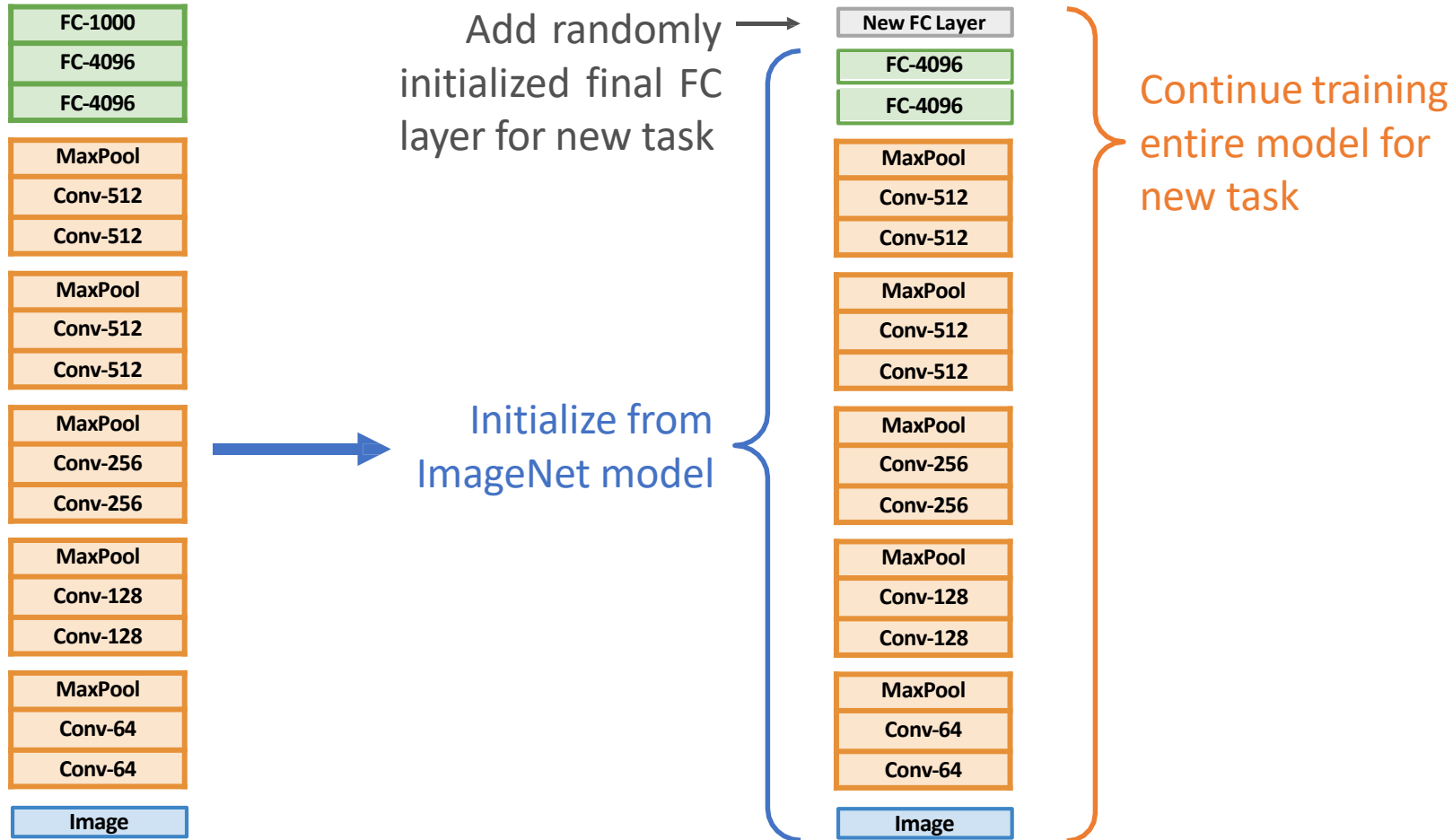


# Lecture 15:

## Object Detectors (II)

# Last Time: Transfer Learning

## 1. Train on ImageNet



# Last Time: Localization Tasks

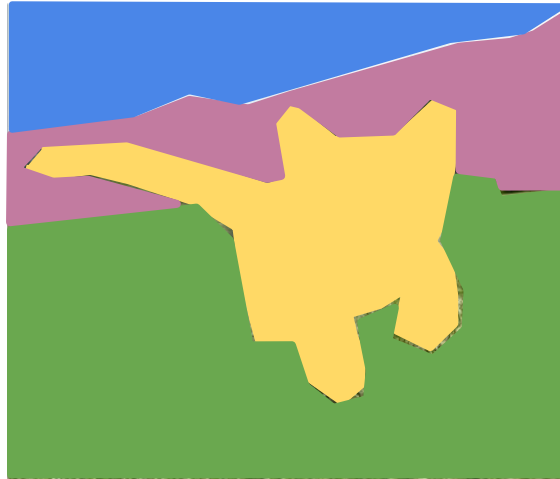
## Classification



**CAT**

No spatial extent

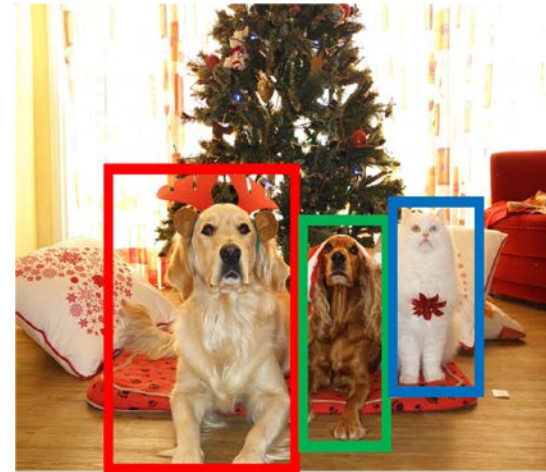
