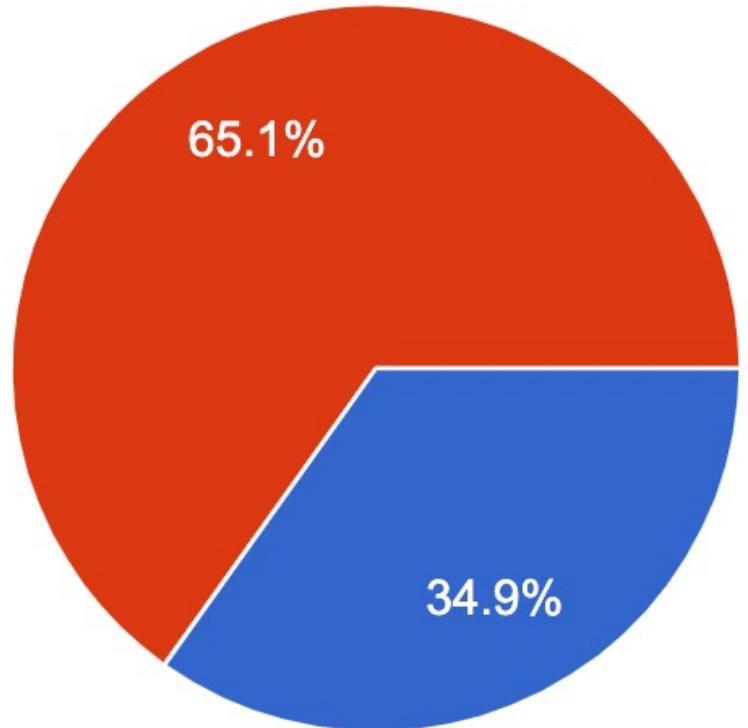


# Lecture 14: Object Detectors

# Poll Results



- Option 1: Keep mini-project, only 1.5 weeks between each of HW4, HW5, HW6, and project
- Option 2: Cancel mini-project, allowing for 2 weeks between each of HW4, HW5, and HW6

Many comments / suggestions in comments and on Piazza:

- Option 2: Want more weight on HW4-6, less on midterm
- Optional project
- Drop one HW assignment
- Extra late days

# Decision

- We will keep 2-week gap between each of HW4-6
- Students can also complete a project if they wish (spec out next week)
- Each student can choose one of the following options:

# Decision

- We will keep 2-week gap between each of HW4-6
- Students can also complete a project if they wish (spec out next week)
- Each student can choose one of the following options:

## **Option A:**

Do all assignments,  
Do not do project.

Grading scheme:

HW1-3: 12%

Midterm: 22%

HW4-6: 14%

# Decision

- We will keep 2-week gap between each of HW4-6
- Students can also complete a project if they wish (spec out next week)
- Each student can choose one of the following options:

## Option A:

Do all assignments,  
Do not do project.

## Option B:

Do 5 or 6 assignments  
Do project

### Grading scheme:

HW1-3: 12%  
Midterm: 22%  
HW4-6: 14%

### Grading scheme (whichever gives you better grade):

HW1-3: 12%  
Midterm: 22%  
HW4-6: 14%  
Project: Replaces lowest HW

Original grading scheme:  
HW1-6: 10%  
Midterm: 20%  
Project: 20%

# Decision

In addition: Everyone gets +3 late days  
(cannot be applied to A6 or project)

- We will keep 2-week gap between each of HW4-6
- Students can also complete a project if they wish (spec out next week)
- Each student can choose one of the following options:

## Option A:

Do all assignments,  
Do not do project.

## Option B:

Do 5 or 6 assignments  
Do project

### Grading scheme:

HW1-3: 12%

Midterm: 22%

HW4-6: 14%

### Grading scheme (whichever gives you better grade):

HW1-3: 12%

Midterm: 22%

HW4-6: 14%

Project: Replaces lowest HW

Original grading scheme:

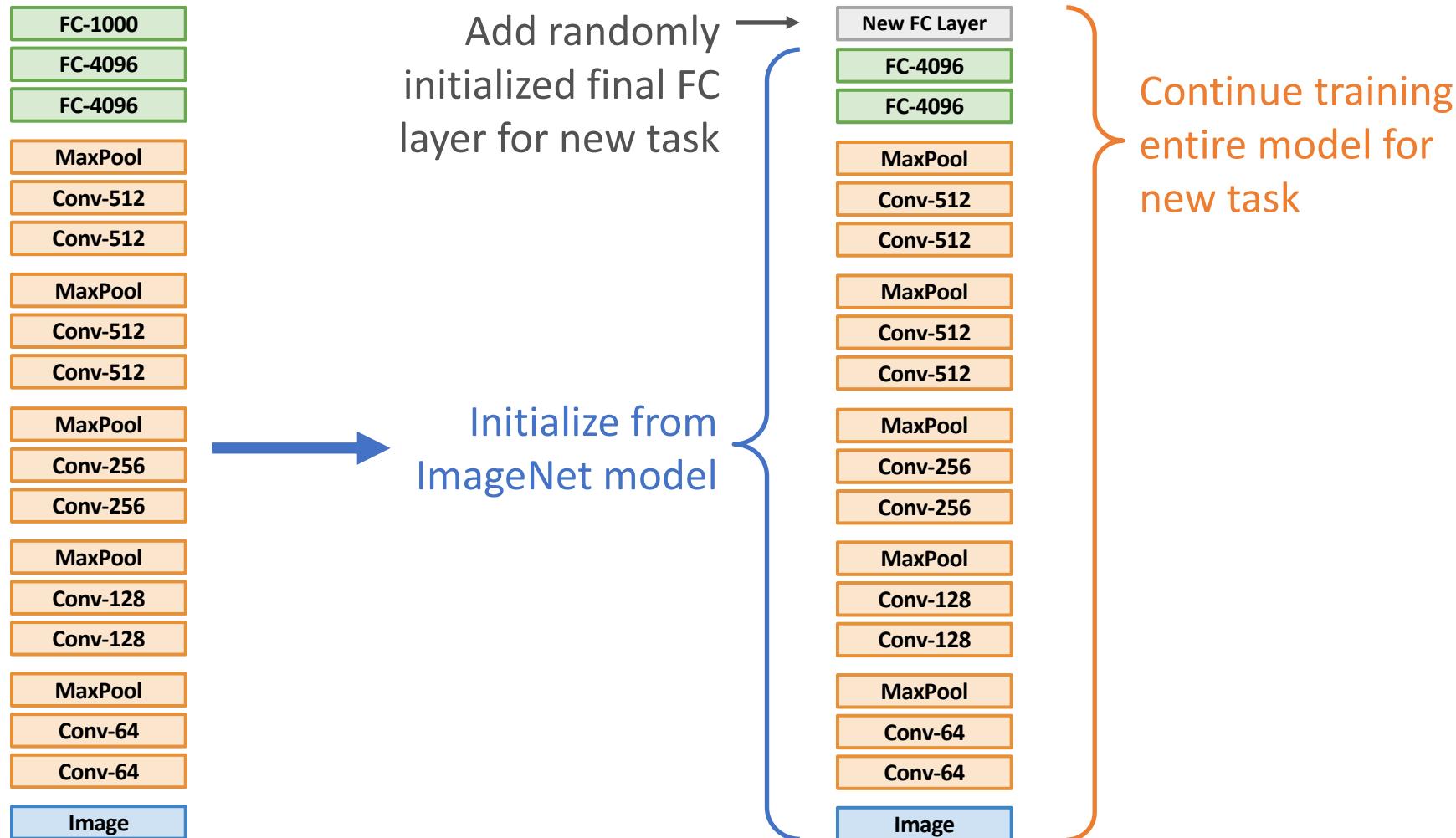
HW1-6: 10%

Midterm: 20%

Project: 20%

# Last Time: Transfer Learning

# 1. Train on ImageNet



# Last Time: Localization Tasks

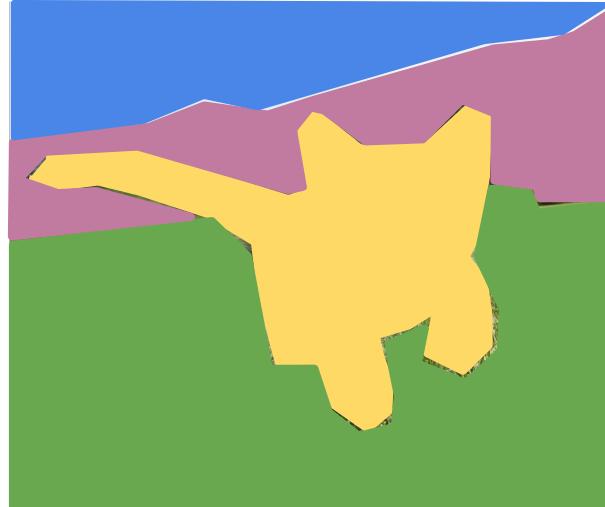
**Classification**



CAT

No spatial extent

**Semantic Segmentation**



GRASS, CAT, TREE,  
SKY

No objects, just pixels

**Object Detection**



DOG, DOG, CAT

Multiple Objects

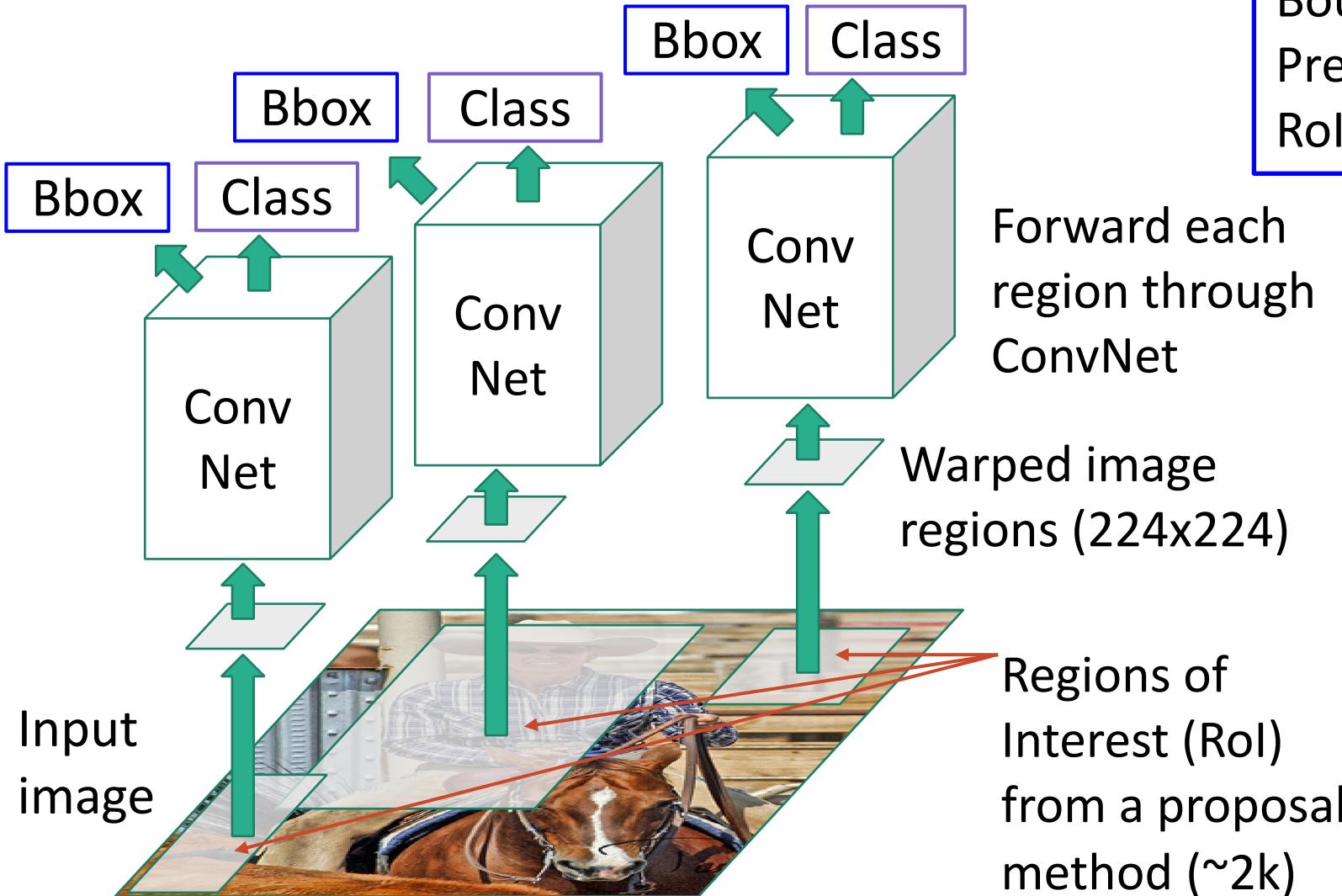
**Instance Segmentation**



DOG, DOG, CAT

[This image is CC0 public domain](#)

# Last Time: R-CNN



Classify each region

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

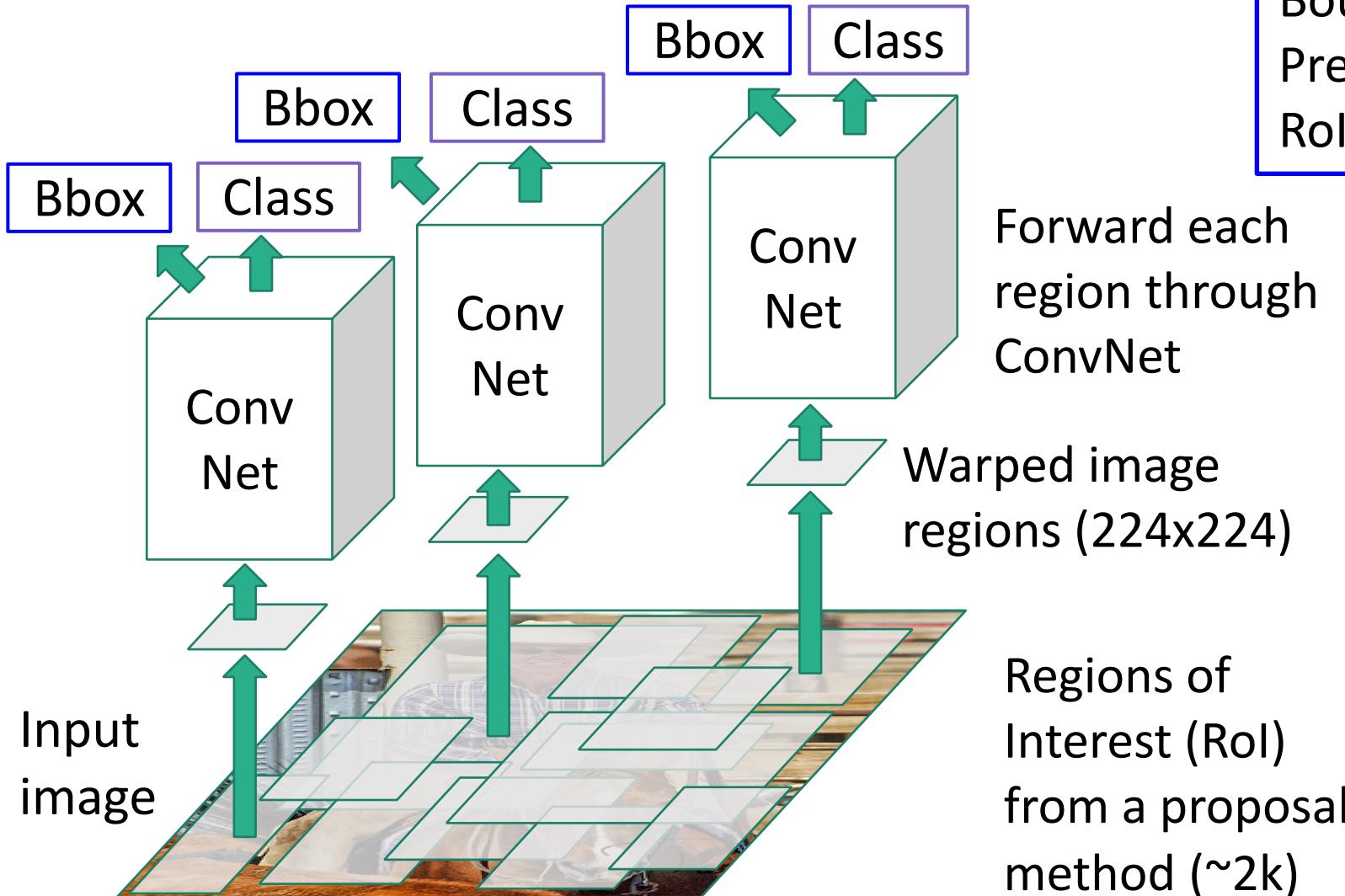
Forward each  
region through  
ConvNet

Warped image  
regions (224x224)

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

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



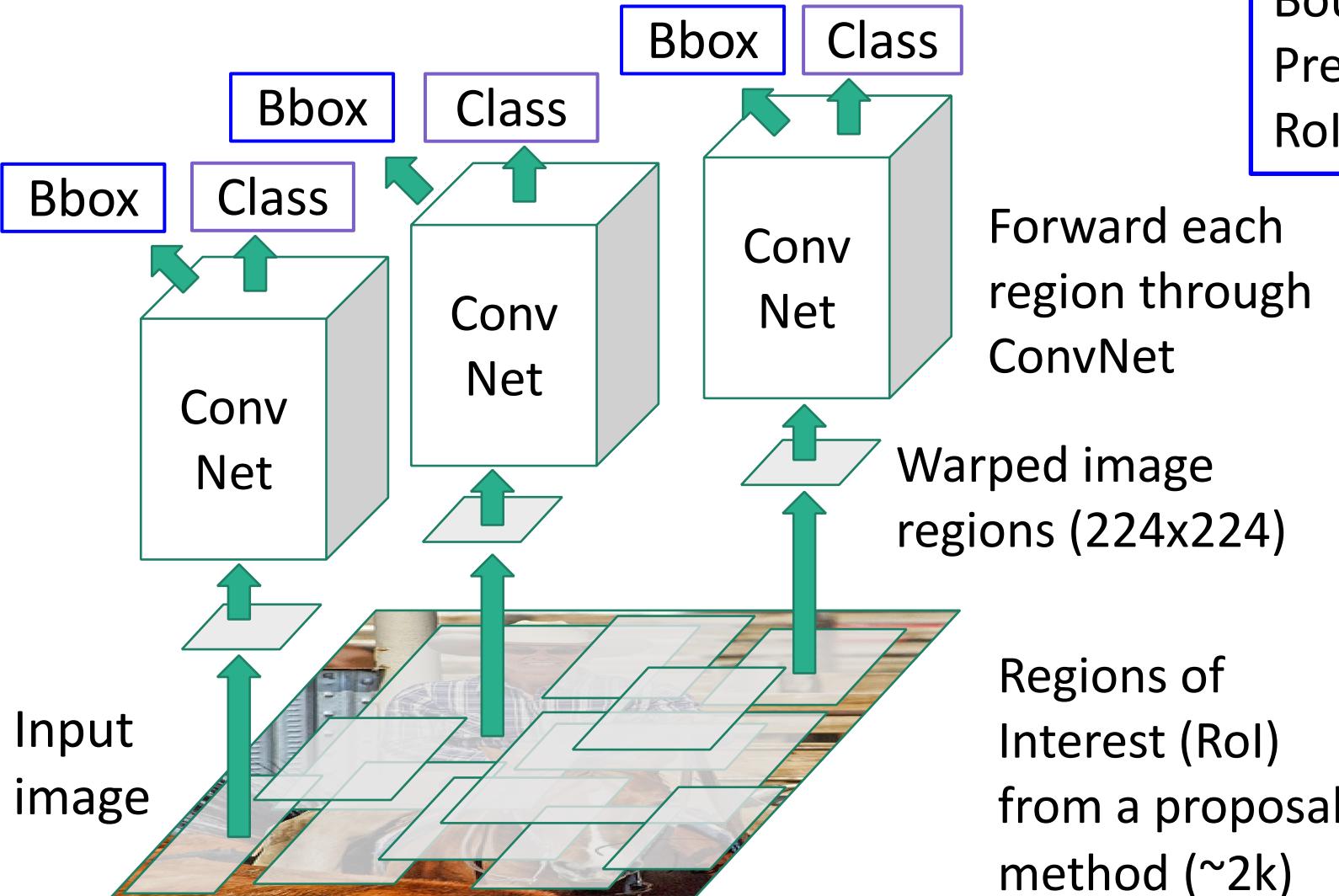
Classify each region

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



Classify each region

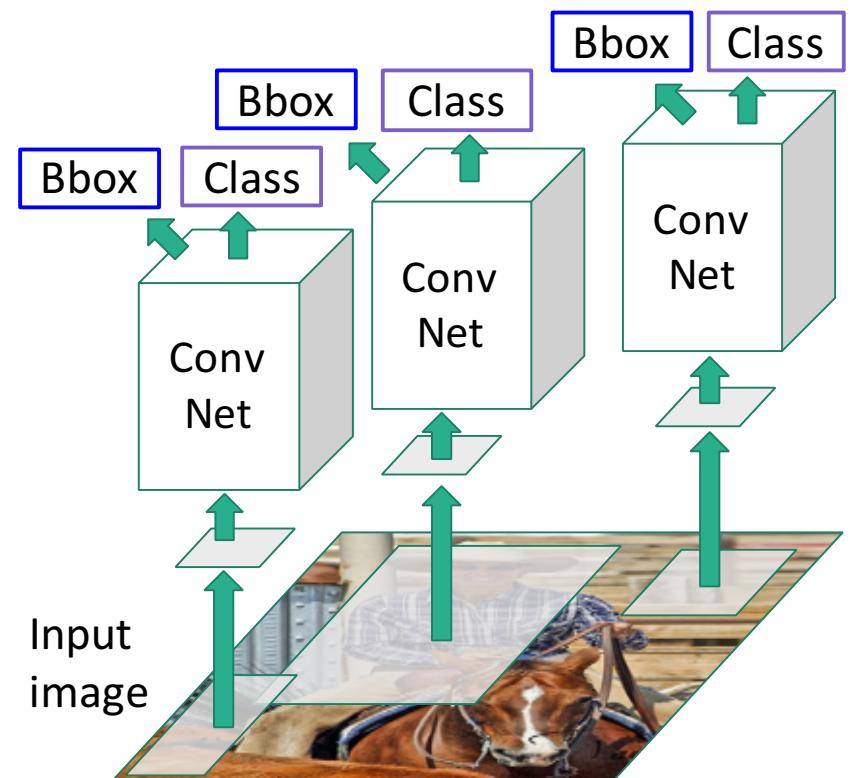
Bounding box regression:  
Predict “transform” to correct the  
RoI: 4 numbers ( $t_x, t_y, t_h, t_w$ )

Problem: Very slow! Need to do 2000 forward passes through CNN per image

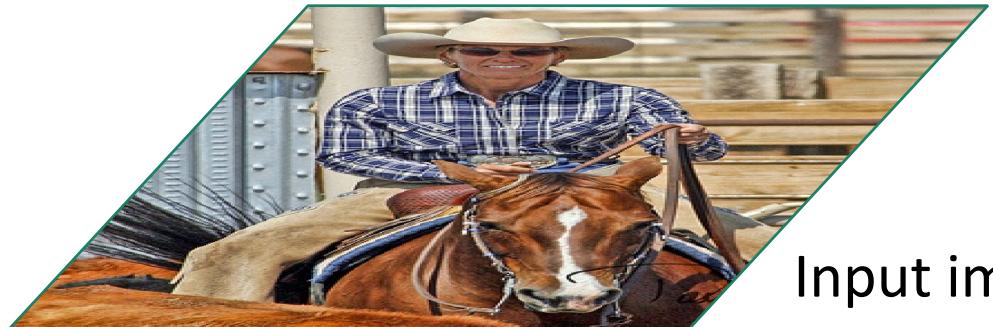
Idea: Overlapping proposals cause a lot of repeated work: same pixels processed many times. Can we avoid this?

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

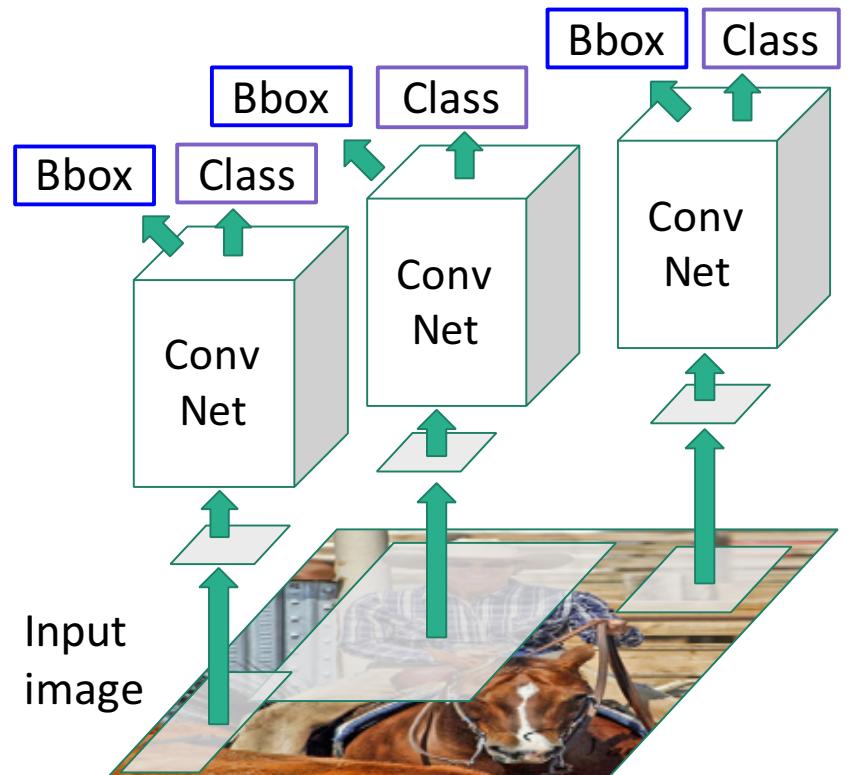


# Fast R-CNN



Input image

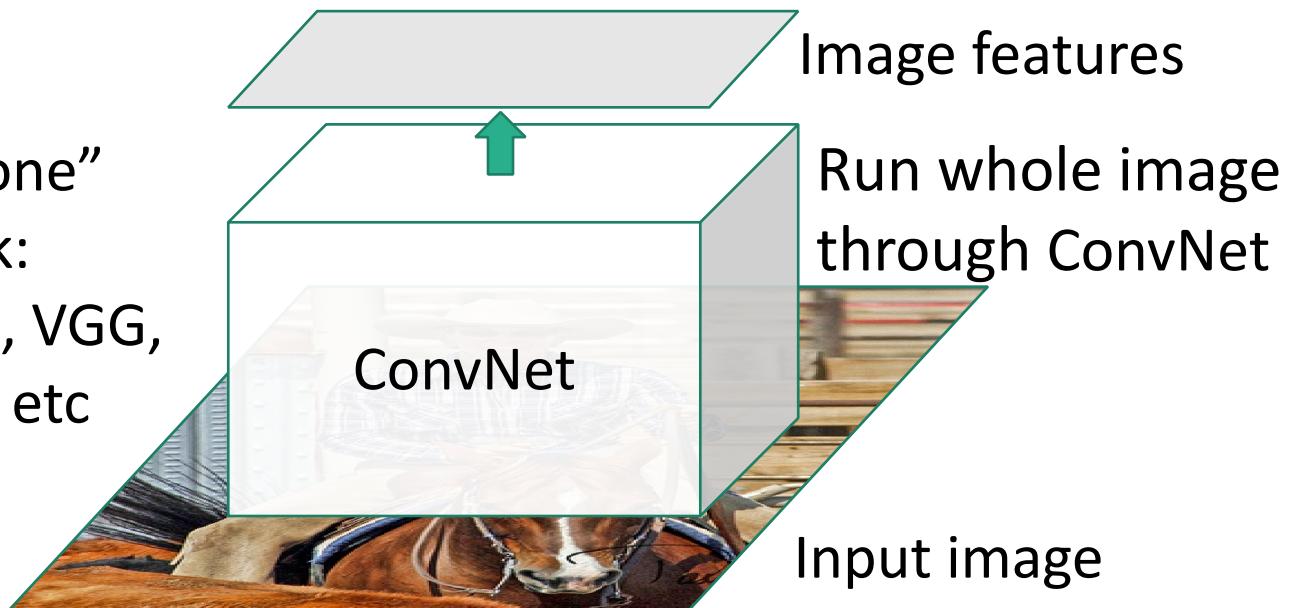
“Slow” R-CNN  
Process each region  
independently



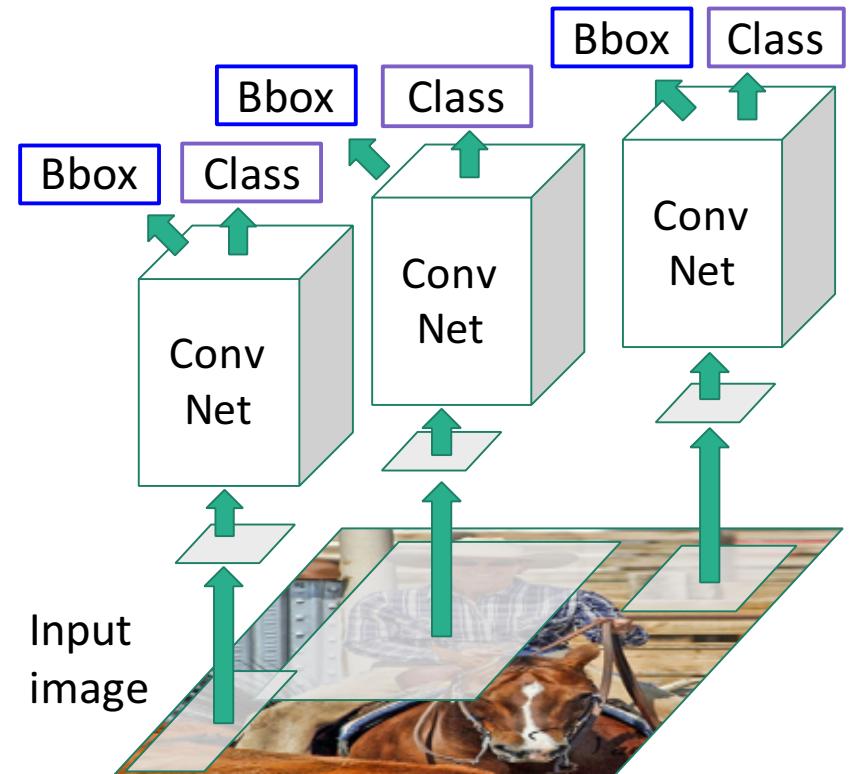
Girshick, “Fast R-CNN”, ICCV 2015. Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.

# Fast R-CNN

“Backbone” network:  
AlexNet, VGG,  
ResNet, etc



“Slow” R-CNN  
Process each region  
independently

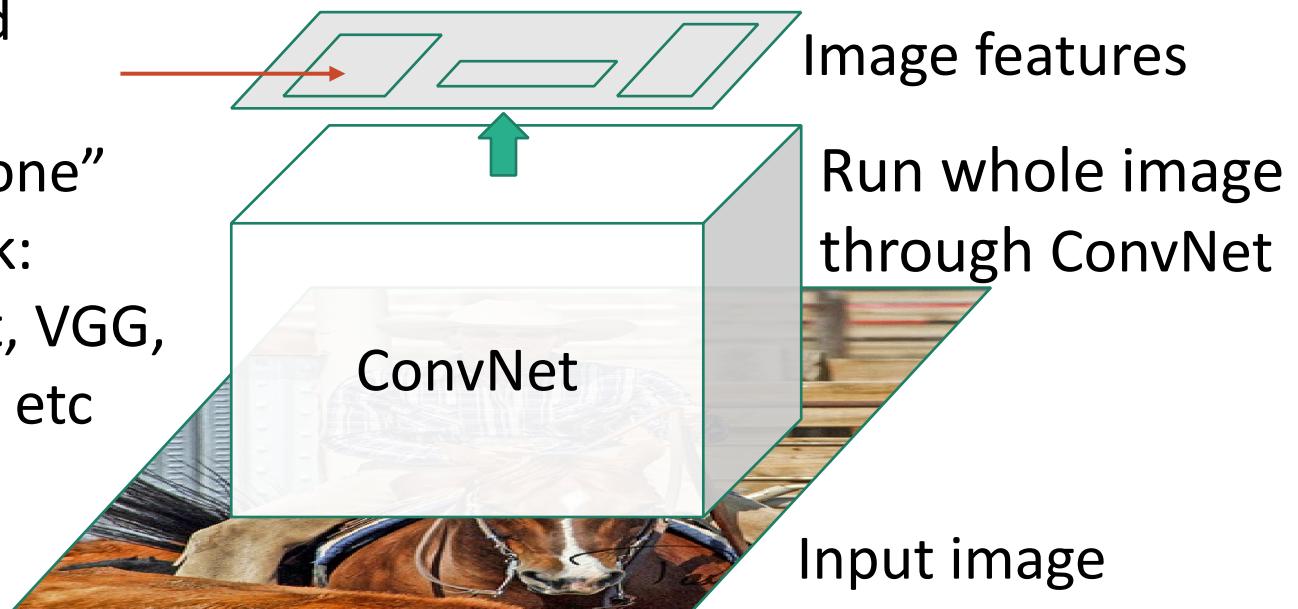


Girshick, “Fast R-CNN”, ICCV 2015. Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.

# Fast R-CNN

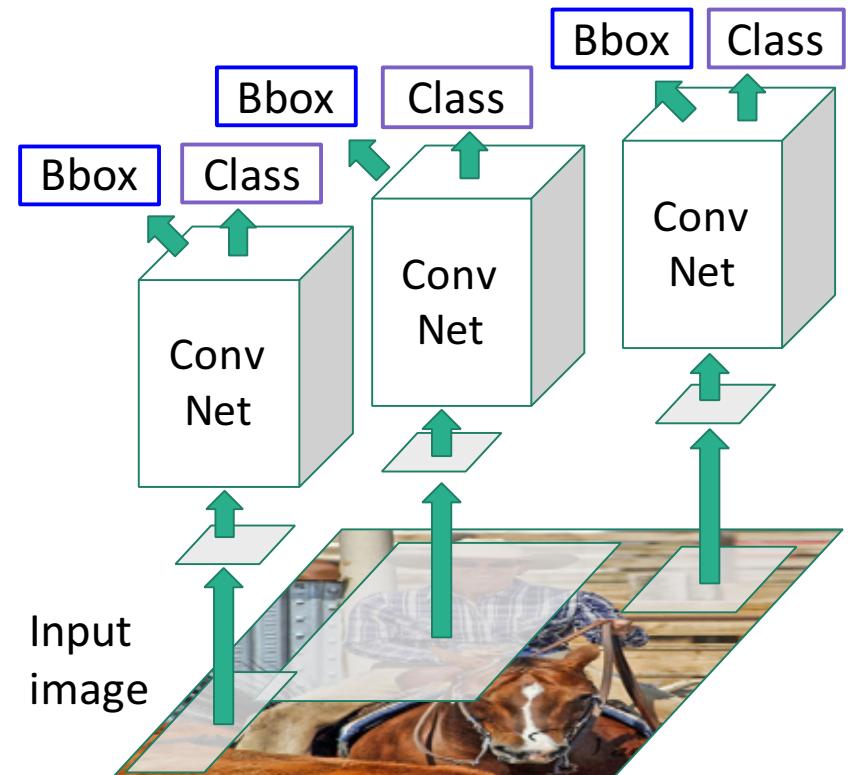
Regions of  
Interest (RoIs)  
from a proposal  
method

“Backbone”  
network:  
AlexNet, VGG,  
ResNet, etc



## “Slow” R-CNN

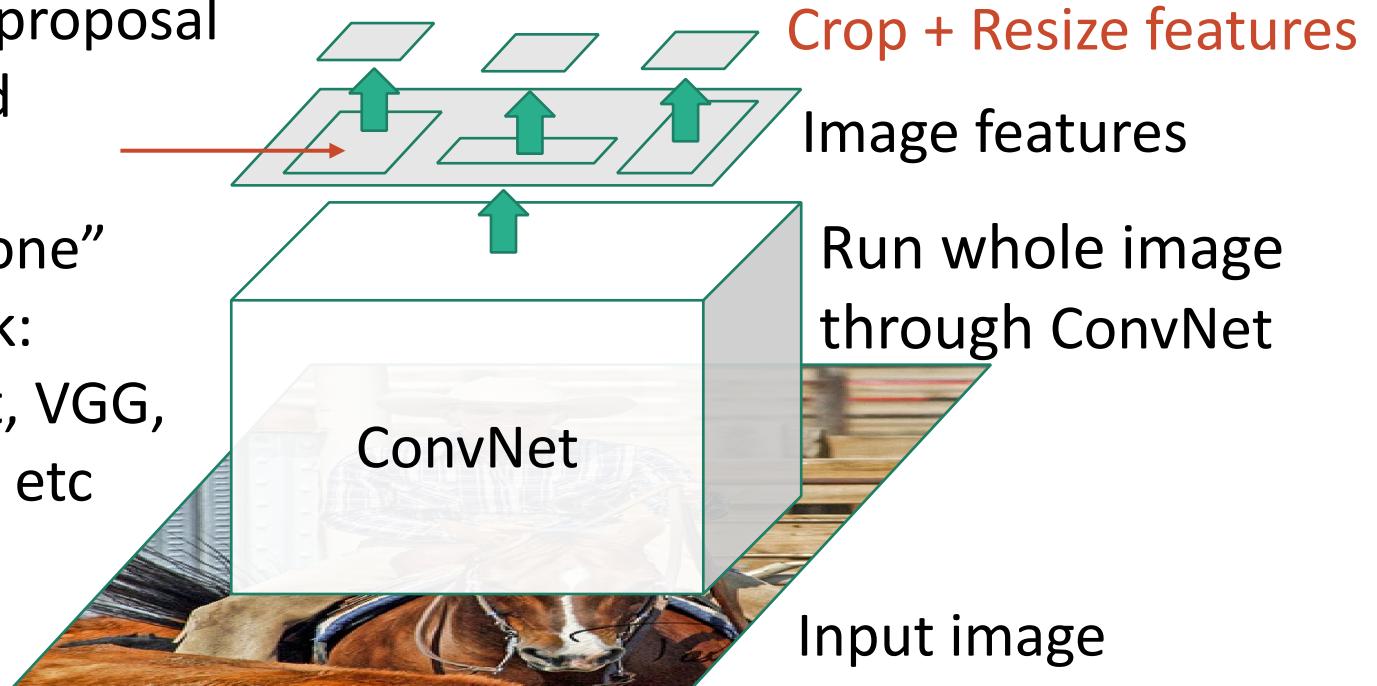
Process each region  
independently



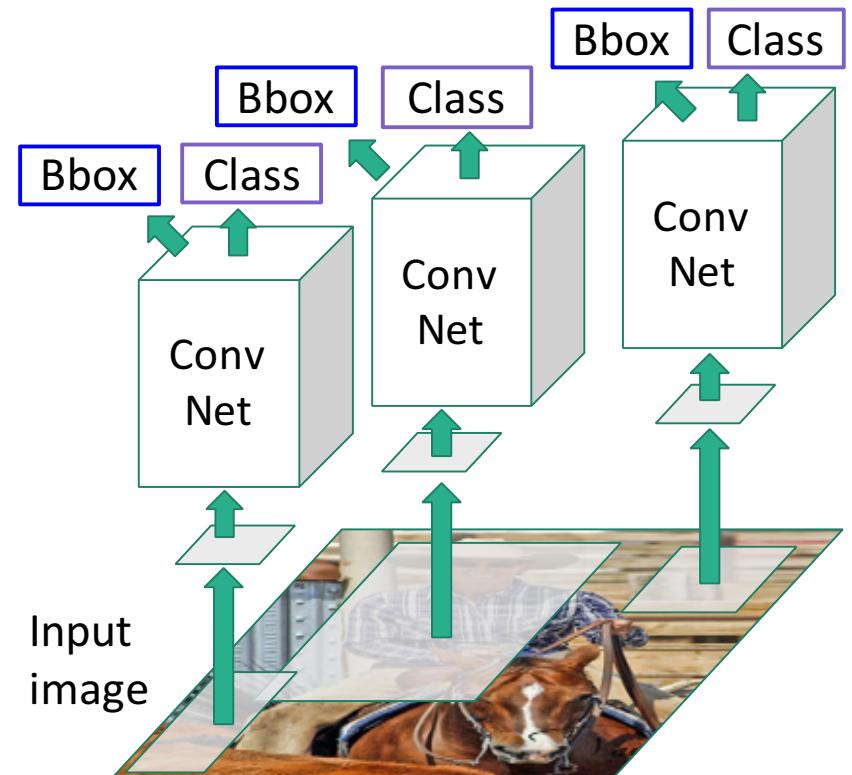
# Fast R-CNN

Regions of  
Interest (RoIs)  
from a proposal  
method

“Backbone”  
network:  
AlexNet, VGG,  
ResNet, etc



“Slow” R-CNN  
Process each region  
independently

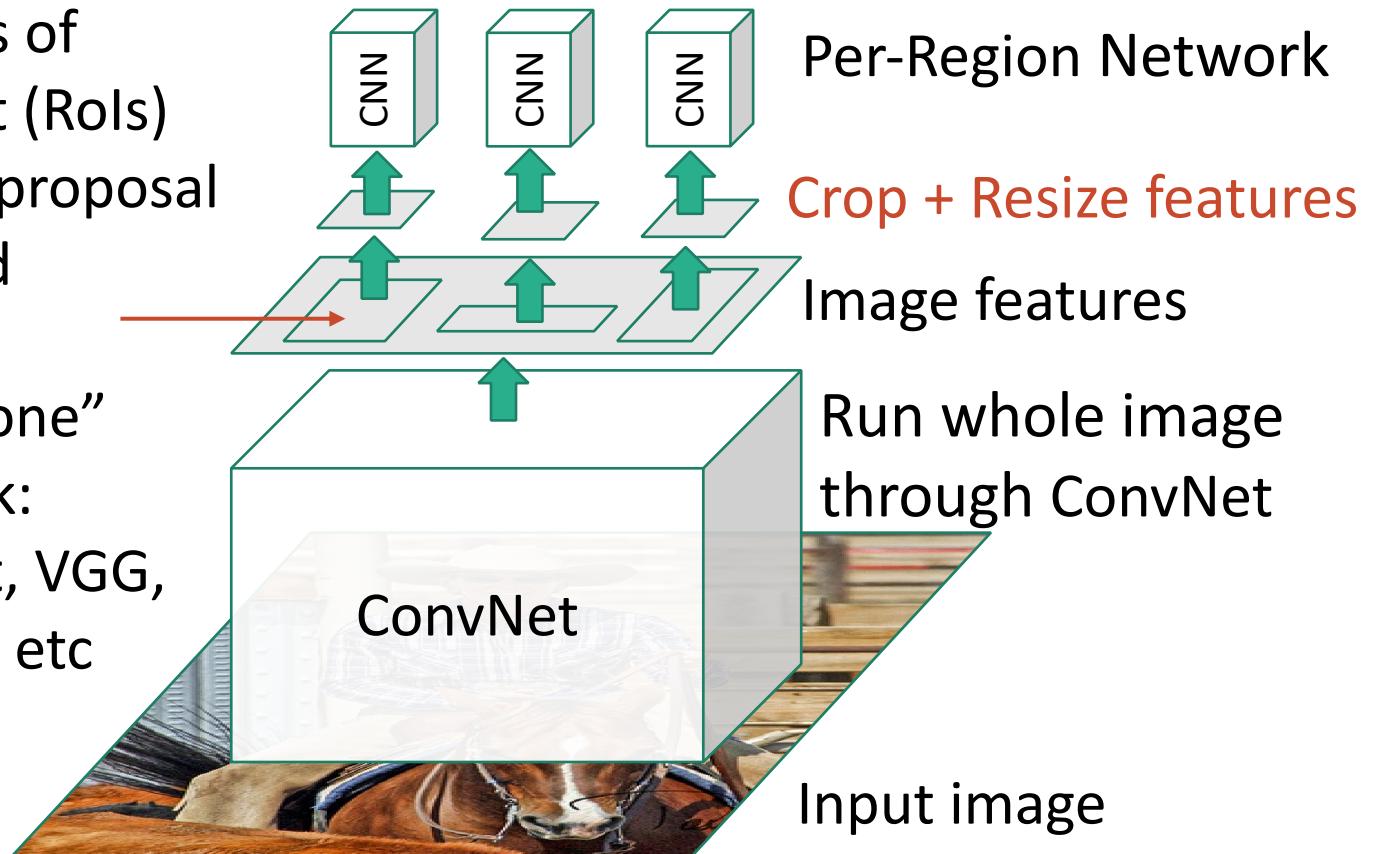


Girshick, “Fast R-CNN”, ICCV 2015. Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.

# Fast R-CNN

Regions of Interest (Rois)  
from a proposal  
method

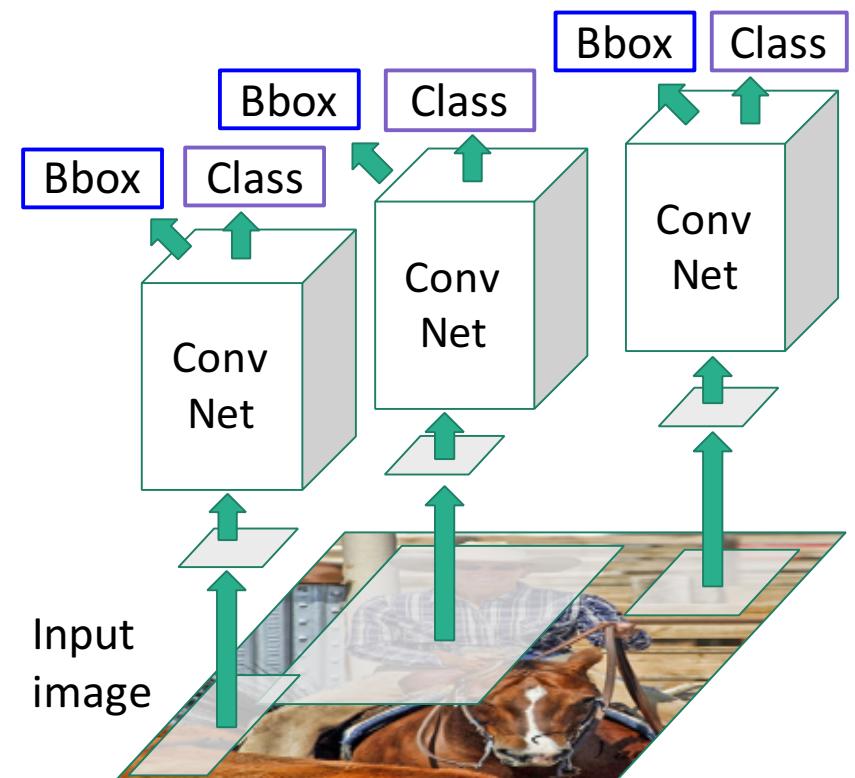
“Backbone”  
network:  
AlexNet, VGG,  
ResNet, etc



Lecture 13 - 17

## “Slow” R-CNN

Process each region  
independently



Girshick, “Fast R-CNN”, ICCV 2015. Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.

