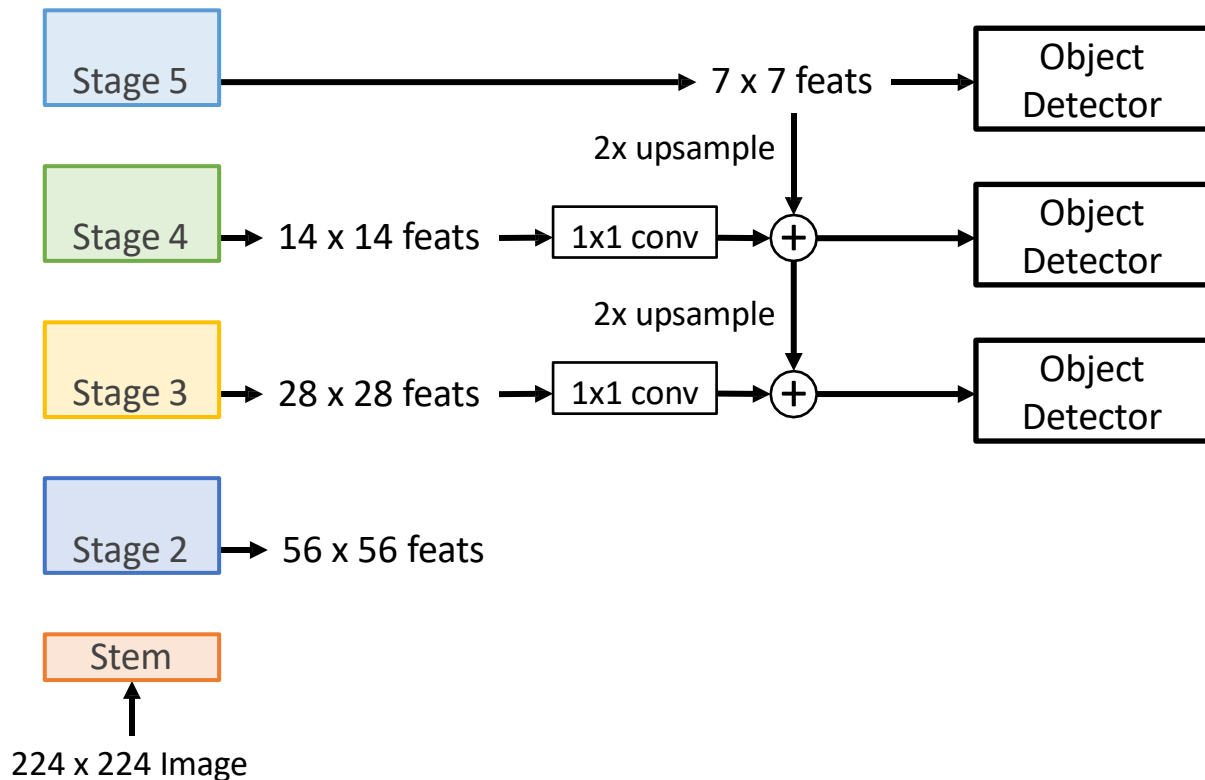
# Dealing with Scale: Feature Pyramid Network

Add *top down connections* that feed information from high level features back down to lower level features