## Semantic Segmentation



**GRASS, CAT, TREE, SKY**

No objects, just pixels

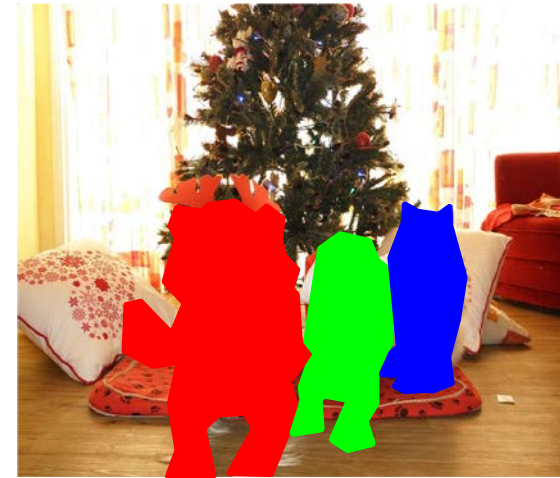
## Object Detection



**DOG, DOG, CAT**

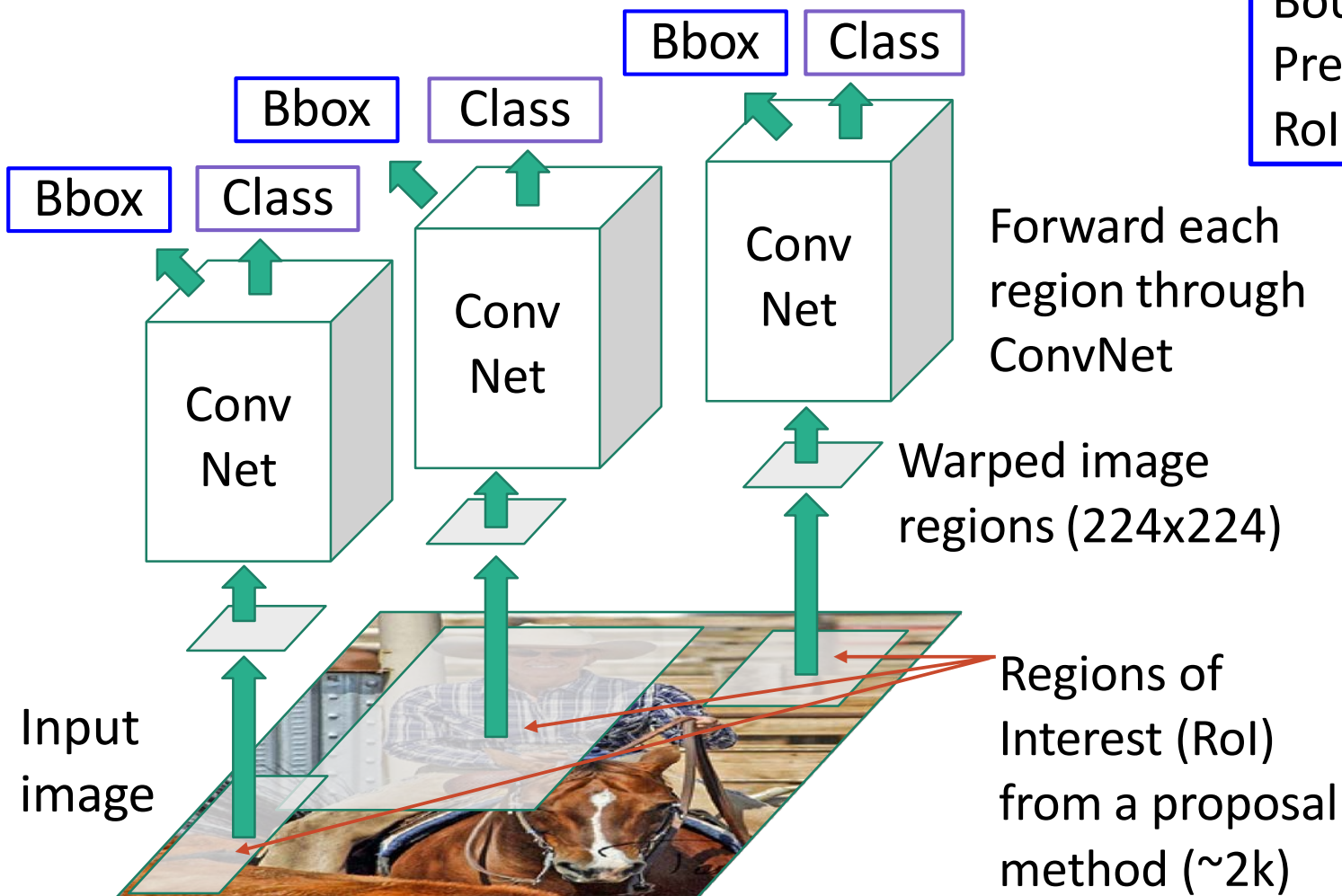
Multiple Objects

## Instance Segmentation



**DOG, DOG, CAT**

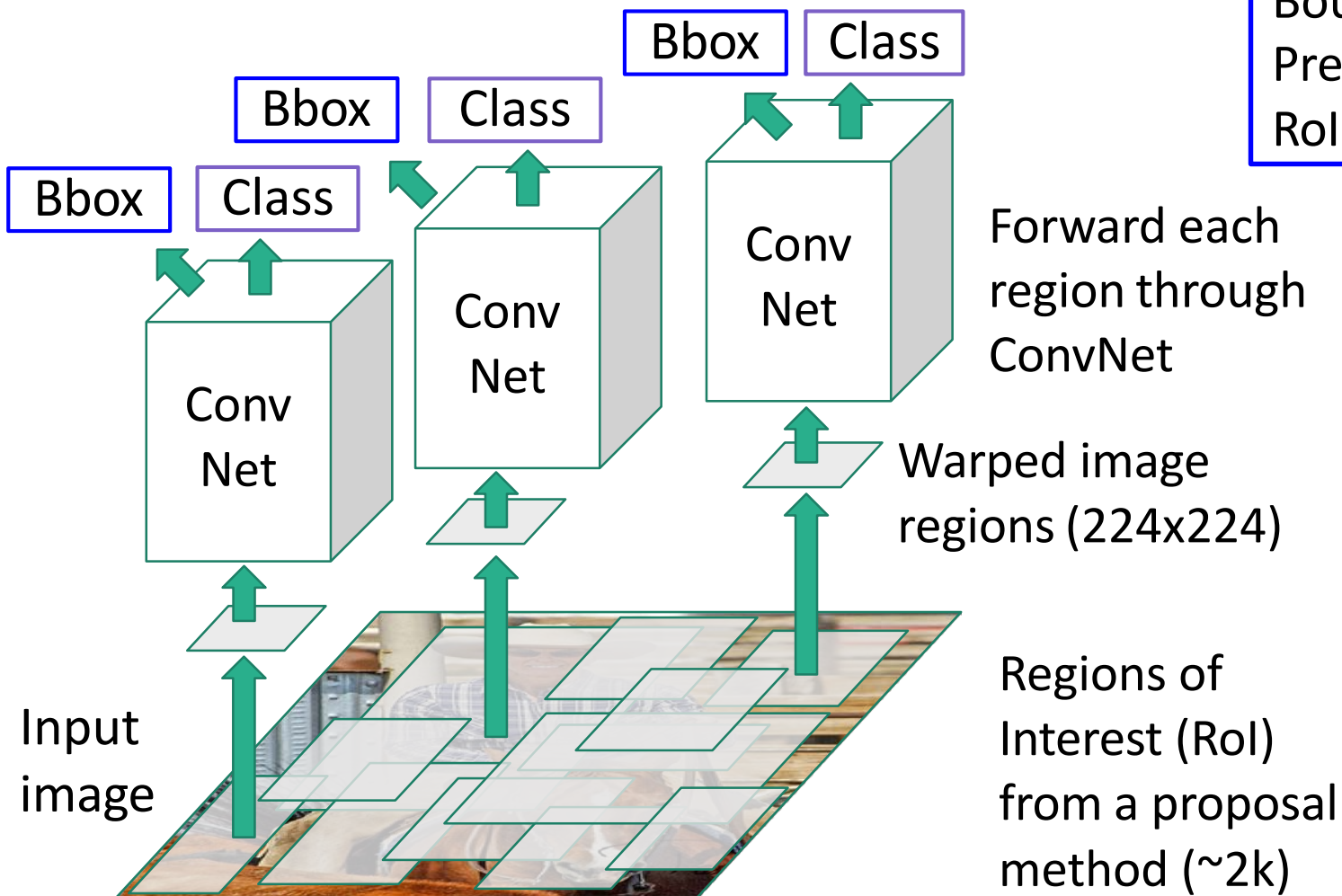
# Last Time: R-CNN



Classify each region