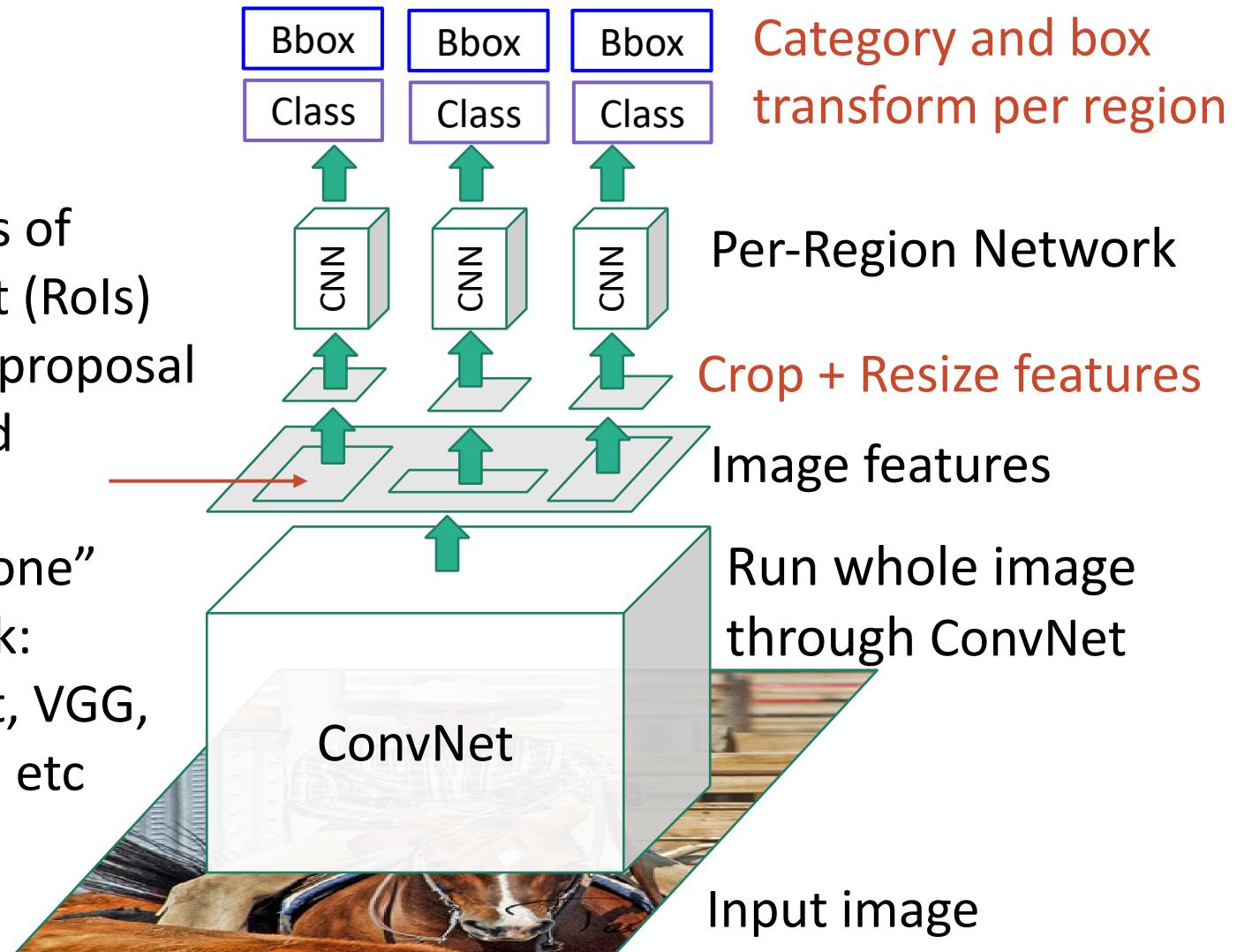
Justin Johnson

March 7, 2022

# Fast R-CNN

Regions of Interest (Rois)  
from a proposal  
method

“Backbone”  
network:  
AlexNet, VGG,  
ResNet, etc



Category and box  
transform per region

Per-Region Network

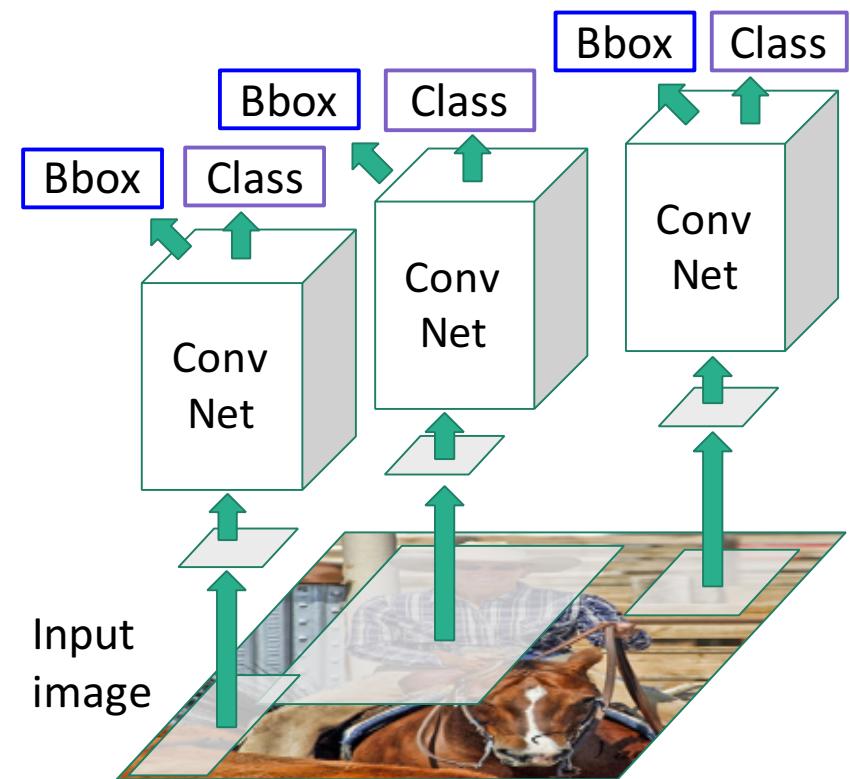
Crop + Resize features

Image features

Run whole image  
through ConvNet

Input image

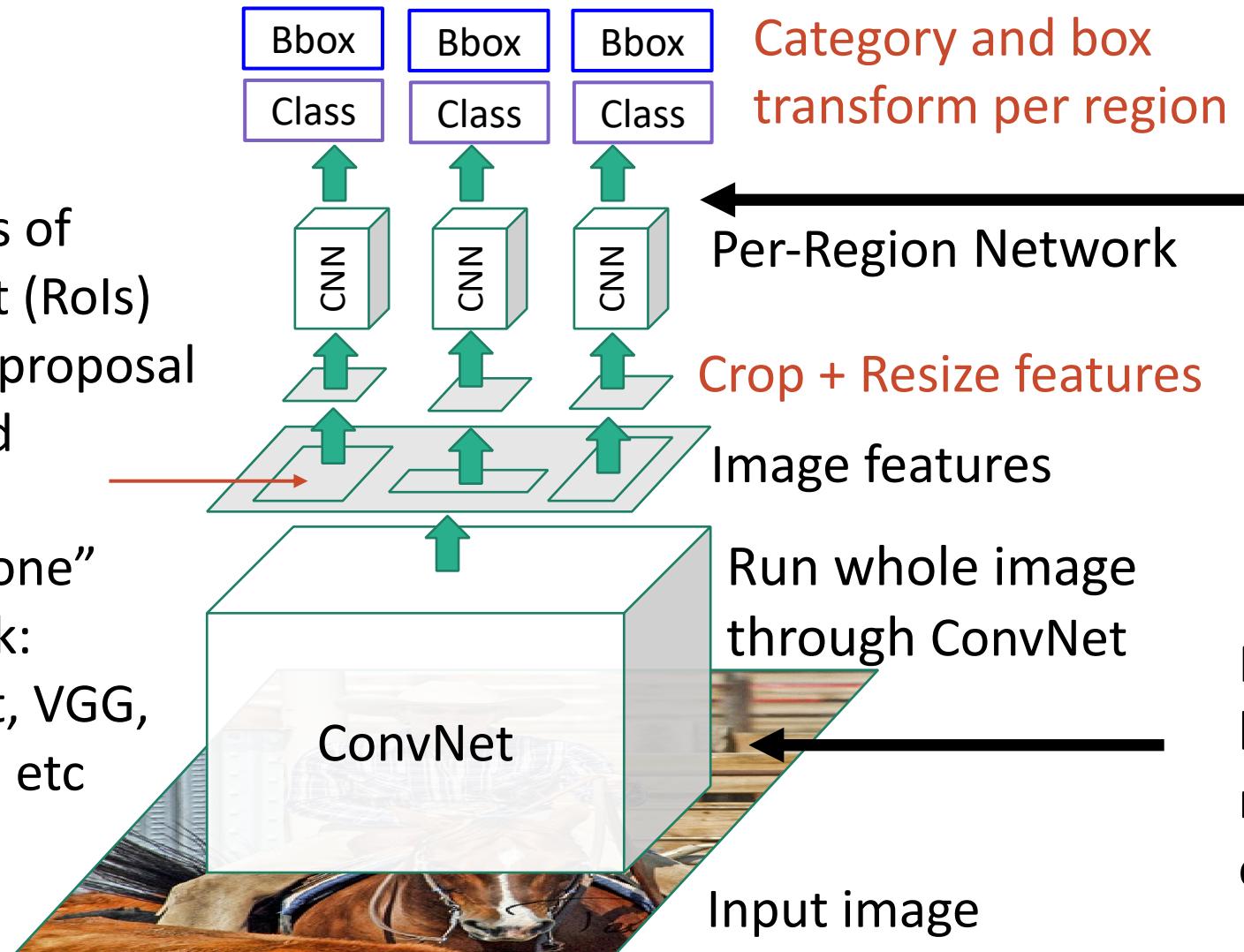
**“Slow” R-CNN**  
Process each region  
independently



# Fast R-CNN

Regions of Interest (Rois)  
from a proposal  
method

“Backbone”  
network:  
AlexNet, VGG,  
ResNet, etc



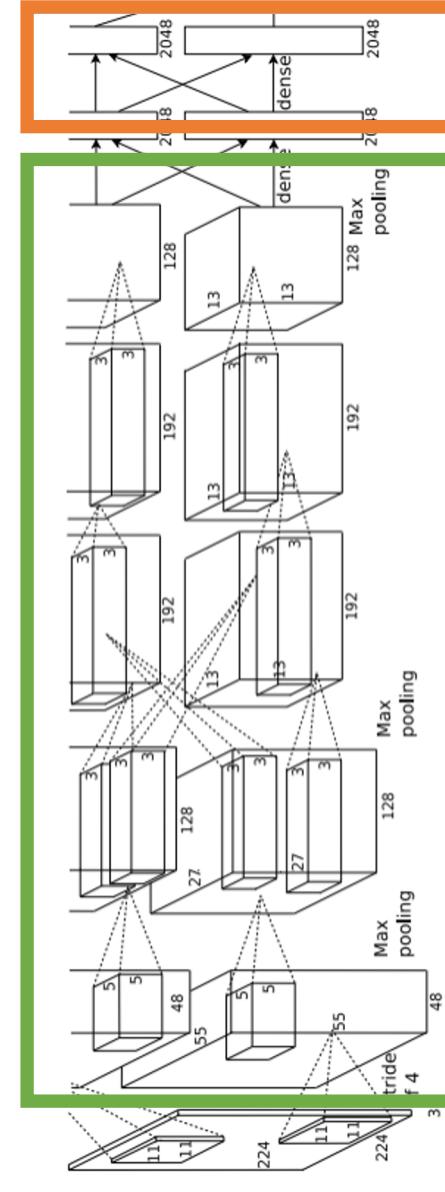
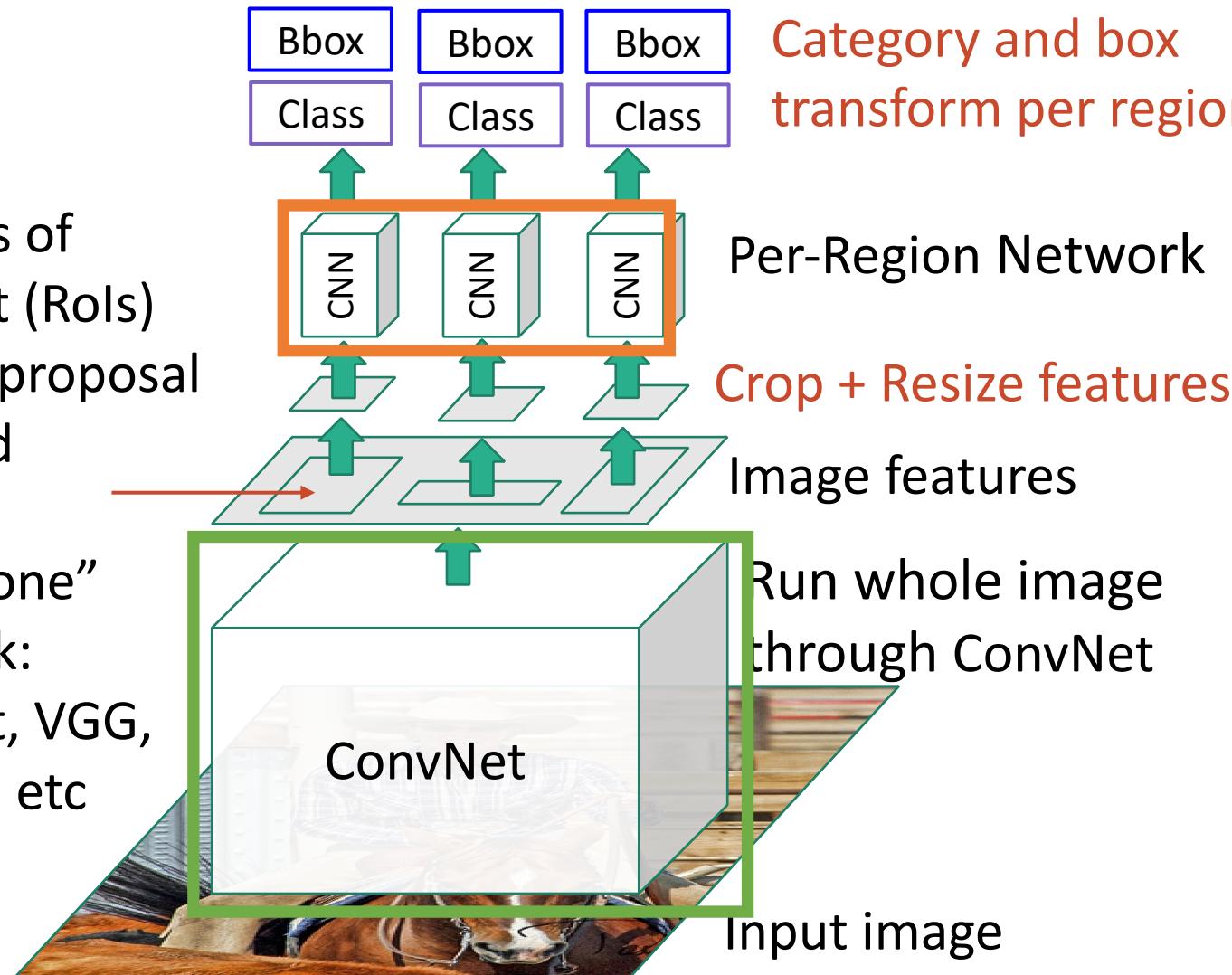
Per-Region network is  
relatively lightweight

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

# Fast R-CNN

Regions of Interest (Rois)  
from a proposal  
method

“Backbone”  
network:  
AlexNet, VGG,  
ResNet, etc

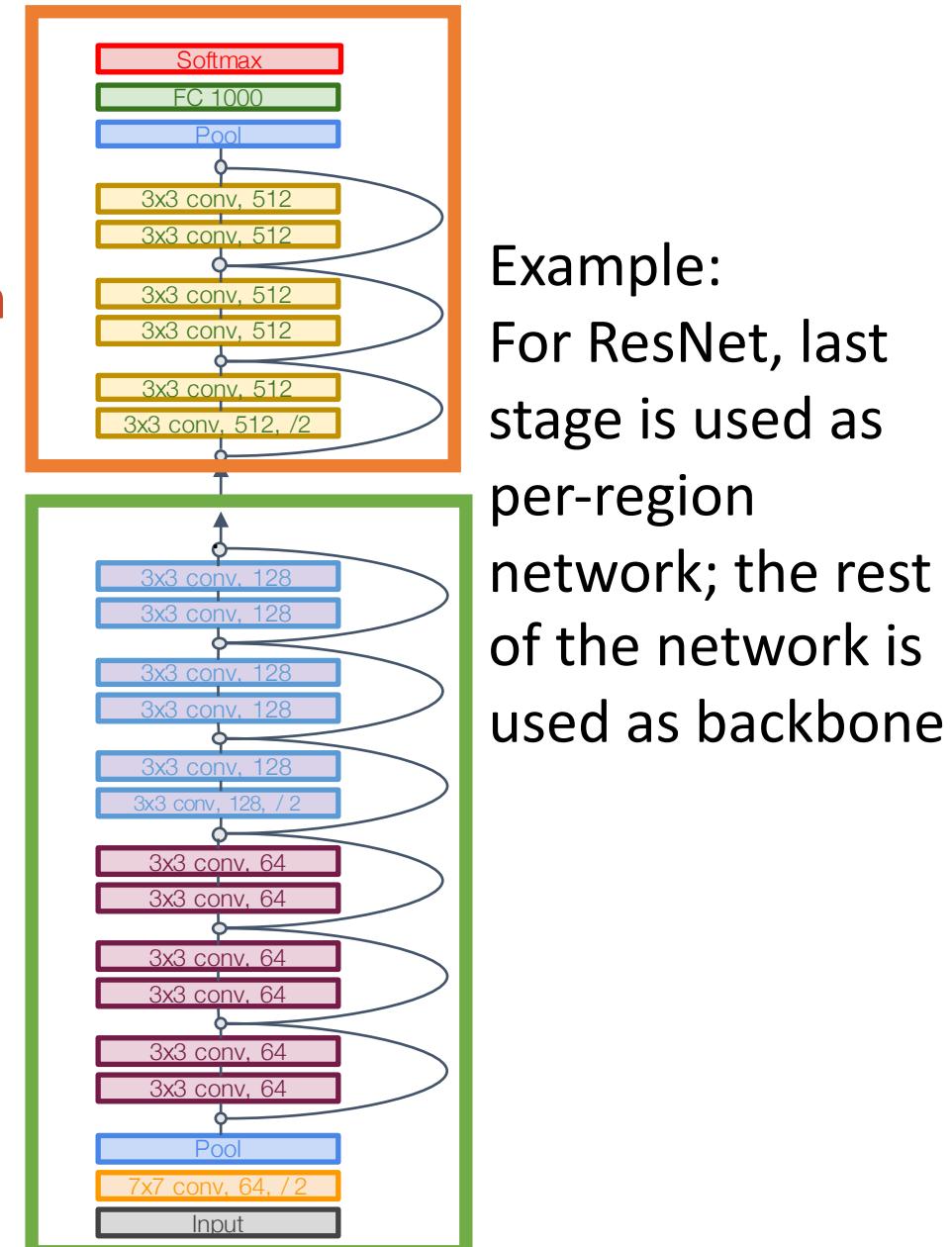
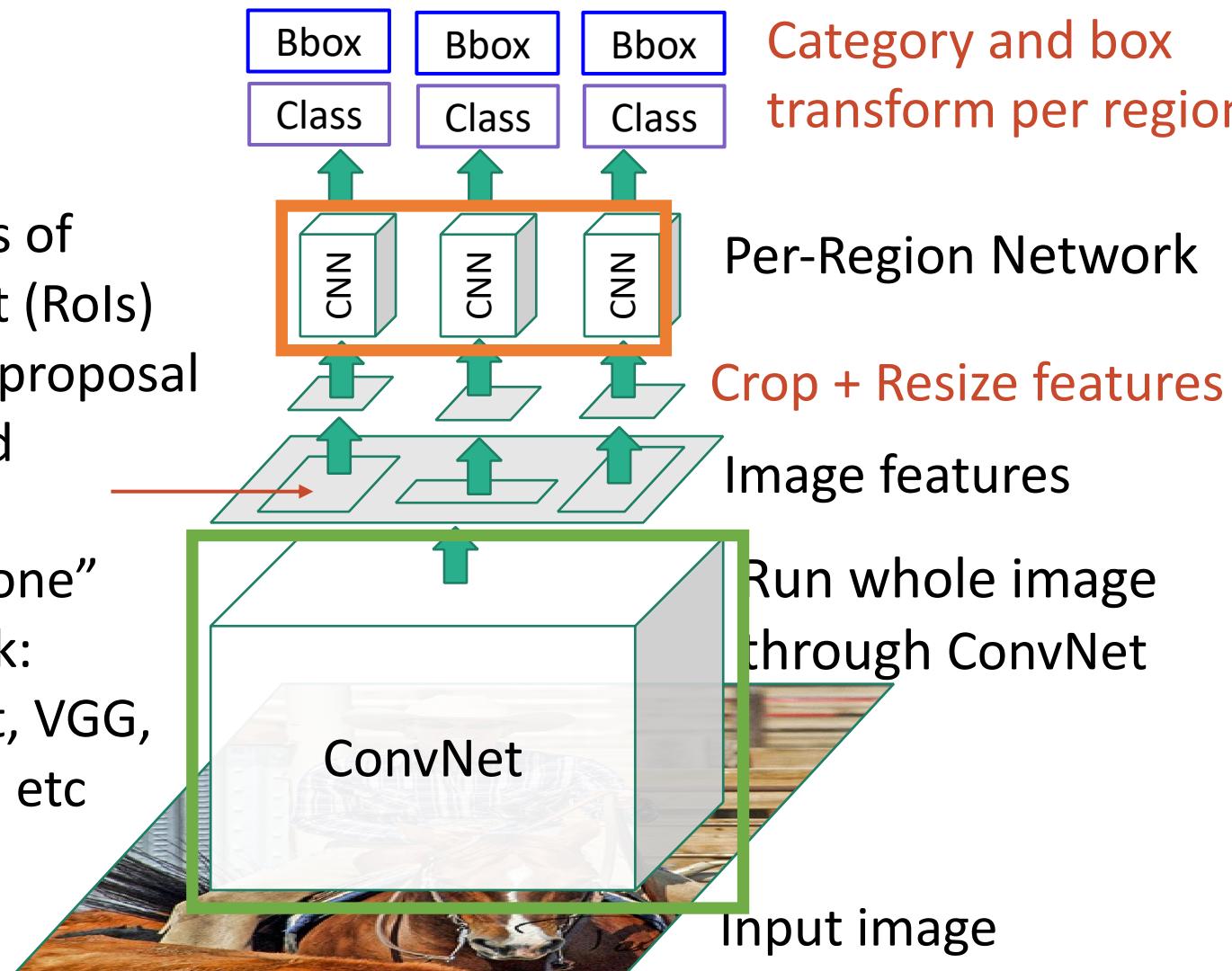


Example:  
When using  
AlexNet for  
detection, five  
conv layers are  
used for  
backbone and  
two FC layers are  
used for per-  
region network

# Fast R-CNN

Regions of Interest (Rois)  
from a proposal  
method

“Backbone”  
network:  
AlexNet, VGG,  
ResNet, etc

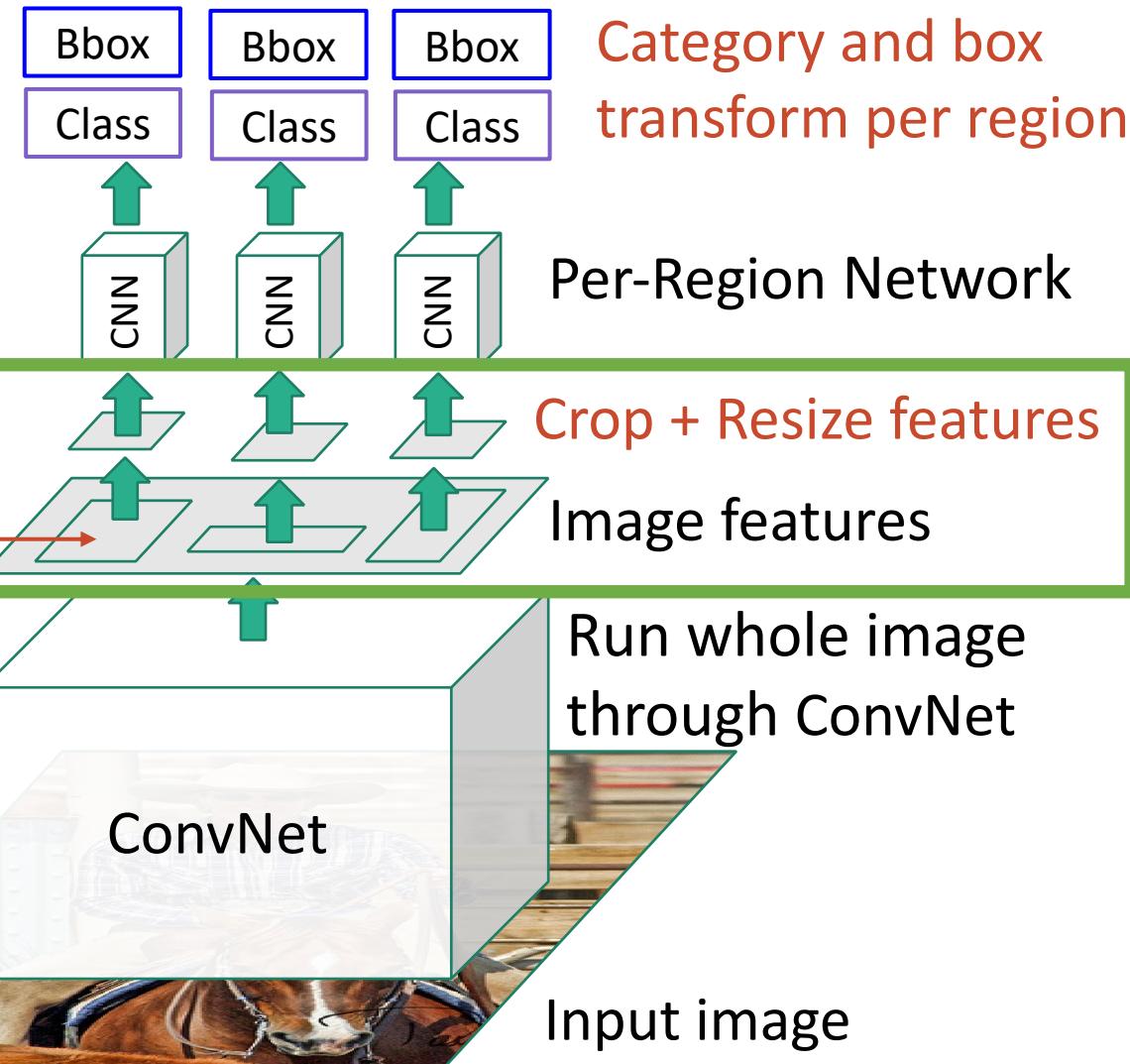


Example:  
For ResNet, last  
stage is used as  
per-region  
network; the rest  
of the network is  
used as backbone

# Fast R-CNN

Regions of Interest (Rois)  
from a proposal  
method

“Backbone”  
network:  
AlexNet, VGG,  
ResNet, etc



Category and box  
transform per region

Per-Region Network

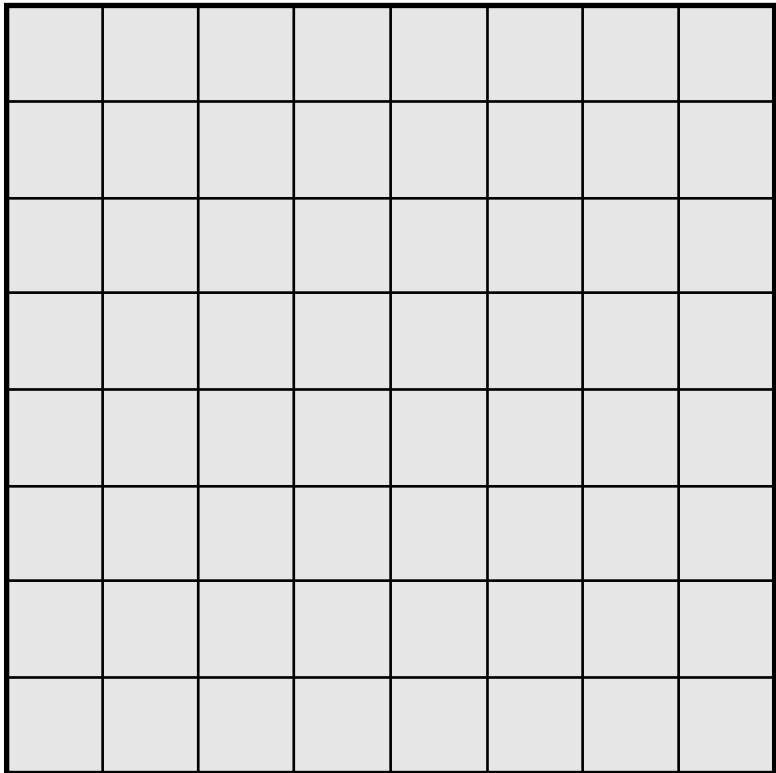
Crop + Resize features

Image features

Run whole image  
through ConvNet

How to crop  
features?

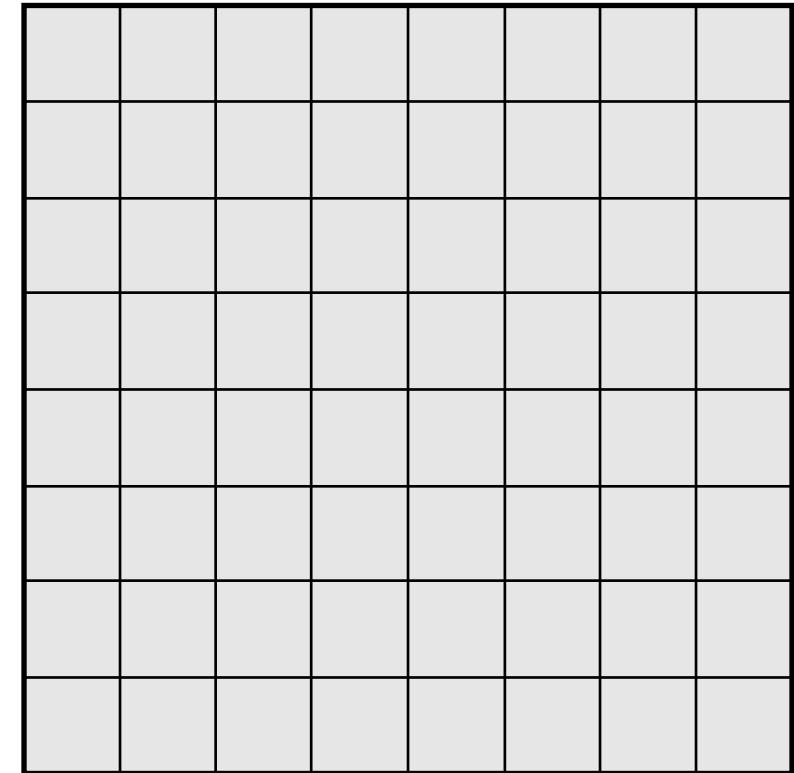
# Recall: Receptive Fields



Input Image: 8 x 8

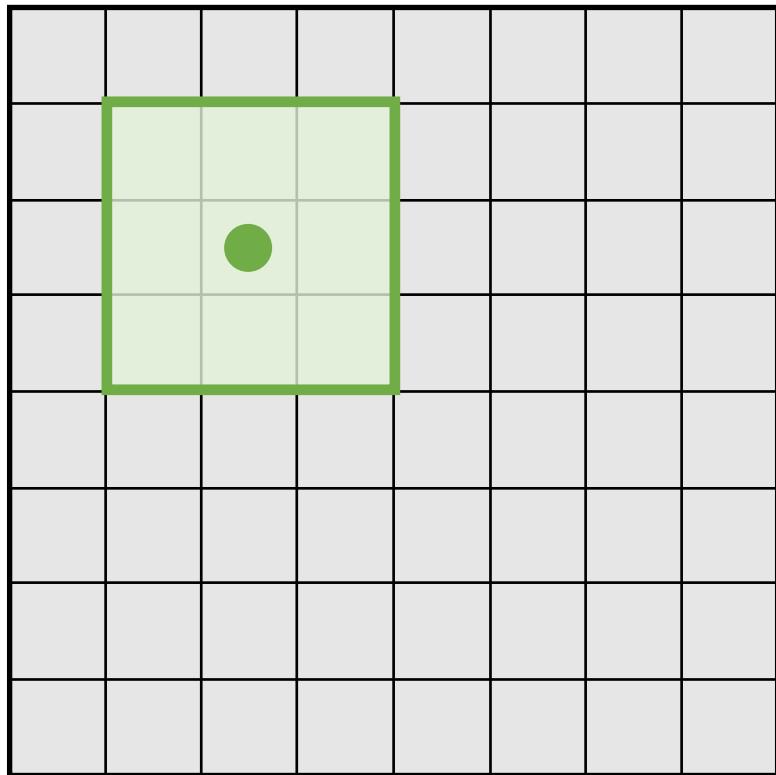
Every position in the output feature map depends on a 3x3 receptive field in the input

3x3 Conv  
Stride 1, pad 1



Output Image: 8 x 8

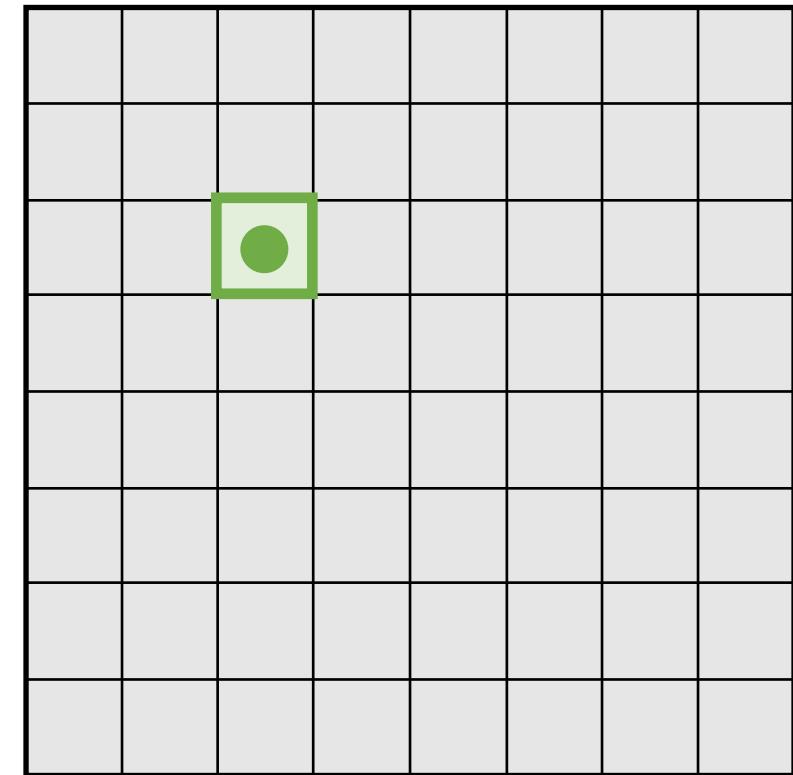
# Recall: Receptive Fields



Input Image: 8 x 8

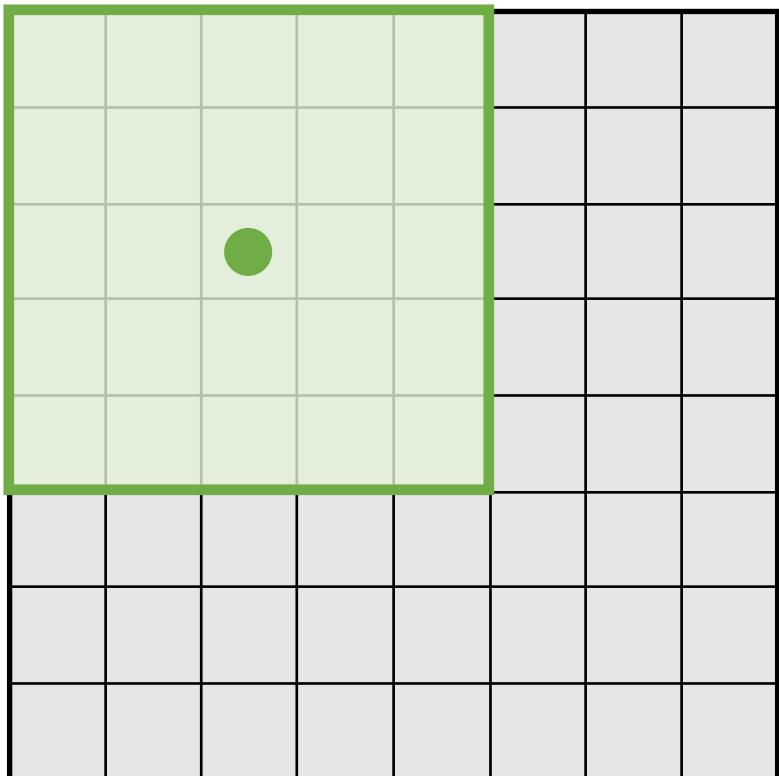
Every position in the output feature map depends on a 3x3 receptive field in the input

3x3 Conv  
Stride 1, pad 1



Output Image: 8 x 8

# Recall: Receptive Fields

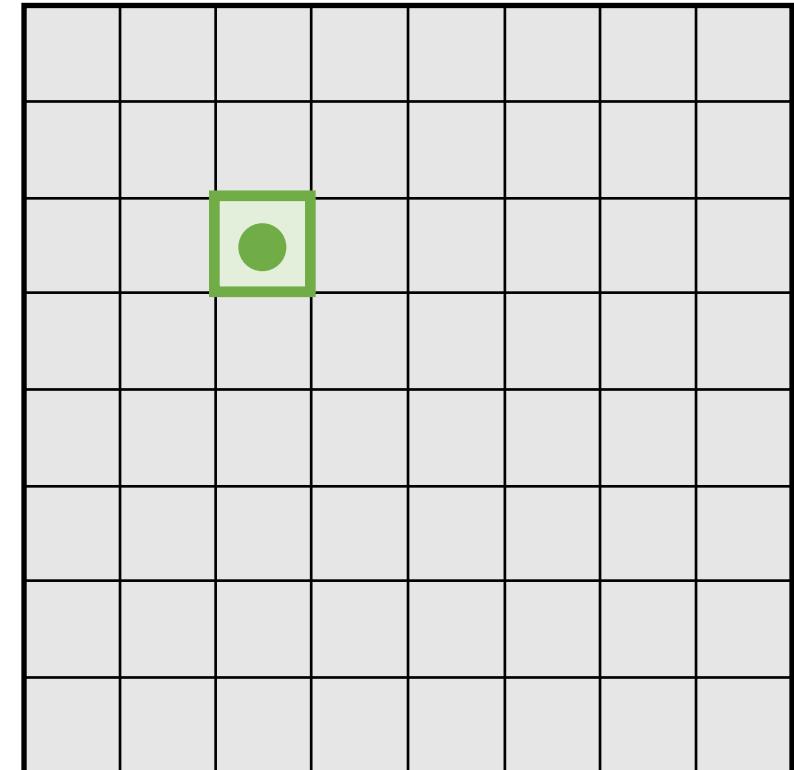


Input Image: 8 x 8

Every position in the output feature map depends on a 5x5 receptive field in the input

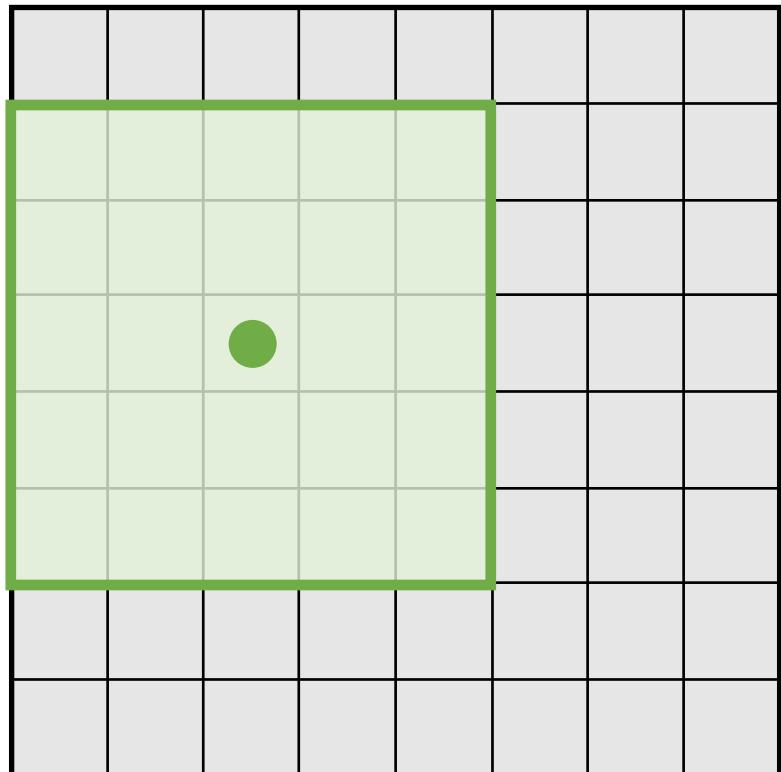
3x3 Conv  
Stride 1, pad 1

3x3 Conv  
Stride 1, pad 1



Output Image: 8 x 8

# Recall: Receptive Fields

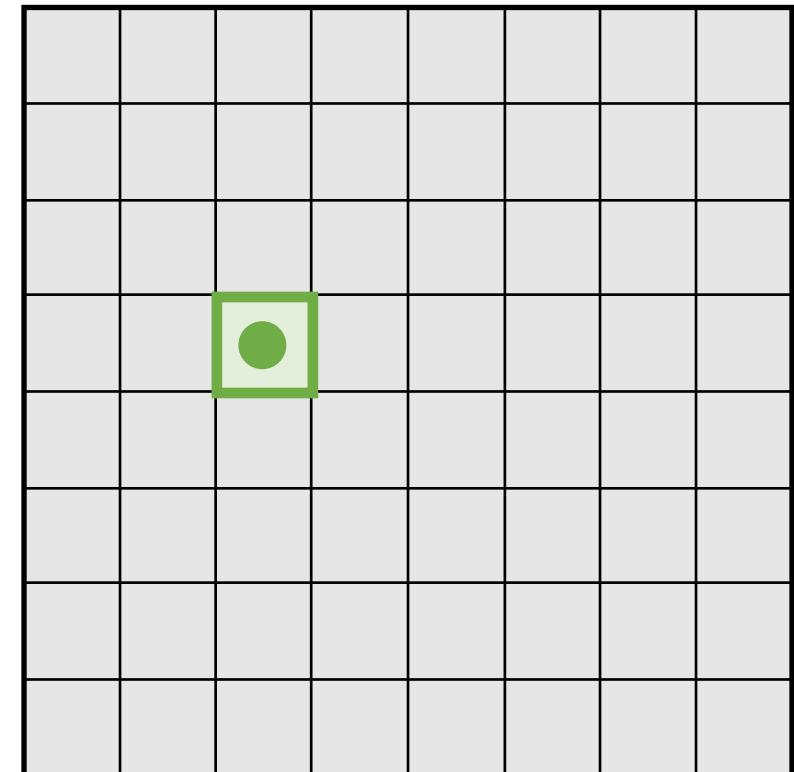


Input Image: 8 x 8

Moving one unit in the output space also moves the receptive field by one

3x3 Conv  
Stride 1, pad 1

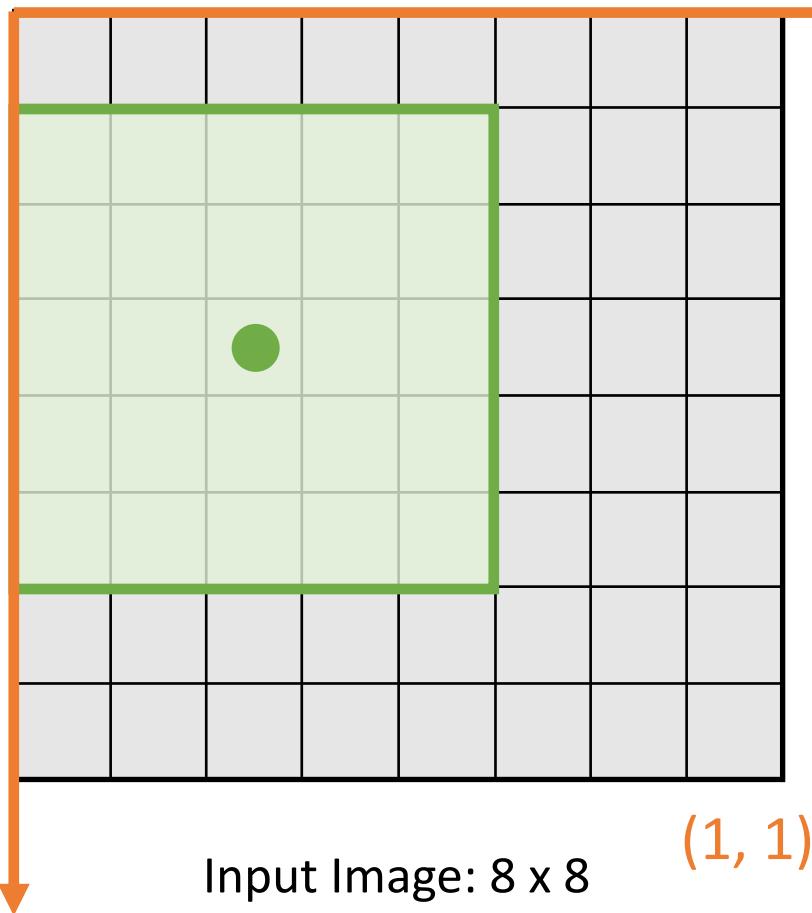
3x3 Conv  
Stride 1, pad 1



Output Image: 8 x 8

# Recall: Receptive Fields

(0, 0)



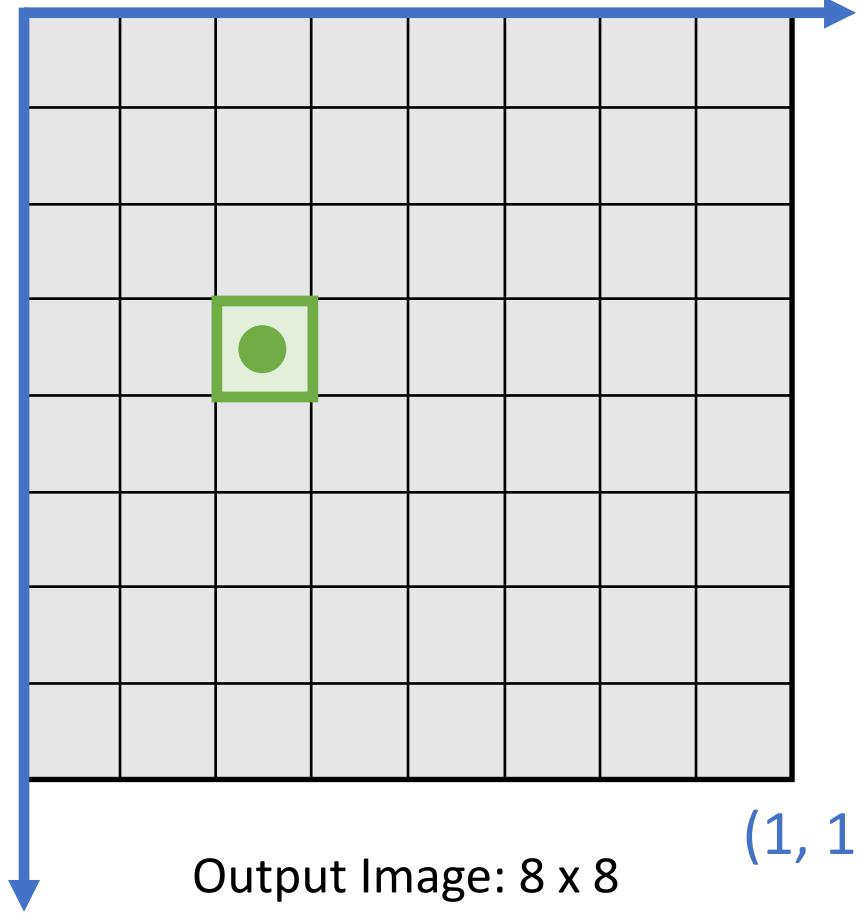
Moving one unit in the output space also moves the receptive field by one

3x3 Conv  
Stride 1, pad 1

3x3 Conv  
Stride 1, pad 1

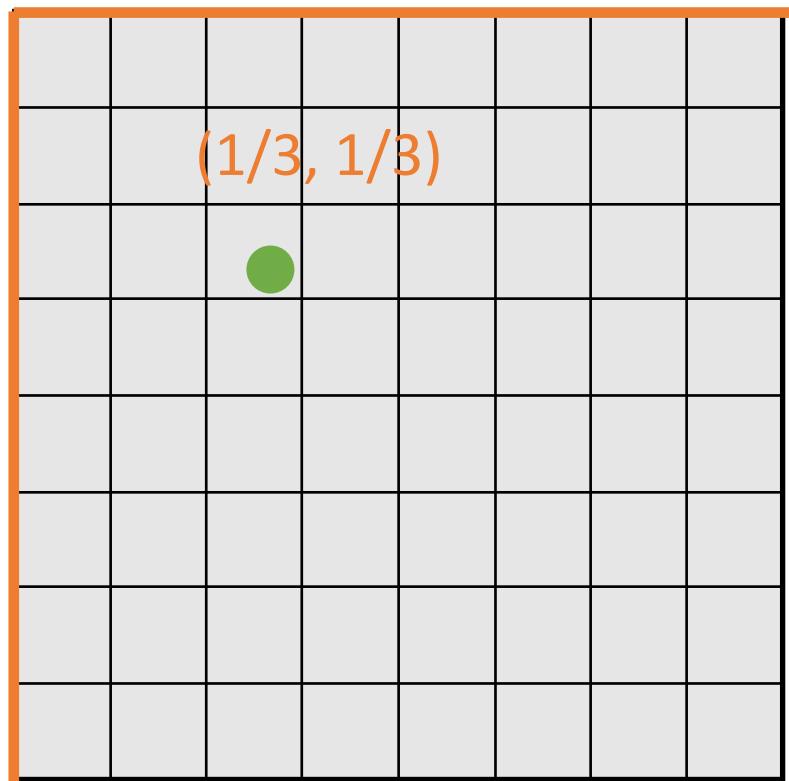
There is a correspondence between the coordinate system of the input and the coordinate system of the output