Lin et al, "Feature Pyramid Networks for Object Detection", ICCV 2017
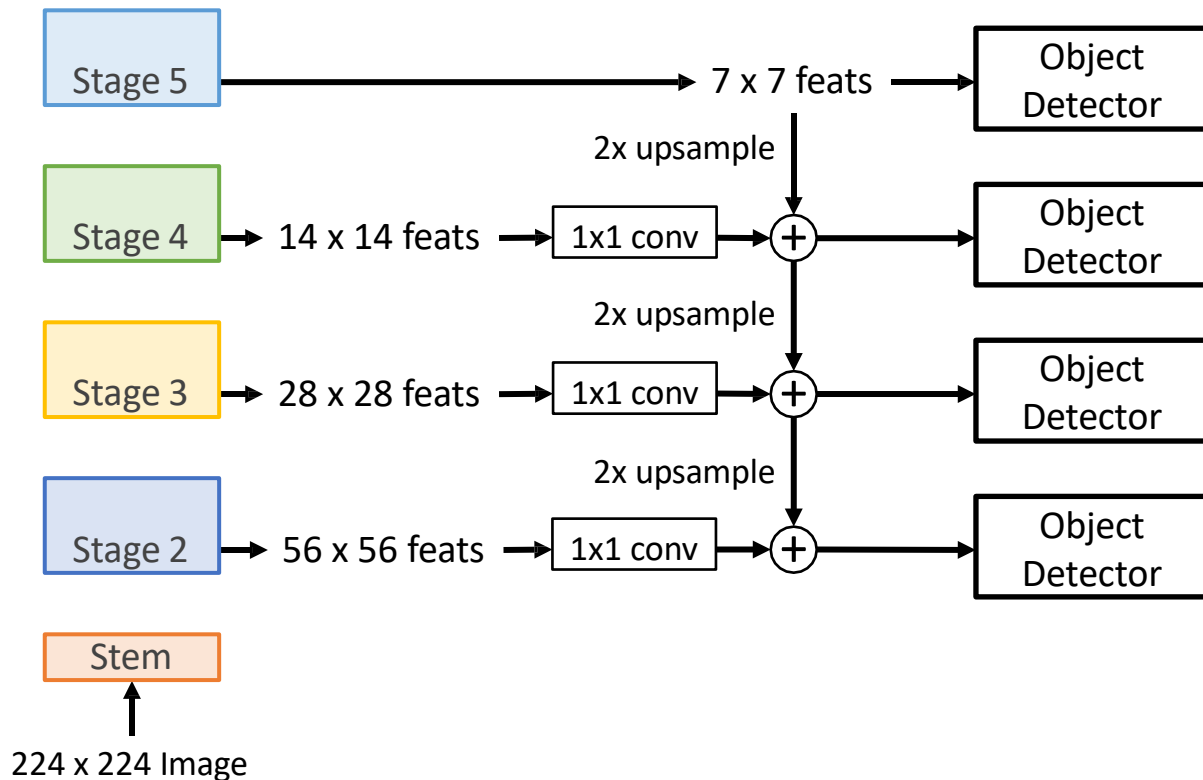
# Dealing with Scale: Feature Pyramid Network

Add *top down connections* that feed information from high level features back down to lower level features

Efficient multiscale features where all levels benefit from the whole backbone! Widely used in practice