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

# Last Time: R-CNN



Classify each region

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

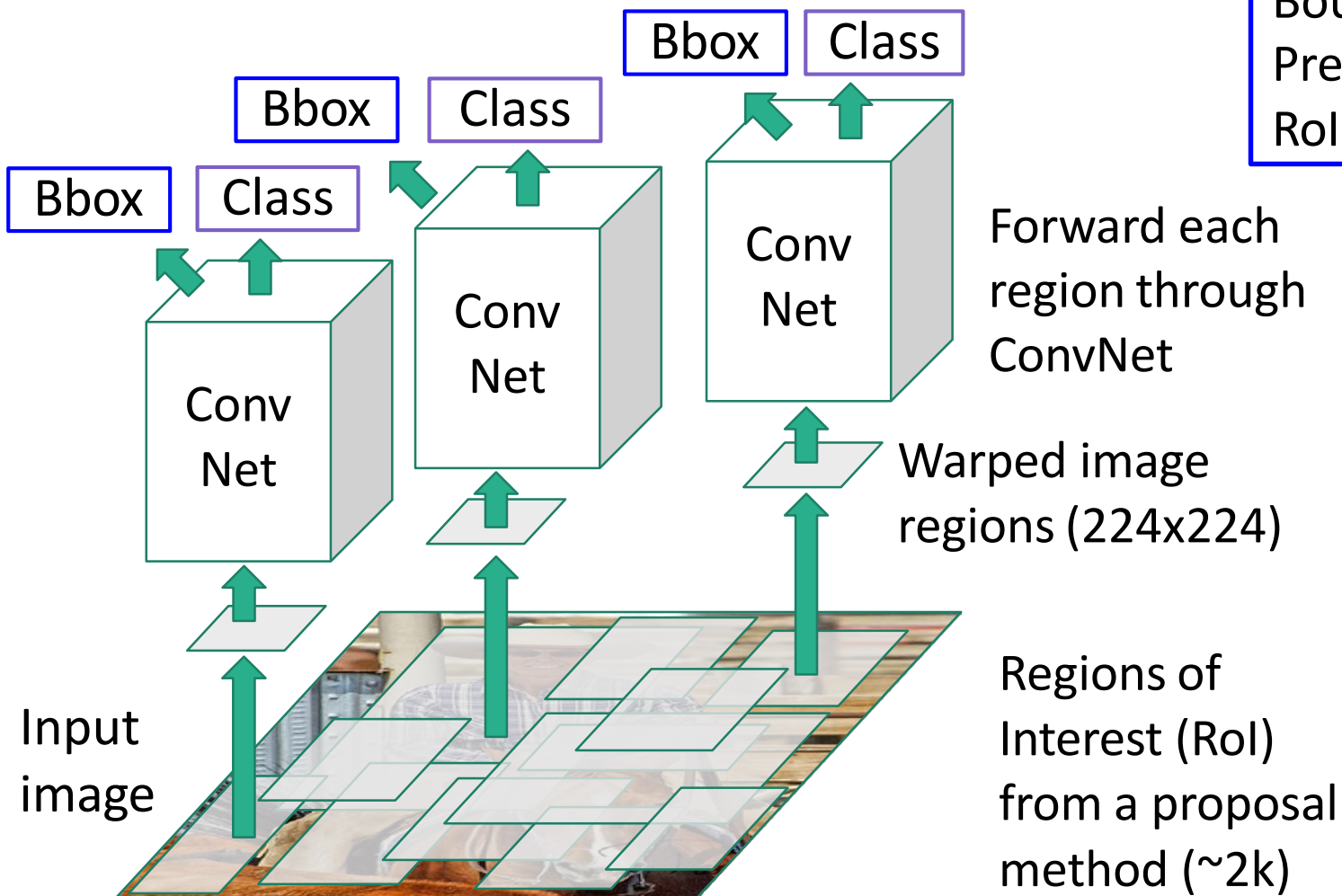
Forward each  
region through  
ConvNet

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

Warped image  
regions (224x224)