(0, 0)



# Projecting Points

(0, 0)



Input Image: 8 x 8

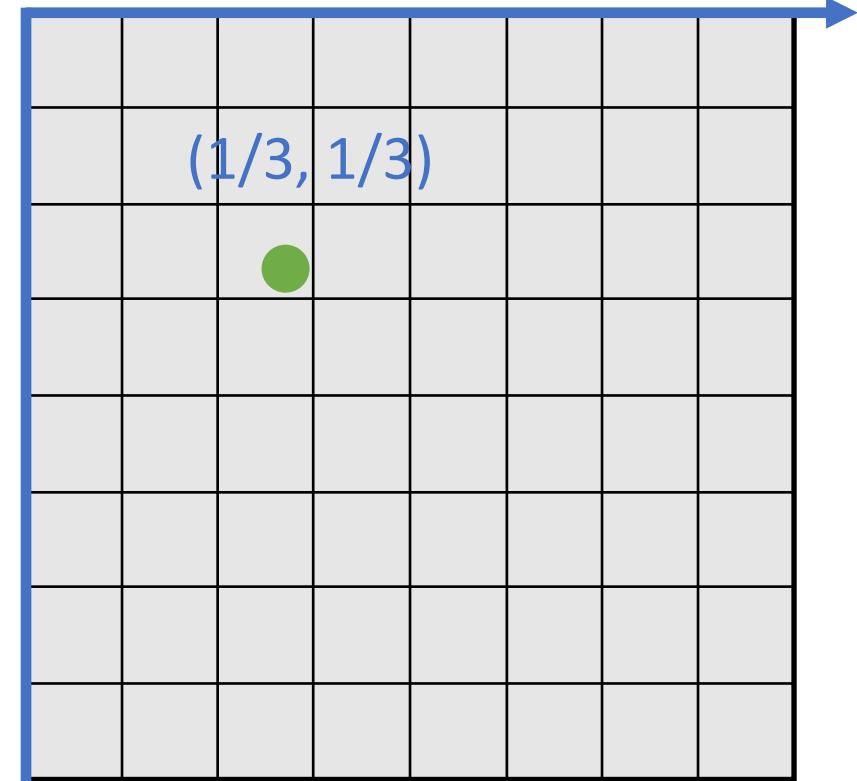
We can align arbitrary points between coordinate system of input and output

3x3 Conv  
Stride 1, pad 1

3x3 Conv  
Stride 1, pad 1

There is a correspondence between the coordinate system of the input and the coordinate system of the output

(0, 0)

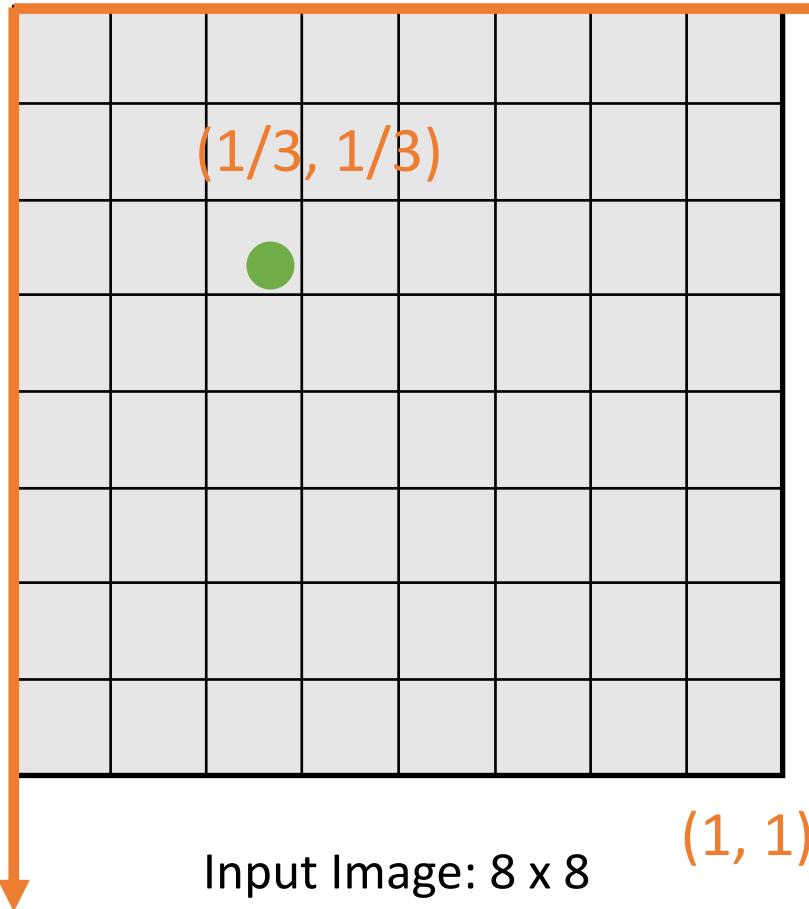


Output Image: 8 x 8

(1, 1)

# Projecting Points

(0, 0)



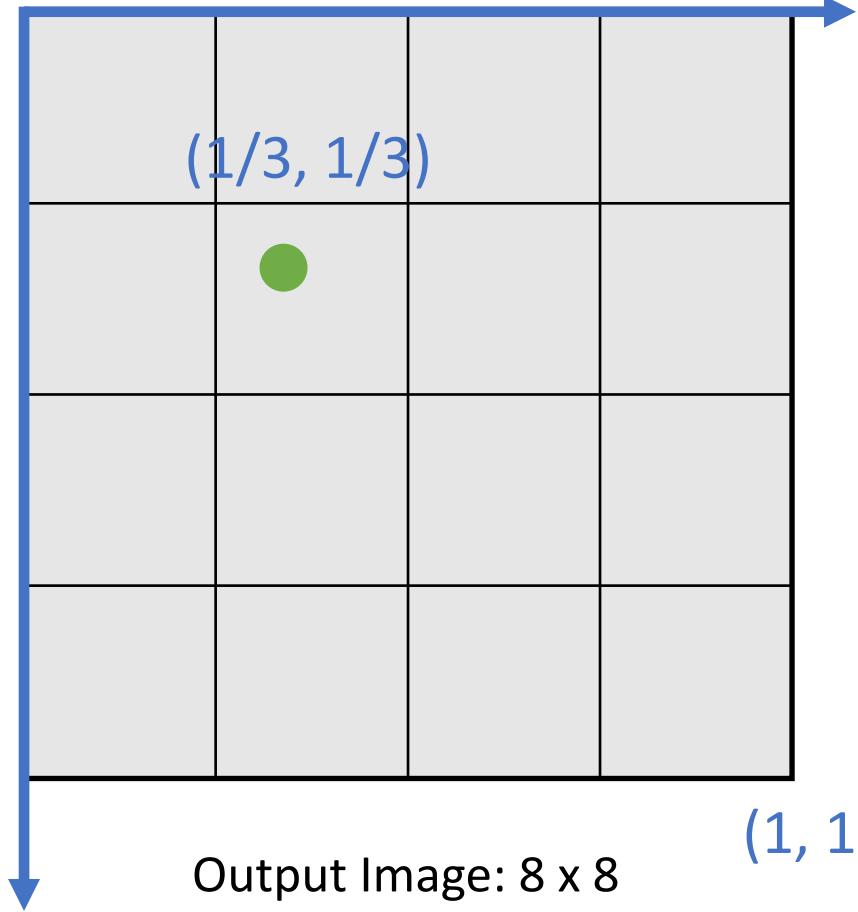
We can align arbitrary points between coordinate system of input and output

3x3 Conv  
Stride 1, pad 1

2x2 MaxPool  
Stride 2

There is a correspondence between the **coordinate system of the input** and the **coordinate system of the output**

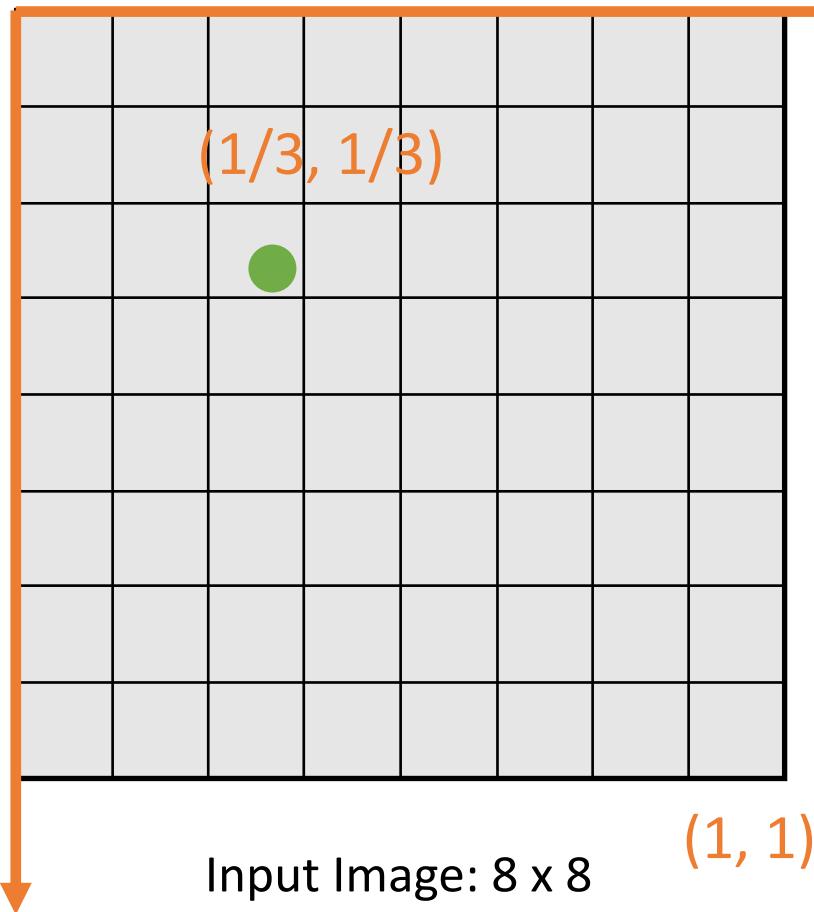
(0, 0)



Same logic holds for more complicated CNNs, even if spatial resolution of input and output are different

# Projecting Points

(0, 0)



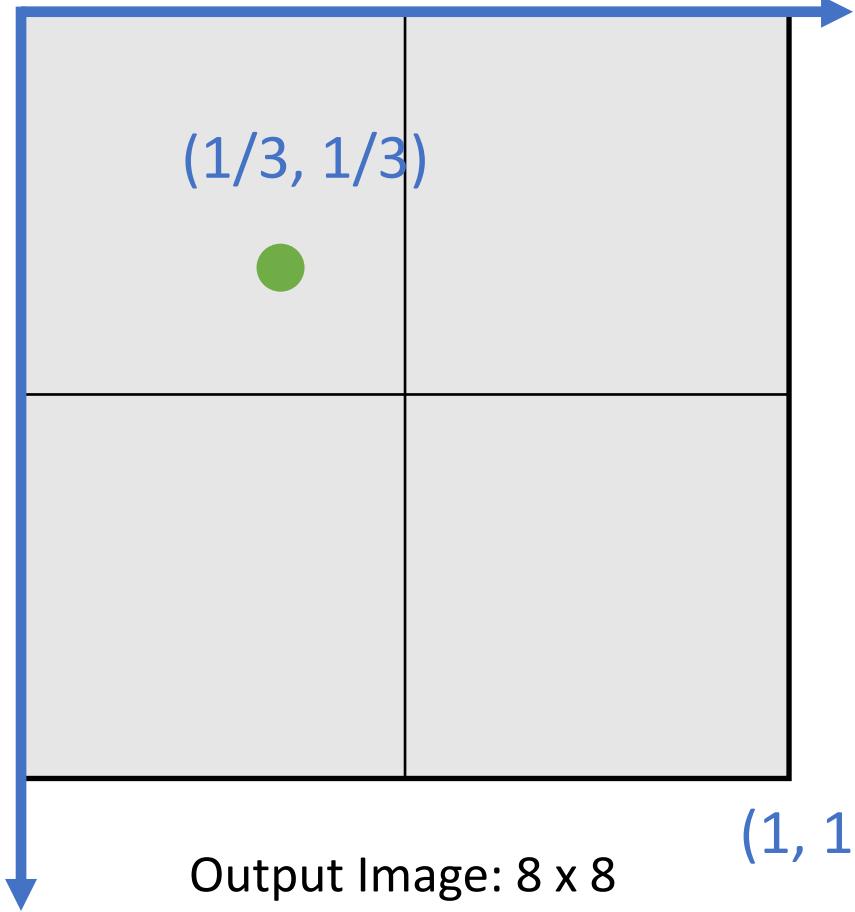
We can align arbitrary points between coordinate system of input and output

3x3 Conv  
Stride 1, pad 1

4x4 MaxPool  
Stride 4

There is a correspondence between the **coordinate system of the input** and the **coordinate system of the output**

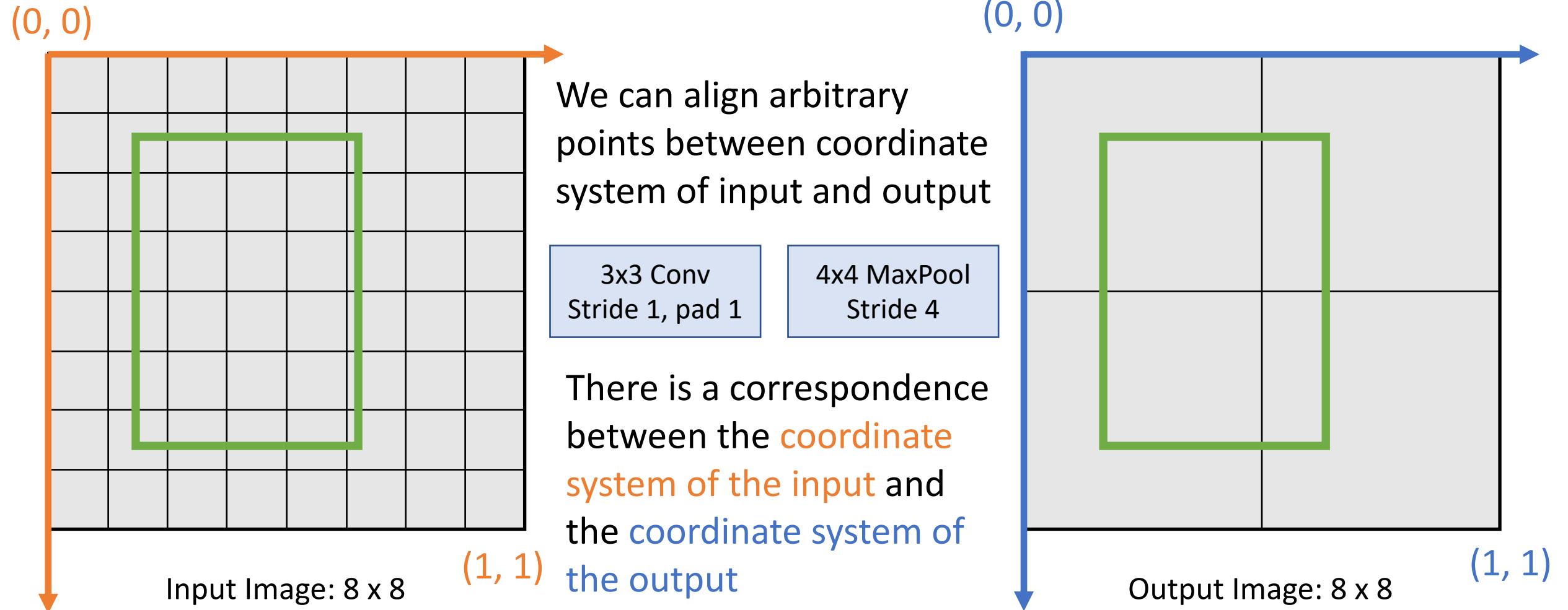
(0, 0)



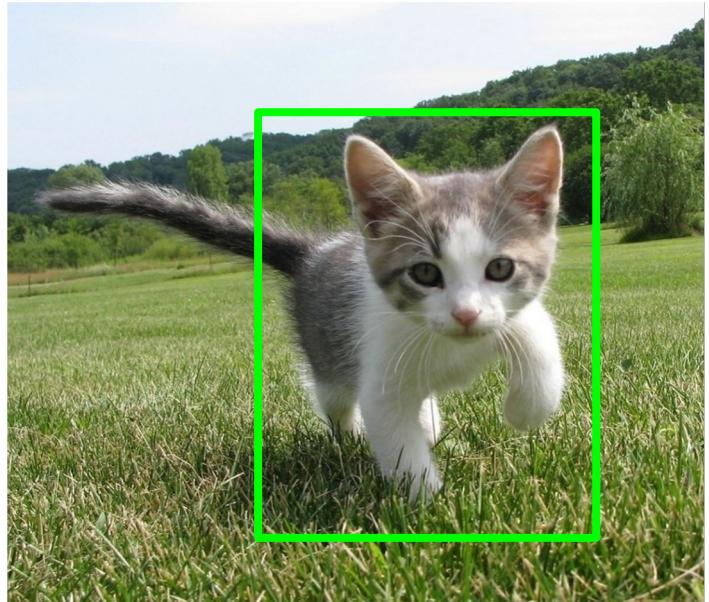
Same logic holds for more complicated CNNs, even if spatial resolution of input and output are different

# Projecting Boxes

We can use this idea to project **bounding boxes** between an input image and a feature map



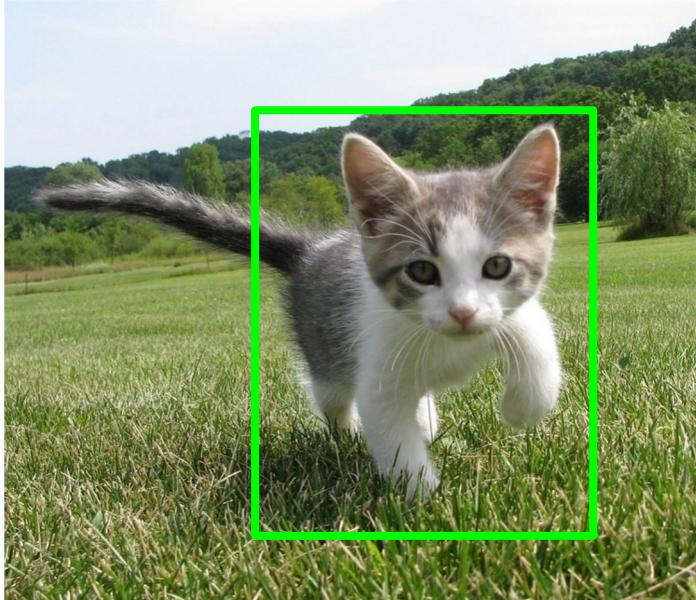
# Cropping Features: RoI Pool



Input Image  
(e.g.  $3 \times 640 \times 480$ )

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

# Cropping Features: RoI Pool



Input Image  
(e.g.  $3 \times 640 \times 480$ )

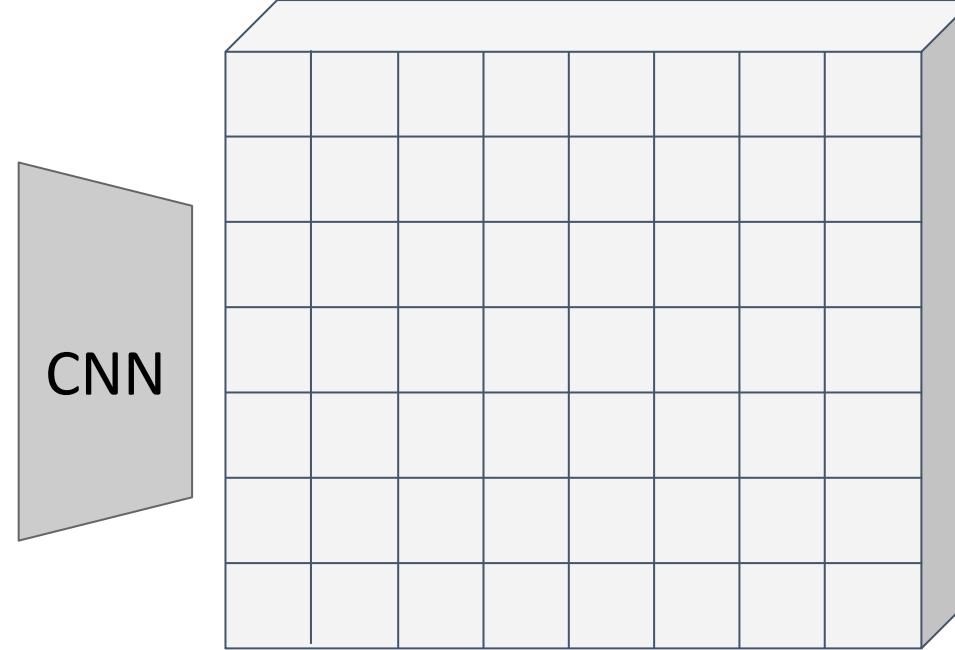
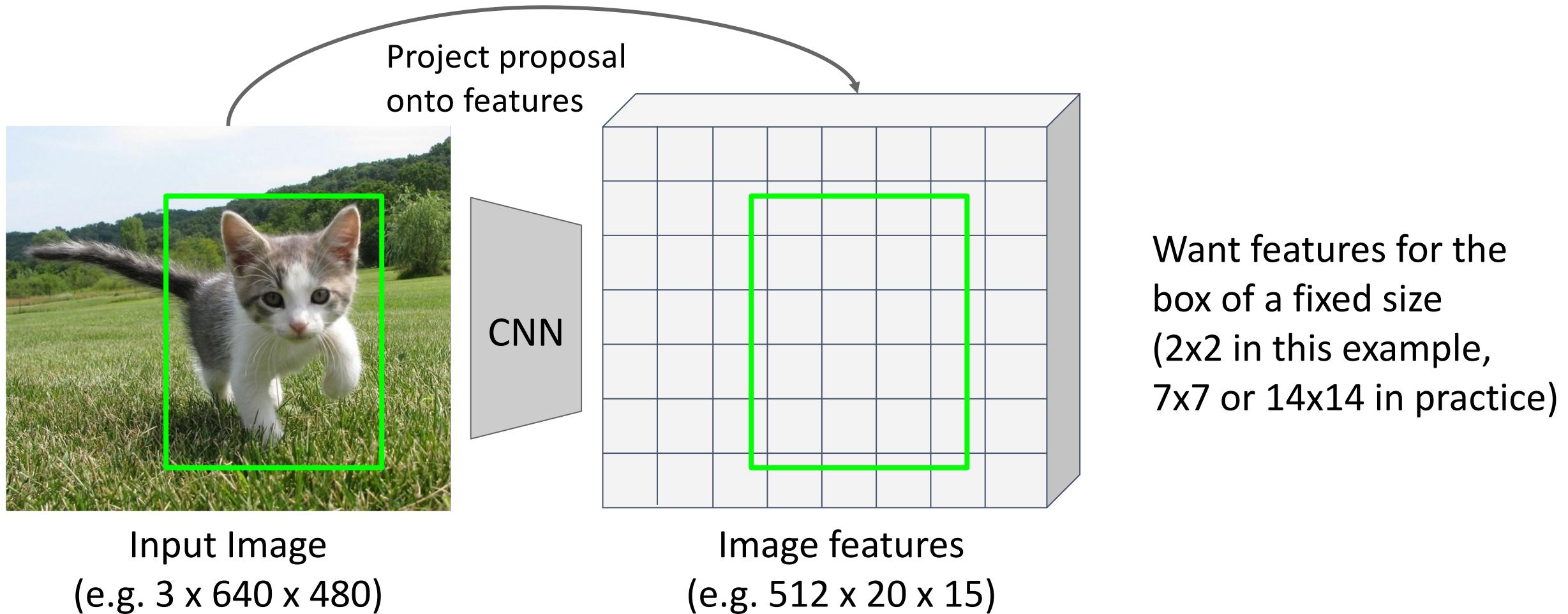


Image features  
(e.g.  $512 \times 20 \times 15$ )

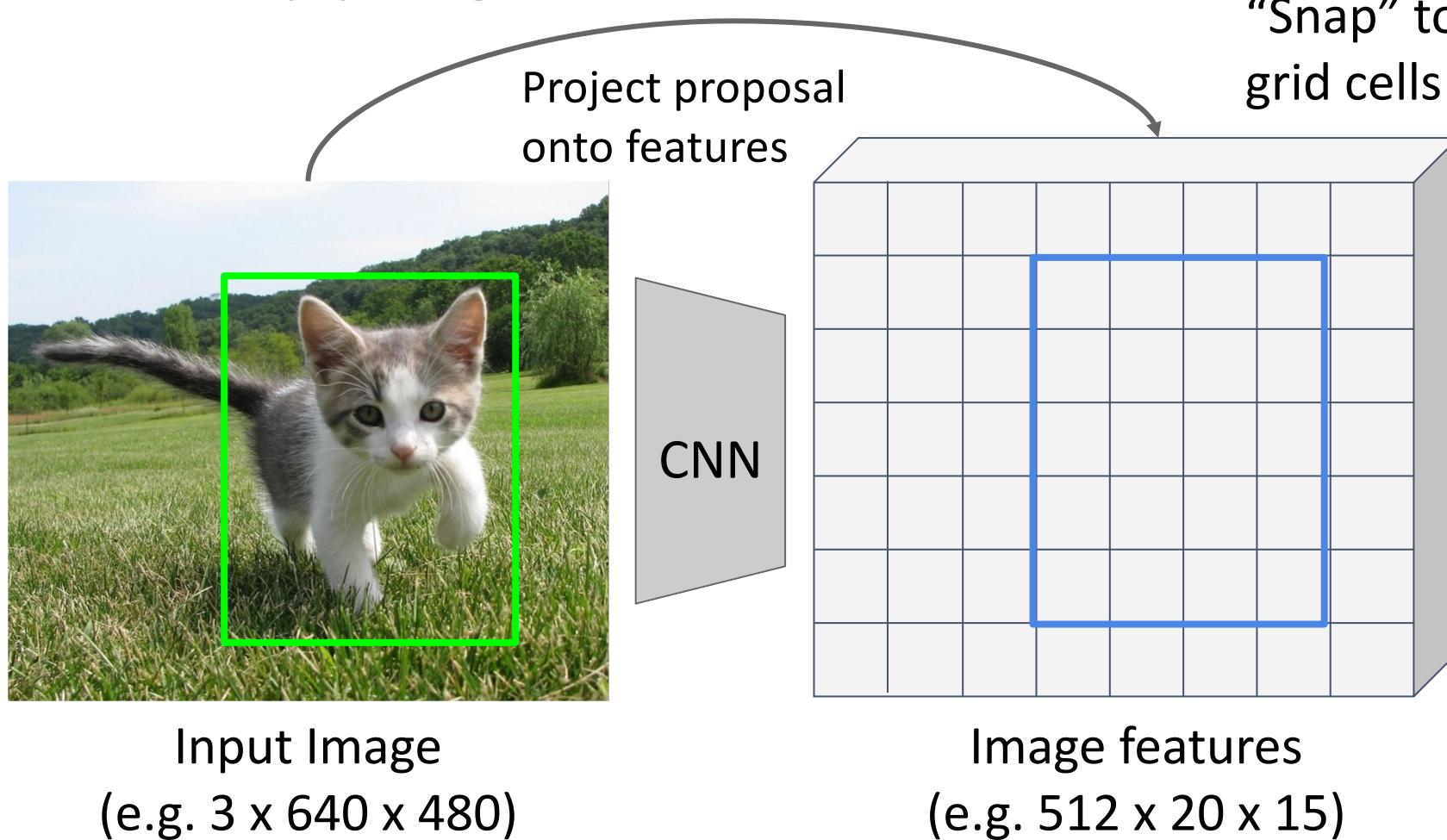
Want features for the  
box of a fixed size  
( $2 \times 2$  in this example,  
 $7 \times 7$  or  $14 \times 14$  in practice)

# Cropping Features: RoI Pool



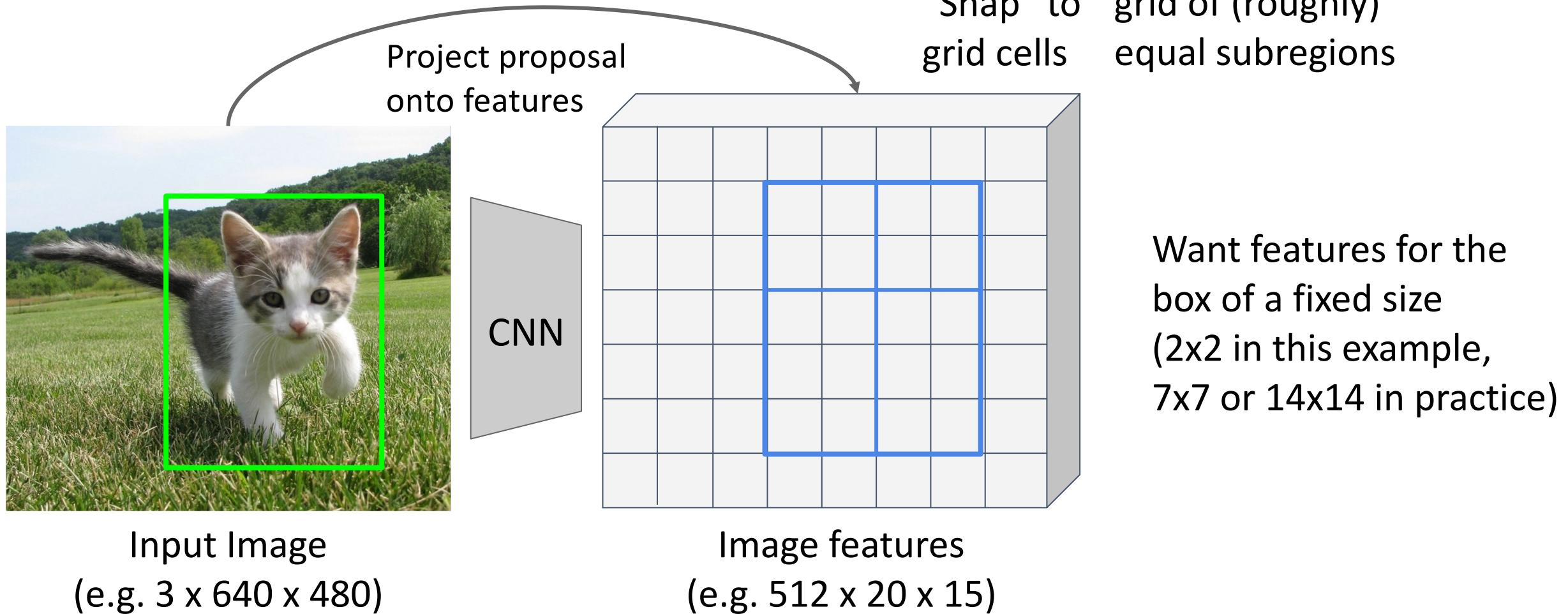
Girshick, "Fast R-CNN", ICCV 2015.

# Cropping Features: RoI Pool



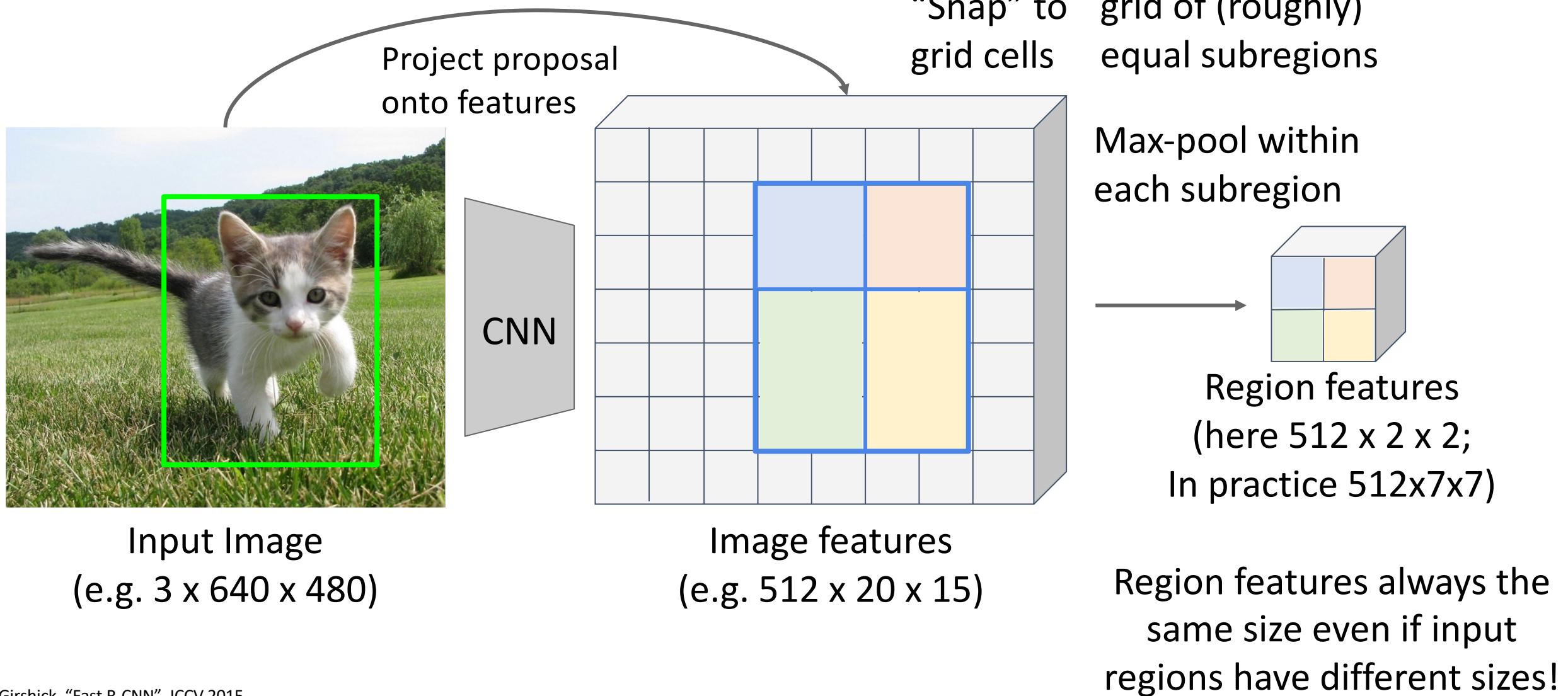
Want features for the  
box of a fixed size  
( $2 \times 2$  in this example,  
 $7 \times 7$  or  $14 \times 14$  in practice)

# Cropping Features: RoI Pool



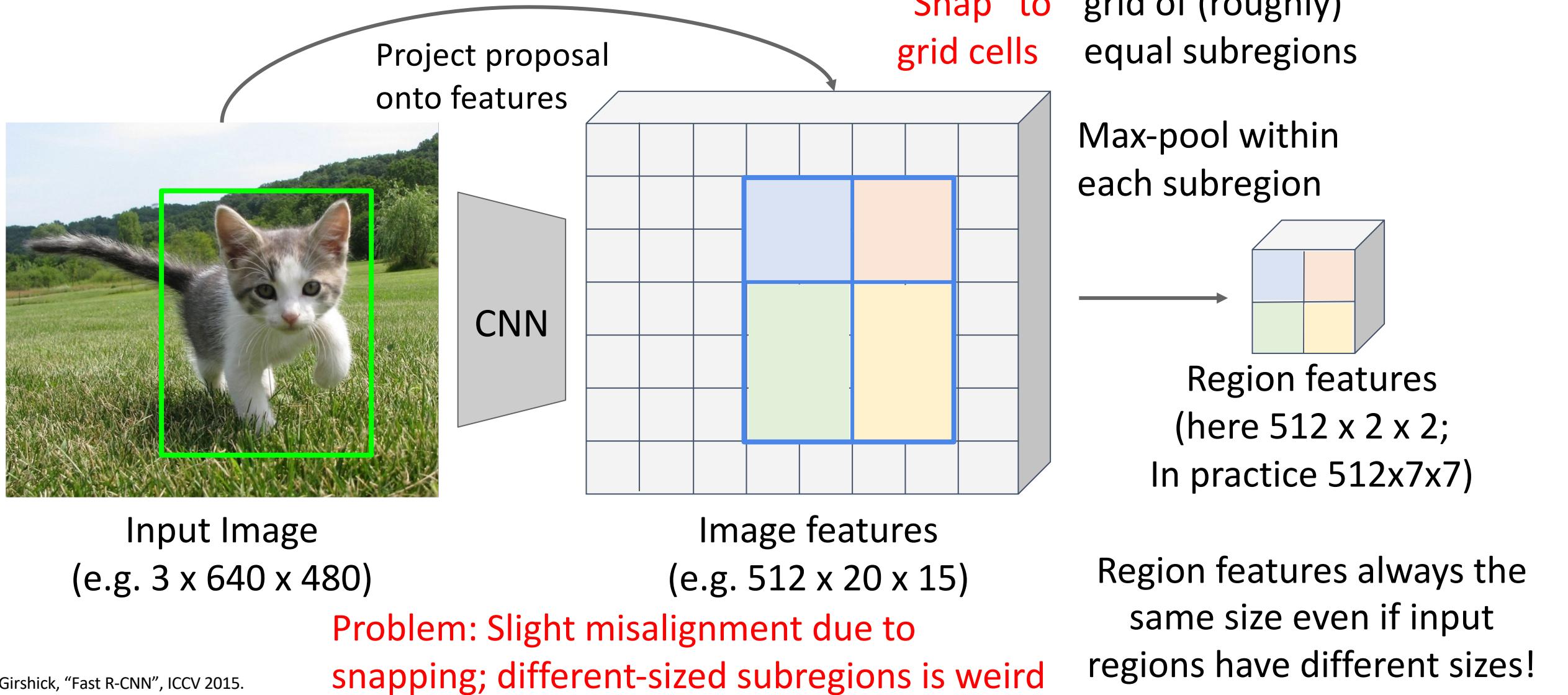
Girshick, "Fast R-CNN", ICCV 2015.

# Cropping Features: RoI Pool



Girshick, “Fast R-CNN”, ICCV 2015.

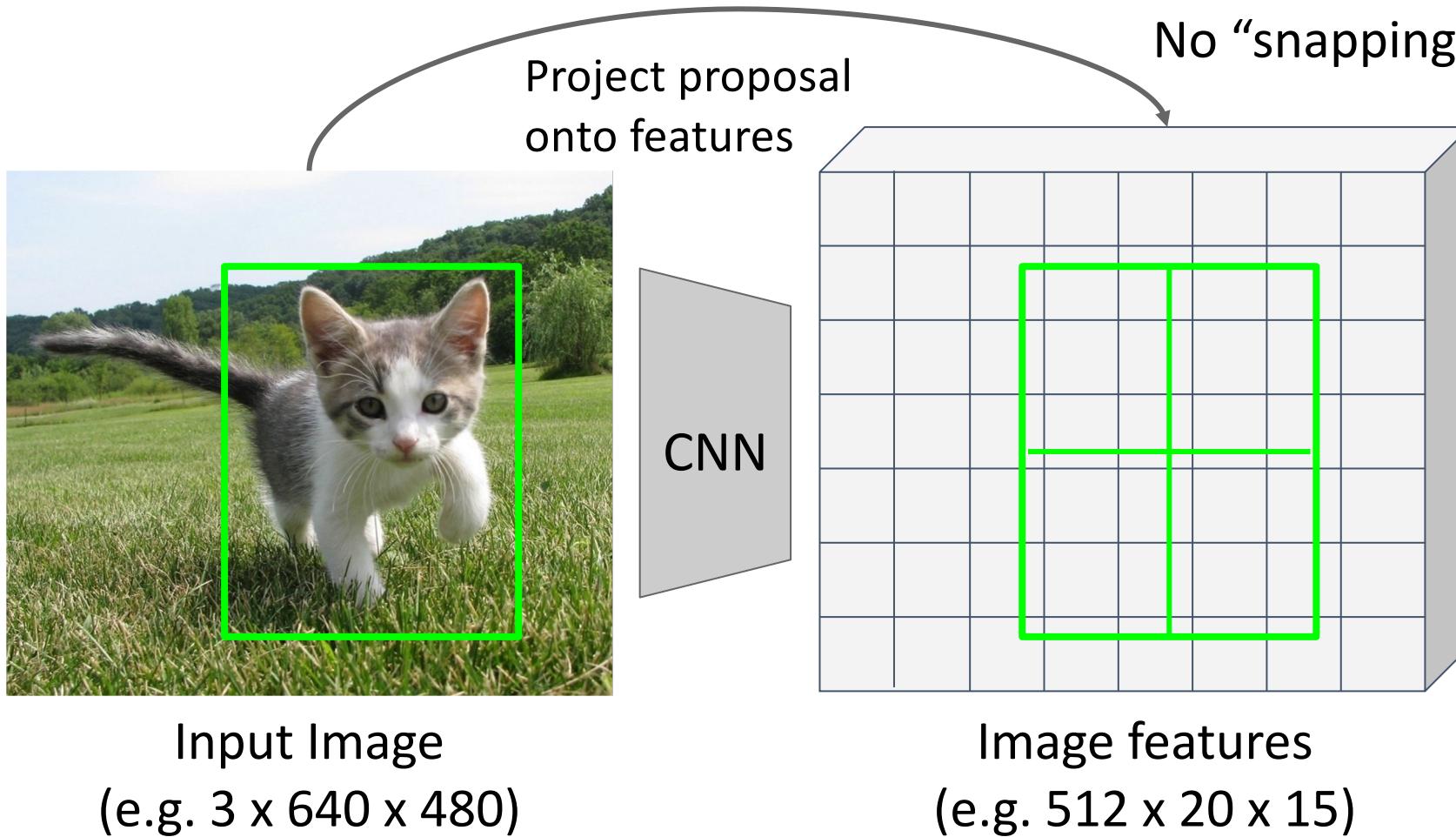
# Cropping Features: RoI Pool



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

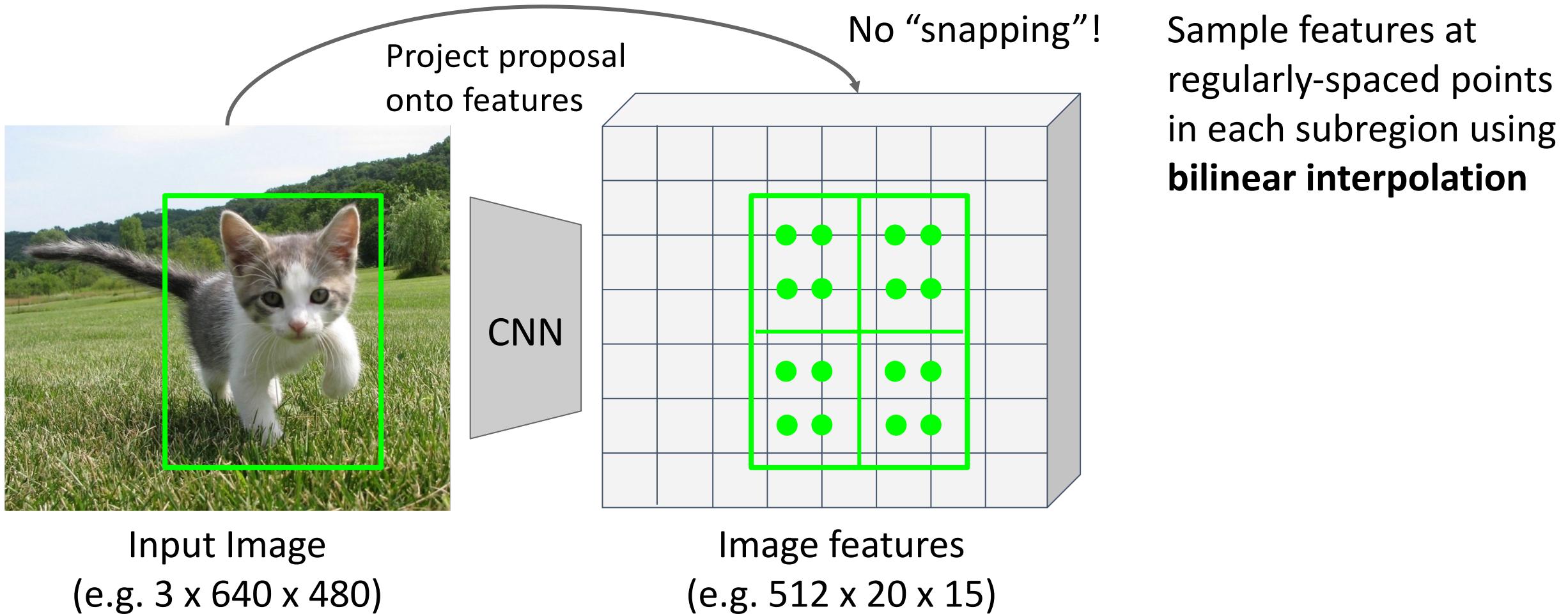
# Cropping Features: RoI Align

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



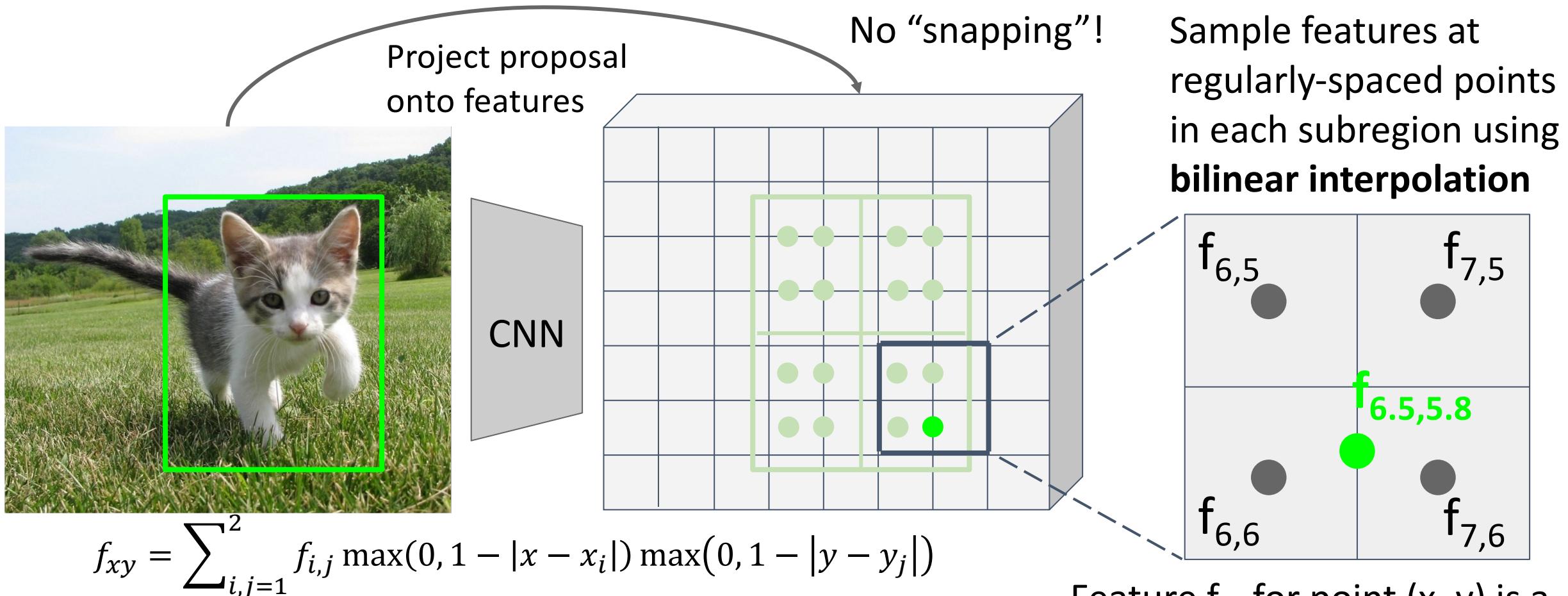
# Cropping Features: RoI Align

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



# Cropping Features: RoI Align

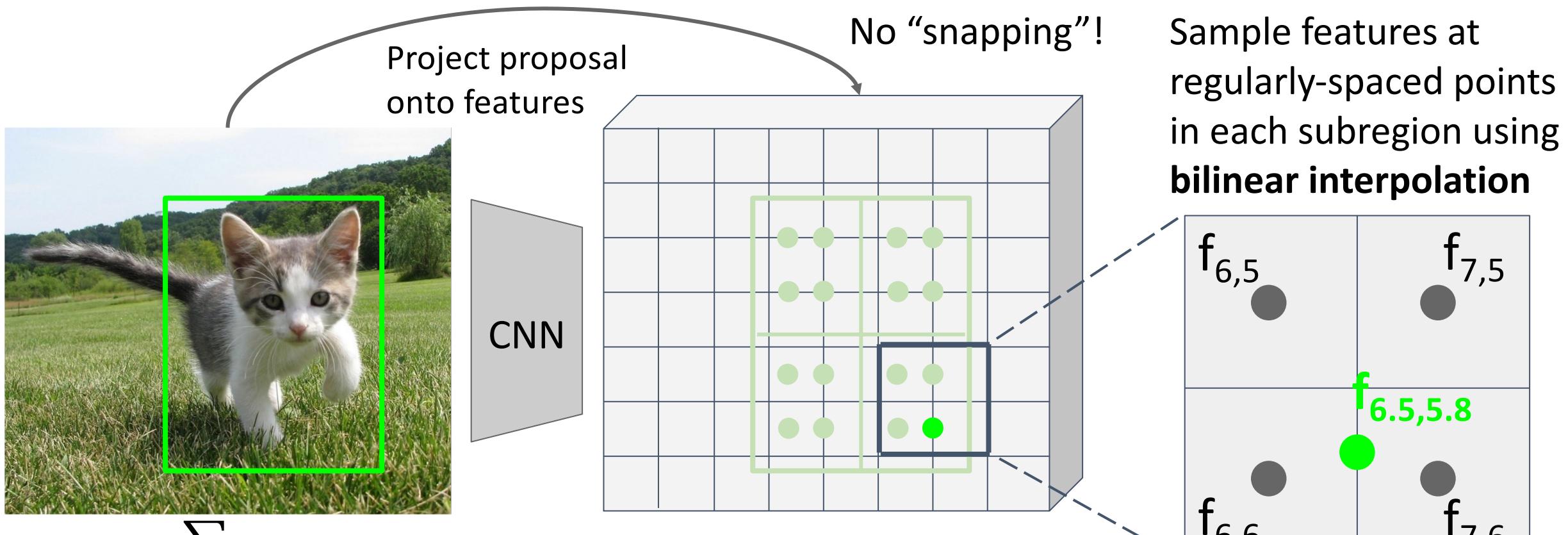
Divide into equal-sized subregions  
(may not be aligned to grid!)