# Dealing with Scale: Feature Pyramid Network

Add *top down connections* that feed information from high level features back down to lower level features
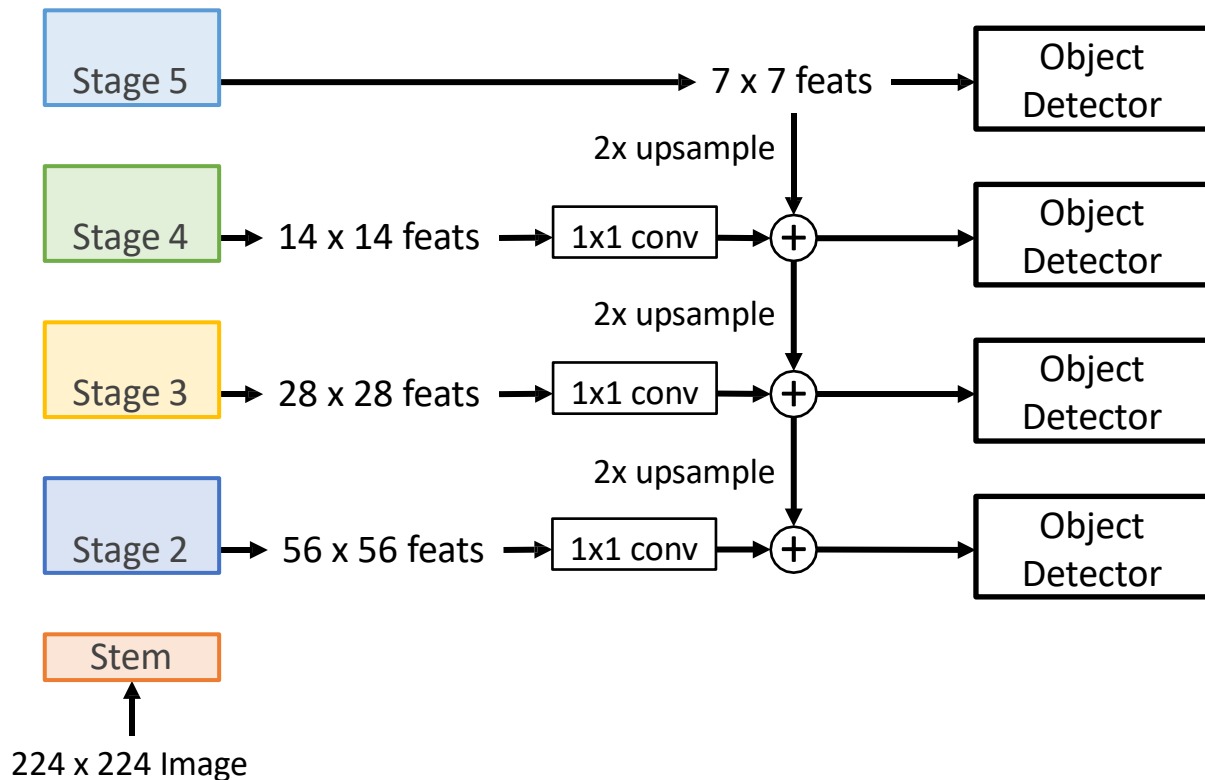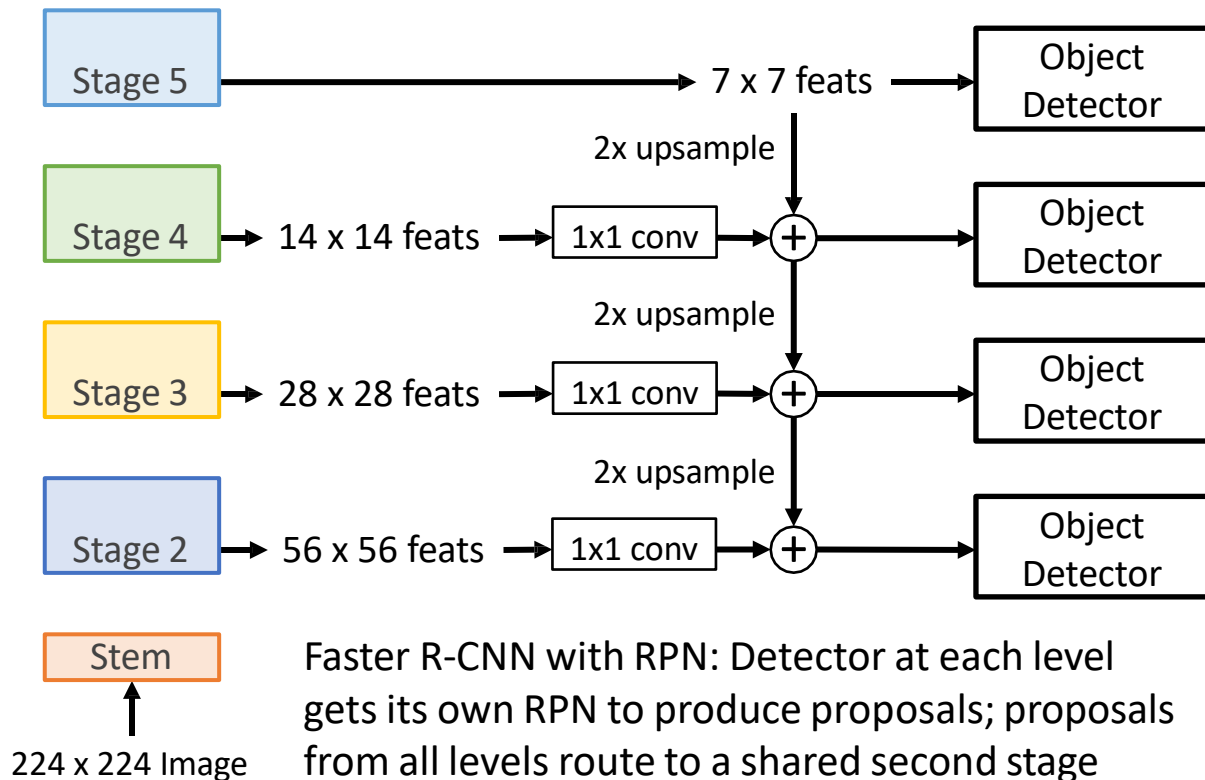
Efficient multiscale features where all levels benefit from the whole backbone! Widely used in practice