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

# Last Time: R-CNN



Classify each region

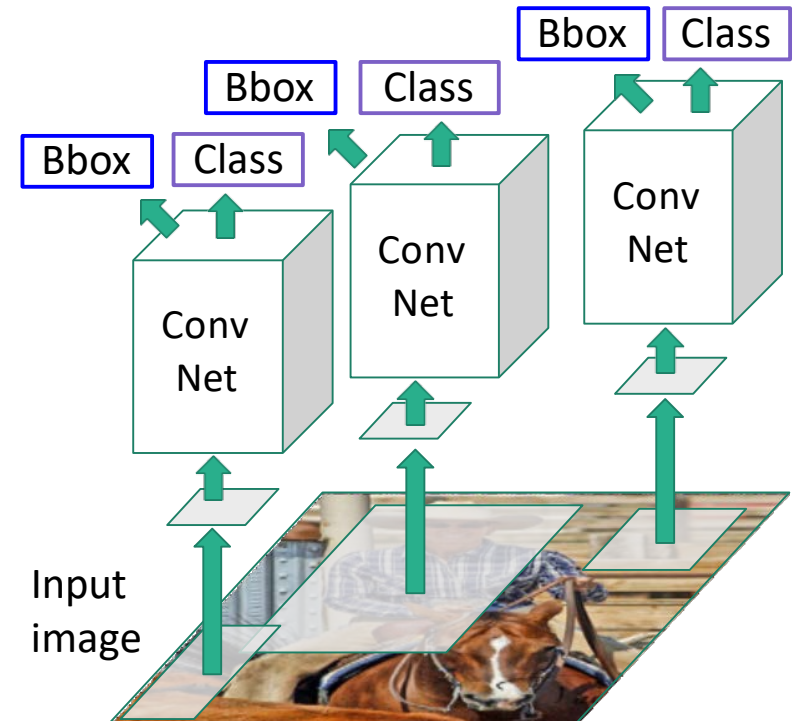
Bounding box regression:  
Predict "transform" to correct the  
RoI: 4 numbers ( $t_x, t_y, t_h, t_w$ )

**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?**

## “Slow” R-CNN

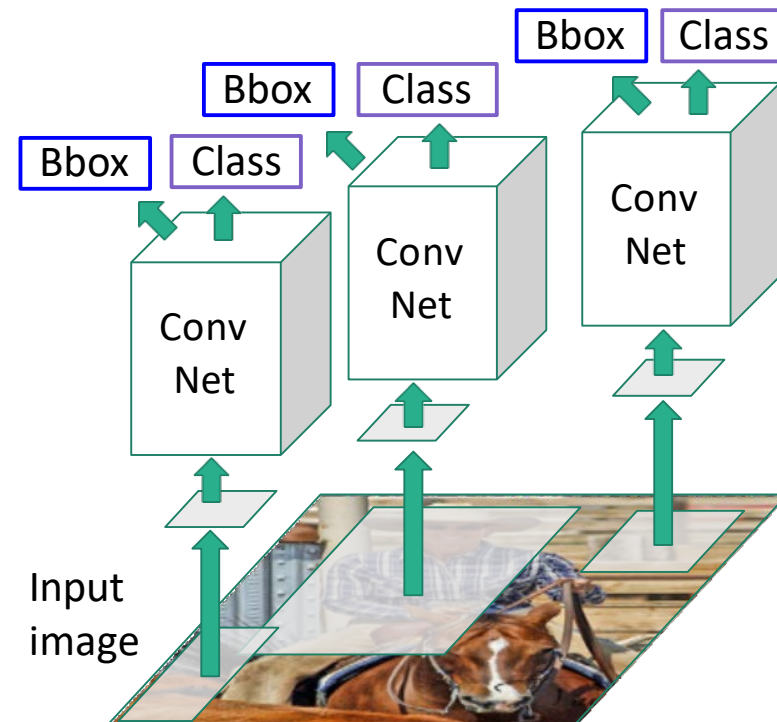
Process each region  
independently



# Fast R-CNN



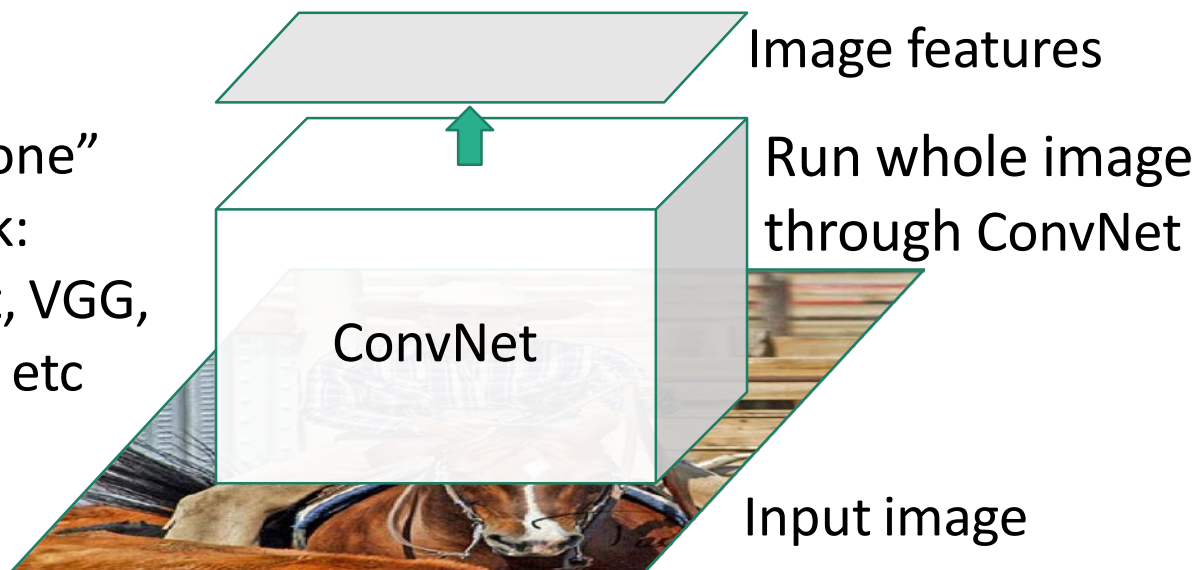
“Slow” R-CNN  
Process each region  
independently