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

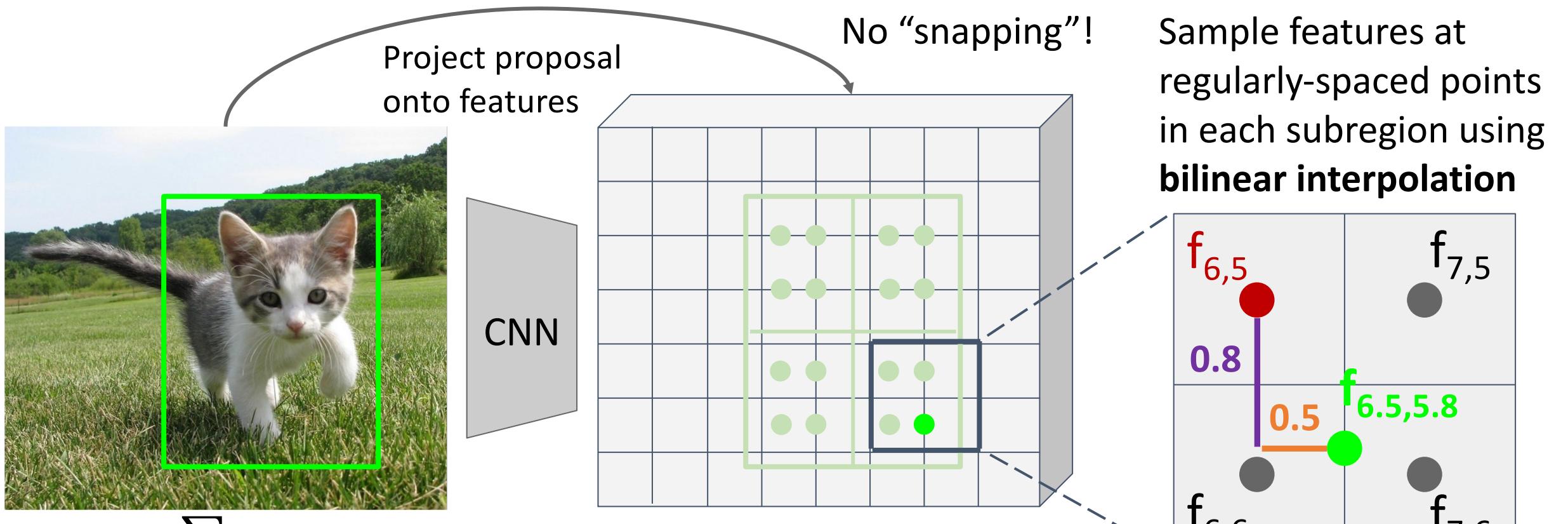
# Cropping Features: RoI Align

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



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

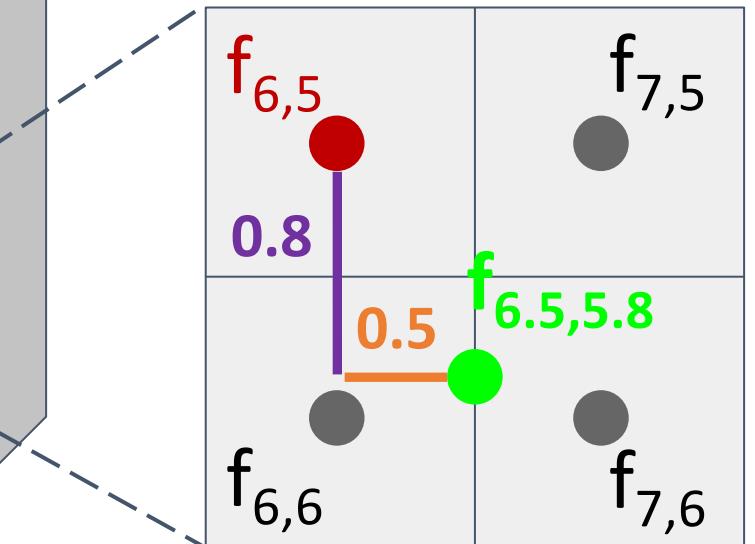
# Cropping Features: RoI Align



$$f_{xy} = \sum_{i,j} f_{i,j} \max(0, 1 - |x - x_i|) \max(0, 1 - |y - y_i|)$$

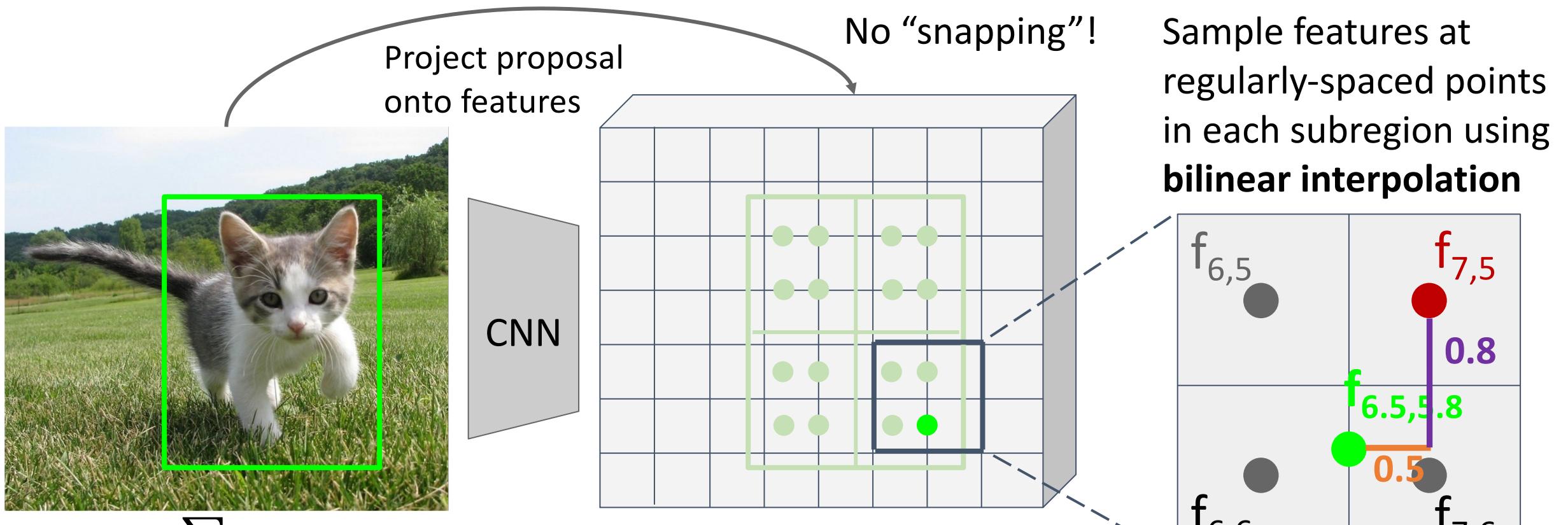
$$\begin{aligned} f_{6.5,5.8} &= (f_{6,5} * 0.5 * 0.2) + (f_{7,5} * 0.5 * 0.2) \\ &\quad + (f_{6,6} * 0.5 * 0.8) + (f_{7,6} * 0.5 * 0.8) \end{aligned}$$

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



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

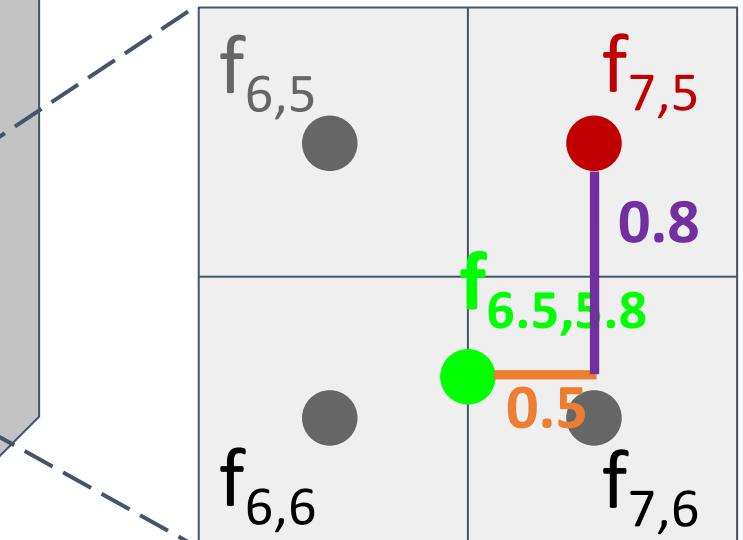
# Cropping Features: RoI Align



$$f_{xy} = \sum_{i,j} f_{i,j} \max(0, 1 - |x - x_i|) \max(0, 1 - |y - y_i|)$$

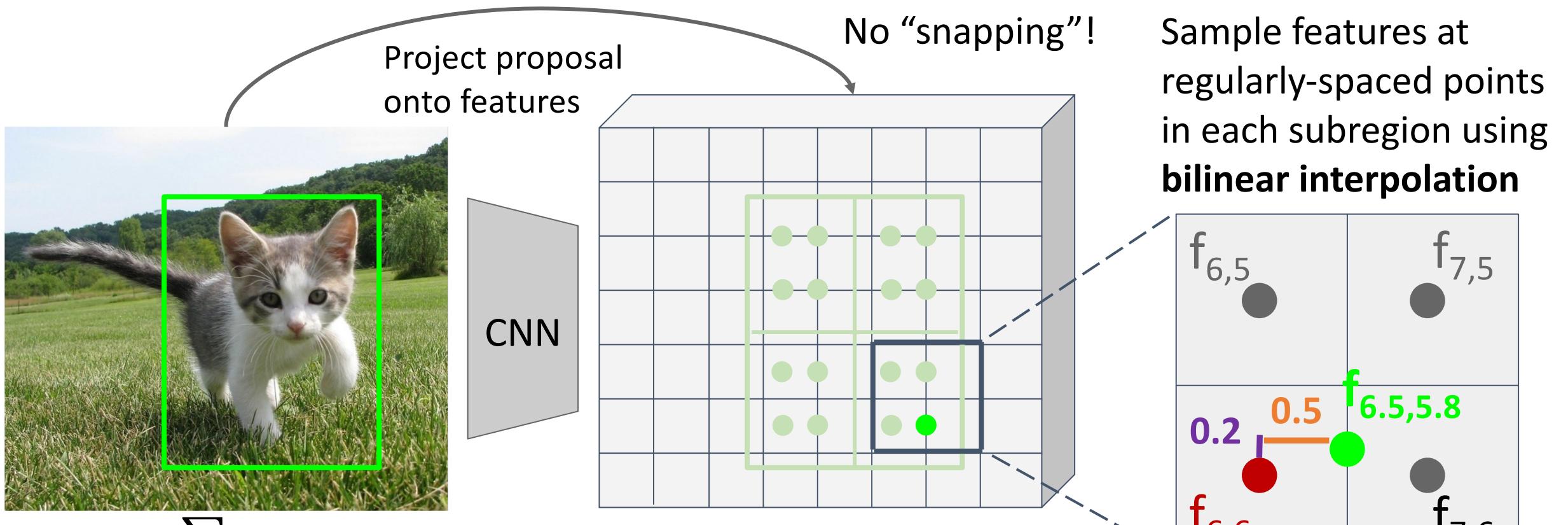
$$\begin{aligned} f_{6.5,5.8} &= (f_{6,5} * 0.5 * 0.2) + (f_{7,5} * 0.5 * 0.2) \\ &\quad + (f_{6,6} * 0.5 * 0.8) + (f_{7,6} * 0.5 * 0.8) \end{aligned}$$

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

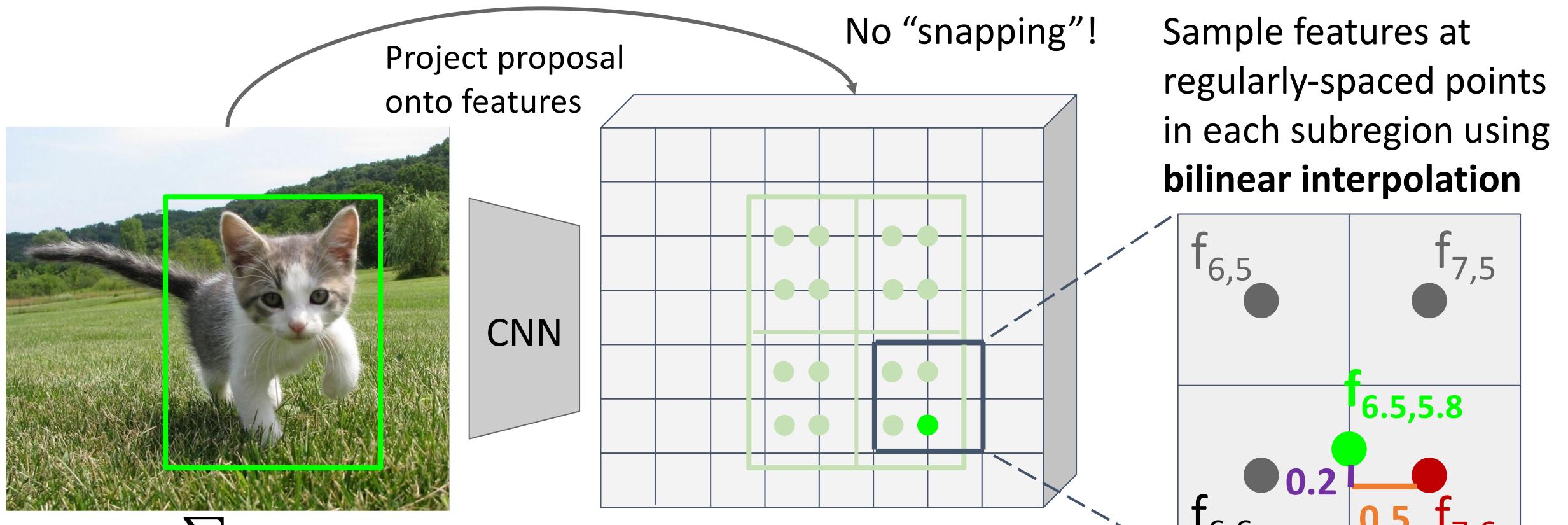


Feature  $f_{xy}$  for point  $(x, y)$  is a linear combination of features at its four neighboring grid cells:

# Cropping Features: RoI Align



# Cropping Features: RoI Align

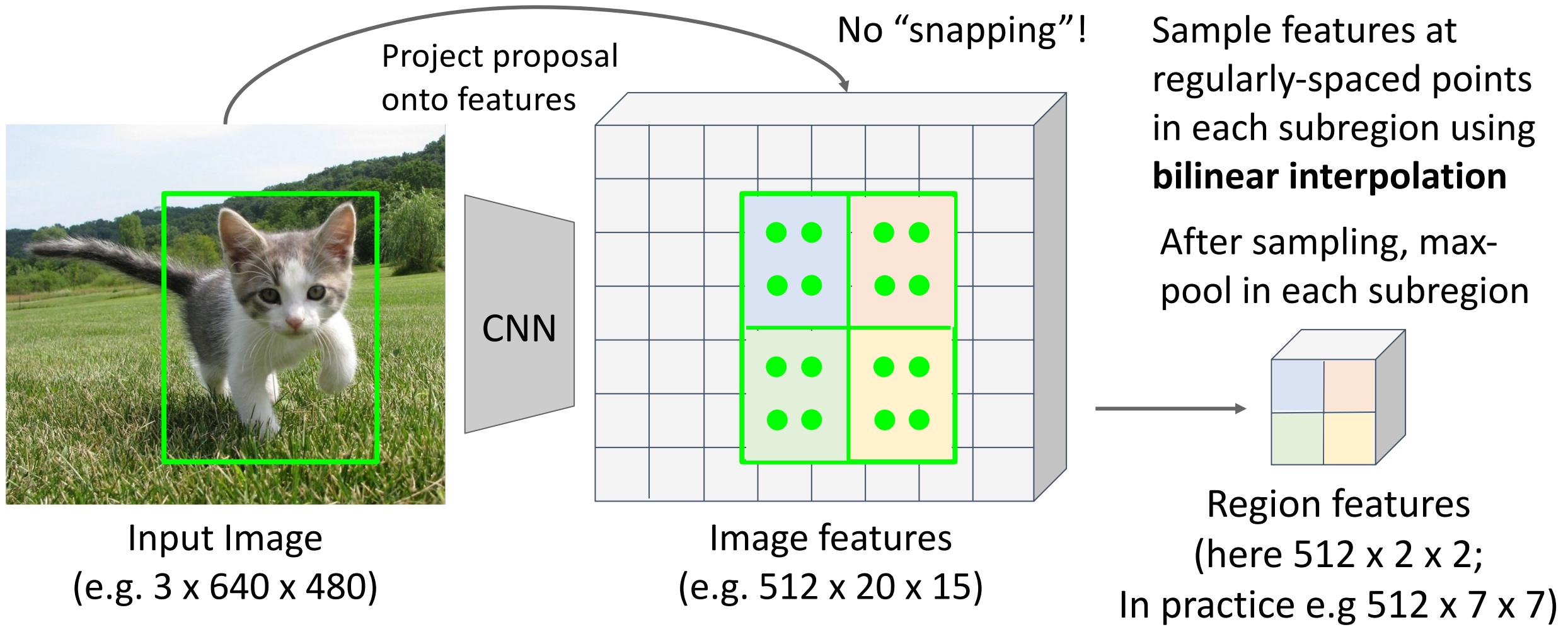


$$f_{xy} = \sum_{i,j} f_{i,j} \max(0, 1 - |x - x_i|) \max(0, 1 - |y - y_i|)$$

$$\begin{aligned} f_{6.5,5.8} &= (f_{6,5} * 0.5 * 0.2) + (f_{7,5} * 0.5 * 0.2) \\ &\quad + (f_{6,6} * 0.5 * 0.8) + (f_{7,6} * 0.5 * 0.8) \end{aligned}$$

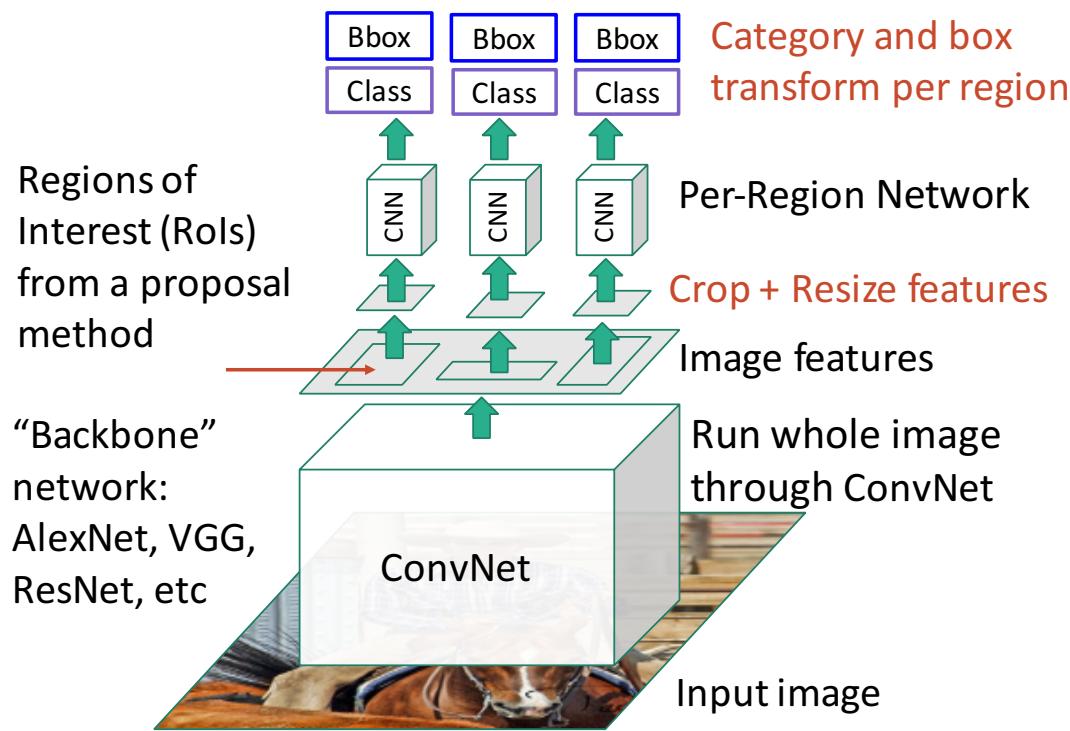
Feature  $f_{xy}$  for point  $(x, y)$  is a linear combination of features at its four neighboring grid cells:

# Cropping Features: RoI Align

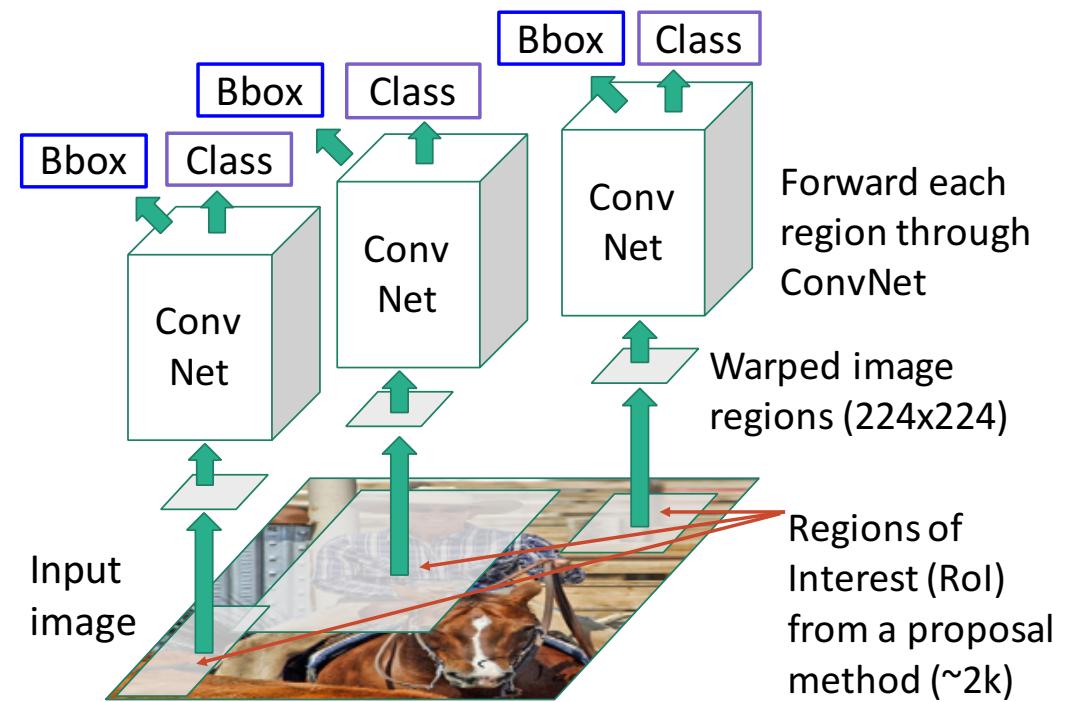


# Fast R-CNN vs “Slow” R-CNN

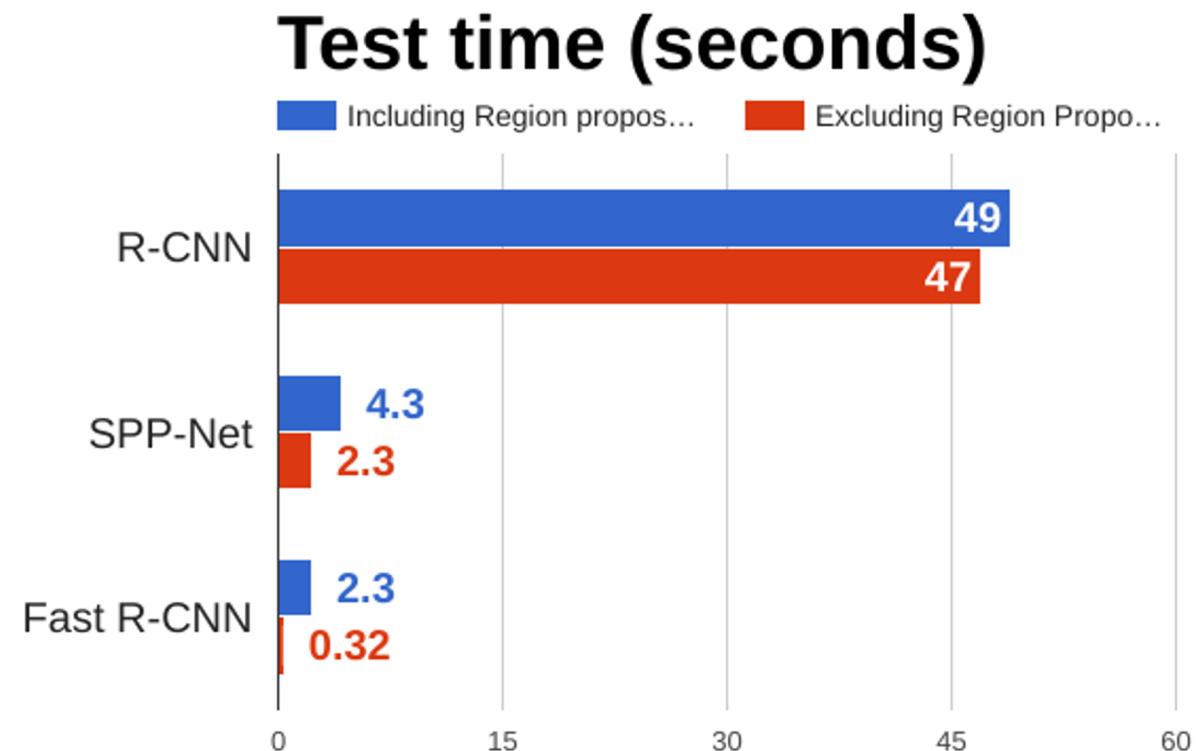
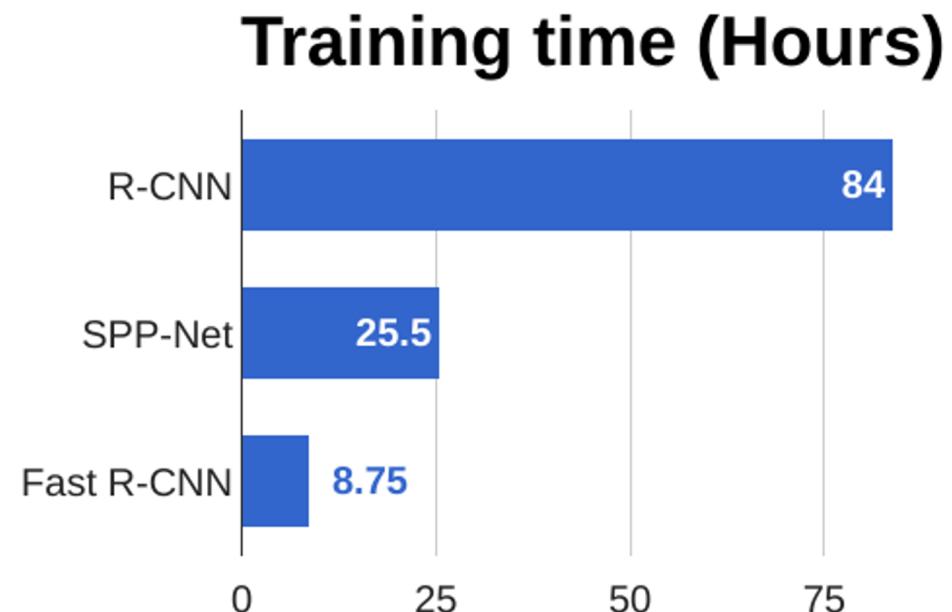
**Fast R-CNN:** Apply differentiable cropping to shared image features



**“Slow” R-CNN:** Apply differentiable cropping to shared image features



# Fast R-CNN vs “Slow” R-CNN



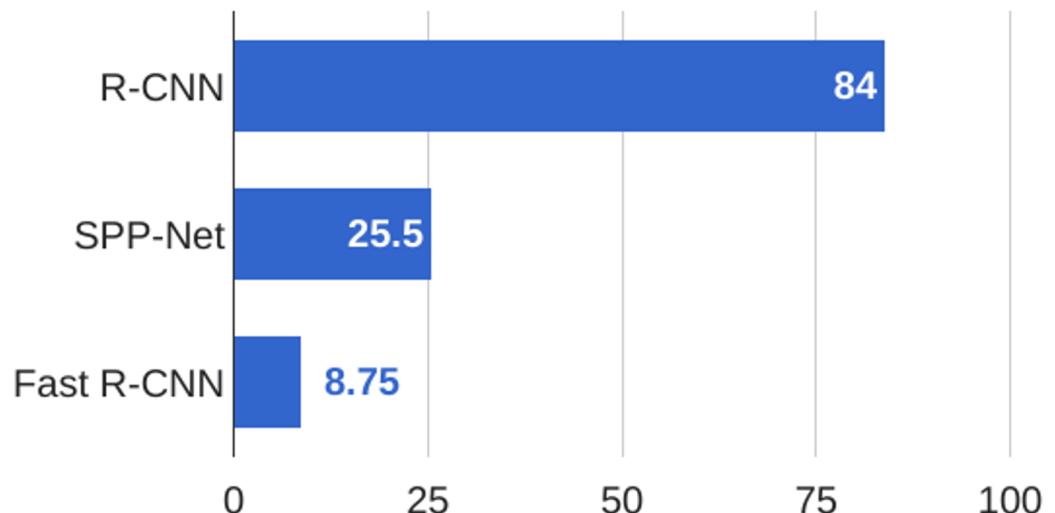
Girshick et al, “Rich feature hierarchies for accurate object detection and semantic segmentation”, CVPR 2014.

He et al, “Spatial pyramid pooling in deep convolutional networks for visual recognition”, ECCV 2014

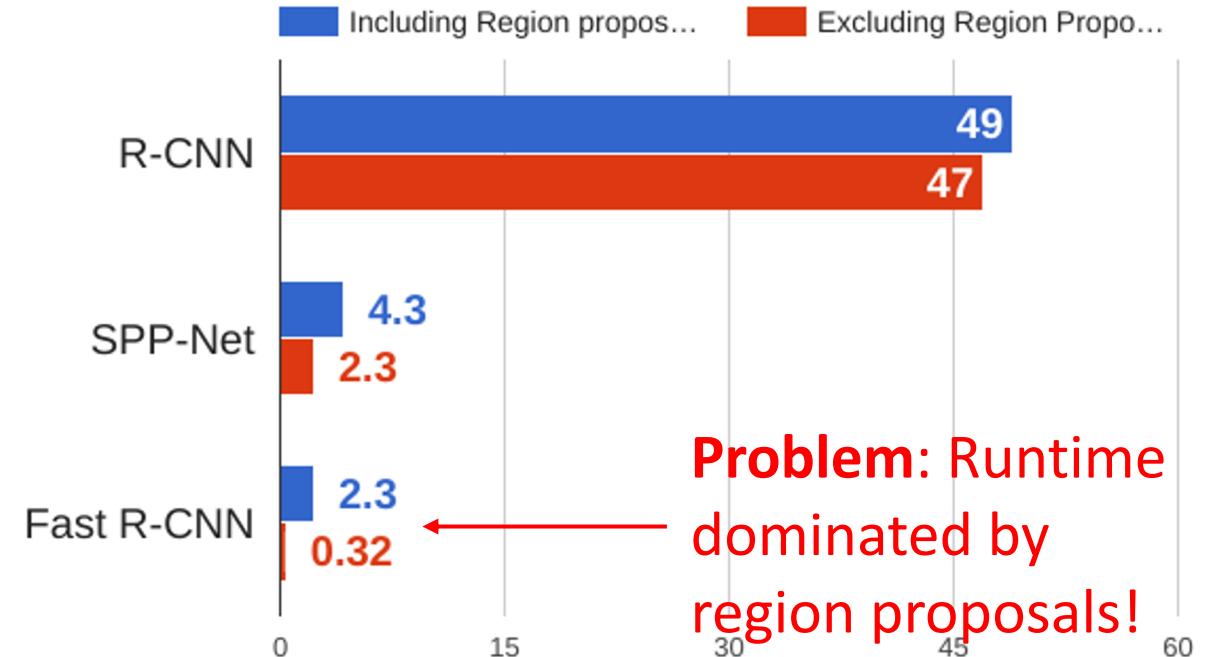
Girshick, “Fast R-CNN”, ICCV 2015

# Fast R-CNN vs “Slow” R-CNN

**Training time (Hours)**



**Test time (seconds)**



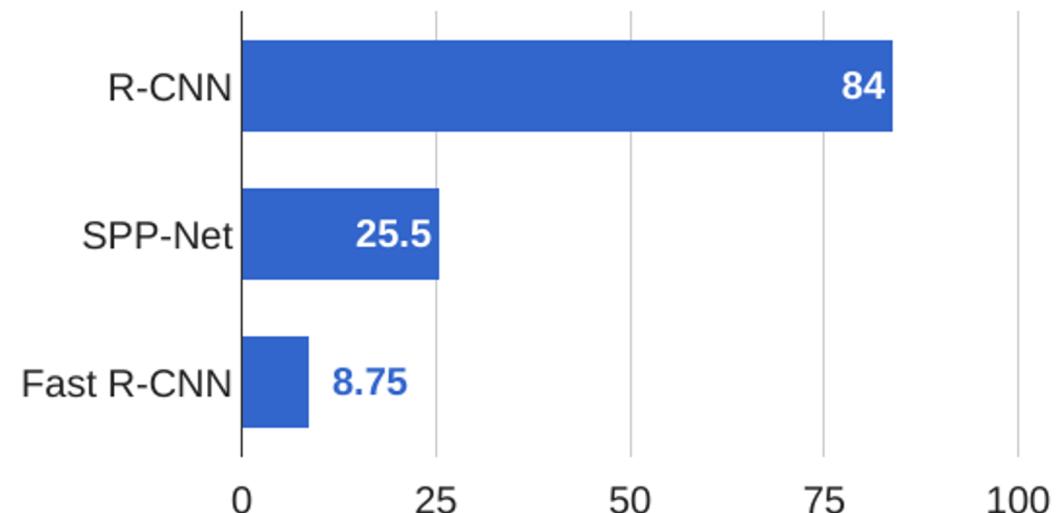
Girshick et al, “Rich feature hierarchies for accurate object detection and semantic segmentation”, CVPR 2014.

He et al, “Spatial pyramid pooling in deep convolutional networks for visual recognition”, ECCV 2014

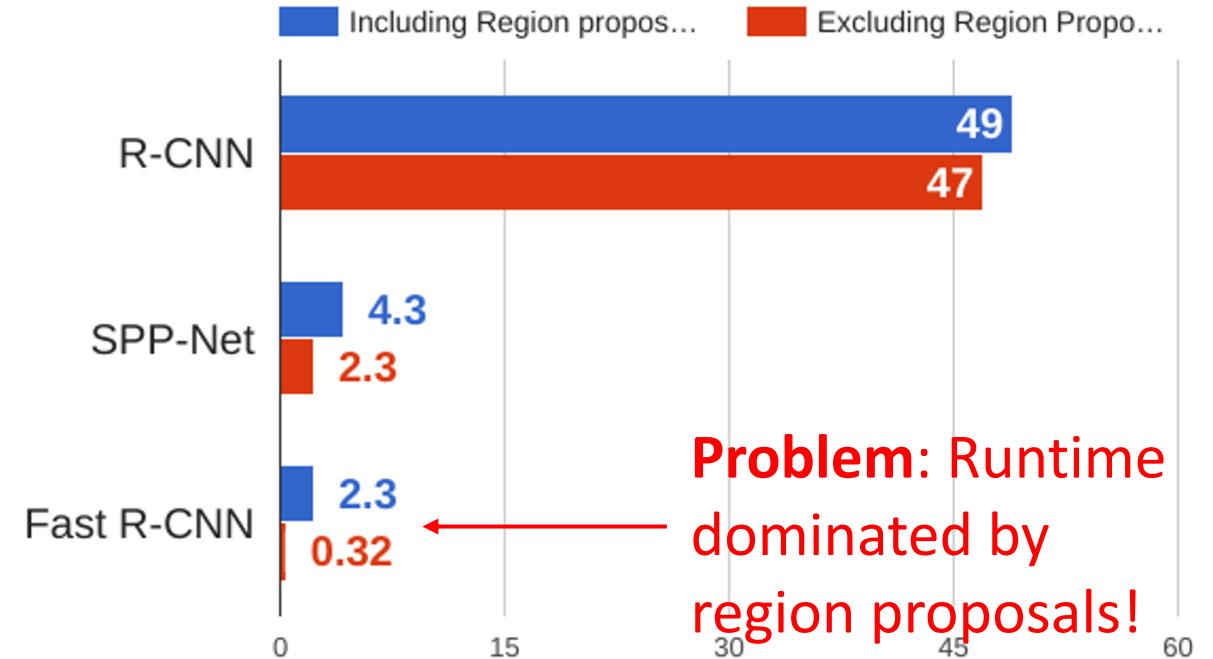
Girshick, “Fast R-CNN”, ICCV 2015

# Fast R-CNN vs “Slow” R-CNN

**Training time (Hours)**



**Test time (seconds)**



**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.

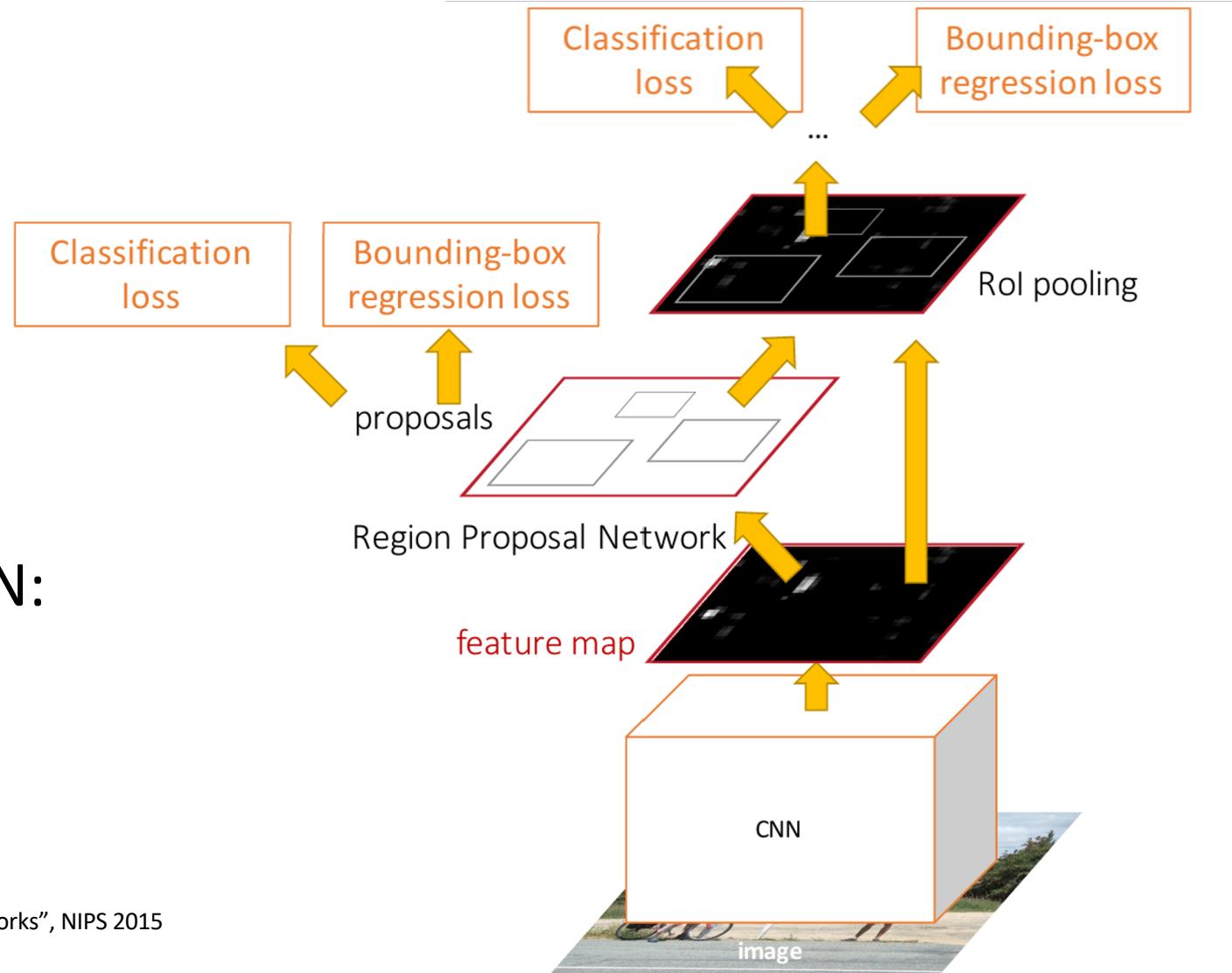
He et al, “Spatial pyramid pooling in deep convolutional networks for visual recognition”, ECCV 2014

Girshick, “Fast R-CNN”, ICCV 2015

# FasterR-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

# Region Proposal Network (RPN)

Run backbone CNN to get  
features aligned to input image



Input Image  
(e.g.  $3 \times 640 \times 480$ )

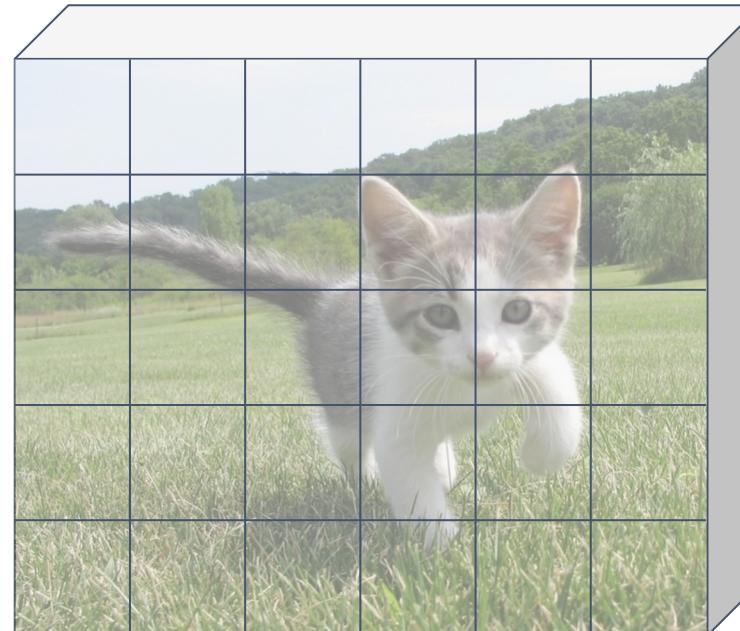


Image features  
(e.g.  $512 \times 5 \times 6$ )

# Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



Input Image  
(e.g.  $3 \times 640 \times 480$ )



Each feature corresponds to a point in the input

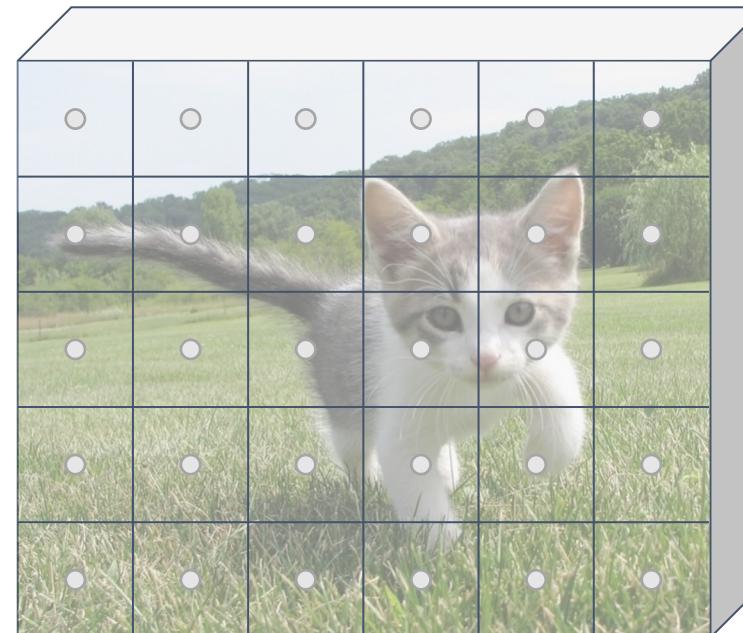
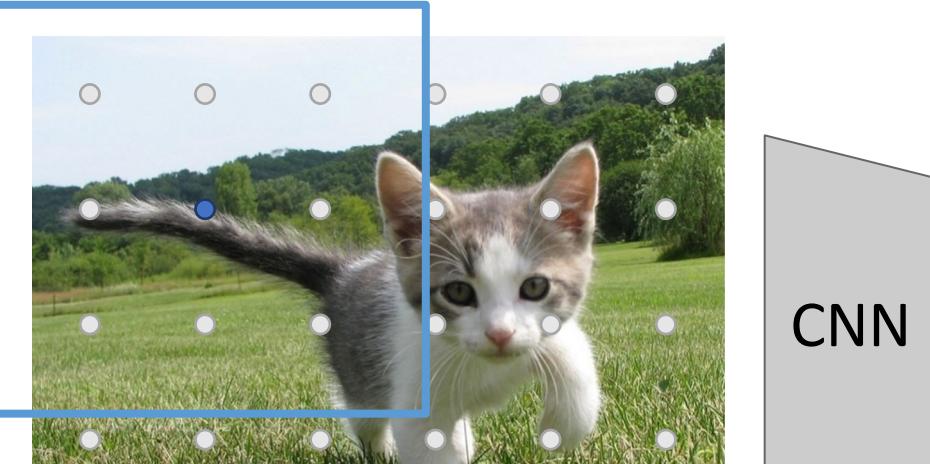


Image features  
(e.g.  $512 \times 5 \times 6$ )

# Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



Input Image  
(e.g.  $3 \times 640 \times 480$ )