Lin et al, "Feature Pyramid Networks for Object Detection", ICCV 2017



Faster R-CNN with RPN: Detector at each level gets its own RPN to produce proposals; proposals from all levels route to a shared second stage

# Faster R-CNN: Learnable Region Proposals

Question: Do we really need the second stage?
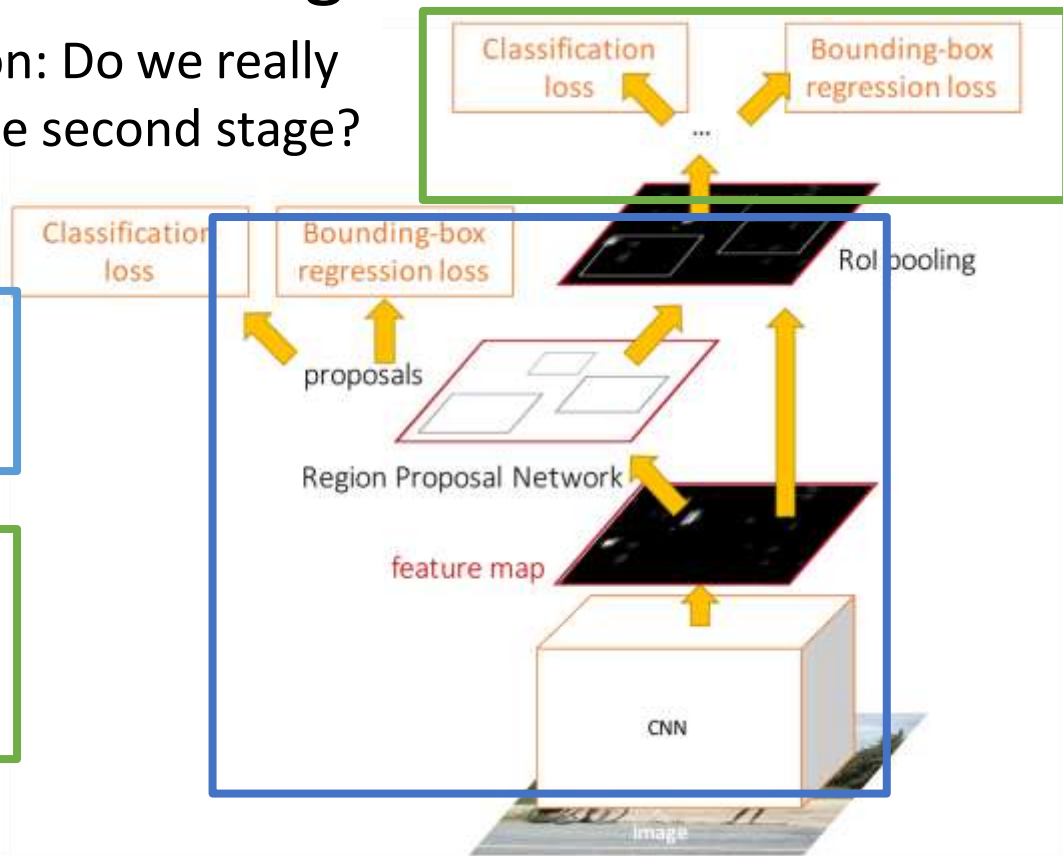
Faster R-CNN is a
**Two-stage object detector**