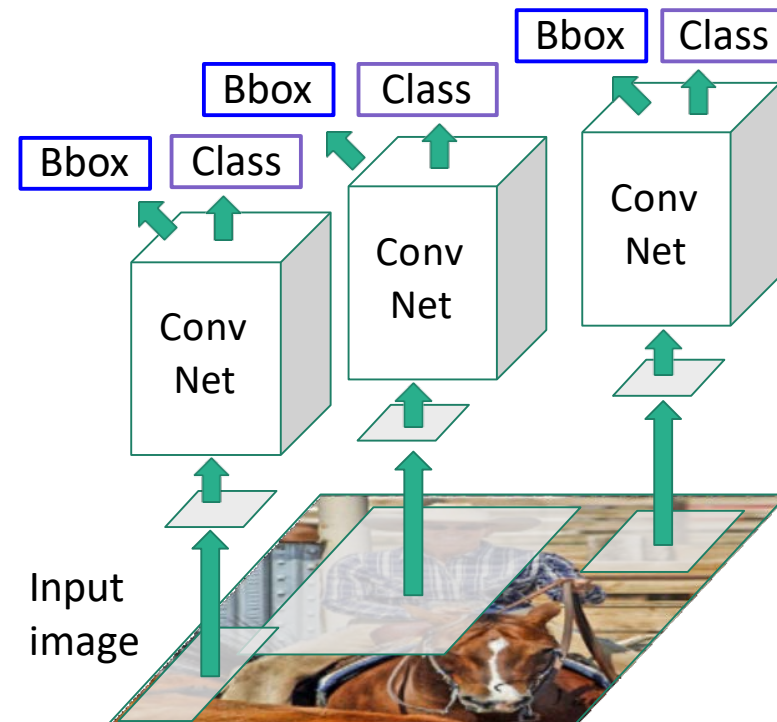


# Fast R-CNN

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



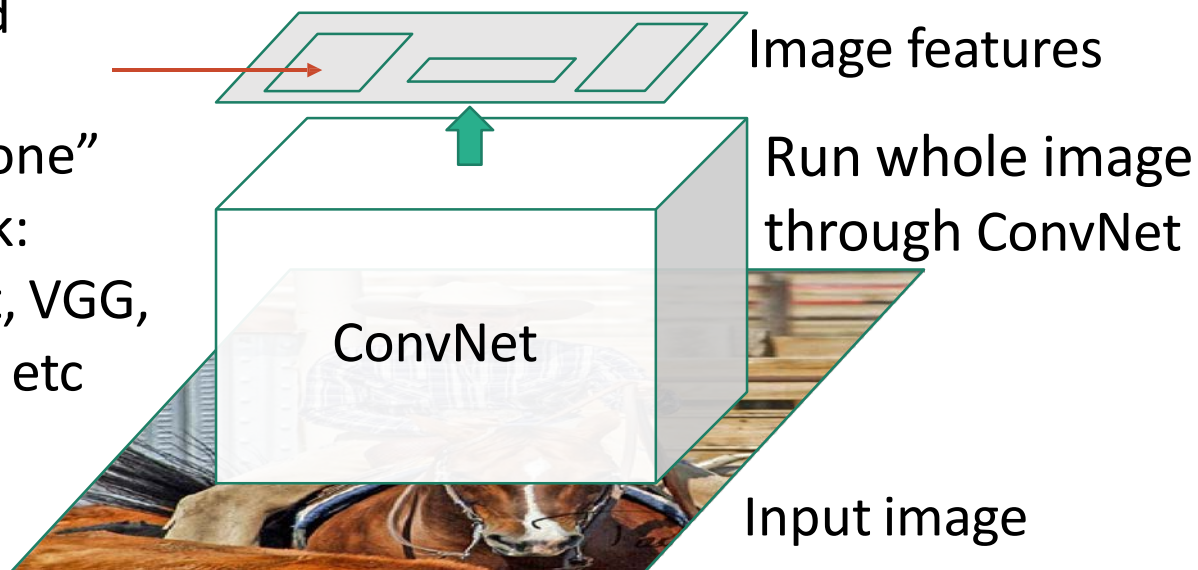
“Slow” R-CNN  
Process each region  
independently



# Fast R-CNN

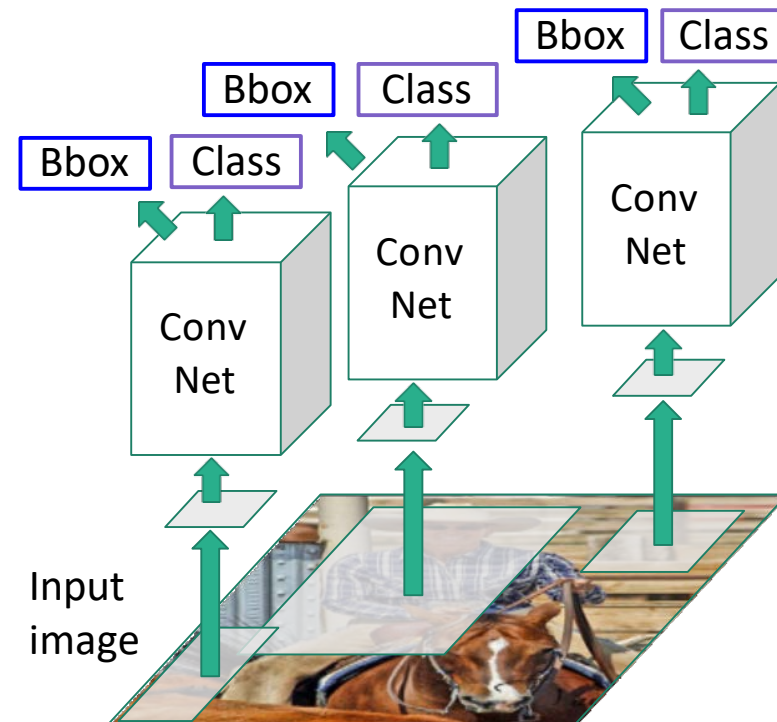
Regions of Interest (Rois) from a proposal method