Each feature corresponds to a point in the input

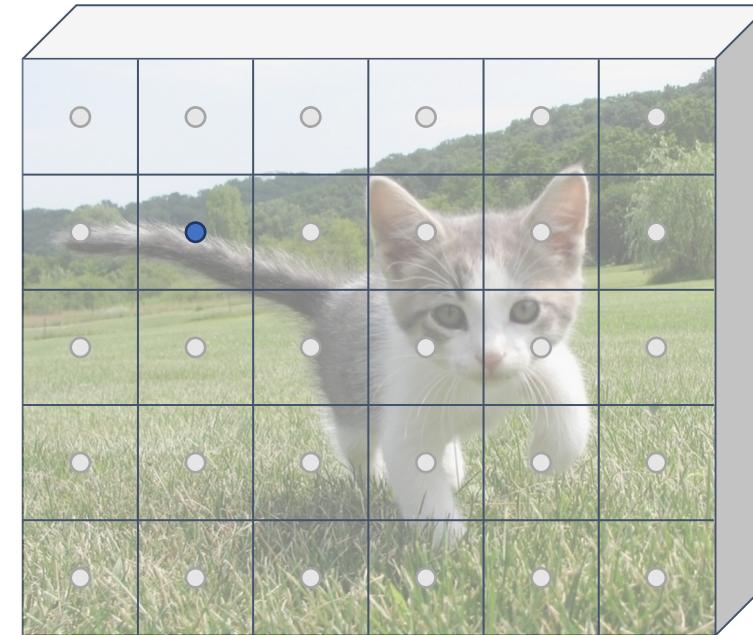
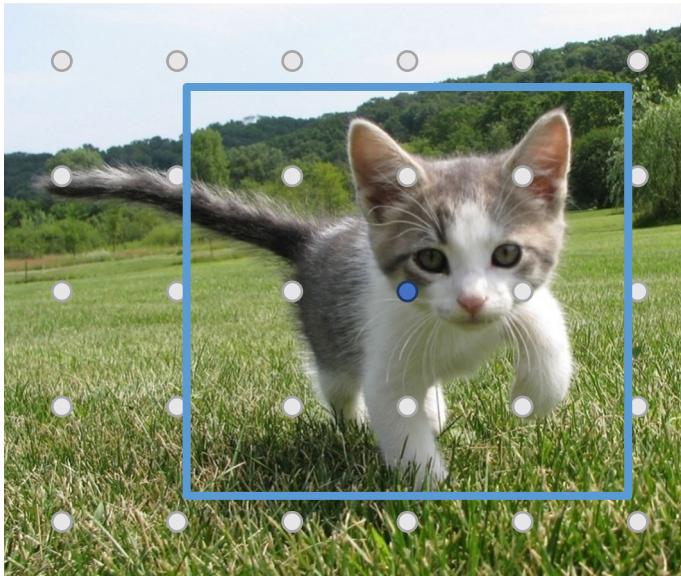


Image features  
(e.g.  $512 \times 5 \times 6$ )

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

# Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



Input Image  
(e.g.  $3 \times 640 \times 480$ )



Each feature corresponds to a point in the input

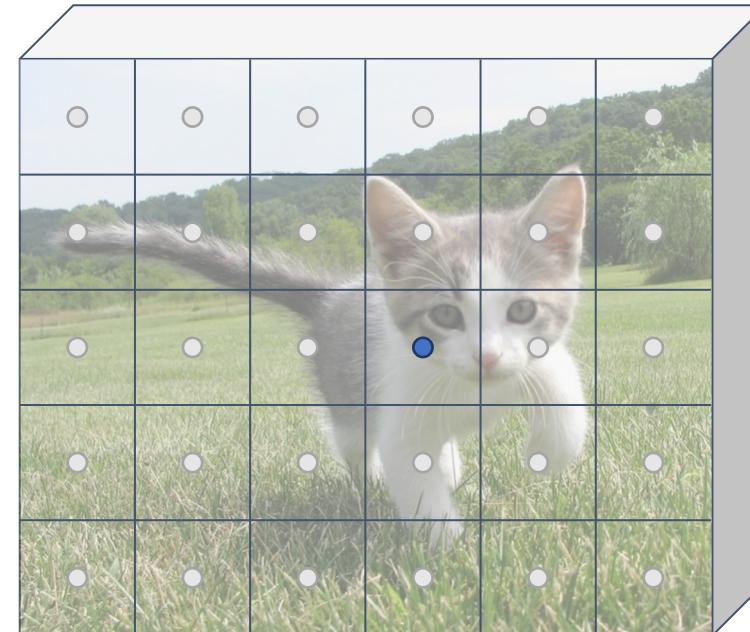
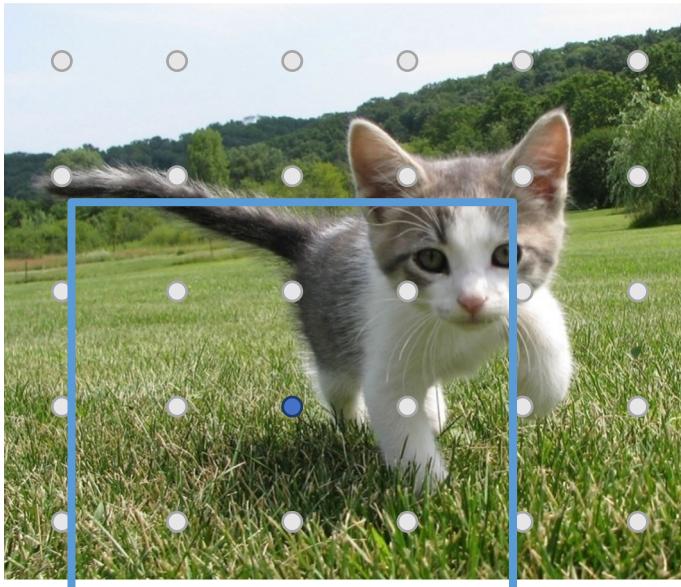


Image features  
(e.g.  $512 \times 5 \times 6$ )

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

# Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



Input Image  
(e.g.  $3 \times 640 \times 480$ )



Each feature corresponds to a point in the input

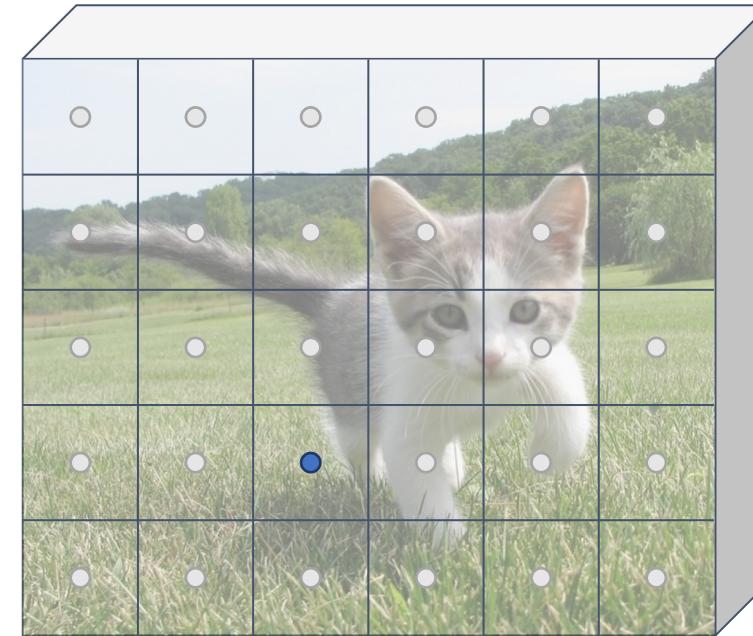
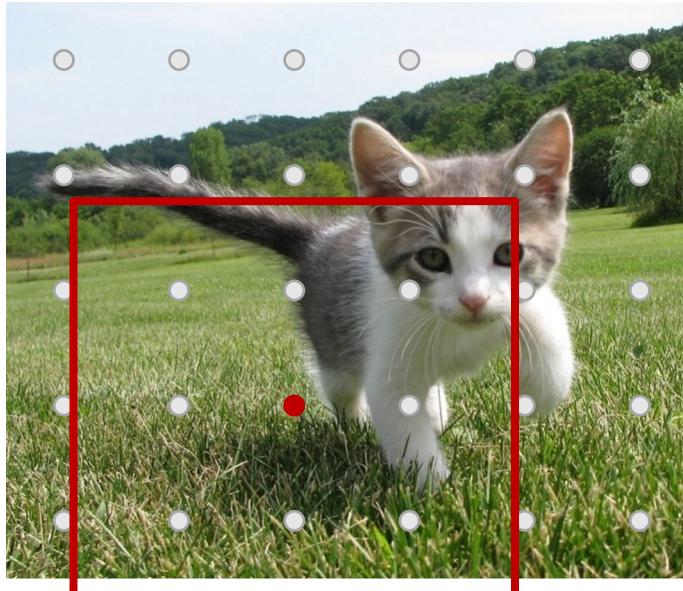


Image features  
(e.g.  $512 \times 5 \times 6$ )

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

# Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



Input Image  
(e.g.  $3 \times 640 \times 480$ )



Each feature corresponds to a point in the input

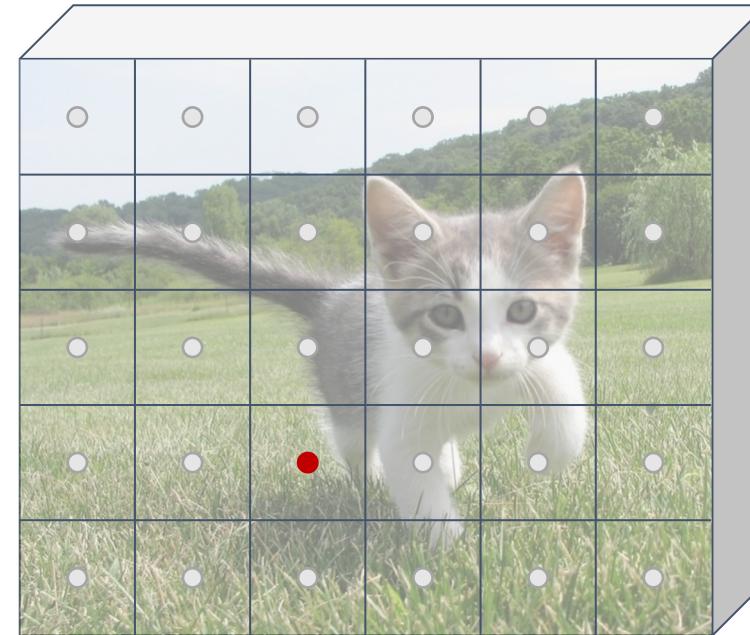


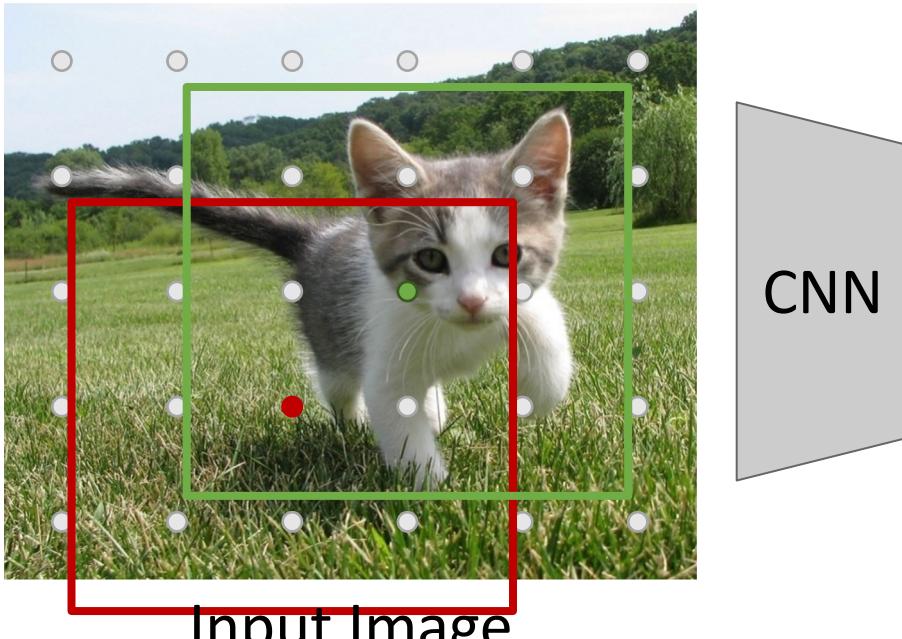
Image features  
(e.g.  $512 \times 5 \times 6$ )

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

Classify each anchor as positive (object) or negative (no object)

# Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



Input Image  
(e.g.  $3 \times 640 \times 480$ )

Each feature corresponds to a point in the input

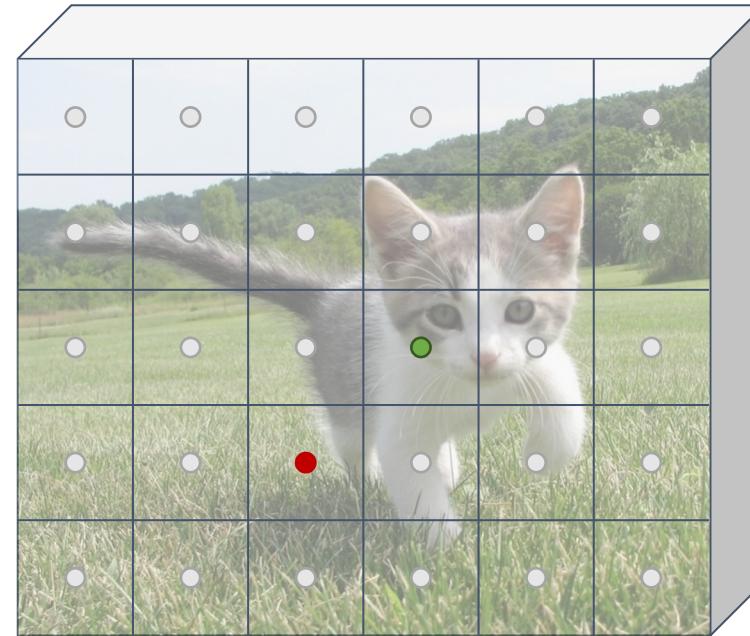


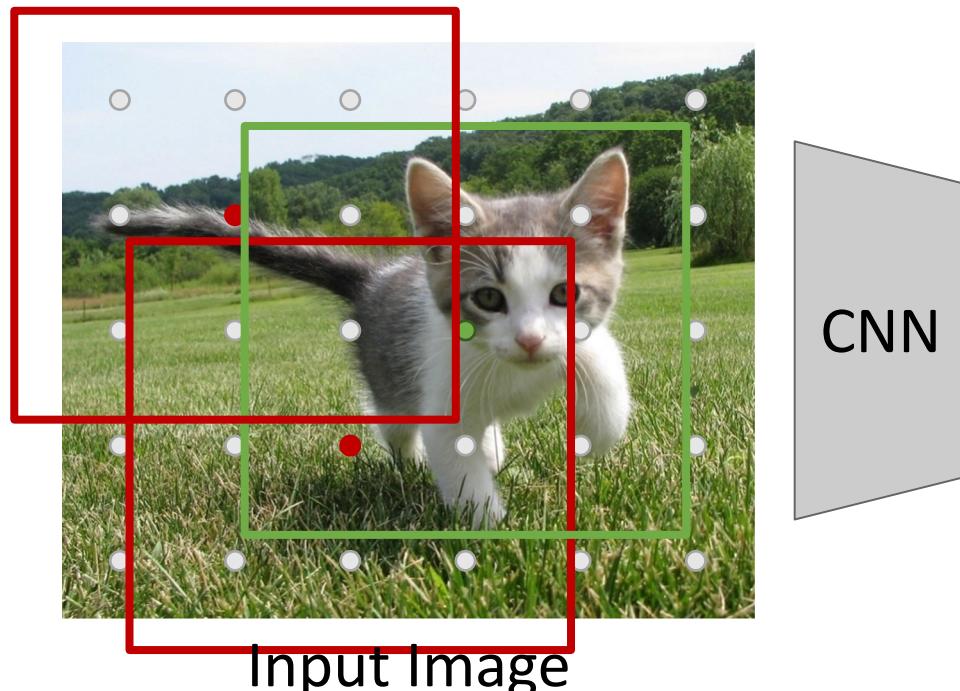
Image features  
(e.g.  $512 \times 5 \times 6$ )

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

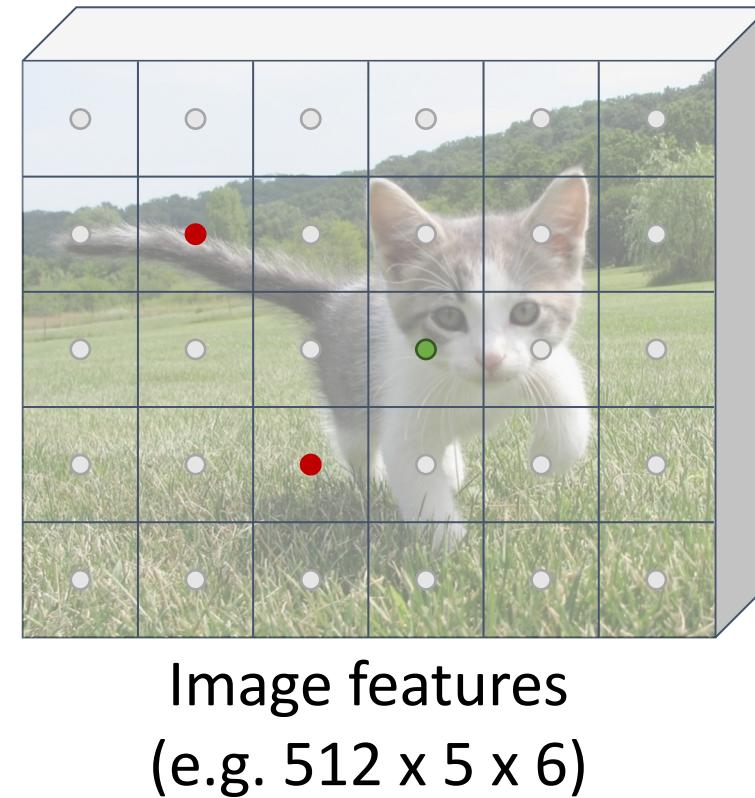
Classify each anchor as positive (object) or negative (no object)

# Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



Each feature corresponds to a point in the input

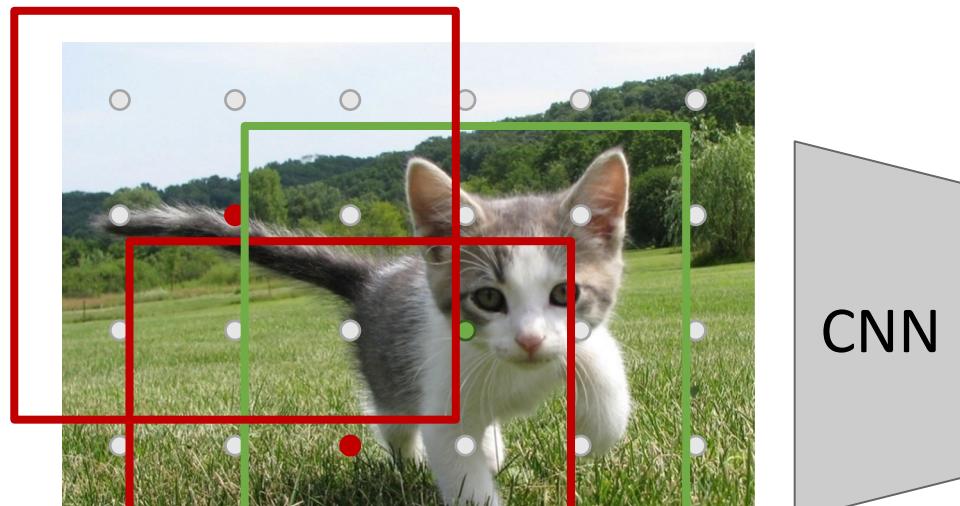


Imagine an **anchor box** of fixed size at each point in the feature map

Classify each anchor as positive (object) or negative (no object)

# Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



Input Image  
(e.g.  $3 \times 640 \times 480$ )

Each feature corresponds to a point in the input

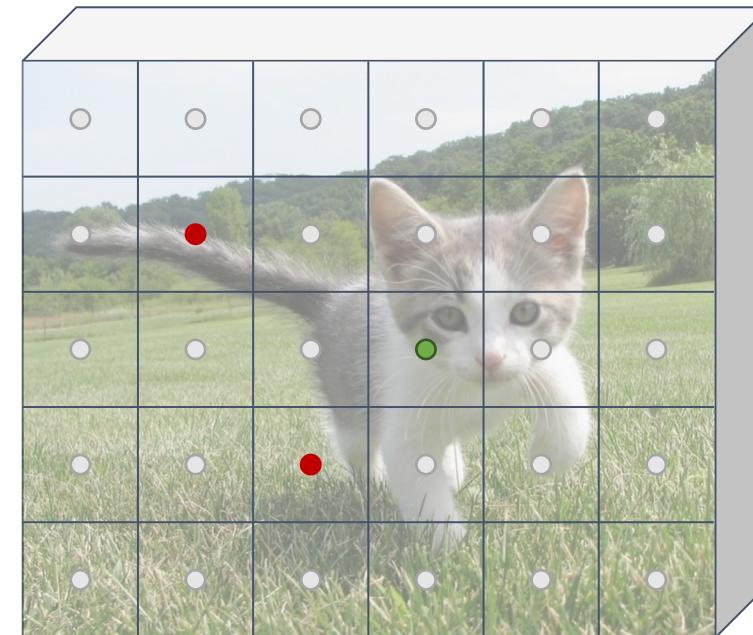
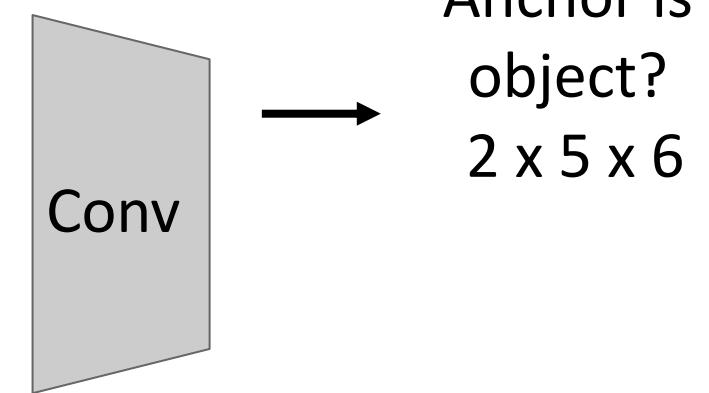


Image features  
(e.g.  $512 \times 5 \times 6$ )

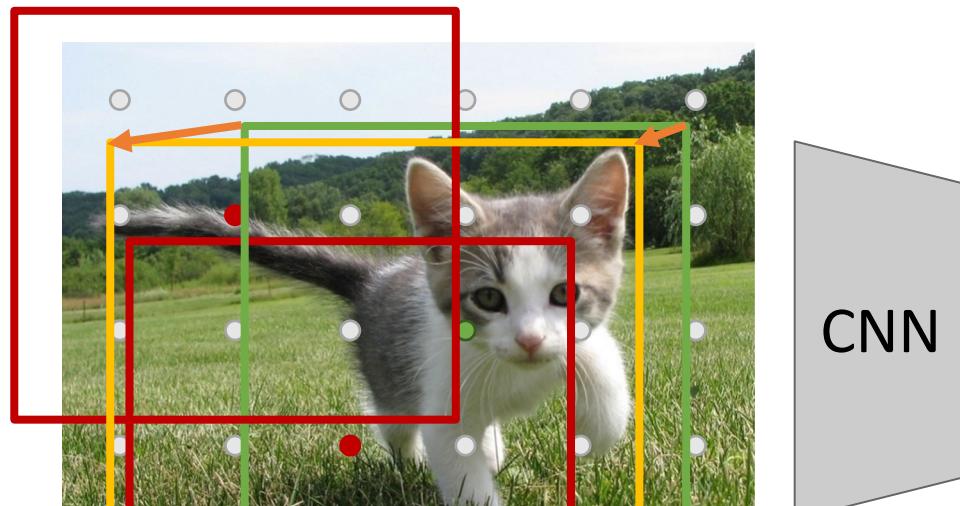
Predict object vs not object scores for all anchors with a conv layer (512 input filters, 2 output filters)



Classify each anchor as positive (object) or negative (no object)

# Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



Input Image  
(e.g.  $3 \times 640 \times 480$ )

Each feature corresponds to a point in the input

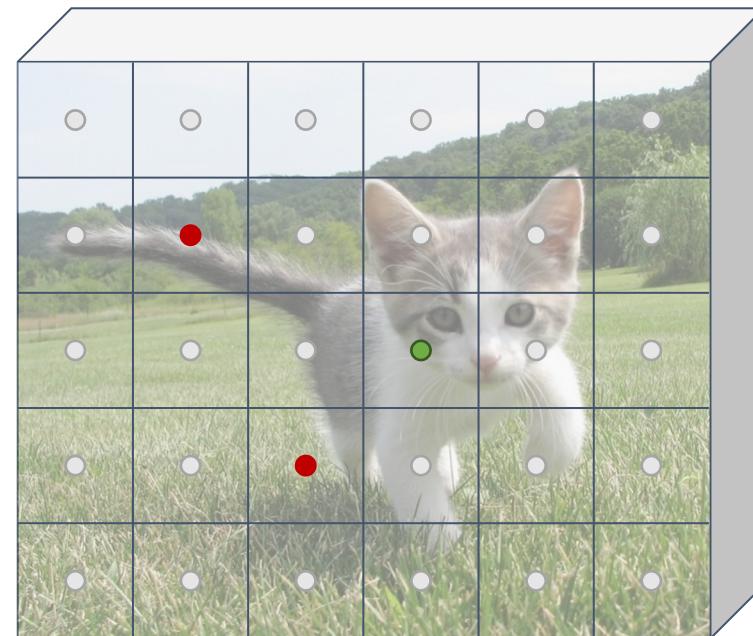
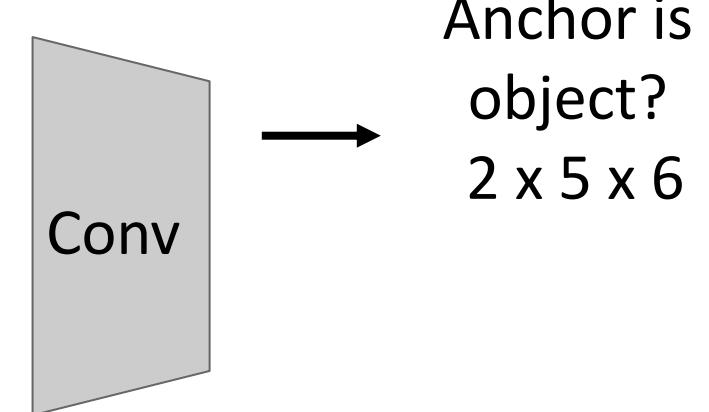


Image features  
(e.g.  $512 \times 5 \times 6$ )

For **positive anchors**, also predict a **transform** that converting the anchor to the **GT box** (like R-CNN)

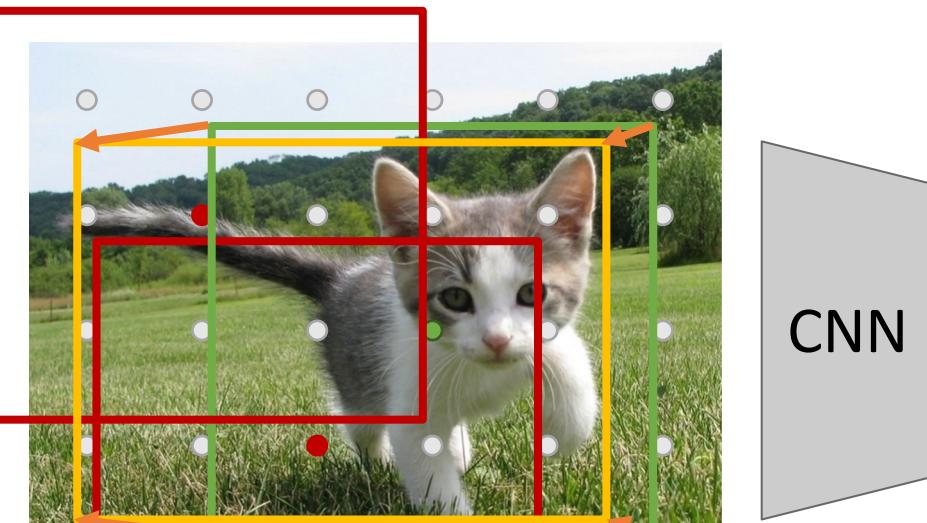


Anchor is object?  
 $2 \times 5 \times 6$

Classify each anchor as **positive (object)** or **negative (no object)**

# Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



Input Image  
(e.g.  $3 \times 640 \times 480$ )

Each feature corresponds to a point in the input

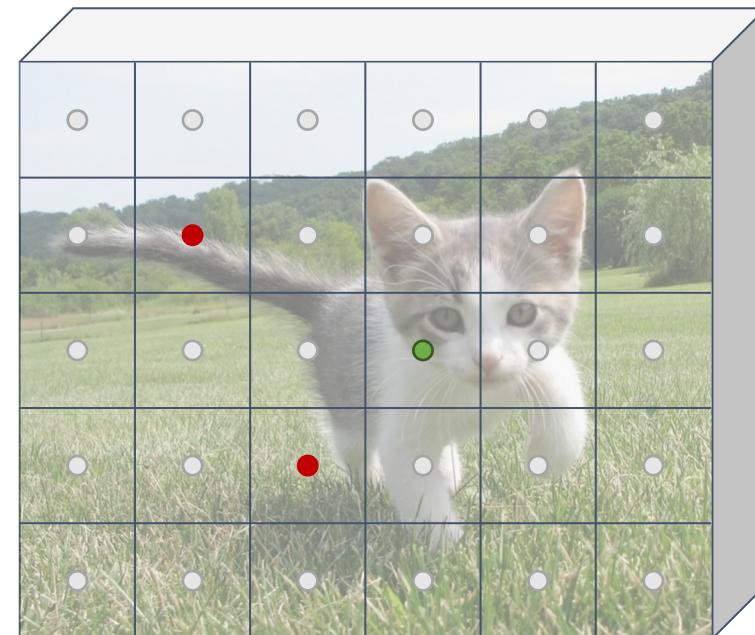
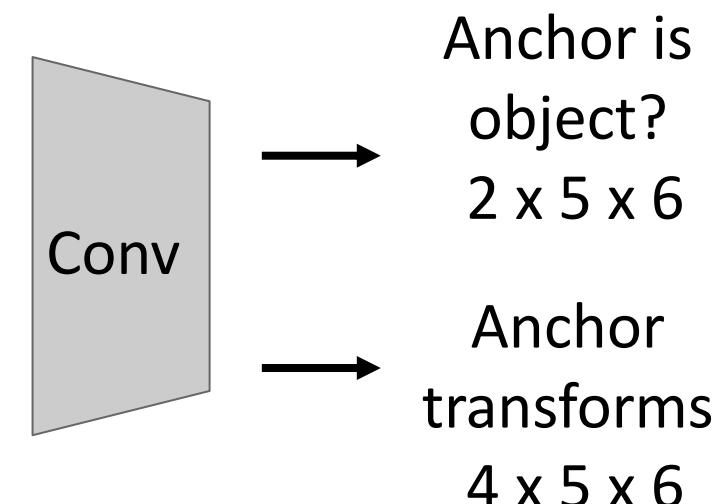


Image features  
(e.g.  $512 \times 5 \times 6$ )

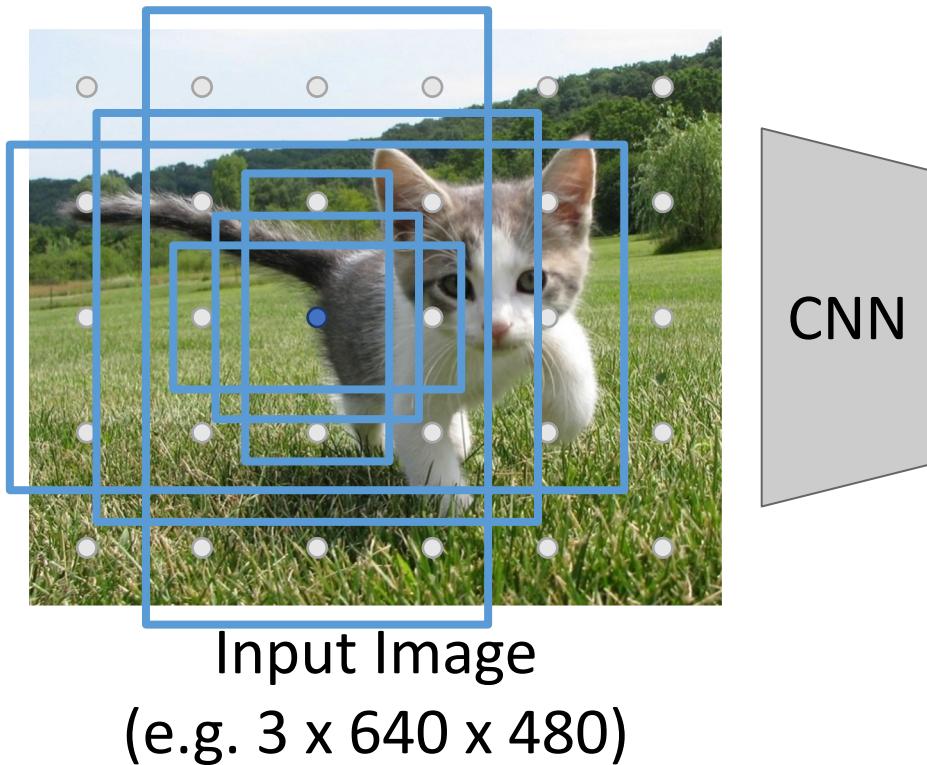
For **positive anchors**, also predict a **transform** that converting the anchor to the **GT box** (like R-CNN)  
Predict transforms with conv



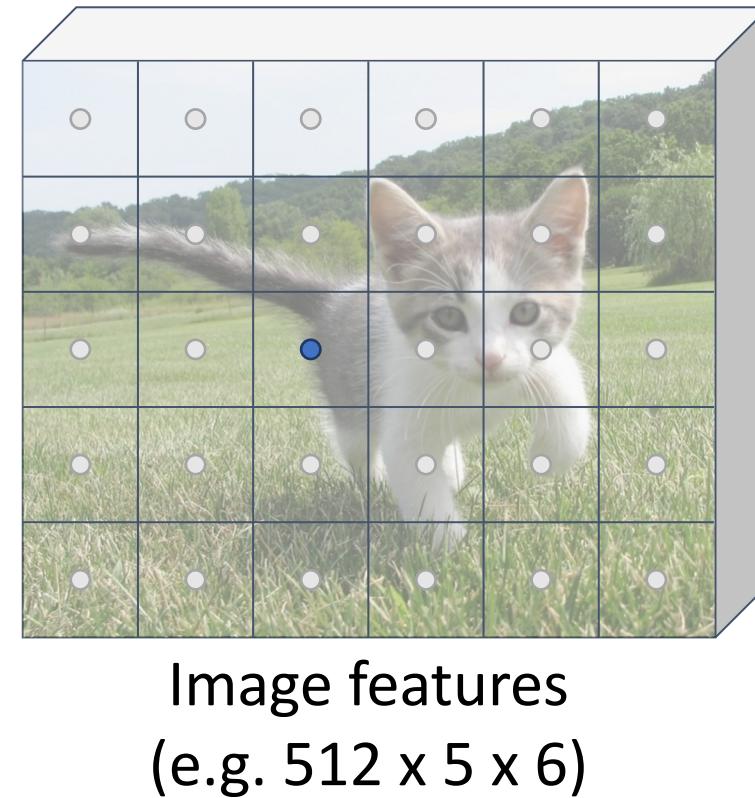
Classify each anchor as **positive (object)** or **negative (no object)**

# Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



Each feature corresponds to a point in the input



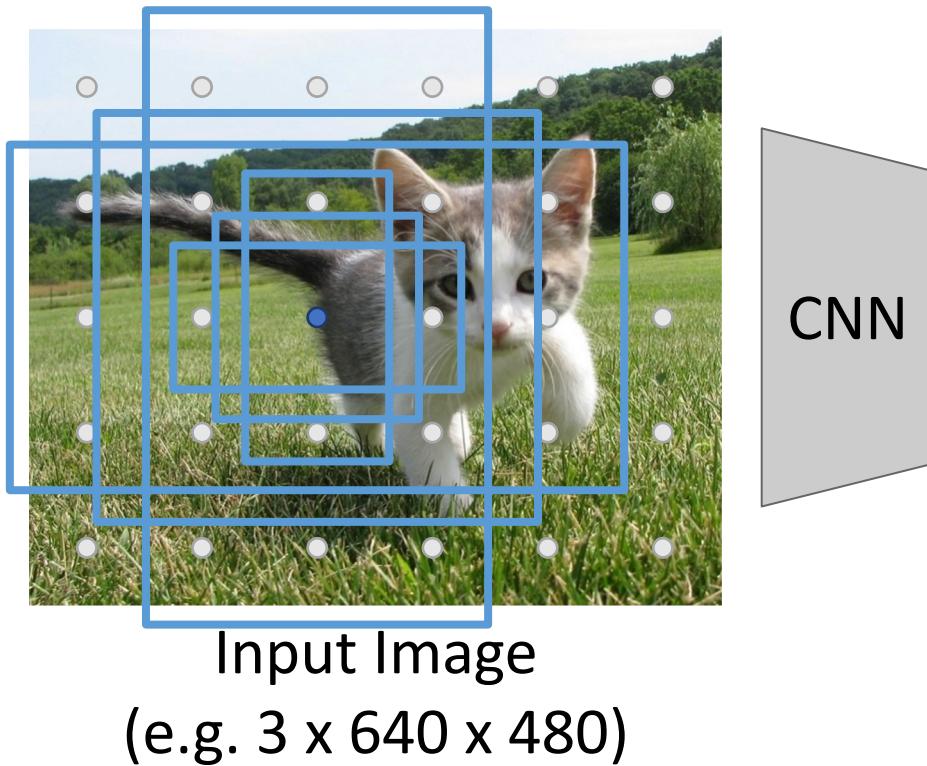
In practice: Rather than using one anchor per point, instead consider K different anchors with different size and scale (here  $K = 6$ )

Anchor is object?  
 $2K \times 5 \times 6$

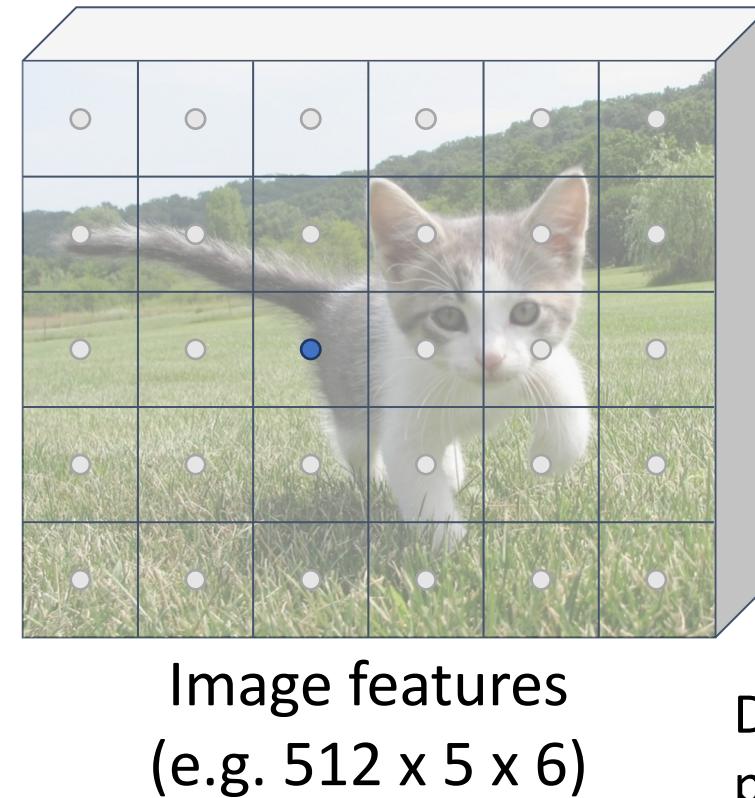
Anchor transforms  
 $4K \times 5 \times 6$

# 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  $K = 6$ )

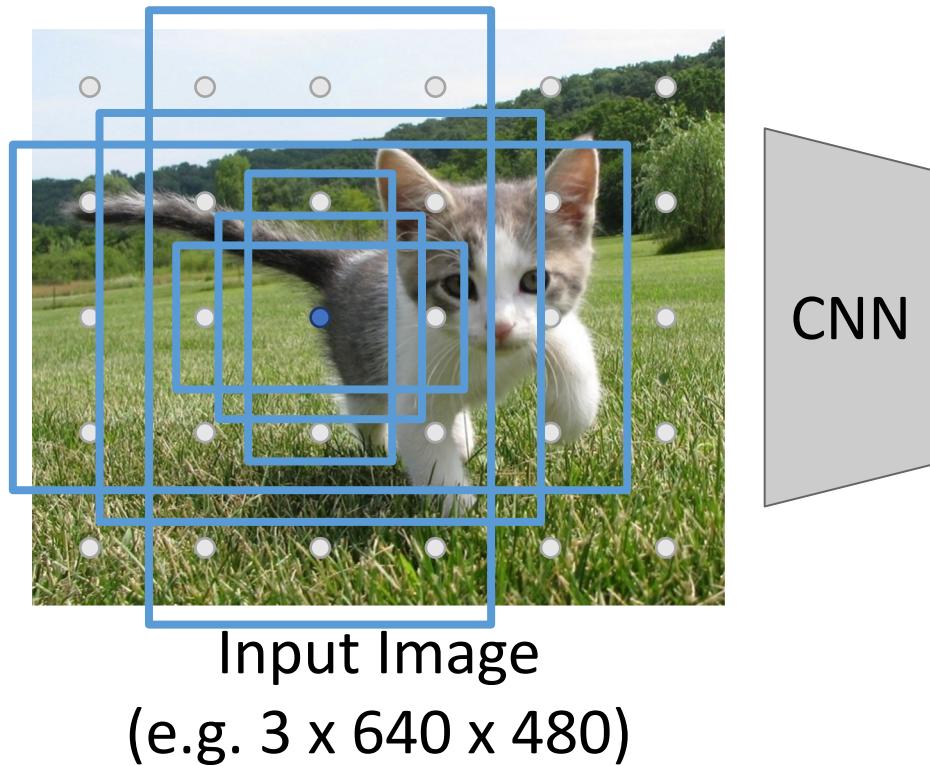
Anchor is object?  
 $2K \times 5 \times 6$

Anchor transforms  
 $4K \times 5 \times 6$

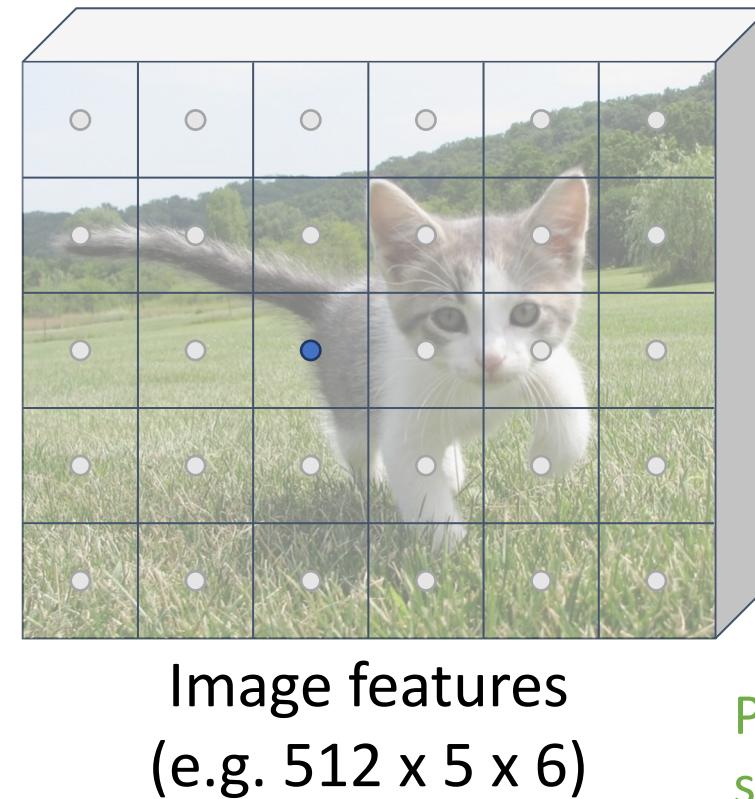
During training, supervised positive / negative anchors and box transforms like R-CNN

# Region Proposal Network (RPN)

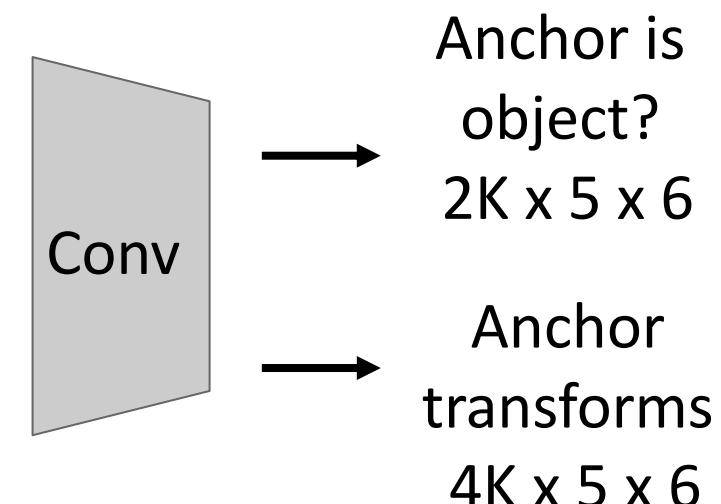
Run backbone CNN to get features aligned to input image



Each feature corresponds to a point in the input



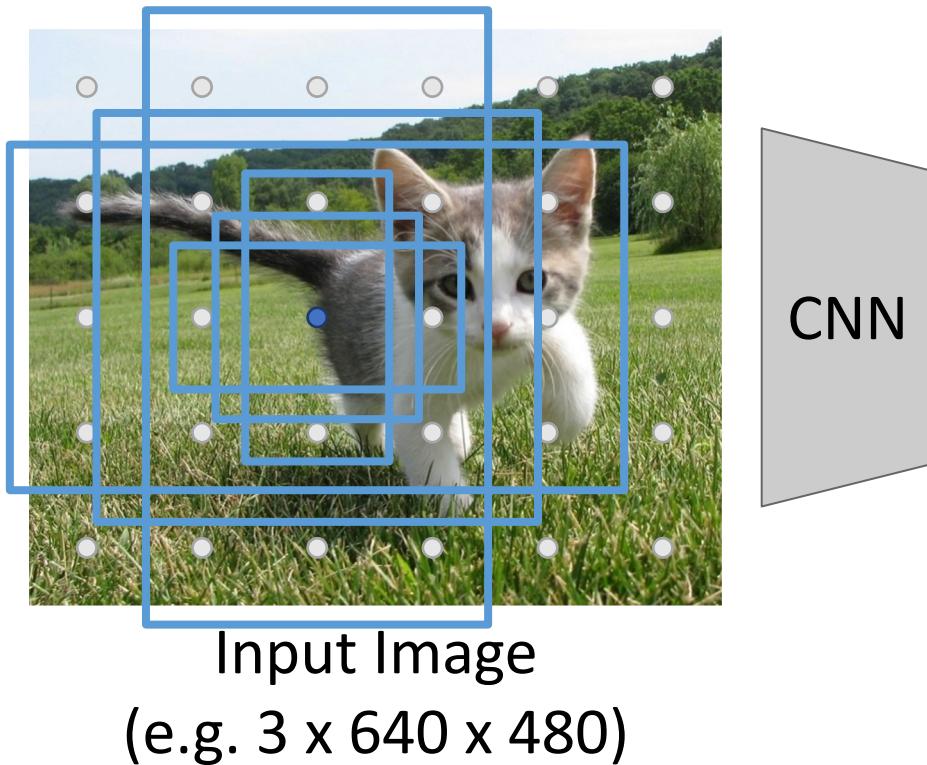
In practice: Rather than using one anchor per point, instead consider K different anchors with different size and scale (here  $K = 6$ )



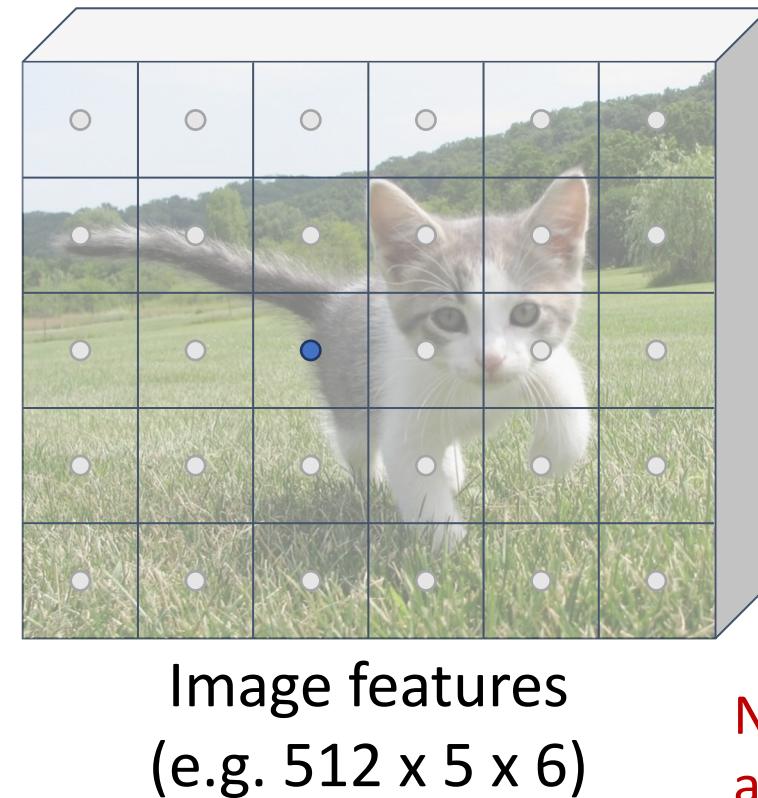
Positive anchors:  $\geq 0.7$  IoU with some GT box (plus highest IoU to each GT)

# 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  $K = 6$ )

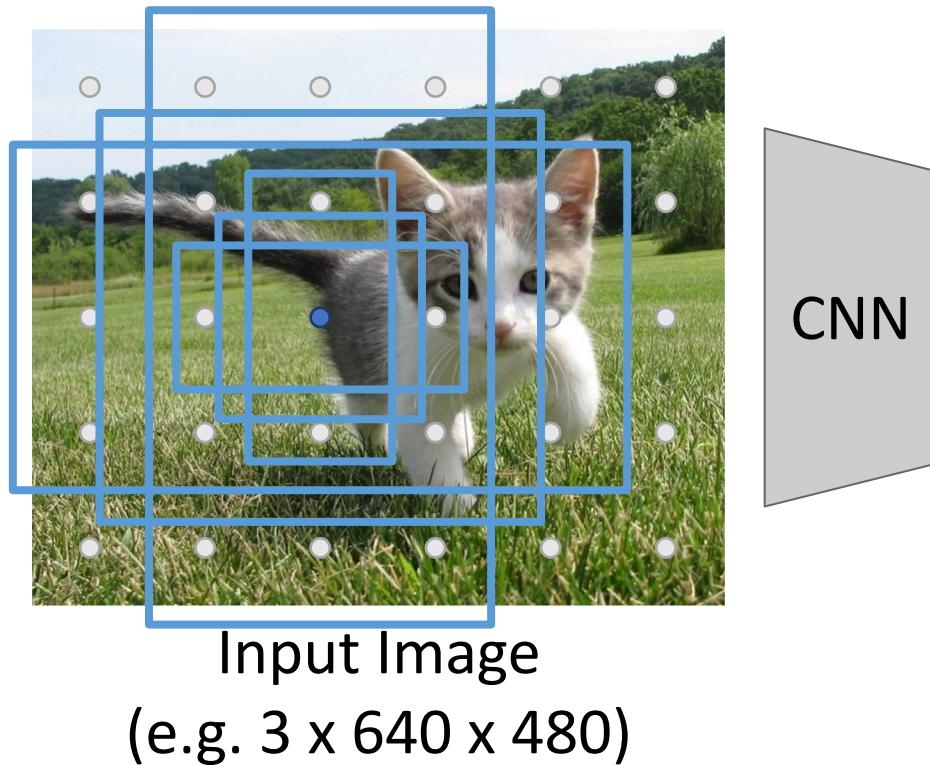
Anchor is object?  
 $2K \times 5 \times 6$

Anchor transforms  
 $4K \times 5 \times 6$

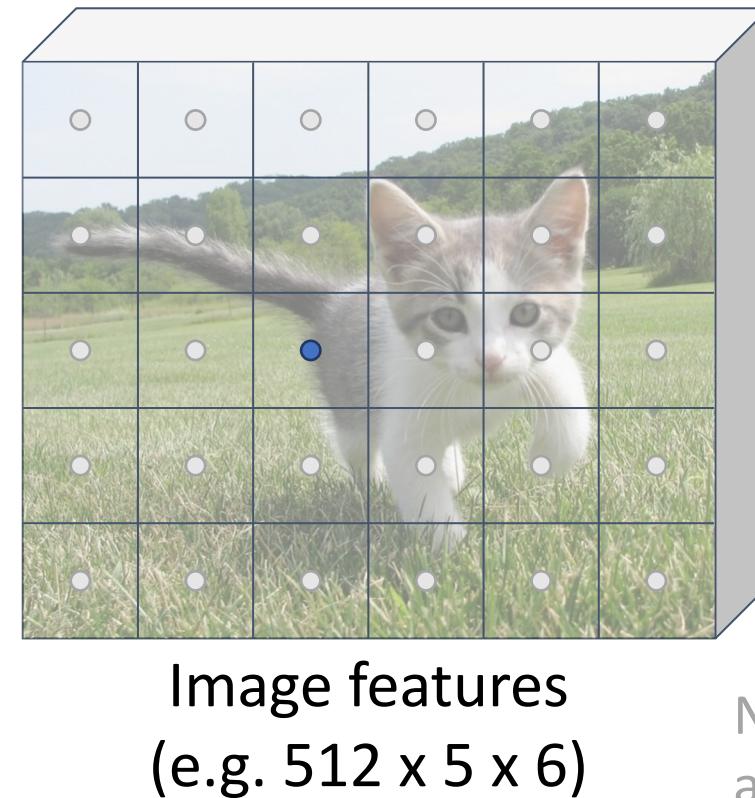
Negative anchors:  $< 0.3$  IoU with all GT boxes. Don't supervise transforms for negative boxes.