First stage: Run once per image
- Backbone network
- Region proposal network
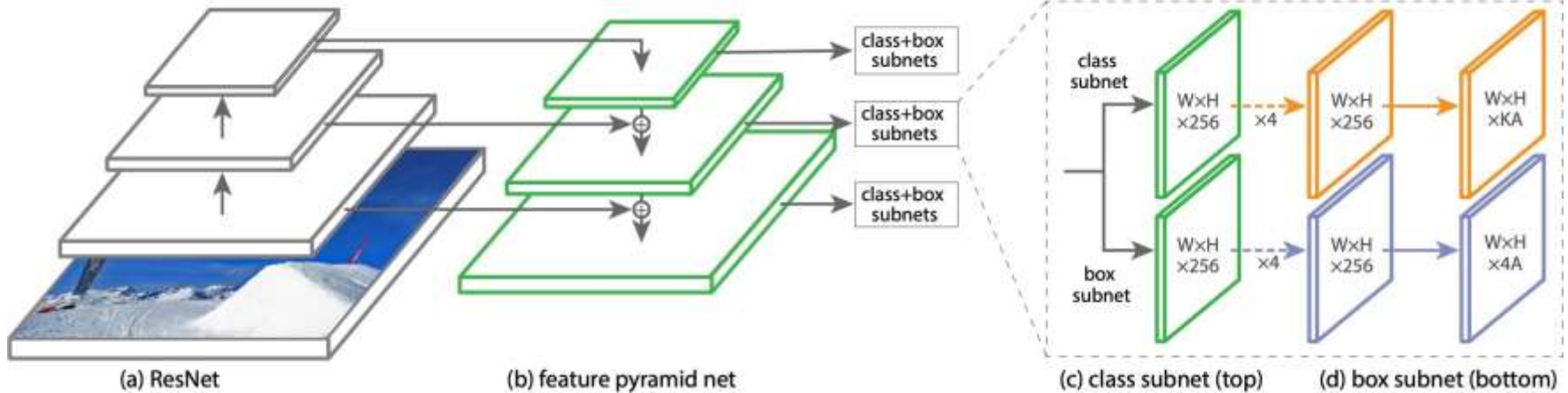
Second stage: Run once per region
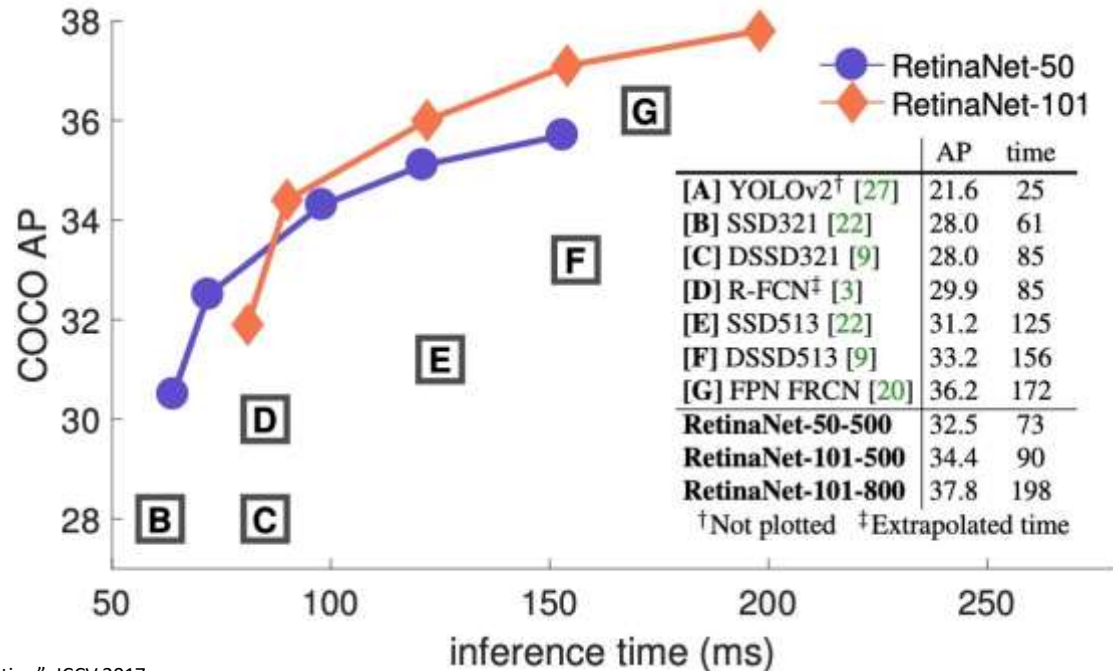- Crop features: RoI pool / align
- Predict object class
- Prediction bbox offset