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



## “Slow” R-CNN

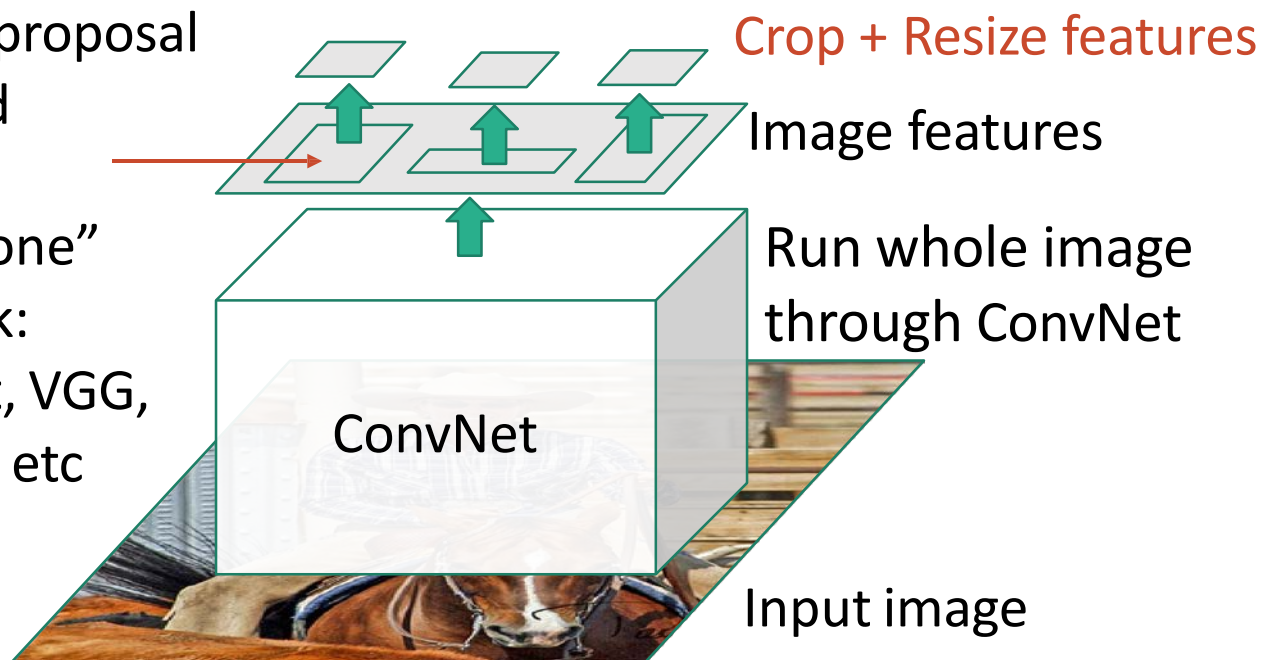
Process each region independently



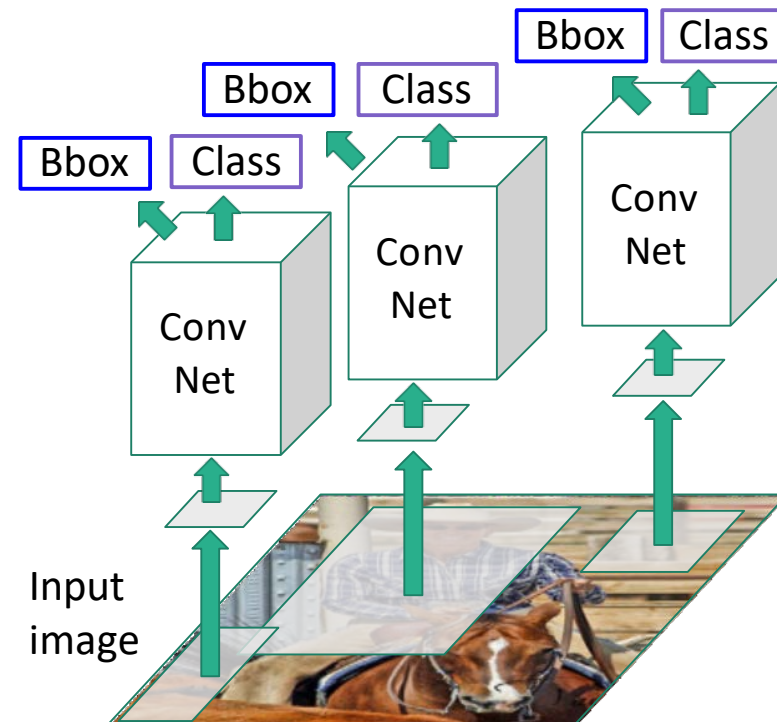
# Fast R-CNN

Regions of Interest (RoIs) from a proposal method

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



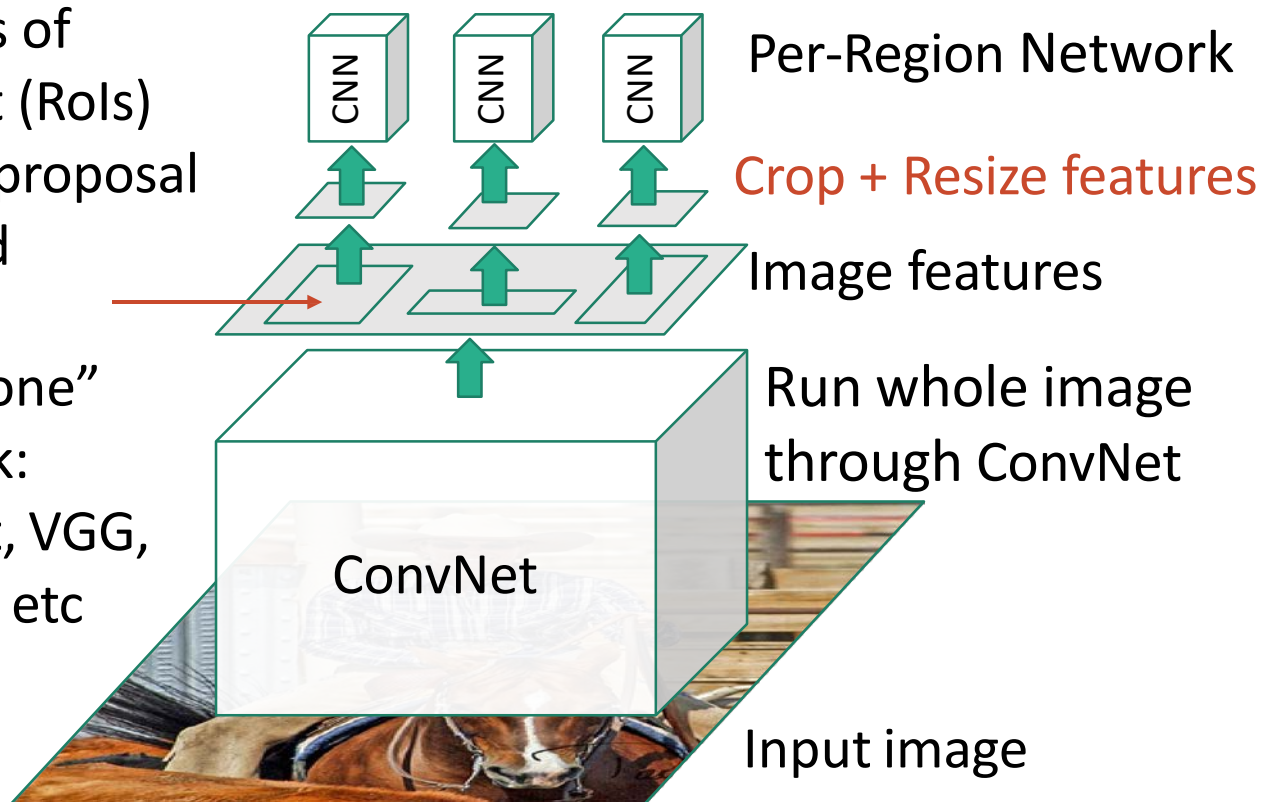
“Slow” R-CNN  
Process each region independently