# Region Proposal Network (RPN)

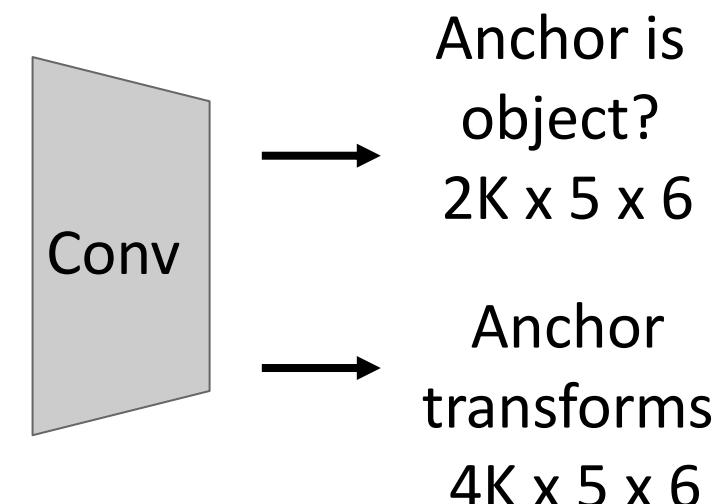
Run backbone CNN to get features aligned to input image



Each feature corresponds to a point in the input



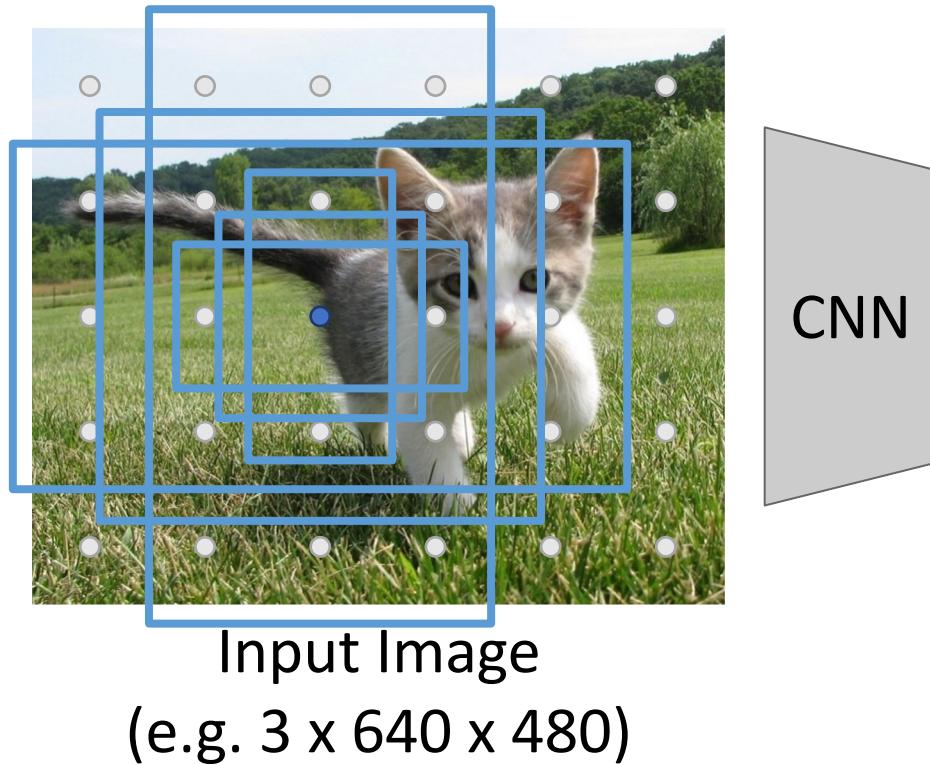
In practice: Rather than using one anchor per point, instead consider K different anchors with different size and scale (here  $K = 6$ )



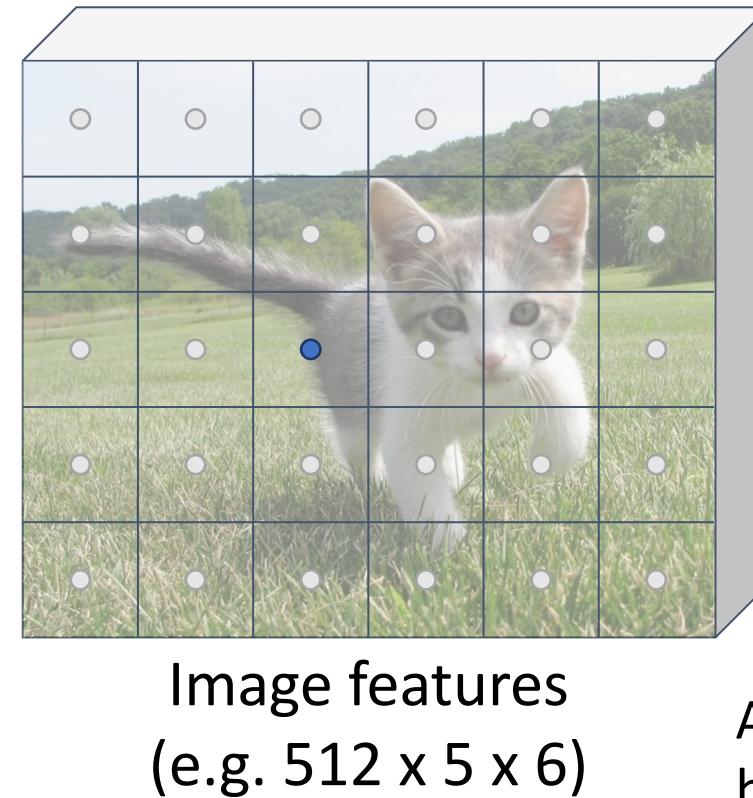
Neutral anchors: between 0.3 and 0.7 IoU with all GT boxes; ignored during training

# Region Proposal Network (RPN)

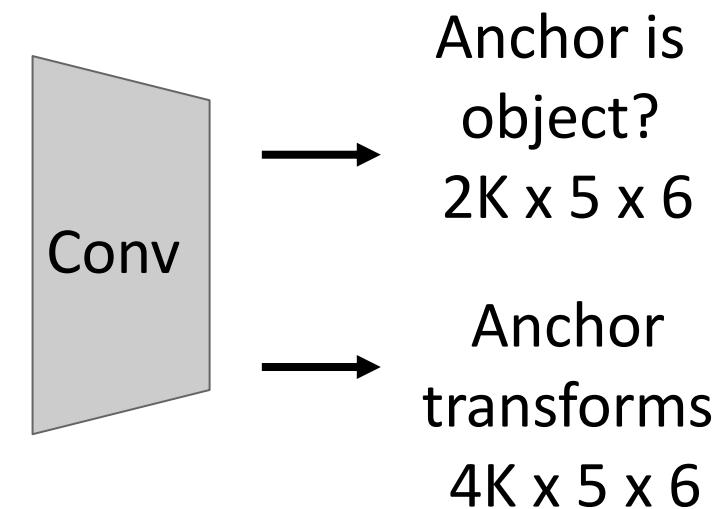
Run backbone CNN to get features aligned to input image



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  $K = 6$ )

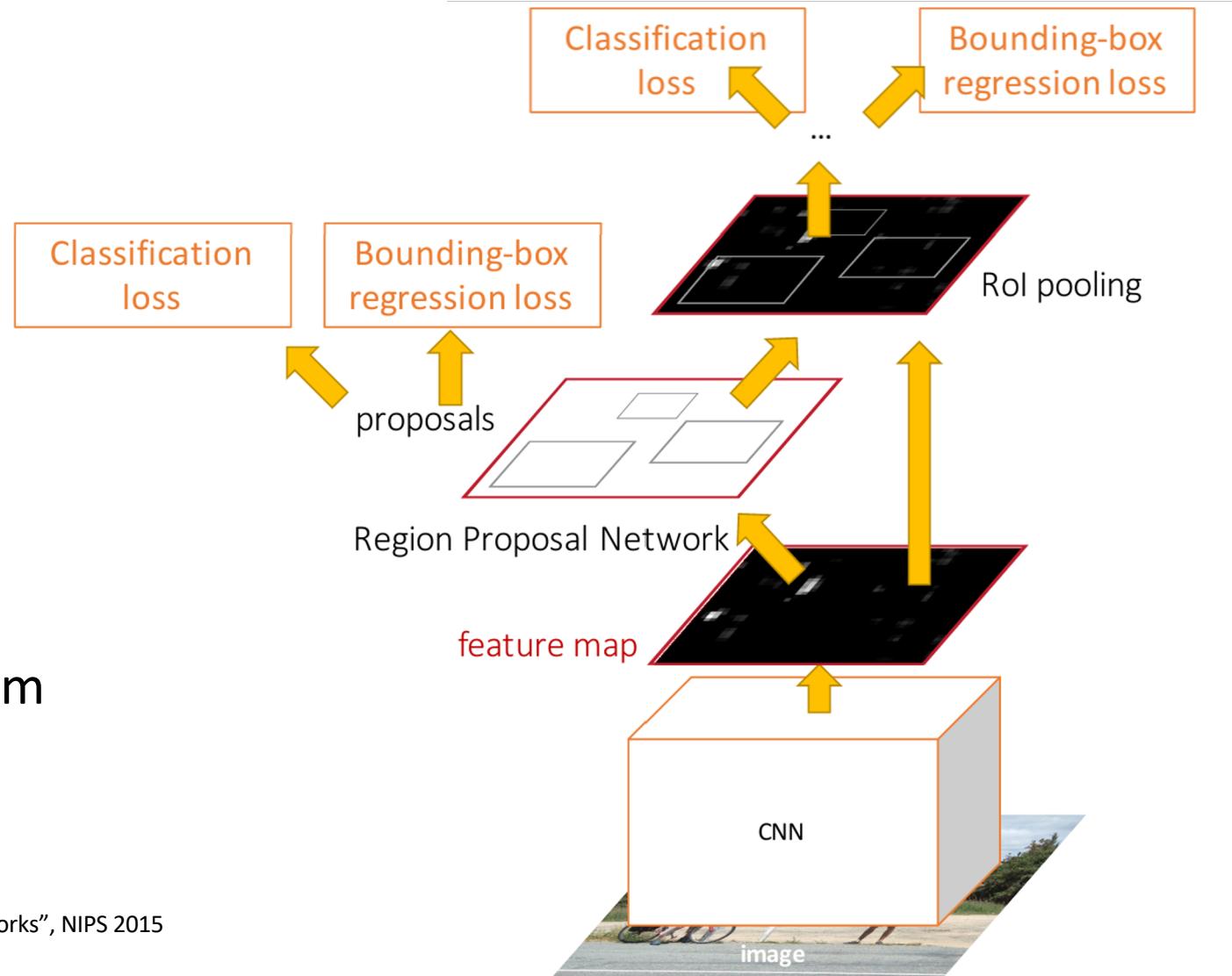


At test-time, sort all  $K^*5^*6$  boxes by their positive score, take top 300 as our region proposals

# FasterR-CNN: Learnable Region Proposals

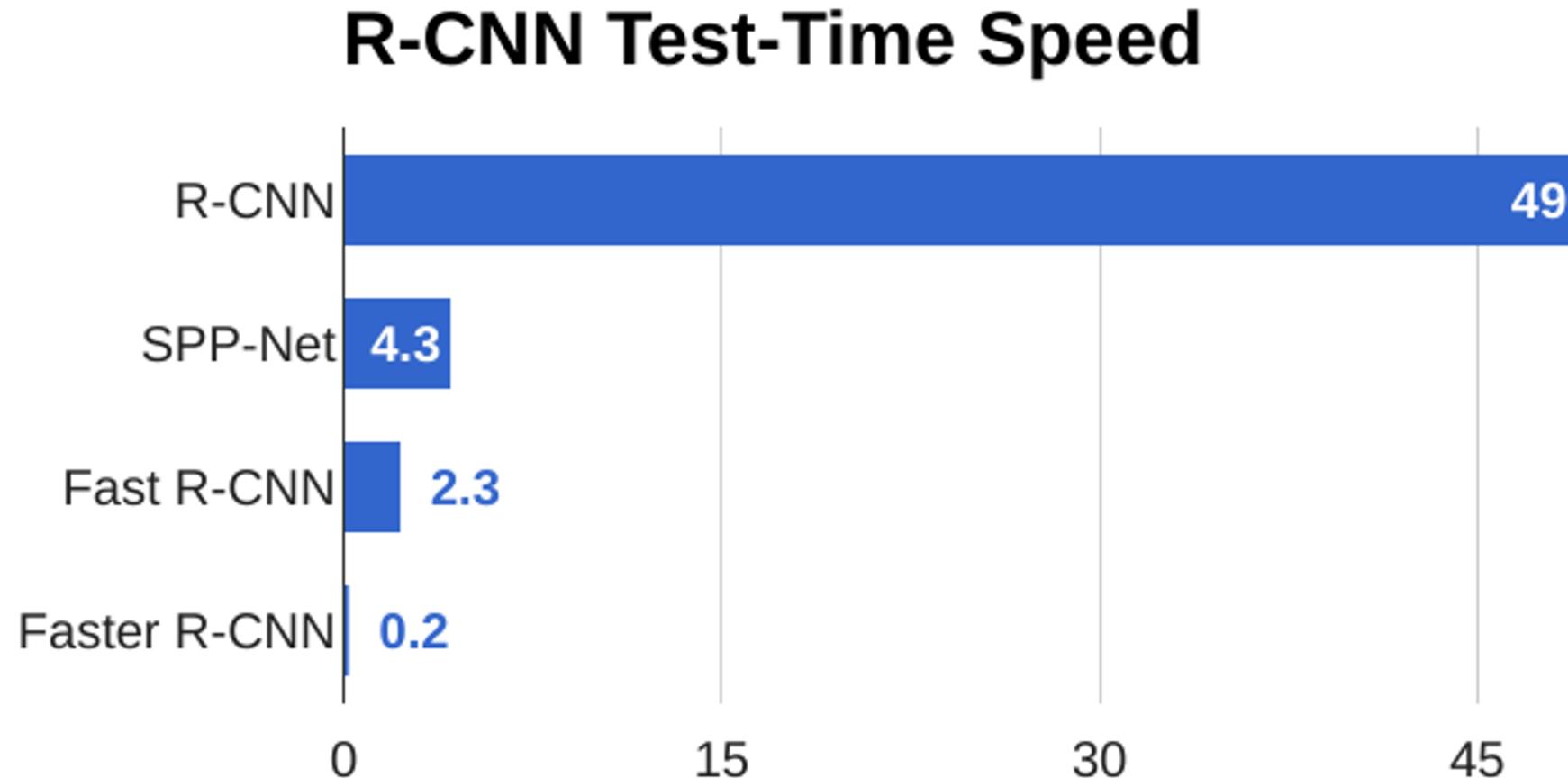
Jointly train with 4 losses:

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

# FasterR-CNN: Learnable Region Proposals



# Faster R-CNN: Learnable Region Proposals

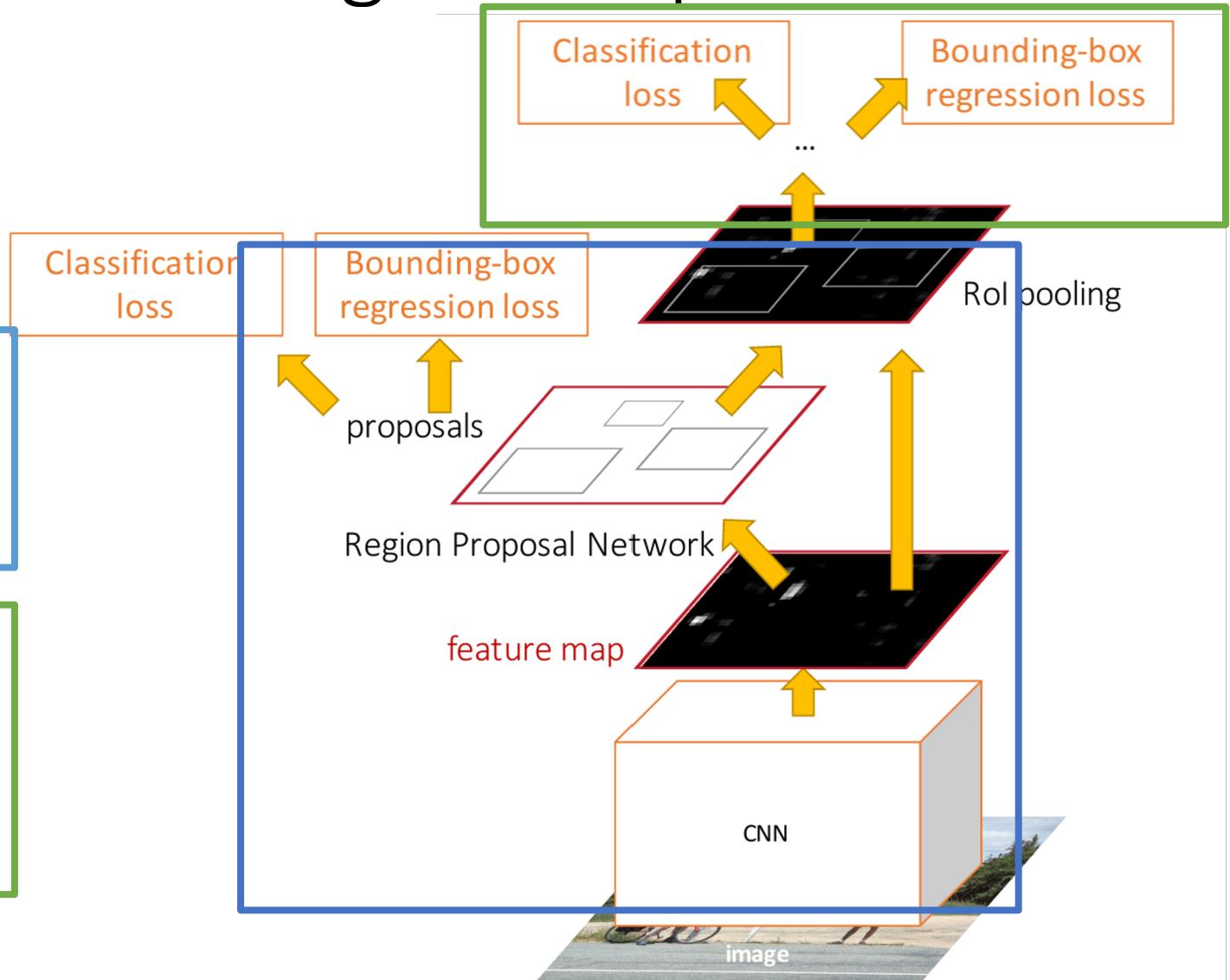
Faster R-CNN is a  
**Two-stage object detector**

First stage: Run once per image

- Backbone network
- Region proposal network

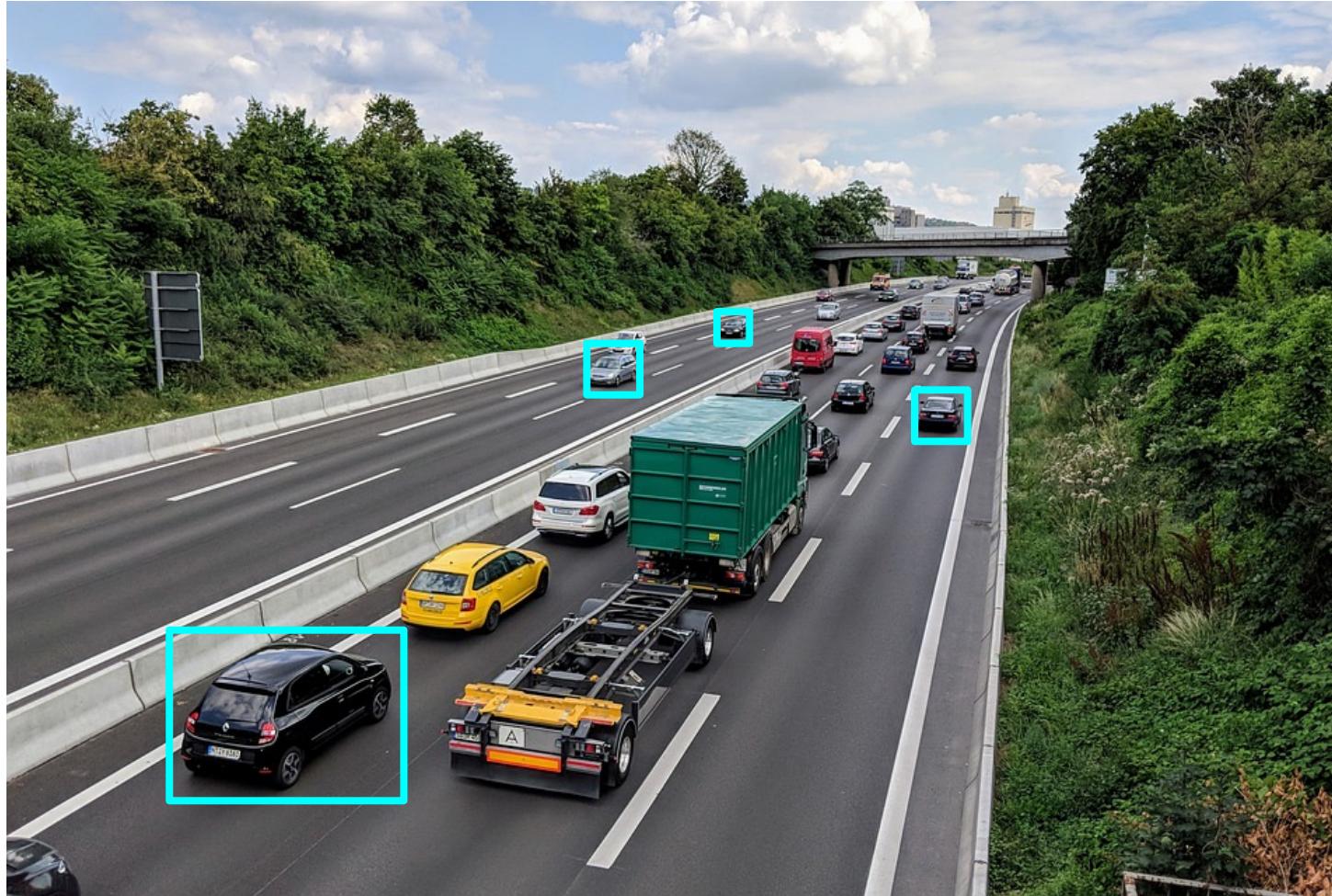
Second stage: Run once per region

- Crop features: RoI pool / align
- Predict object class
- Prediction bbox offset



# Dealing with Scale

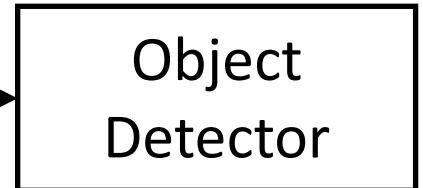
We need to detect objects of many different scales.  
How to improve *scale invariance* of the detector?



This image is free for commercial  
use under the [Pixabay license](#)

# Dealing with Scale: Image Pyramid

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

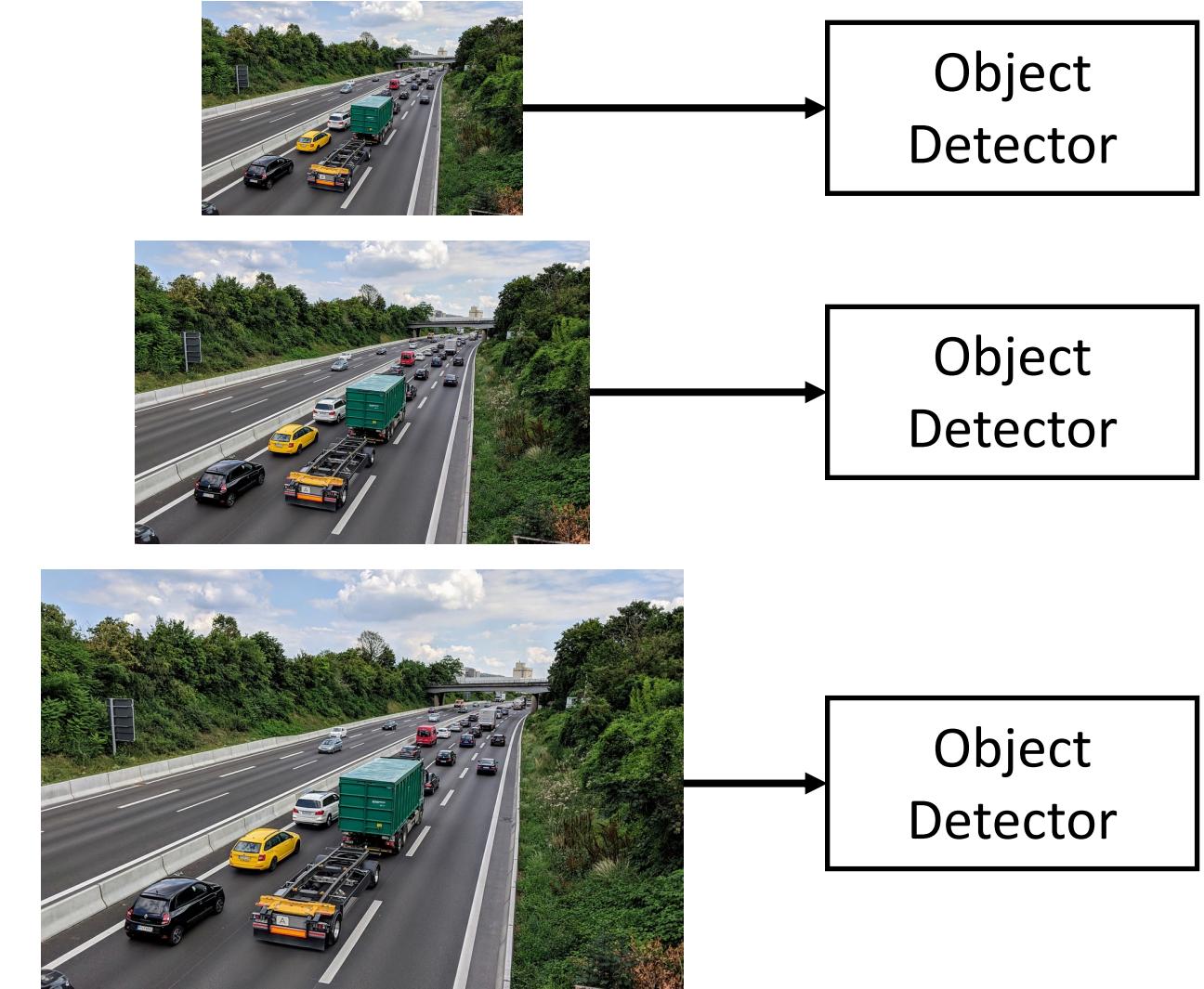


Lin et al, "Feature Pyramid Networks for Object Detection", ICCV 2017

# Dealing with Scale: Image Pyramid

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

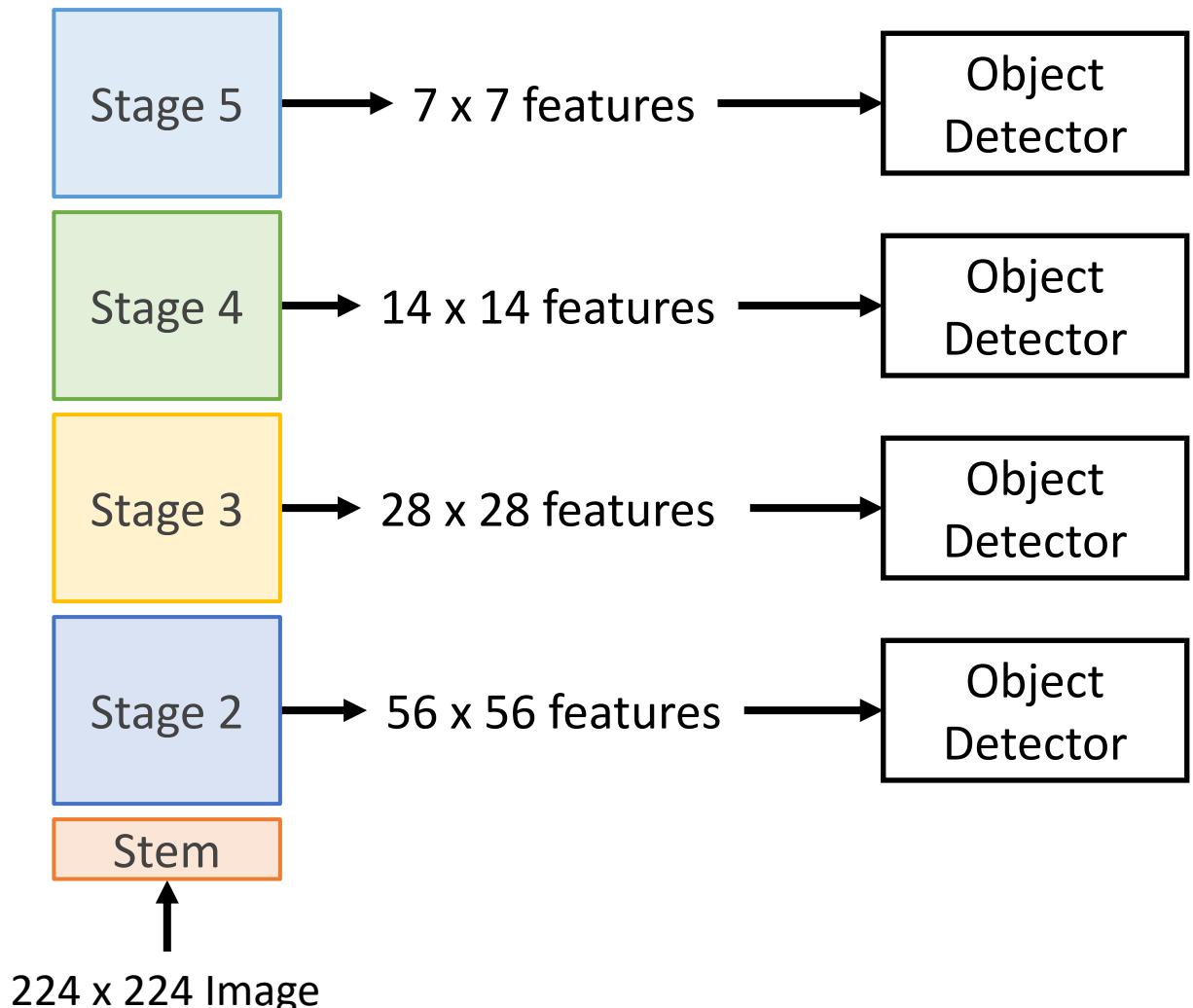
Problem: Expensive! Don't share any computation between scales



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

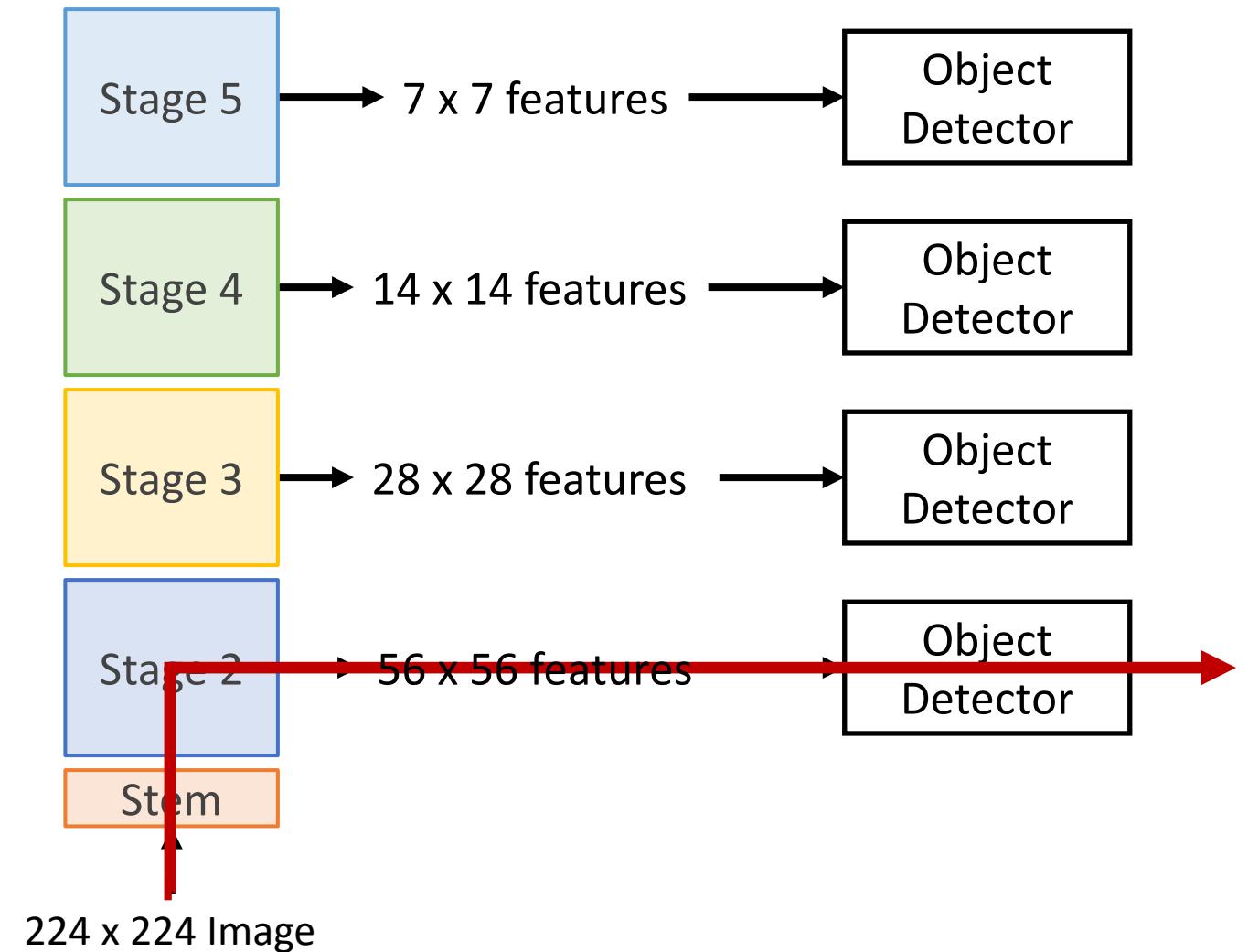


Lin et al, "Feature Pyramid Networks for Object Detection", ICCV 2017

# Dealing with Scale: Multiscale Features

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

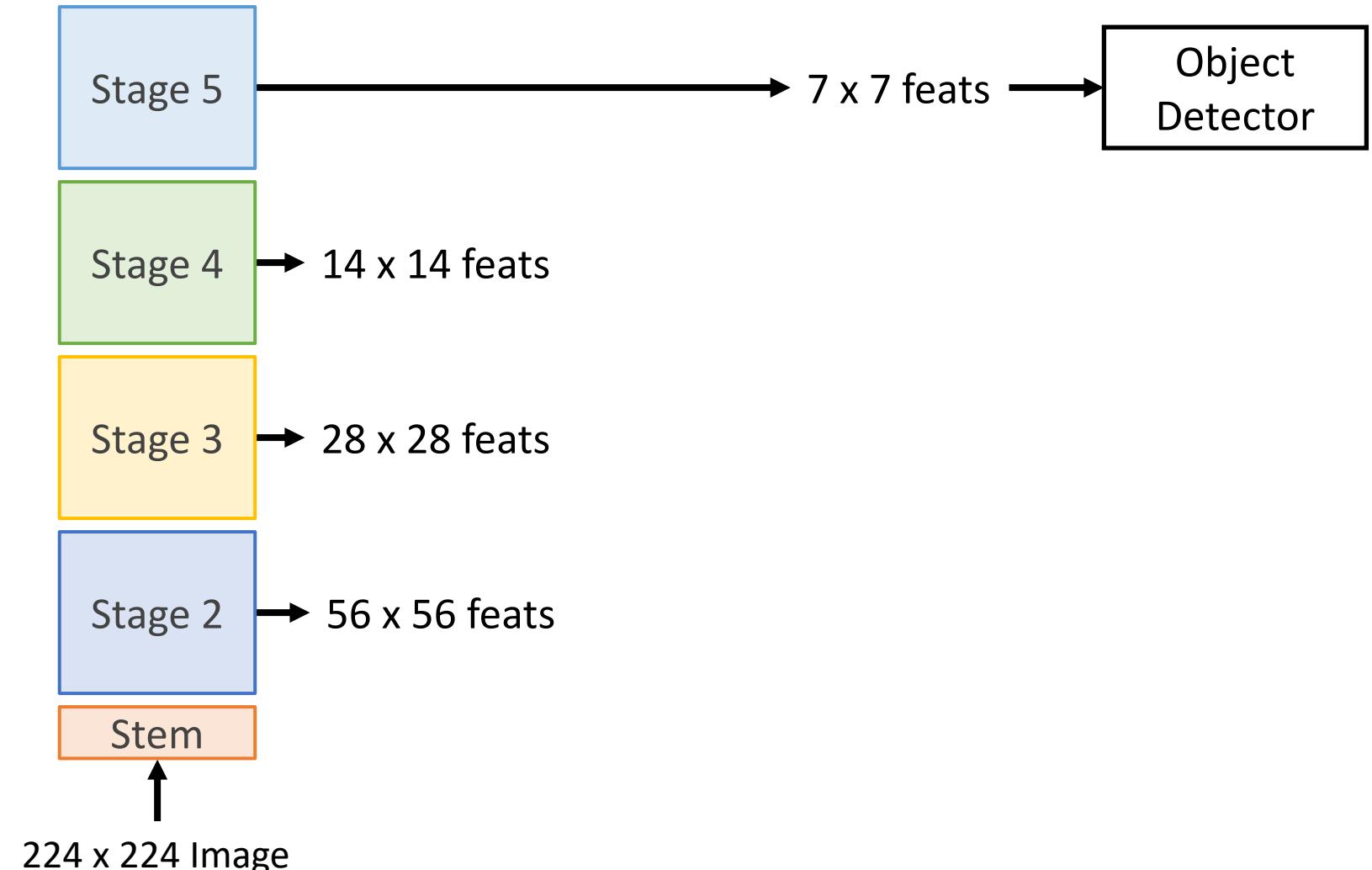
Problem: detector on early features doesn't make use of the entire backbone; doesn't get access to high-level features



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

# Dealing with Scale: Feature Pyramid Network

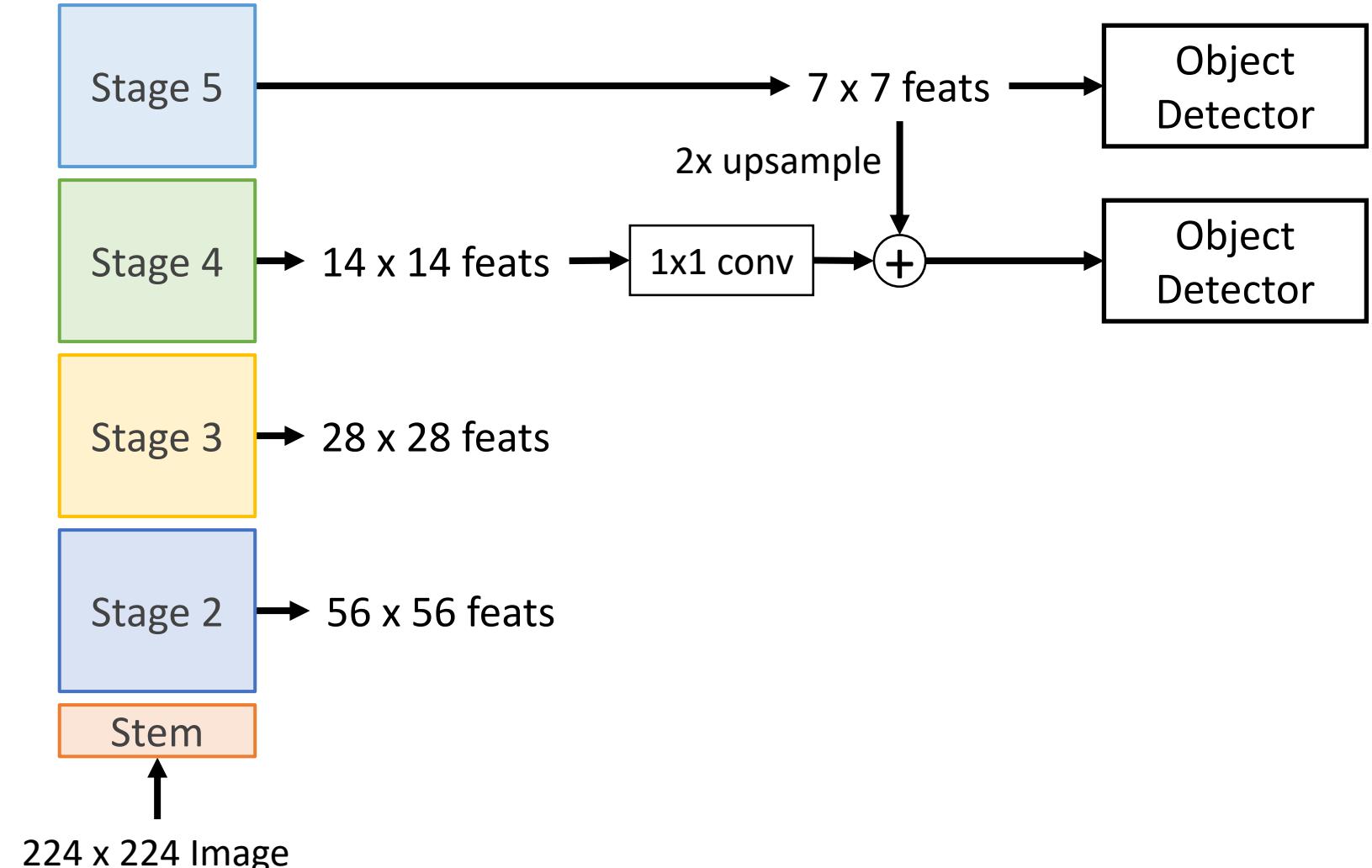
Add *top down connections* that feed information from high level features back down to lower level features