# Single-Stage Detectors: RetinaNet

In practice, RetinaNet also uses Feature Pyramid Network to handle multiscale



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

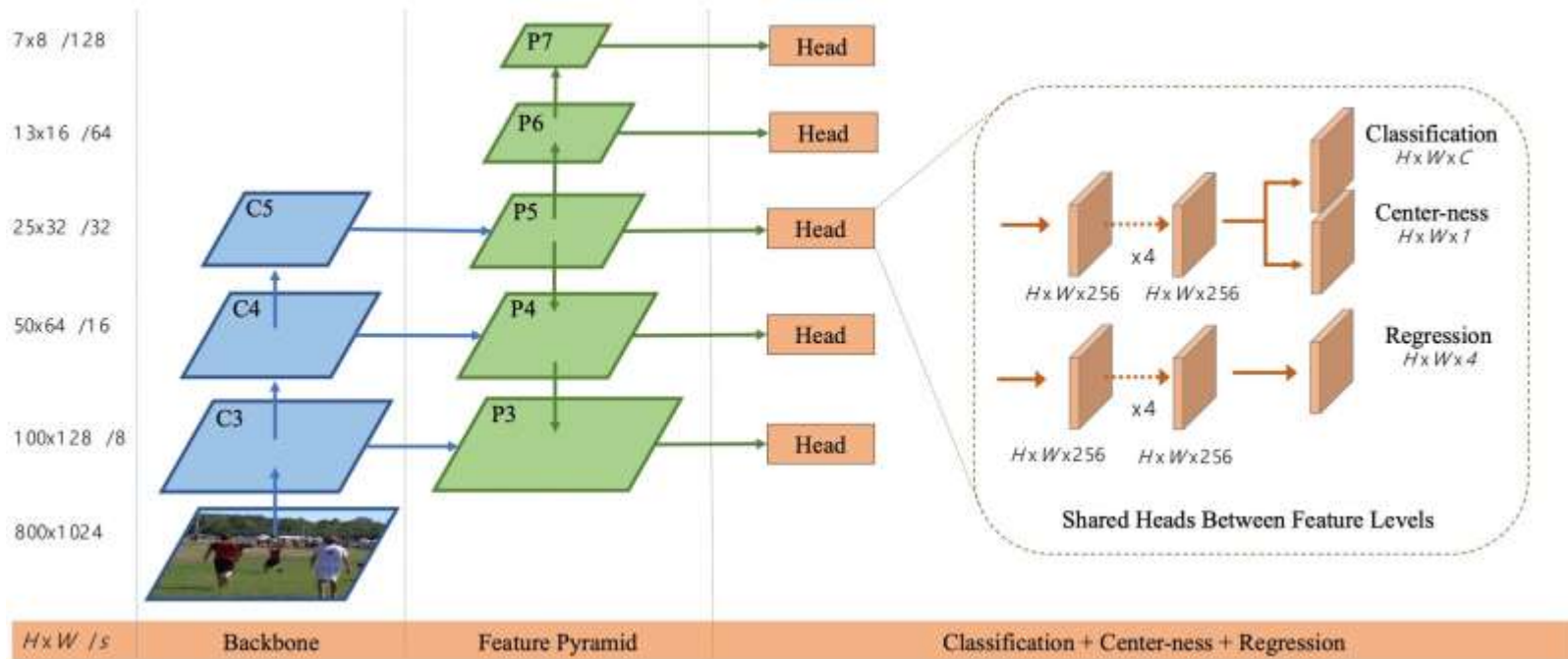Lin et al, "Focal Loss for Dense Object Detection", ICCV 2017

Figure credit: Lin et al, ICCV 2017

# Single-Stage Detectors: RetinaNet

Single-Stage detectors can be much faster than two-stage detectors



Lin et al, "Focal Loss for Dense Object Detection", ICCV 2017

Figure credit: Lin et al, ICCV 2017

# Single-Stage Detectors: FCOS

FCOS also uses a Feature Pyramid Network with heads shared across stages



Tian et al, "FCOS: Fully Convolutional One-Stage Object Detection", ICCV 2019

# Summary

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



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)

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



Bbox Class

Bbox Class

Bbox Class

CNN CNN CNN

Regions of Interest (RoIs) from a proposal method

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

ConvNet

Input image

Category and box transform per region

Per-Region Network

Crop + Resize features

Image features

Run whole image through ConvNet

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



**Single-Stage**: Fully convolutional detector



With anchors: RetinaNet
Anchor-Free: FCOS