# Fast R-CNN

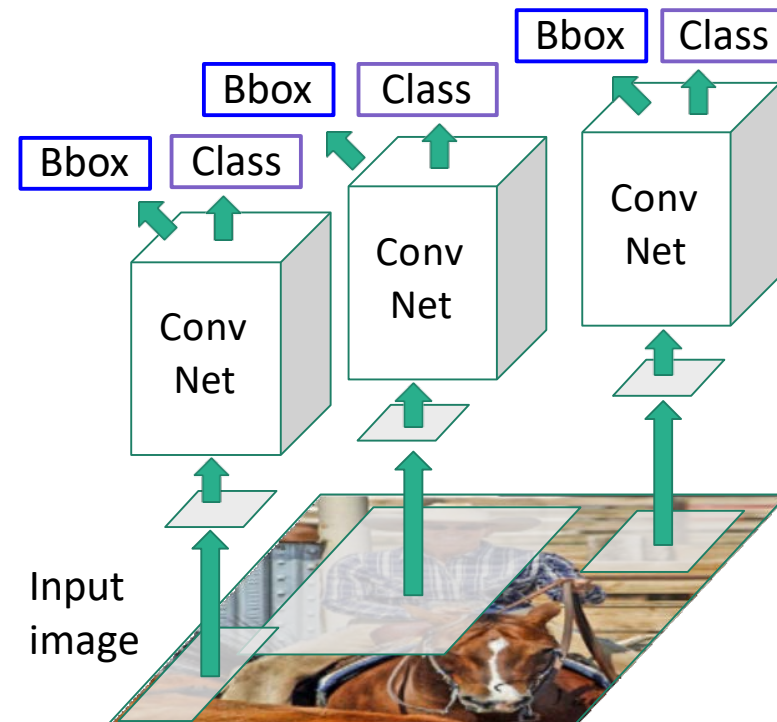
Regions of Interest (RoIs) from a proposal method