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

# Dealing with Scale: Feature Pyramid Network

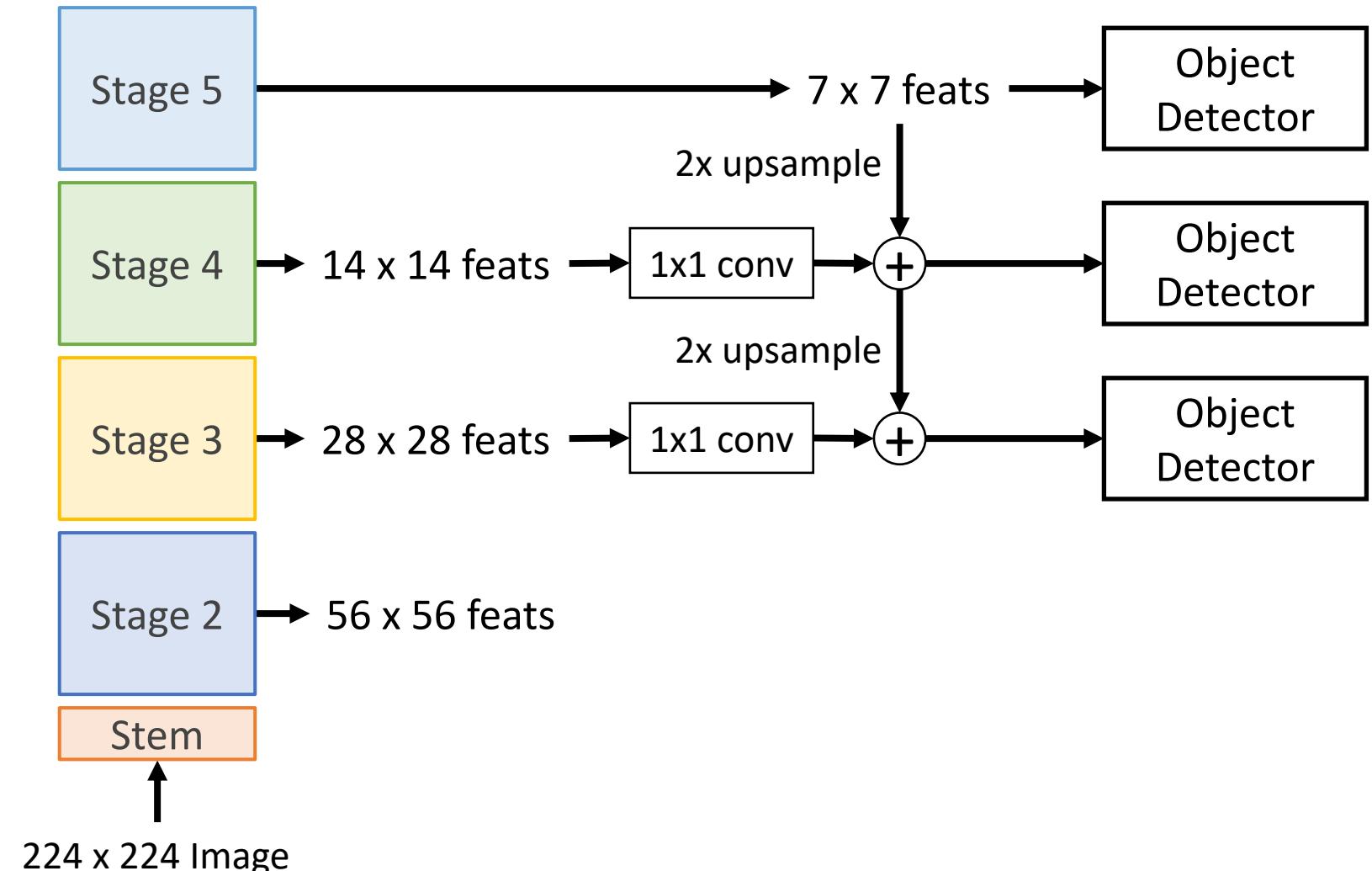
Add *top down connections* that feed information from high level features back down to lower level features



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

# Dealing with Scale: Feature Pyramid Network

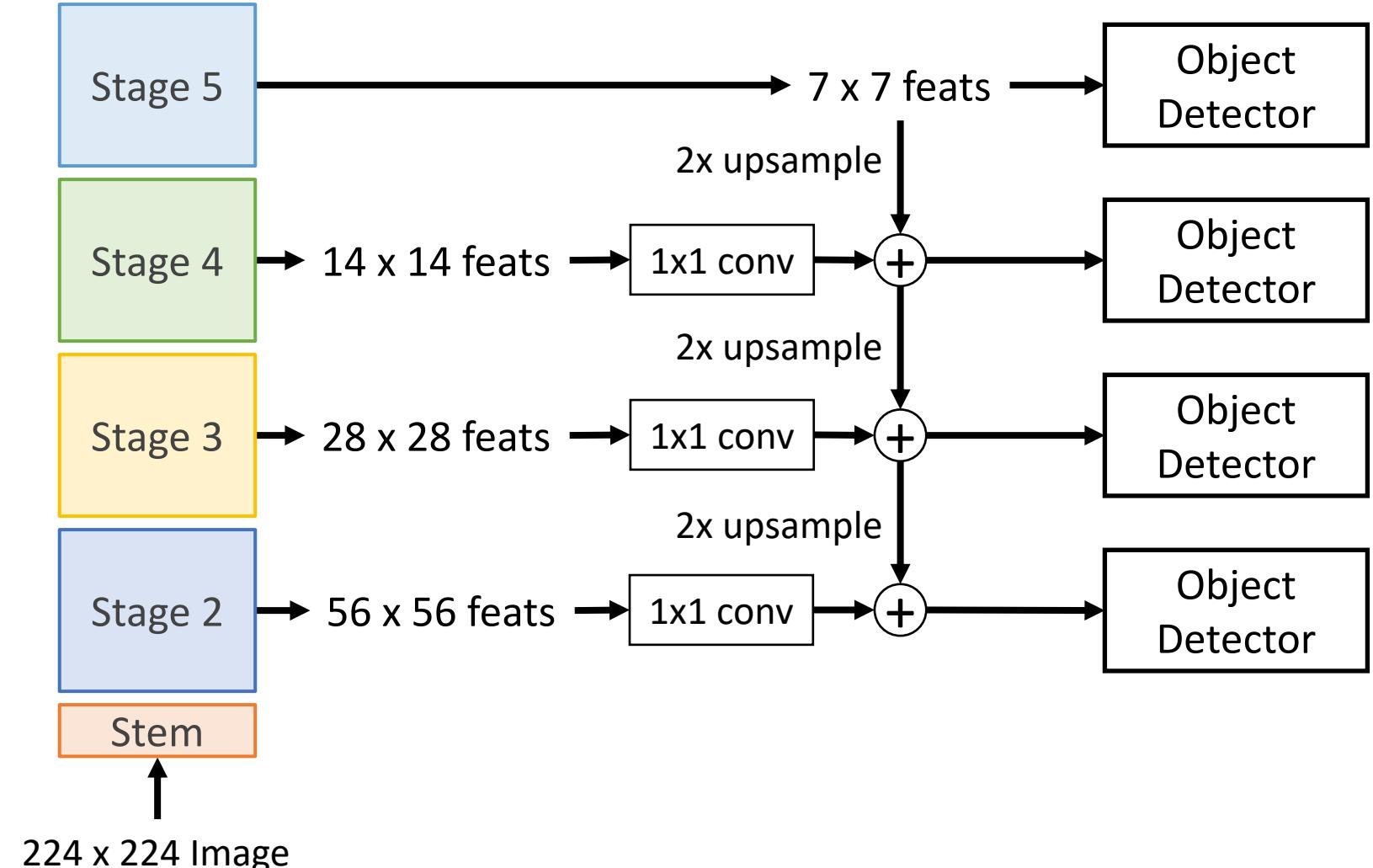
Add *top down connections* that feed information from high level features back down to lower level features



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

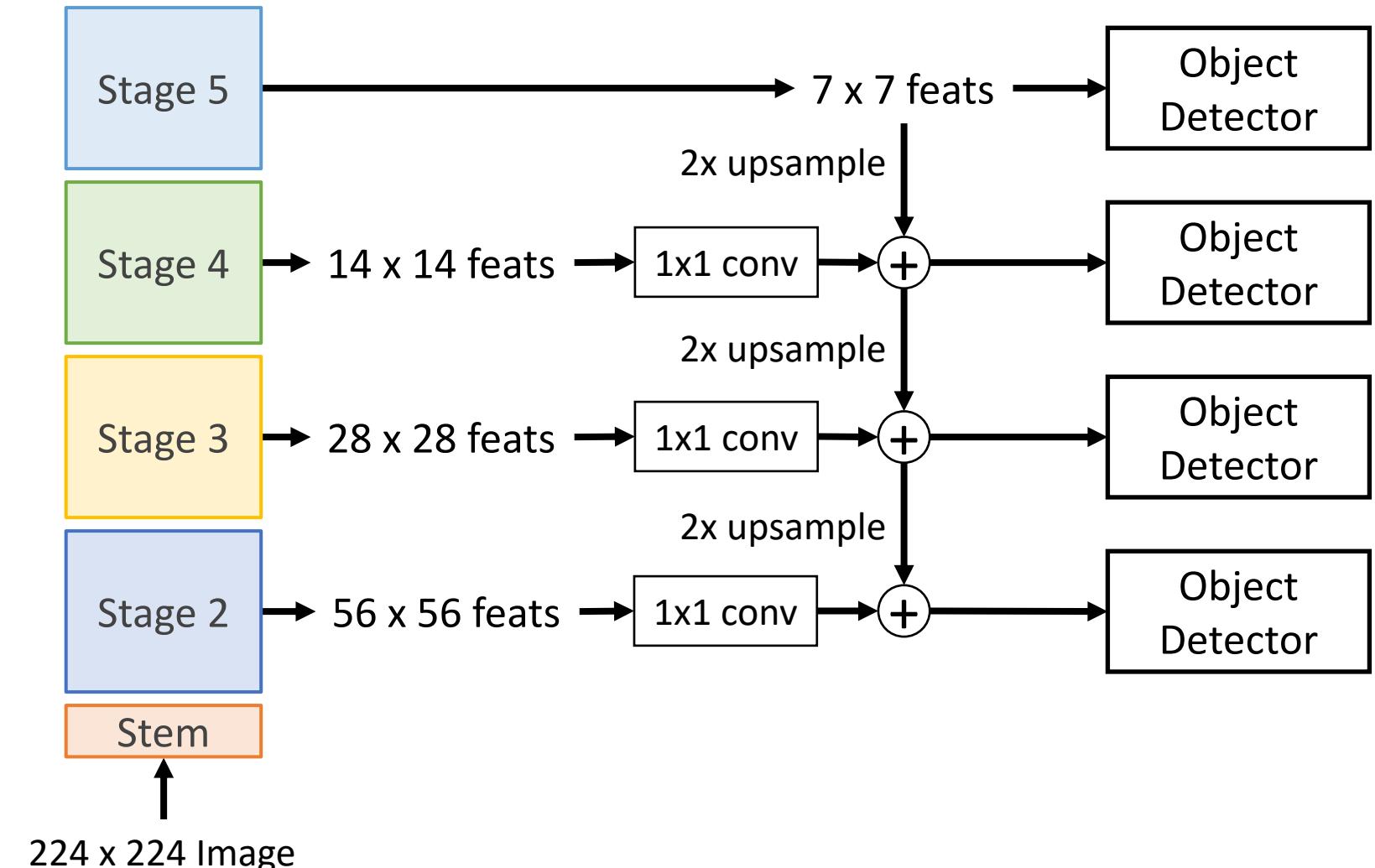


Lin et al, "Feature Pyramid Networks for Object Detection", ICCV 2017

# 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

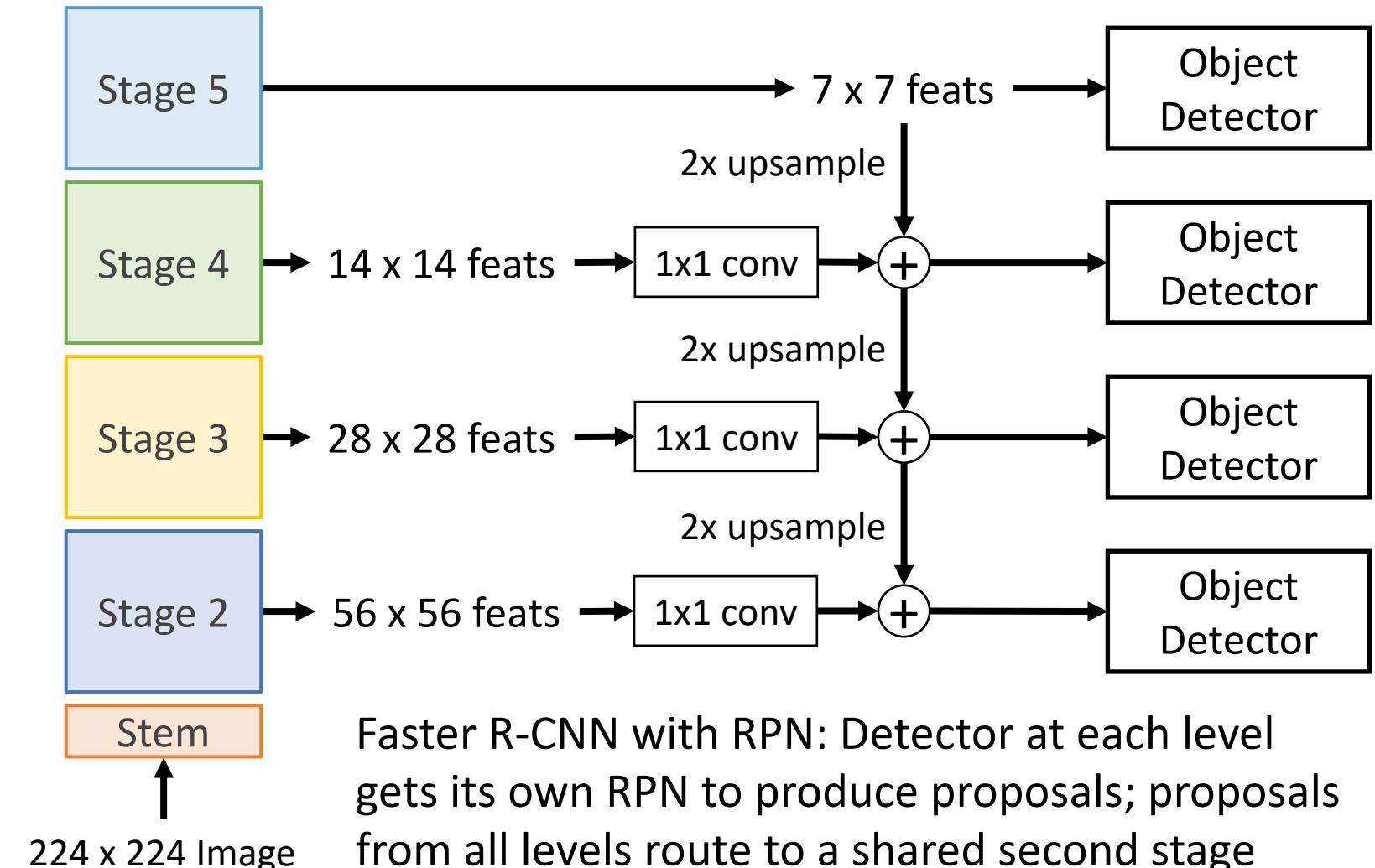


Lin et al, "Feature Pyramid Networks for Object Detection", ICCV 2017

# 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



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

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

# Faster R-CNN: Learnable Region Proposals

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

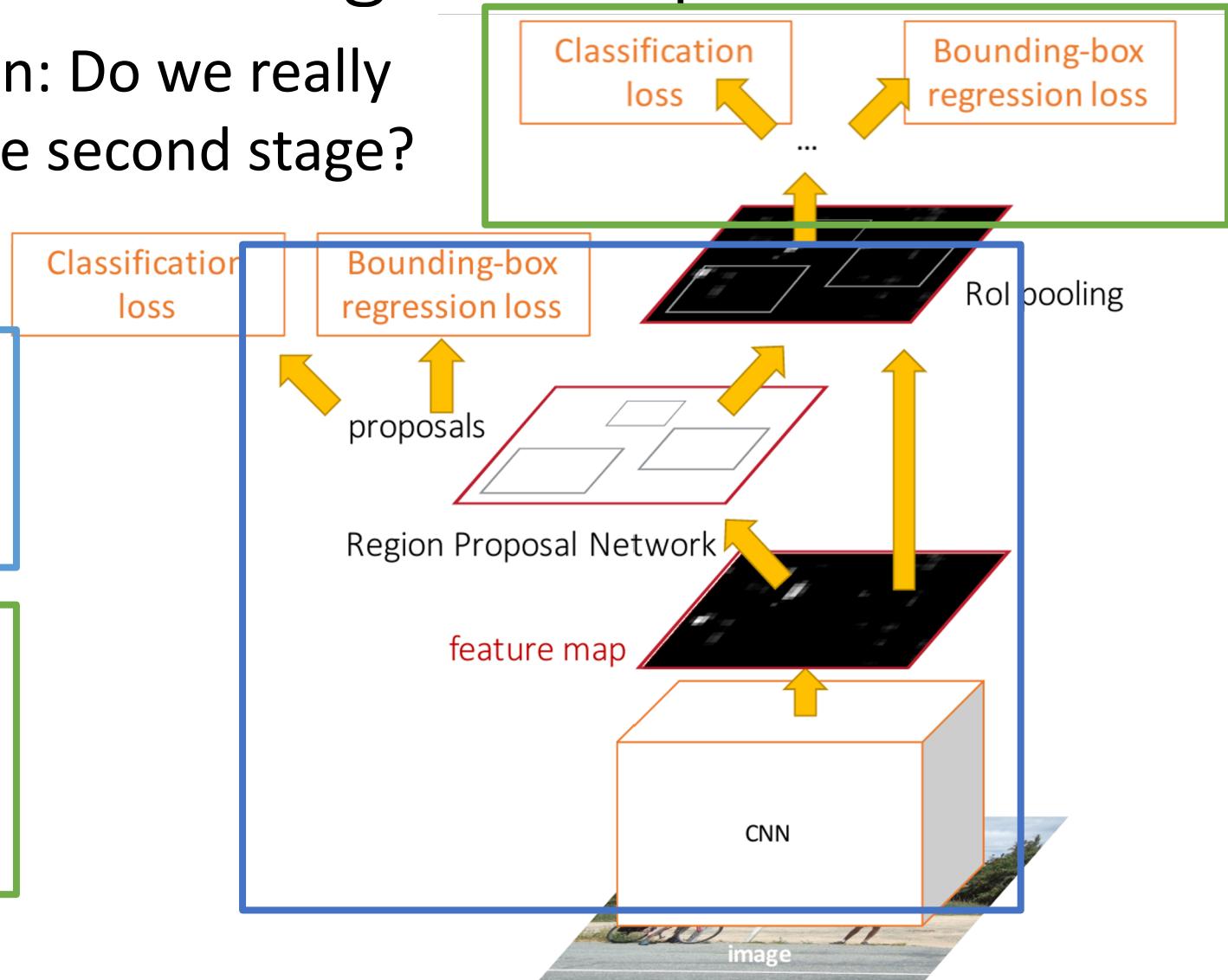
First stage: Run once per image

- Backbone network
- Region proposal network

Second stage: Run once per region

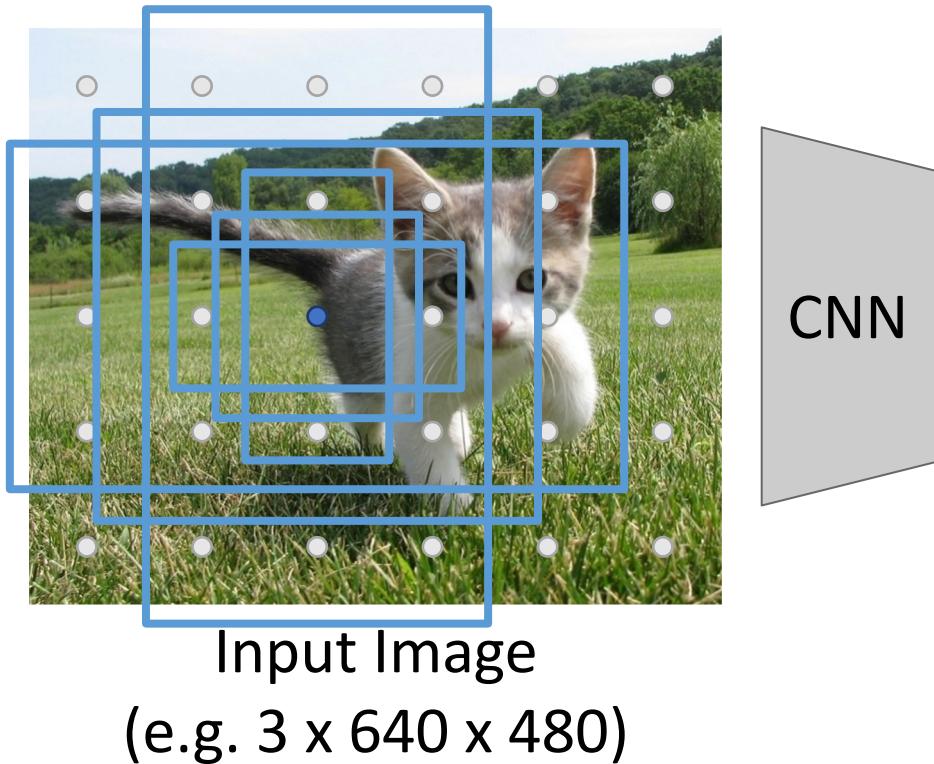
- Crop features: RoI pool / align
- Predict object class
- Prediction bbox offset

Question: Do we really  
need the second stage?

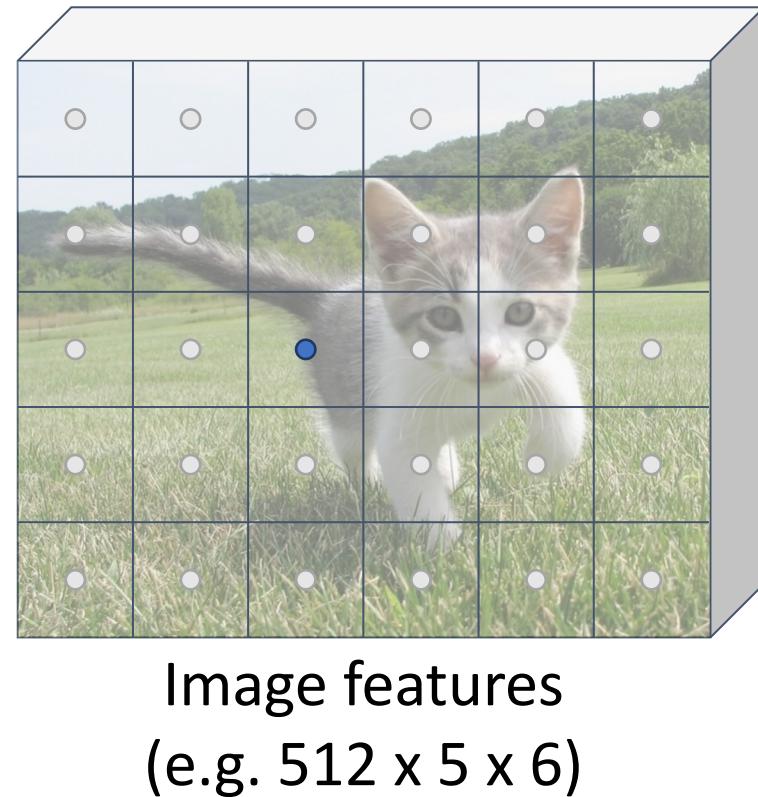


# Single-Stage Detectors: RetinaNet

Run backbone CNN to get features aligned to input image



Each feature corresponds to a point in the input



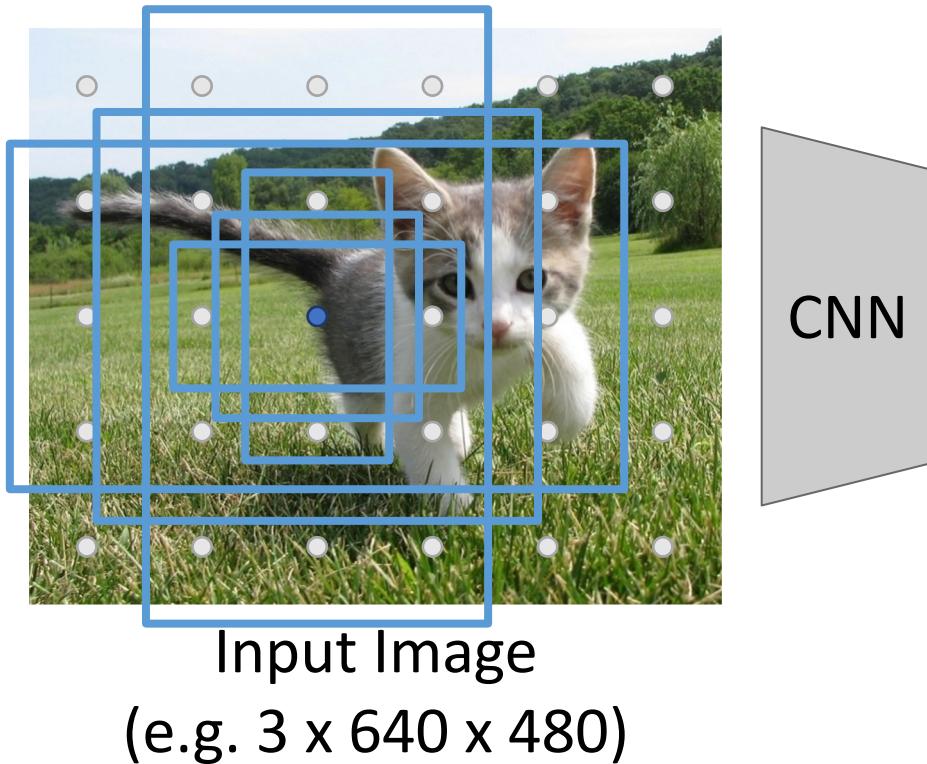
Similar to RPN – but rather than classify anchors as object/no object, directly predict object category (among  $C$  categories) or background

Anchor classification  
 $2K*(C+1) \times 5 \times 6$

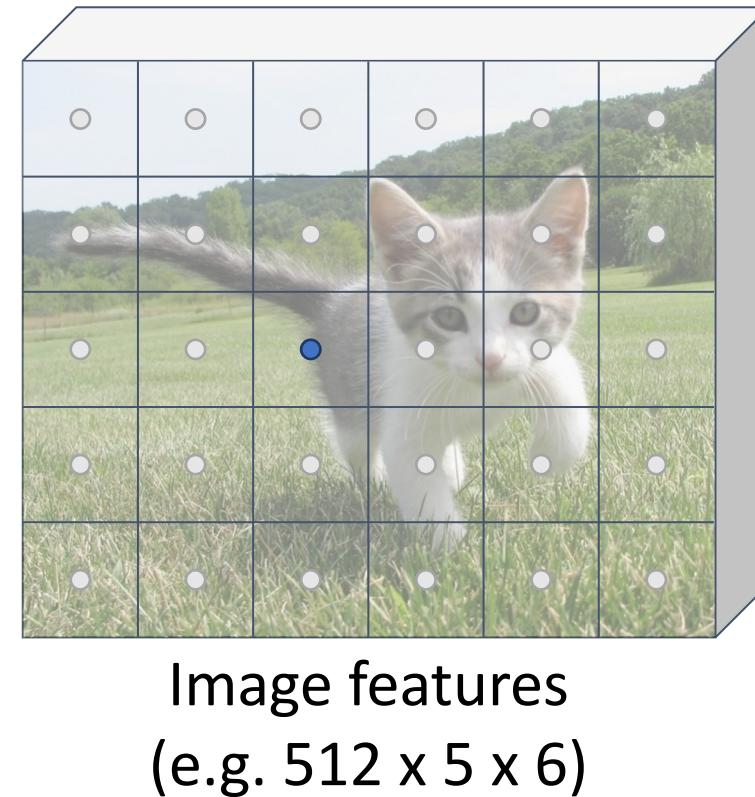
Anchor transforms  
 $4K \times 5 \times 6$

# Single-Stage Detectors: RetinaNet

Run backbone CNN to get features aligned to input image



Each feature corresponds to a point in the input



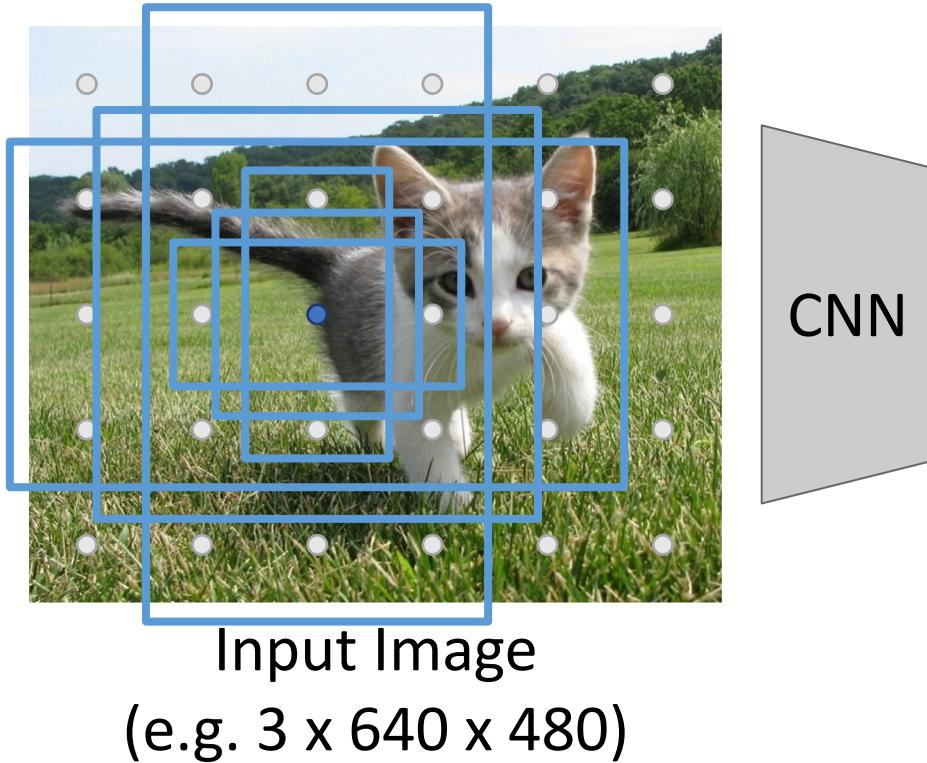
Problem: class imbalance – many more background anchors vs non-background

Anchor classification  
 $2K*(C+1) \times 5 \times 6$

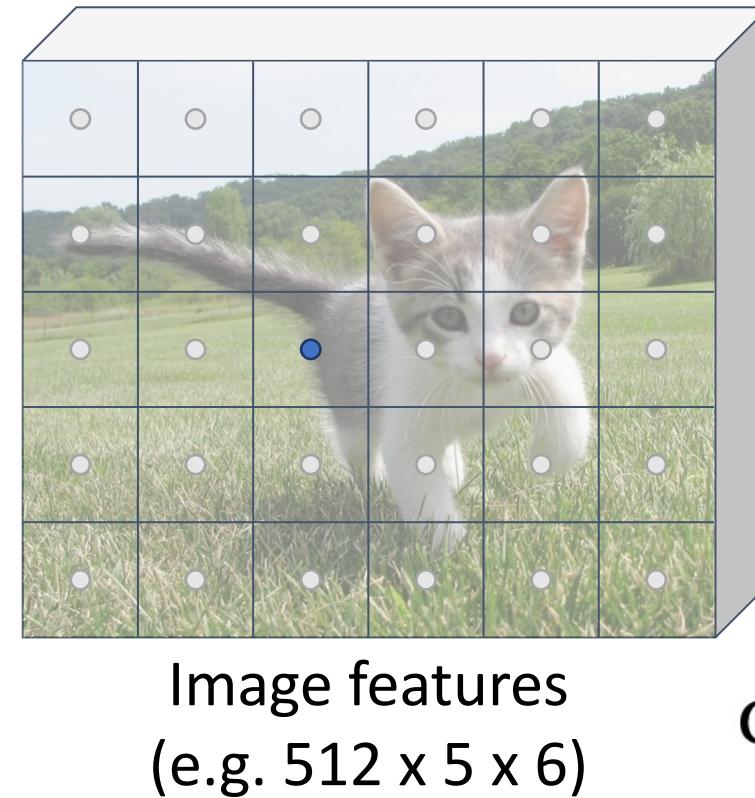
Anchor transforms  
 $4K \times 5 \times 6$

# Single-Stage Detectors: RetinaNet

Run backbone CNN to get features aligned to input image



Each feature corresponds to a point in the input



Problem: class imbalance – many more background anchors vs non-background

Solution: new loss function (Focal Loss); see paper

Anchor classification  
 $2K*(C+1) \times 5 \times 6$

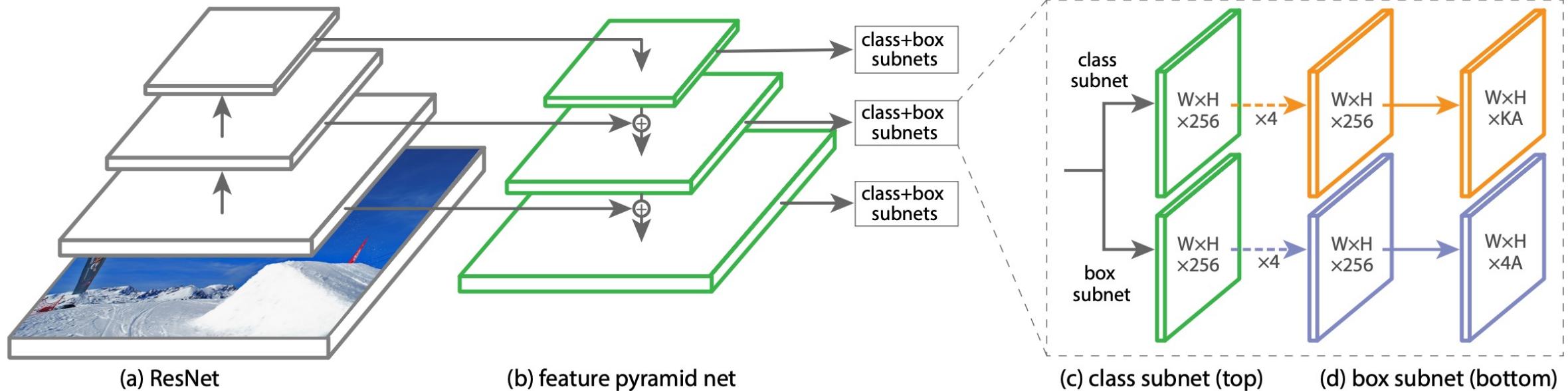
Anchor transforms  
 $4K \times 5 \times 6$

$$\text{CE}(p_t) = -\log(p_t)$$

$$\text{FL}(p_t) = -(1 - p_t)^\gamma \log(p_t)$$

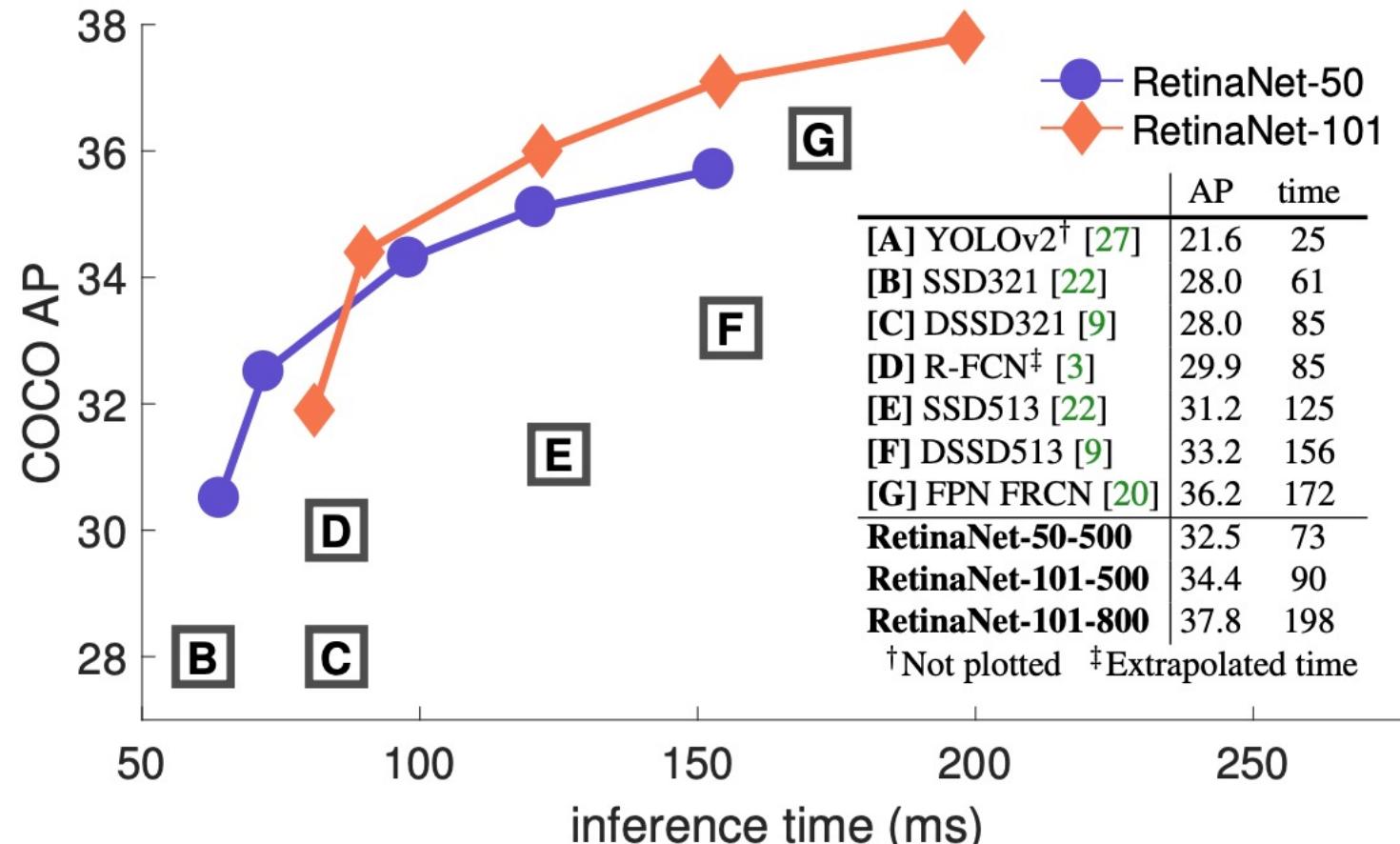
# Single-Stage Detectors: RetinaNet

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



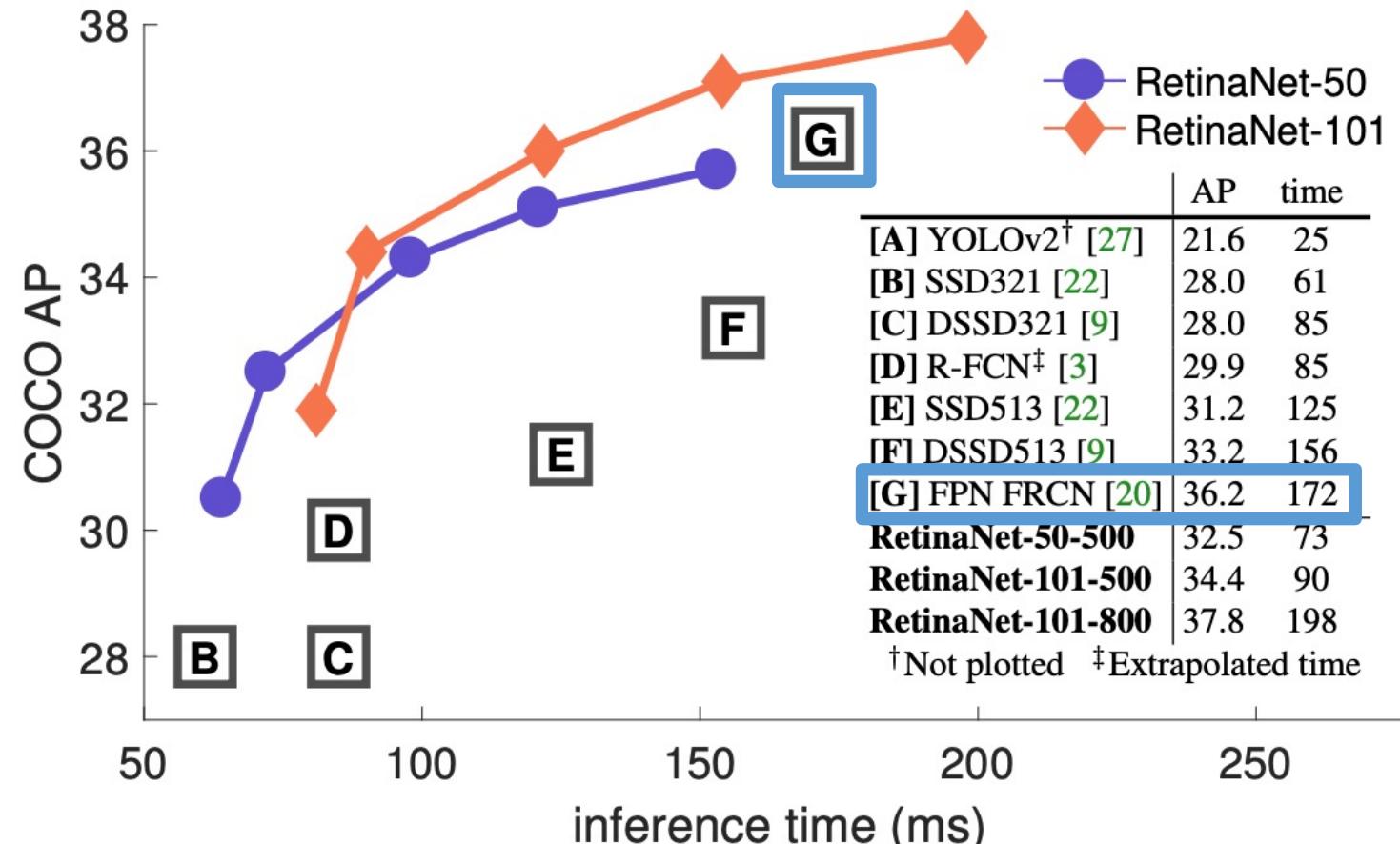
# Single-Stage Detectors: RetinaNet

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



# Single-Stage Detectors: RetinaNet

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

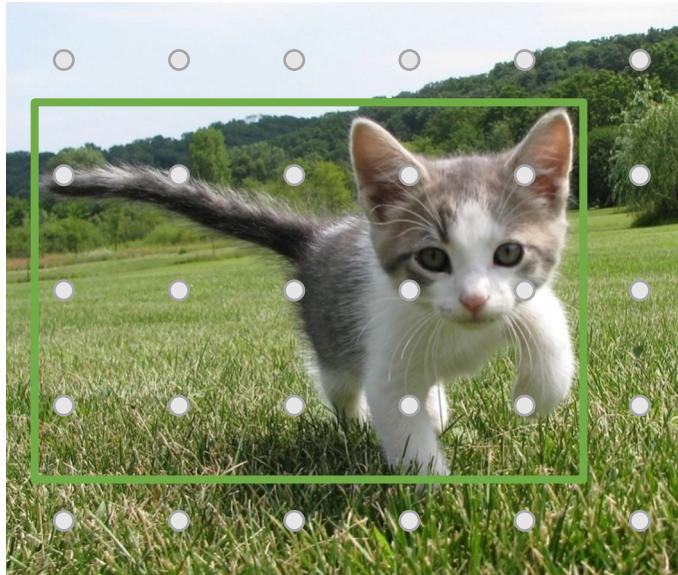


Faster R-CNN  
with Feature  
Pyramid  
Network

# Single-Stage Detectors: FCOS

“Anchor-free” detector

Run backbone CNN to get features aligned to input image



Input Image  
(e.g.  $3 \times 640 \times 480$ )

Each feature corresponds to a point in the input

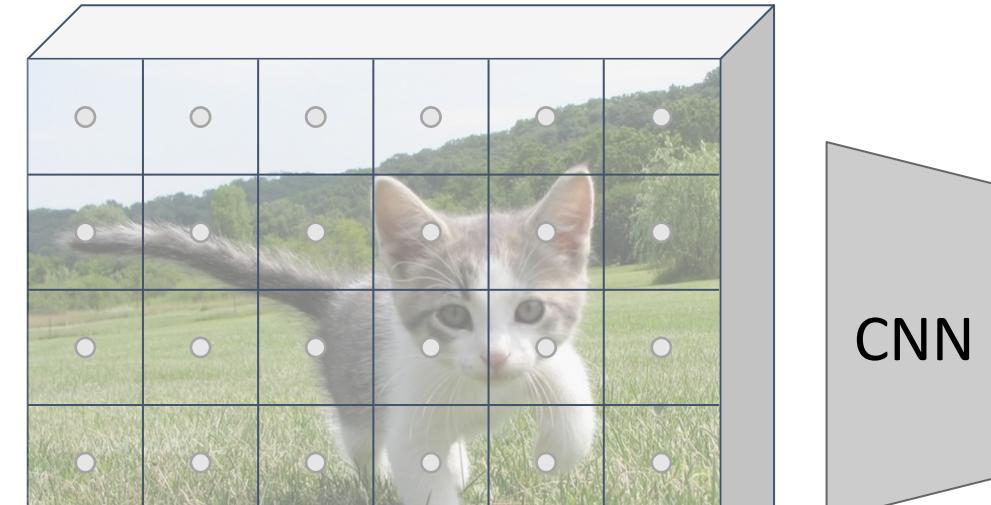
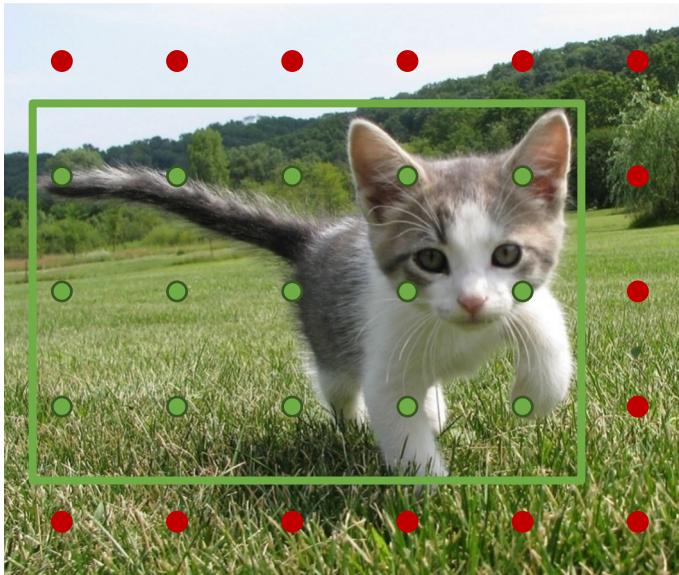


Image features  
(e.g.  $512 \times 5 \times 6$ )

# Single-Stage Detectors: FCOS

Run backbone CNN to get features aligned to input image



Input Image  
(e.g.  $3 \times 640 \times 480$ )



Each feature corresponds to a point in the input

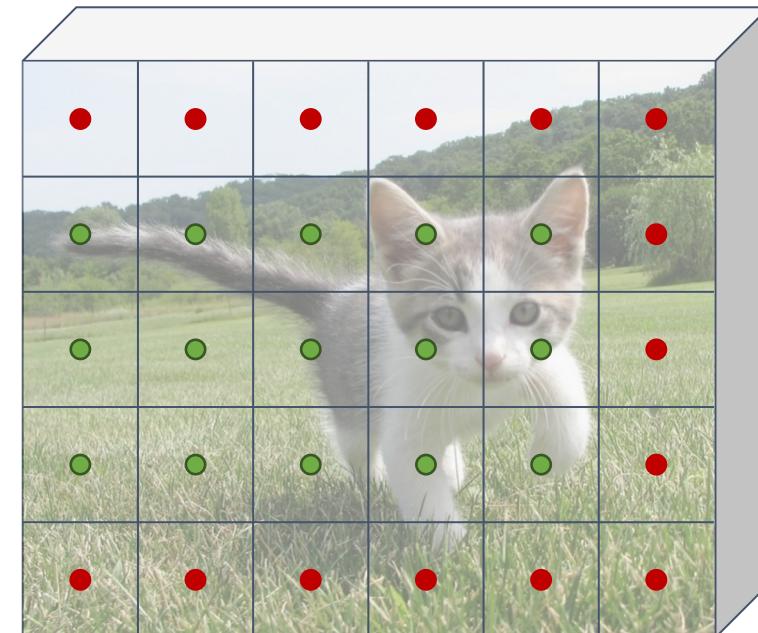


Image features  
(e.g.  $512 \times 5 \times 6$ )

“Anchor-free” detector

Classify points as positive if they fall into a GT box, or negative if they don't

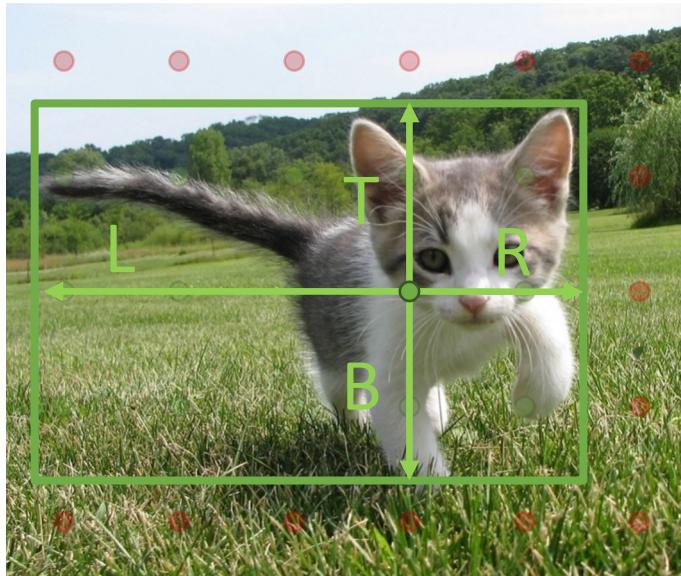
Train independent per-category logistic regressors

→ Class scores  
 $C \times 5 \times 6$



# Single-Stage Detectors: FCOS

Run backbone CNN to get features aligned to input image



Input Image  
(e.g.  $3 \times 640 \times 480$ )

Each feature corresponds to a point in the input

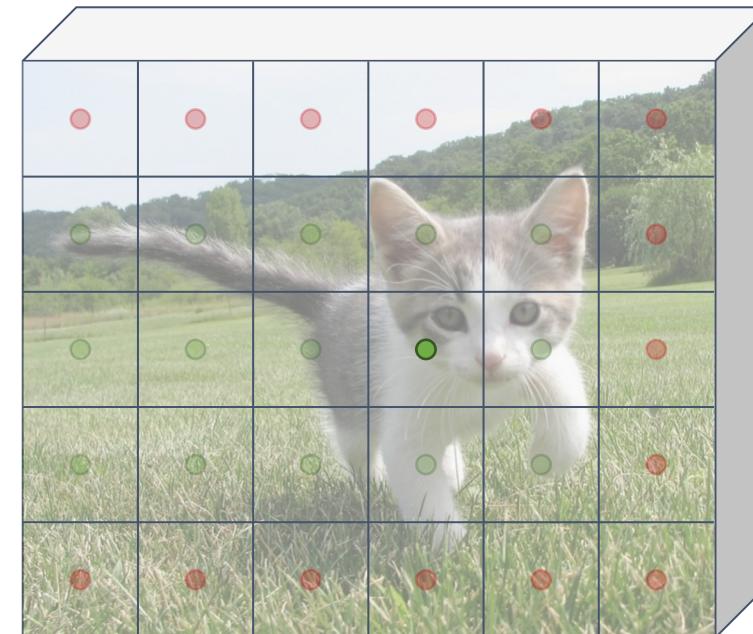
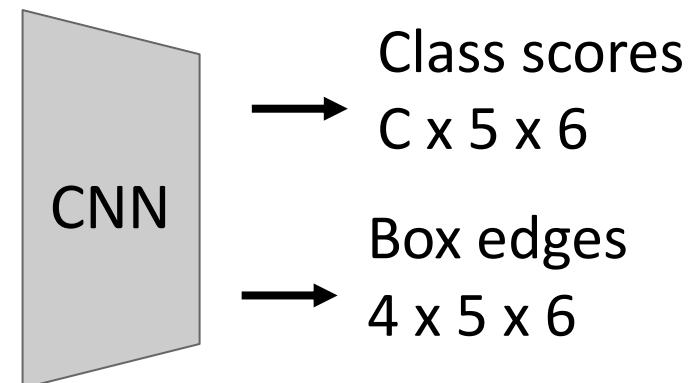


Image features  
(e.g.  $512 \times 5 \times 6$ )

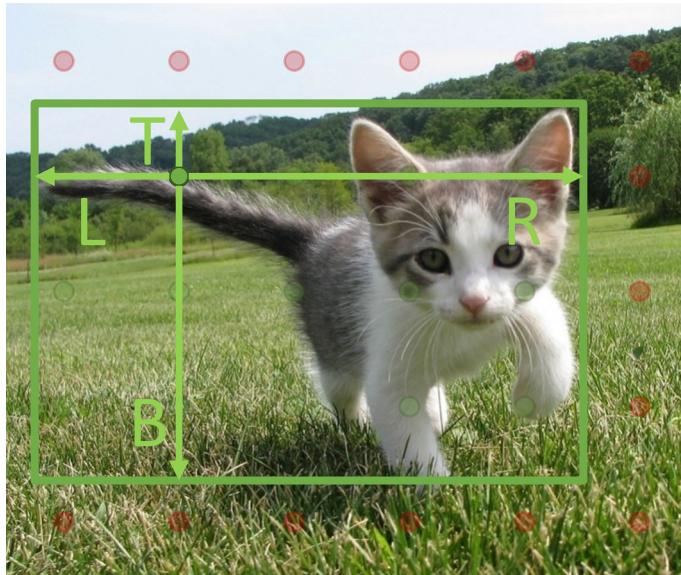
“Anchor-free” detector

For positive points, also regress distance to left, right, top, and bottom of ground-truth box (with L2 loss)



# Single-Stage Detectors: FCOS

Run backbone CNN to get features aligned to input image



Input Image  
(e.g.  $3 \times 640 \times 480$ )

Each feature corresponds to a point in the input

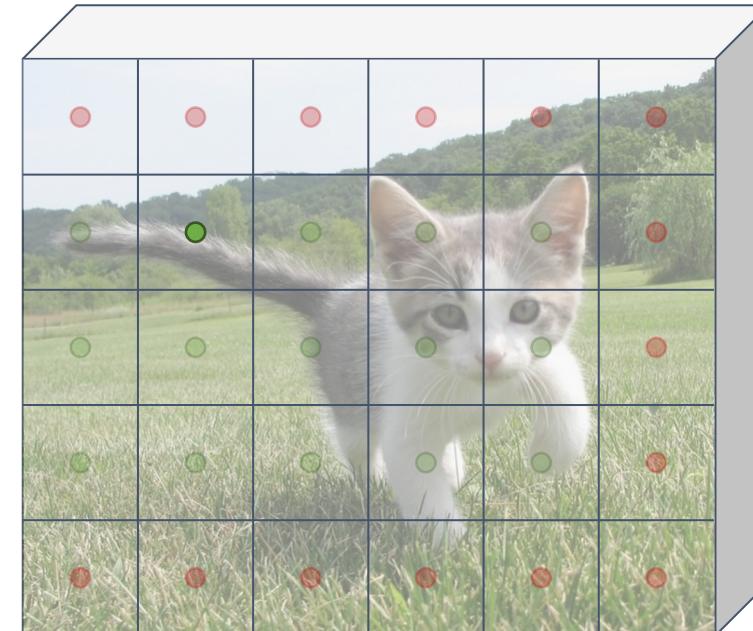
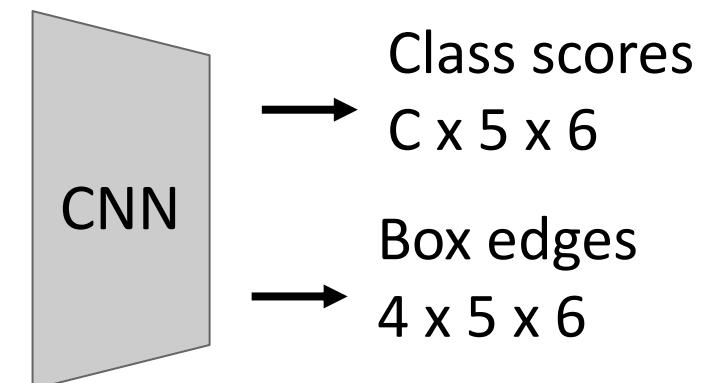


Image features  
(e.g.  $512 \times 5 \times 6$ )

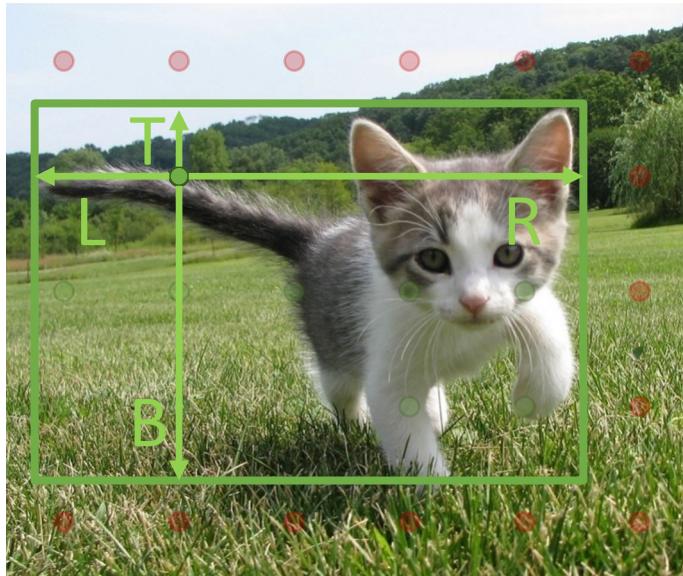
“Anchor-free” detector

For positive points, also regress distance to left, right, top, and bottom of ground-truth box (with L2 loss)



# Single-Stage Detectors: FCOS

Run backbone CNN to get features aligned to input image



Input Image  
(e.g.  $3 \times 640 \times 480$ )



Each feature corresponds to a point in the input

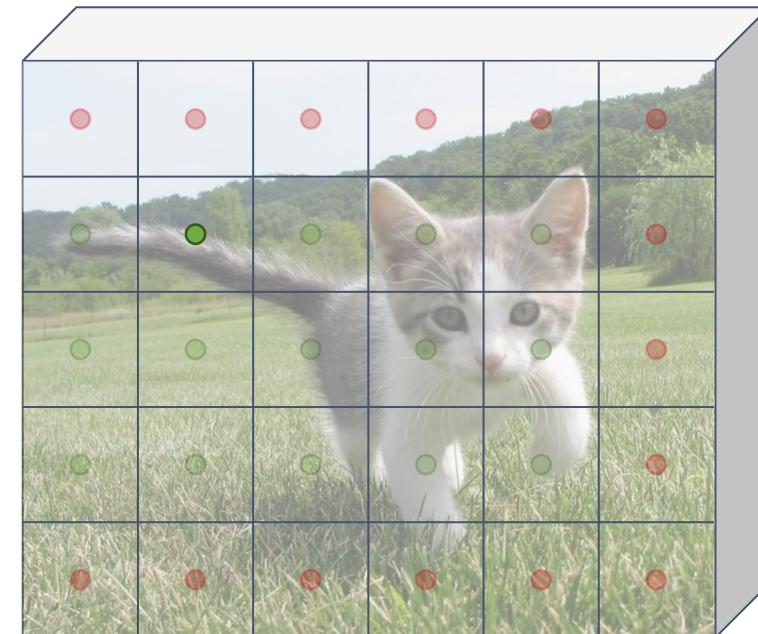
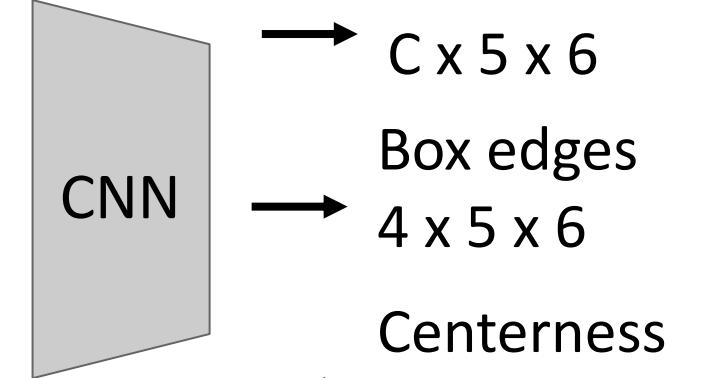


Image features  
(e.g.  $512 \times 5 \times 6$ )

“Anchor-free” detector

Finally, predict “centerness” for all positive points (using logistic regression loss)

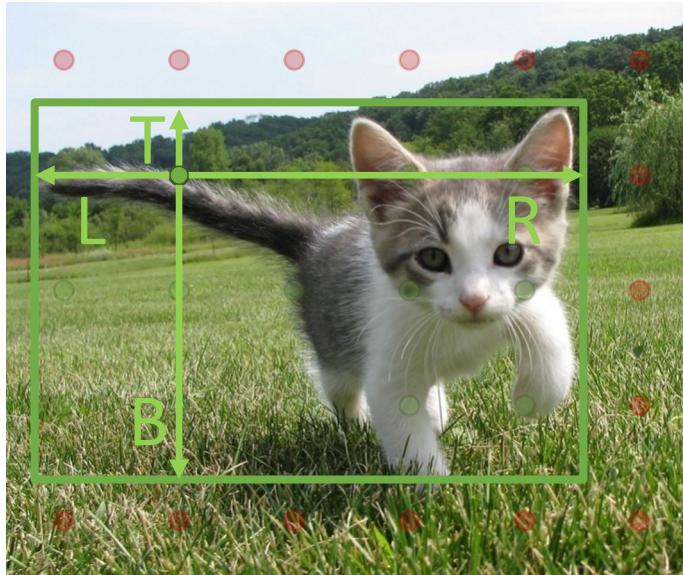


$$\text{centerness} = \sqrt{\frac{\min(L, R)}{\max(L, R)} \cdot \frac{\min(T, B)}{\max(T, B)}}$$

Ranges from 1 at box center to 0 at box edge

# Single-Stage Detectors: FCOS

Run backbone CNN to get features aligned to input image



Input Image  
(e.g.  $3 \times 640 \times 480$ )



Each feature corresponds to a point in the input

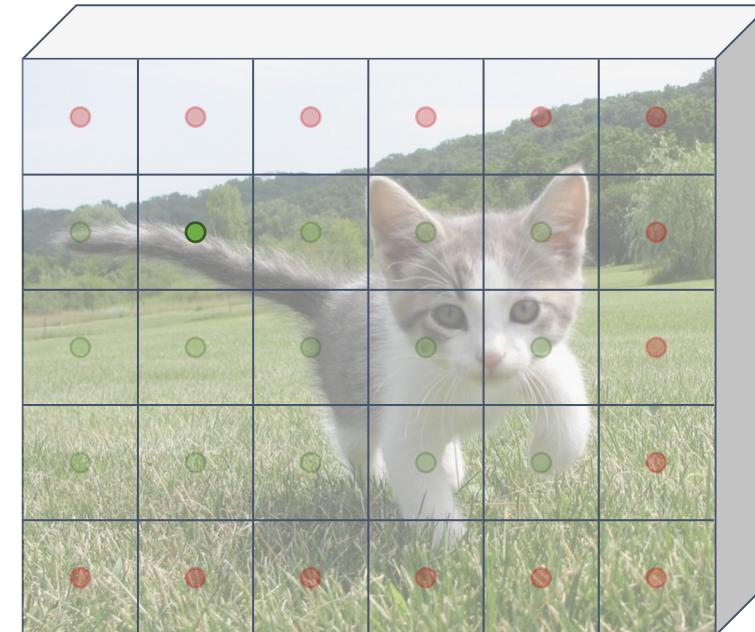
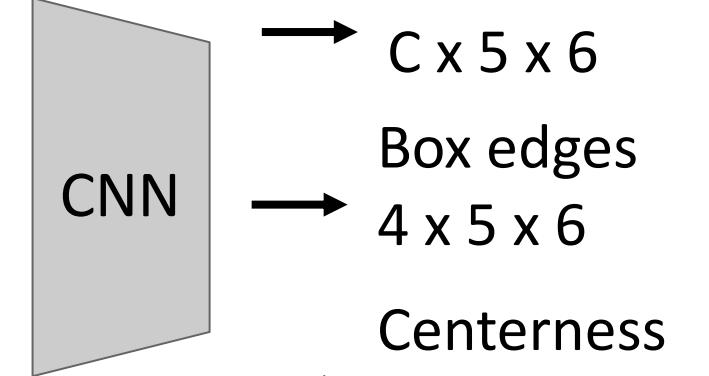


Image features  
(e.g.  $512 \times 5 \times 6$ )

“Anchor-free” detector

Test-time: predicted “confidence” for the box from each point is product of its class score and centerness



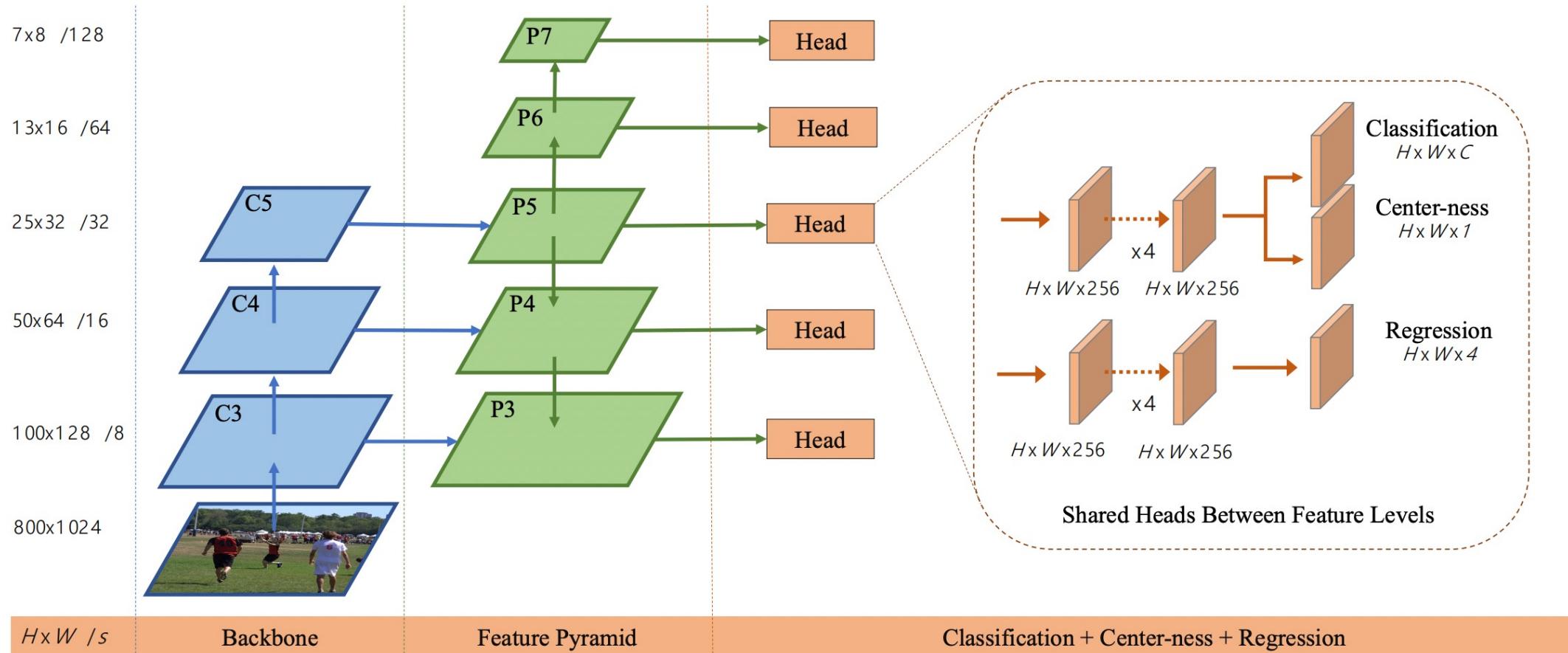
$$\text{centerness} = \sqrt{\frac{\min(L, R)}{\max(L, R)} \cdot \frac{\min(T, B)}{\max(T, B)}}$$

Ranges from 1 at box center to 0 at box edge

“Anchor-free” detector

# Single-Stage Detectors: FCOS

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



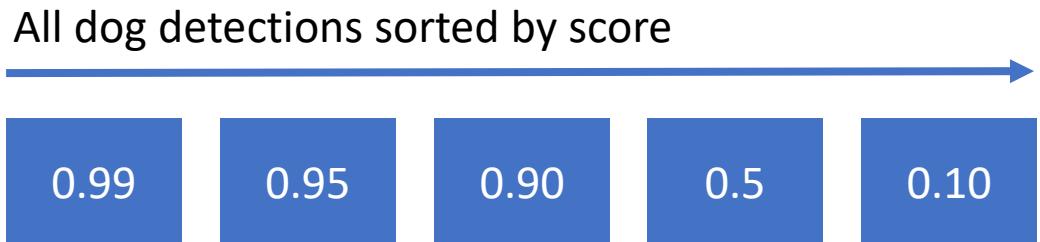
Tian et al, “FCOS: Fully Convolutional One-Stage Object Detection”, ICCV 2019

# 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)

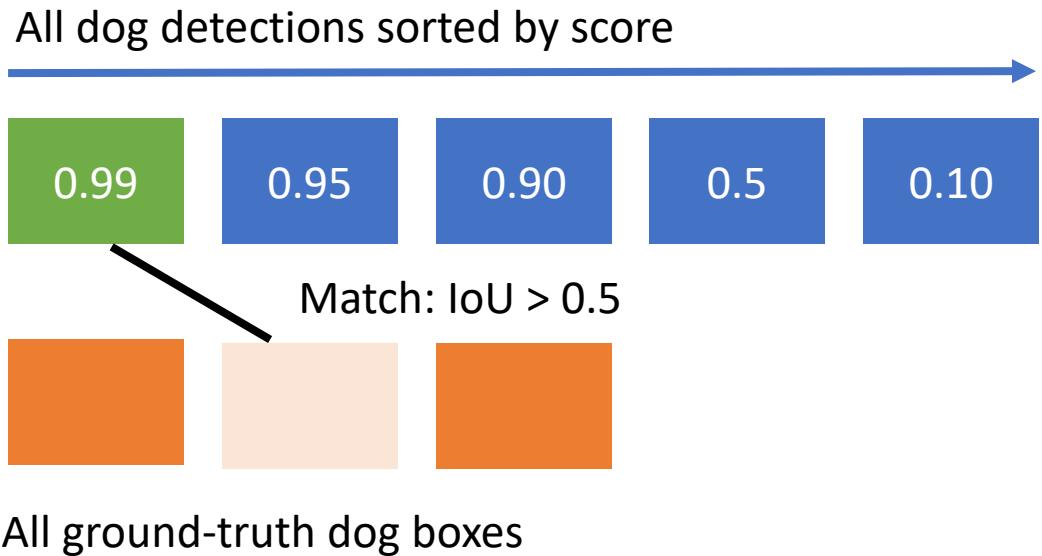
1. Run object detector on all test images (with NMS)
2. For each category, compute Average Precision (AP) = area under Precision vs Recall Curve
  1. For each detection (highest score to lowest score)



All ground-truth dog boxes

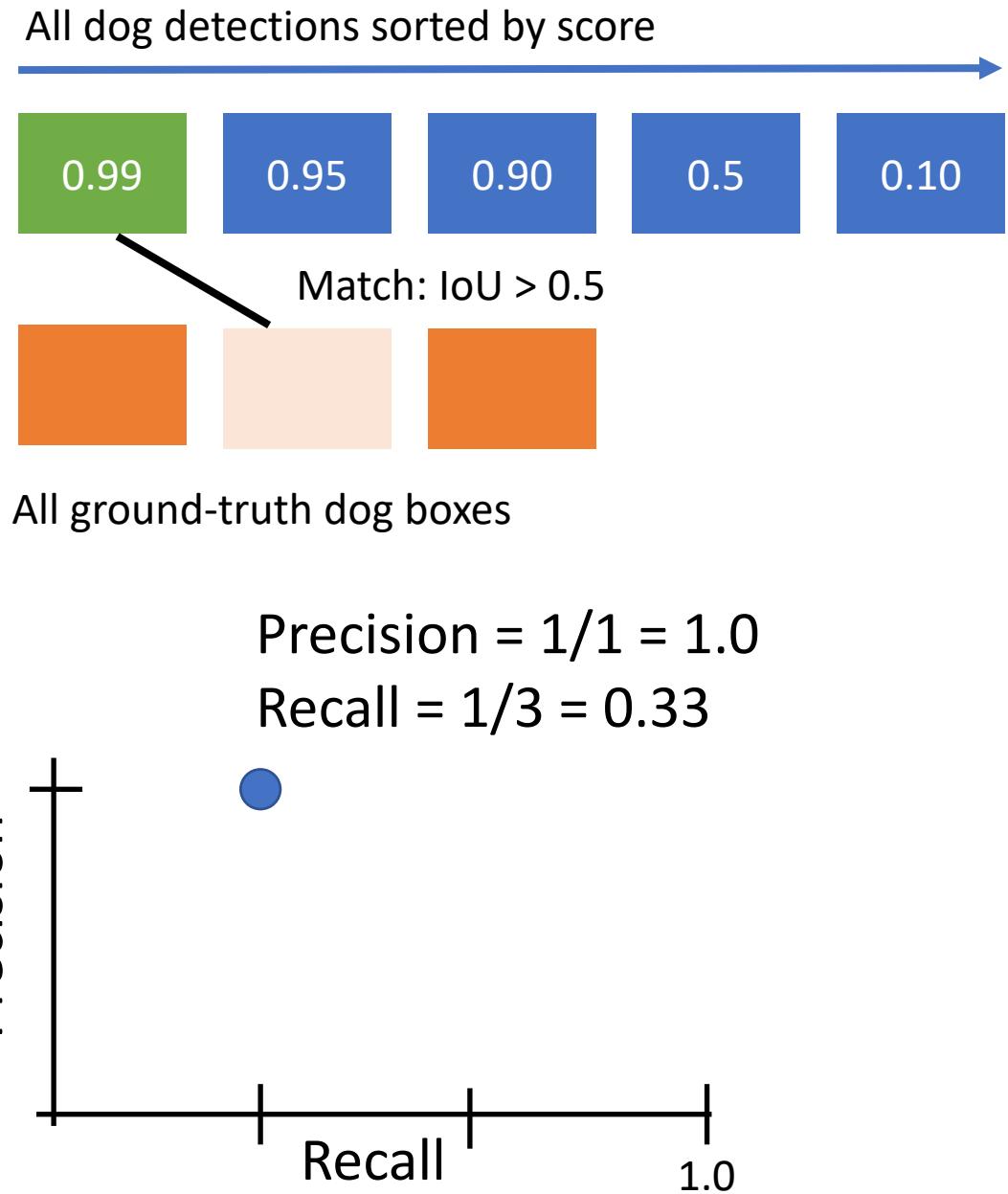
# Evaluating Object Detectors: Mean Average Precision (mAP)

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  $\text{IoU} > 0.5$ , mark it as positive and eliminate the GT
    2. Otherwise mark it as negative



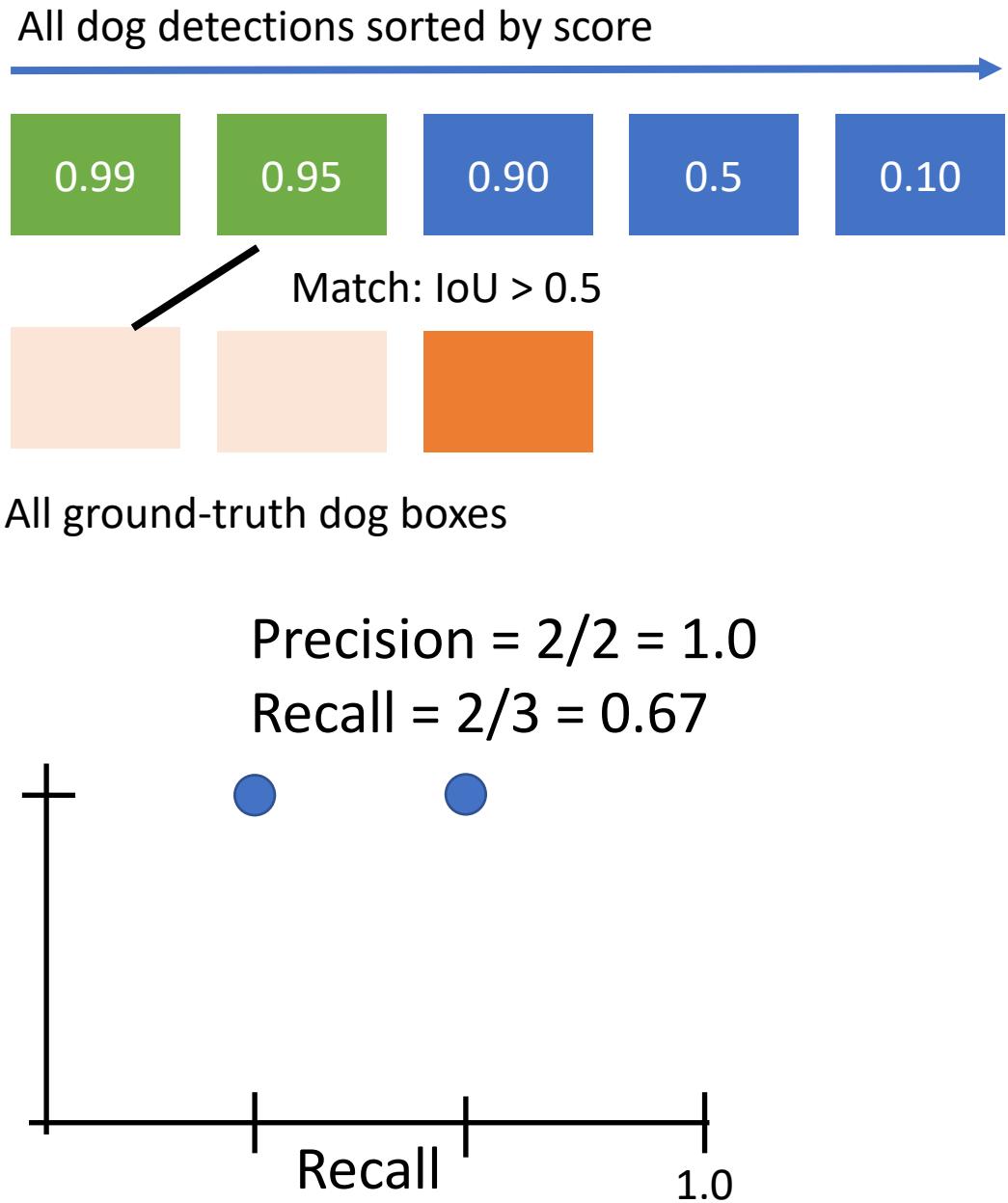
# Evaluating Object Detectors: Mean Average Precision (mAP)

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  $\text{IoU} > 0.5$ , mark it as positive and eliminate the GT
    2. Otherwise mark it as negative
    3. Plot a point on PR Curve



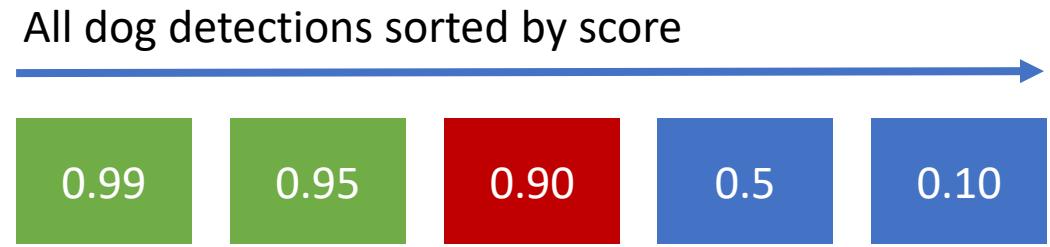
# Evaluating Object Detectors: Mean Average Precision (mAP)

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)
    1. If it matches some GT box with  $\text{IoU} > 0.5$ , mark it as positive and eliminate the GT
    2. Otherwise mark it as negative
    3. Plot a point on PR 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)
    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



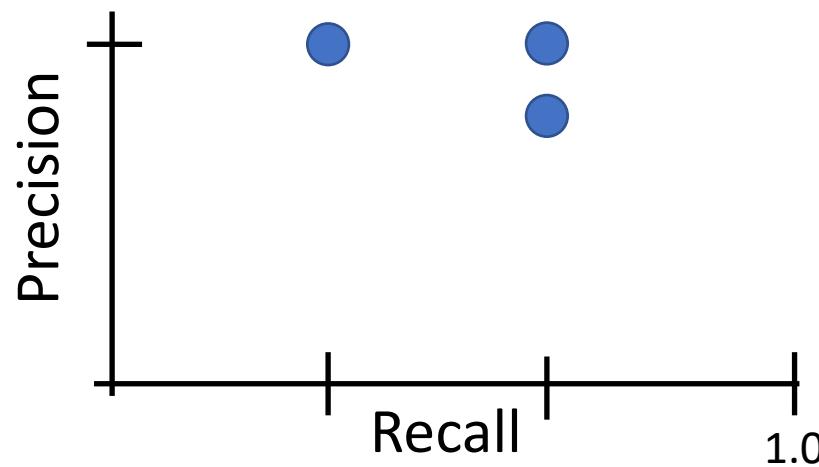
No match > 0.5 IoU with GT



All ground-truth dog boxes

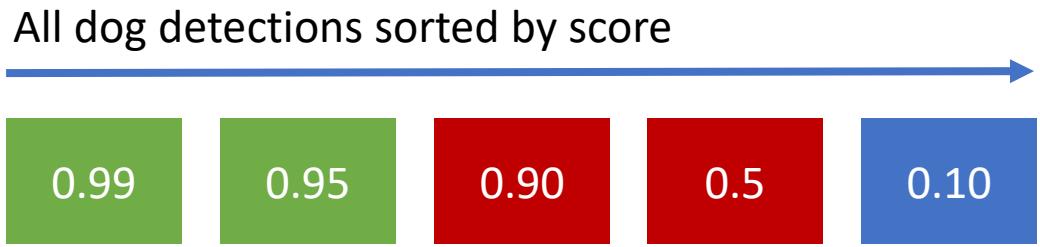
$$\text{Precision} = 2/3 = 0.67$$

$$\text{Recall} = 2/3 = 0.67$$

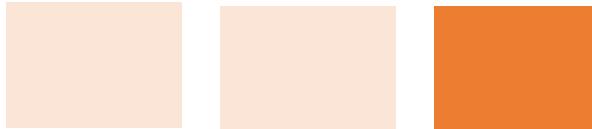


# Evaluating Object Detectors: Mean Average Precision (mAP)

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

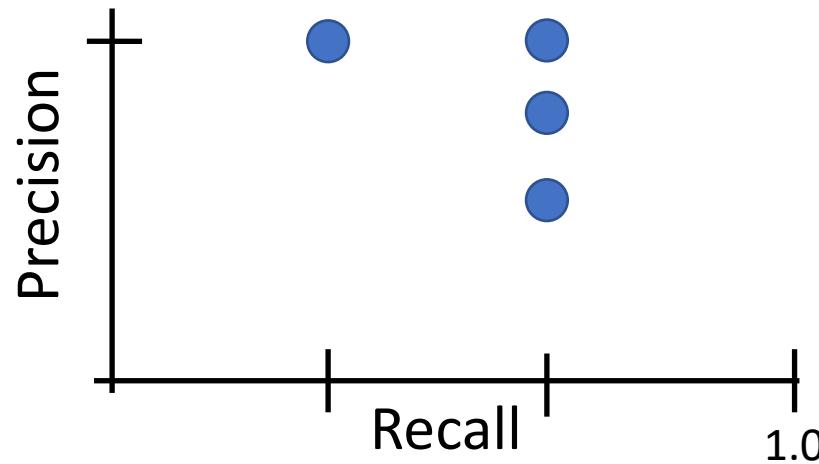


No match > 0.5 IoU with GT



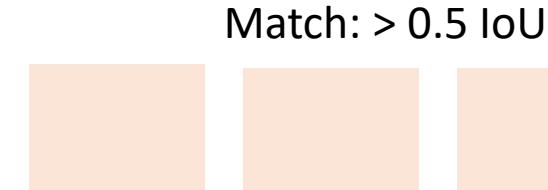
All ground-truth dog boxes

$$\text{Precision} = 2/4 = 0.5$$
$$\text{Recall} = 2/3 = 0.67$$



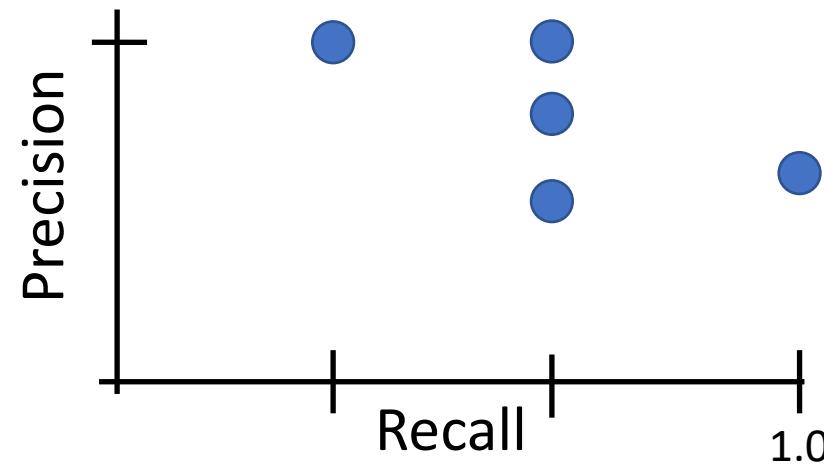
# Evaluating Object Detectors: Mean Average Precision (mAP)

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



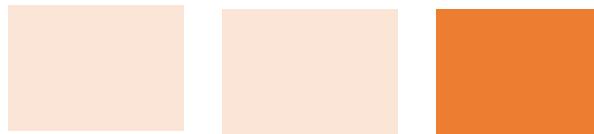
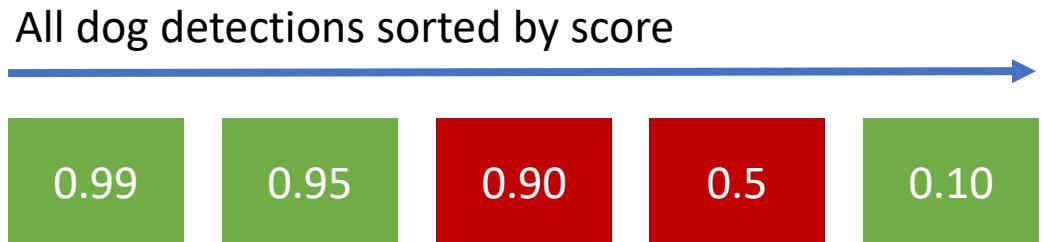
All ground-truth dog boxes

$$\text{Precision} = 3/5 = 0.6$$
$$\text{Recall} = 3/3 = 1.0$$

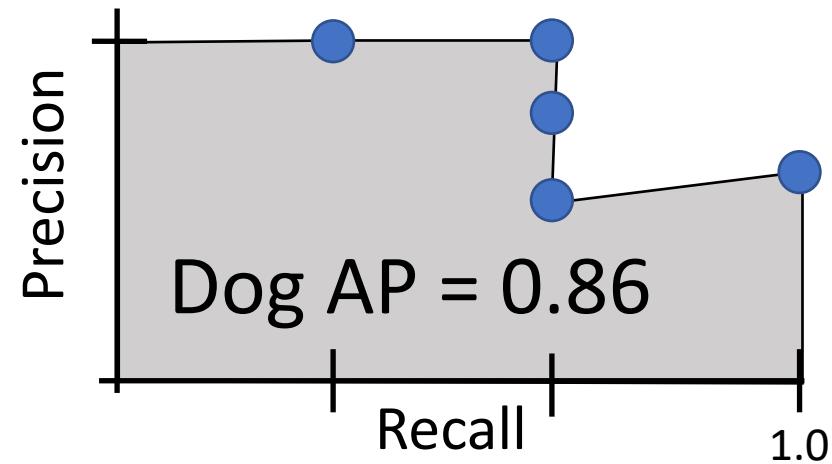


# Evaluating Object Detectors: Mean Average Precision (mAP)

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



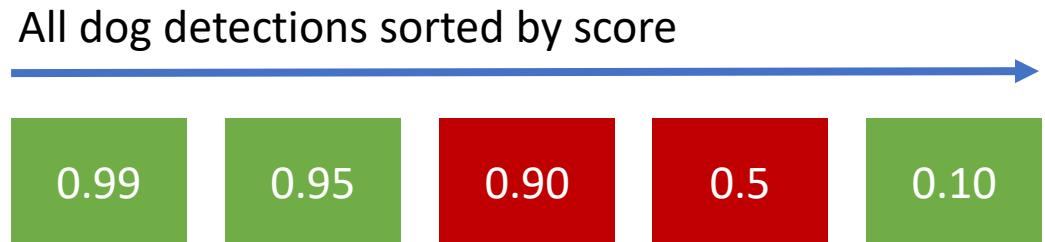
All ground-truth dog boxes



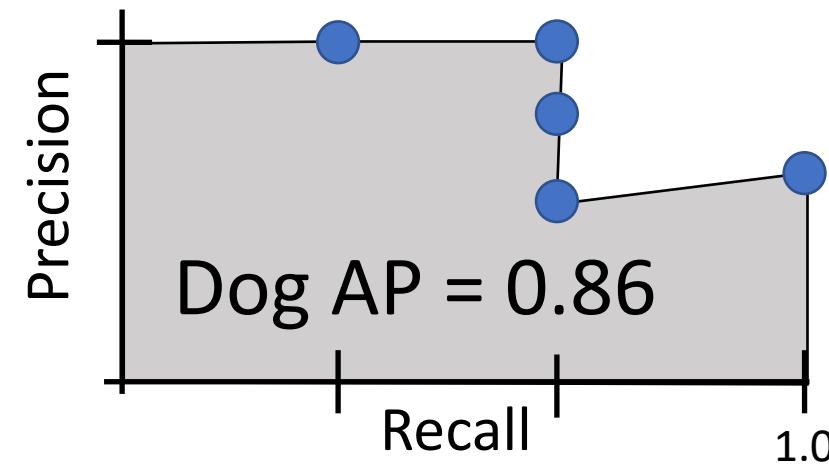
# Evaluating Object Detectors: Mean Average Precision (mAP)

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  $\text{IoU} > 0.5$ , 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

**How to get AP = 1.0: Hit all GT boxes with  $\text{IoU} > 0.5$ , and have no “false positive” detections ranked above any “true positives”**



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  $\text{IoU} > 0.5$ , 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

# 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  $\text{IoU} > 0.5$ , 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
4. For “COCO mAP”: Compute mAP@thresh for each IoU threshold (0.5, 0.55, 0.6, ..., 0.95) and take average

$\text{mAP}@0.5 = 0.77$

$\text{mAP}@0.55 = 0.71$

$\text{mAP}@0.60 = 0.65$

...

$\text{mAP}@0.95 = 0.2$

COCO mAP = 0.4

# Summary: Beyond Image Classification

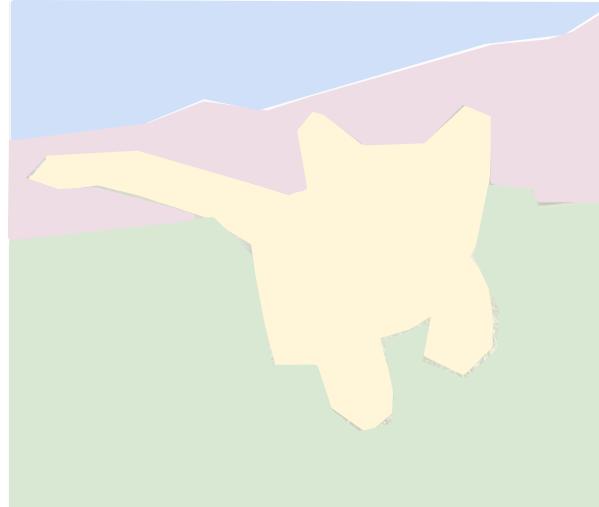
## Classification



CAT

No spatial extent

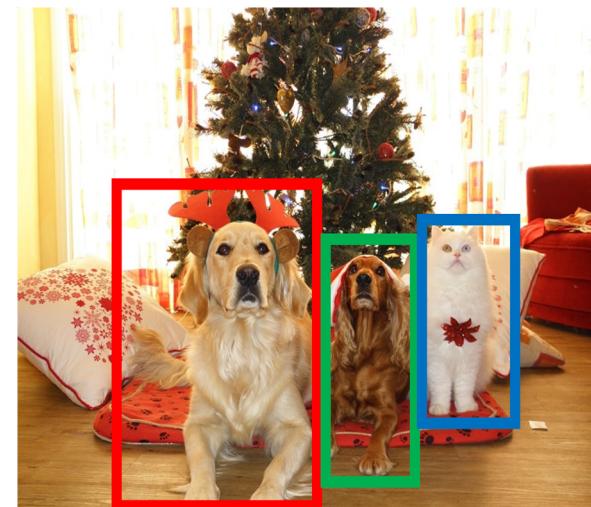
## Semantic Segmentation



GRASS, CAT, TREE,  
SKY

No objects, just pixels

## Object Detection



DOG, DOG, CAT

Multiple Objects

## Instance Segmentation

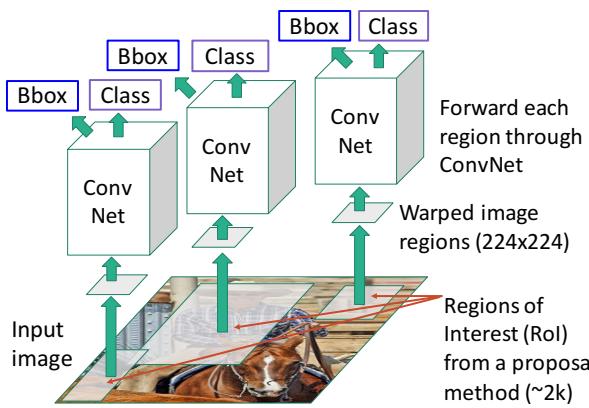


DOG, DOG, CAT

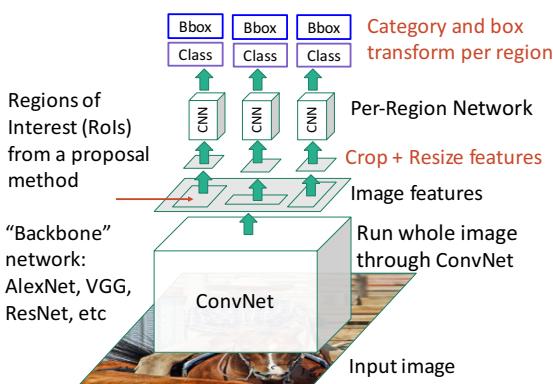
[This image](#) is CCO public domain

# Summary

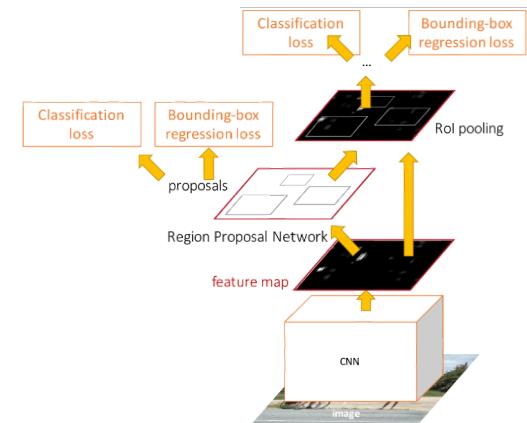
**“Slow” R-CNN:** Run CNN independently for each region



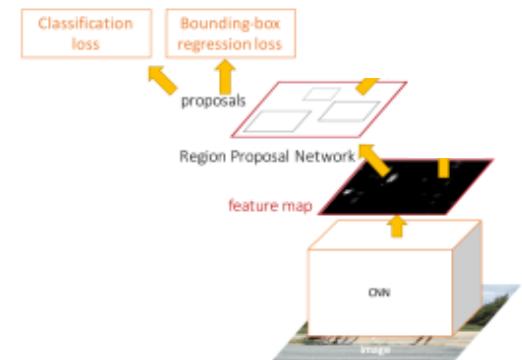
**Fast R-CNN:** Apply differentiable cropping to shared image features



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



**Single-Stage:** Fully convolutional detector



With anchors: RetinaNet  
Anchor-Free: FCOS

Next time:  
Image and Instance Segmentation