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

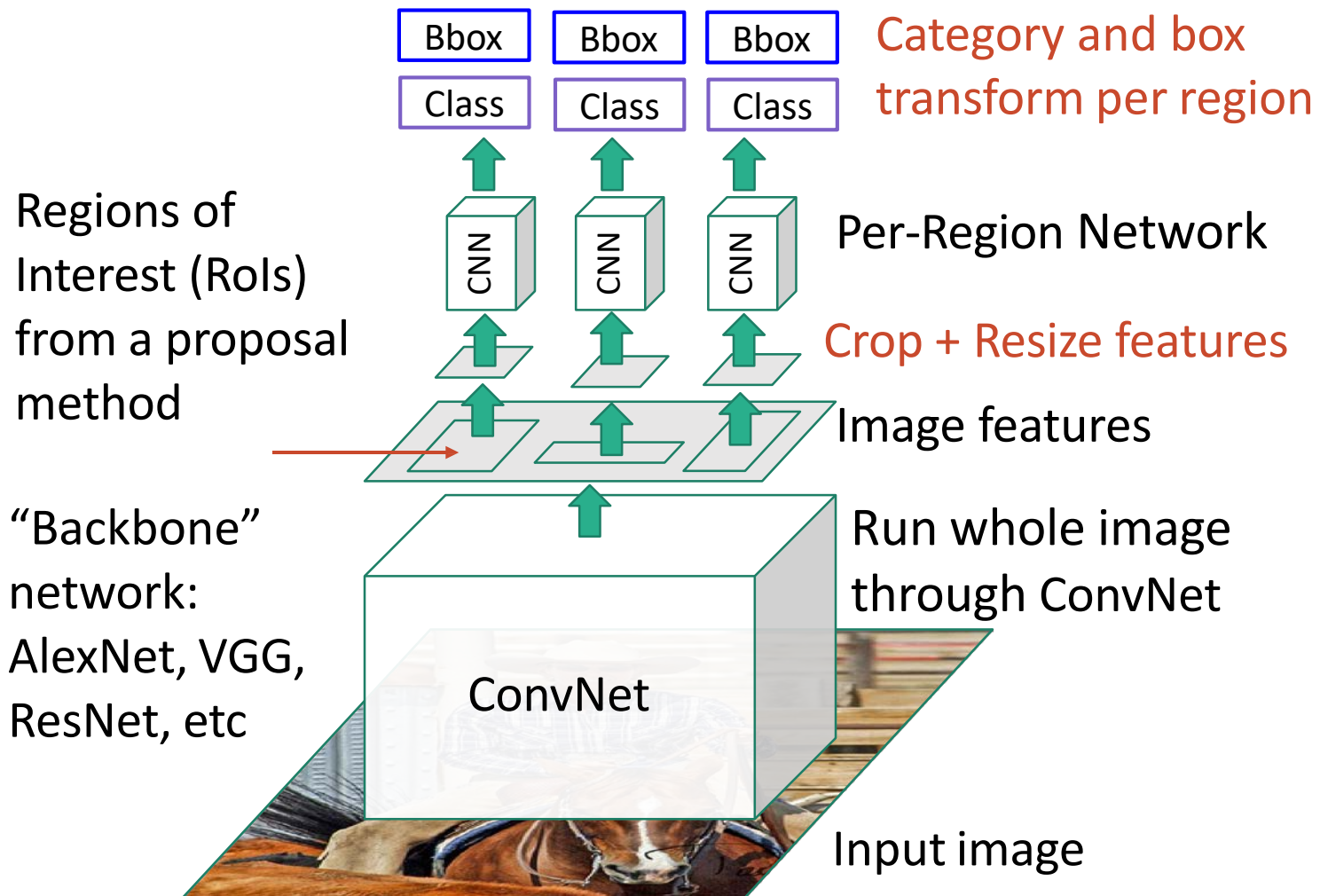


## “Slow” R-CNN

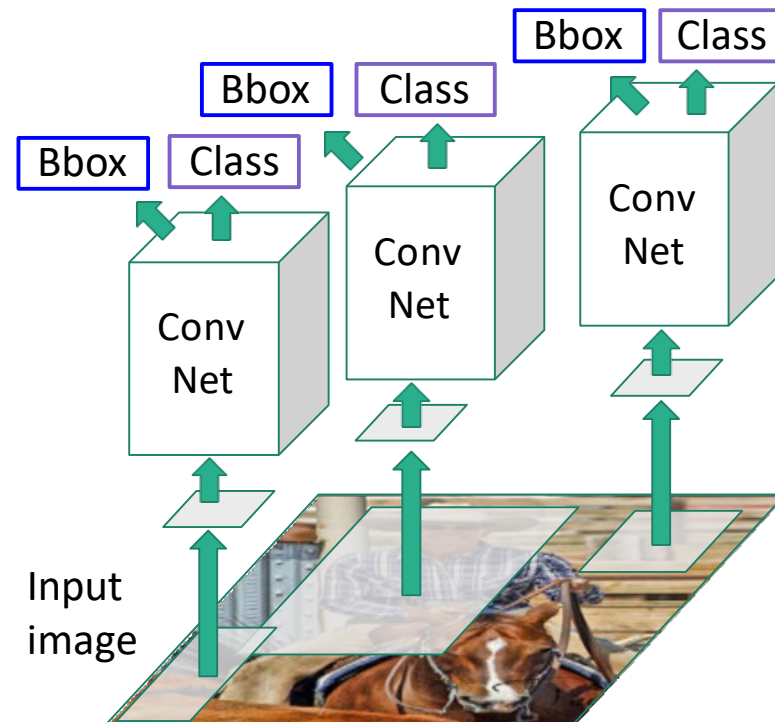
Process each region independently



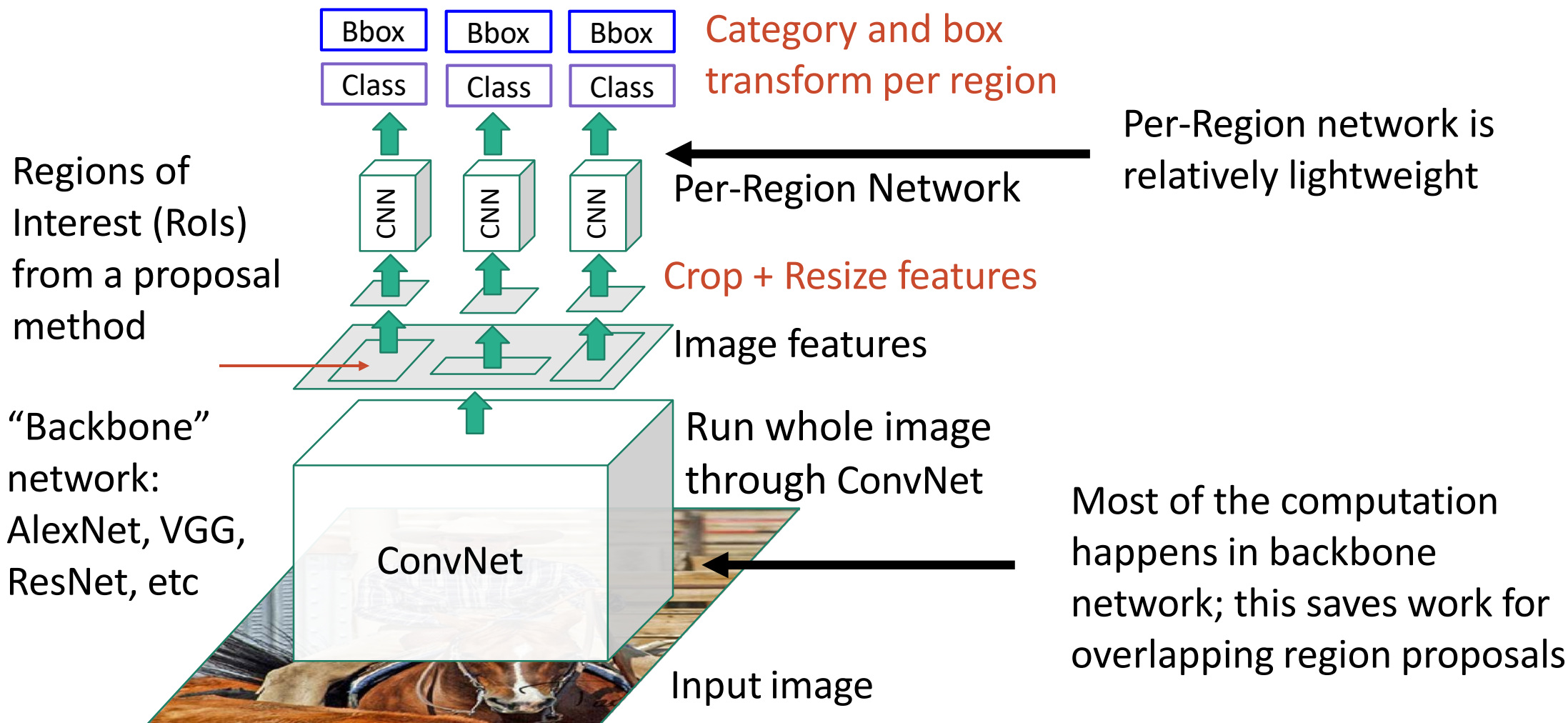
# Fast R-CNN



“Slow” R-CNN  
Process each region independently

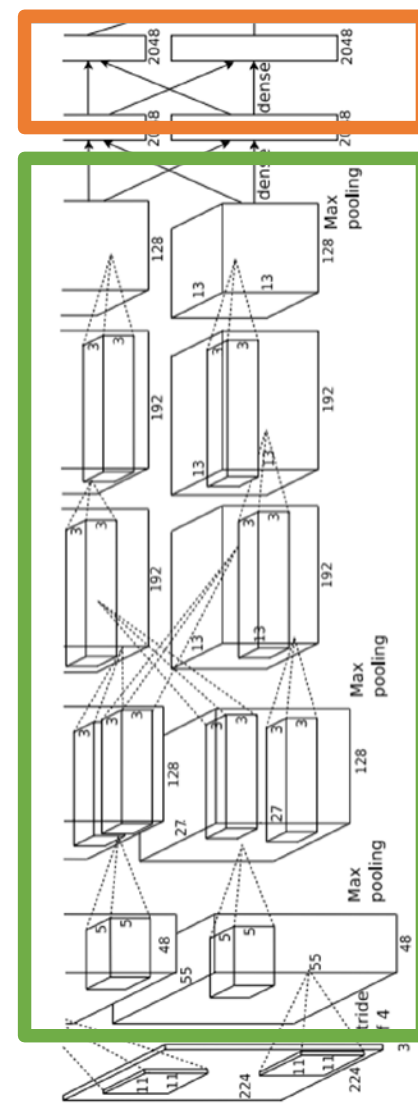
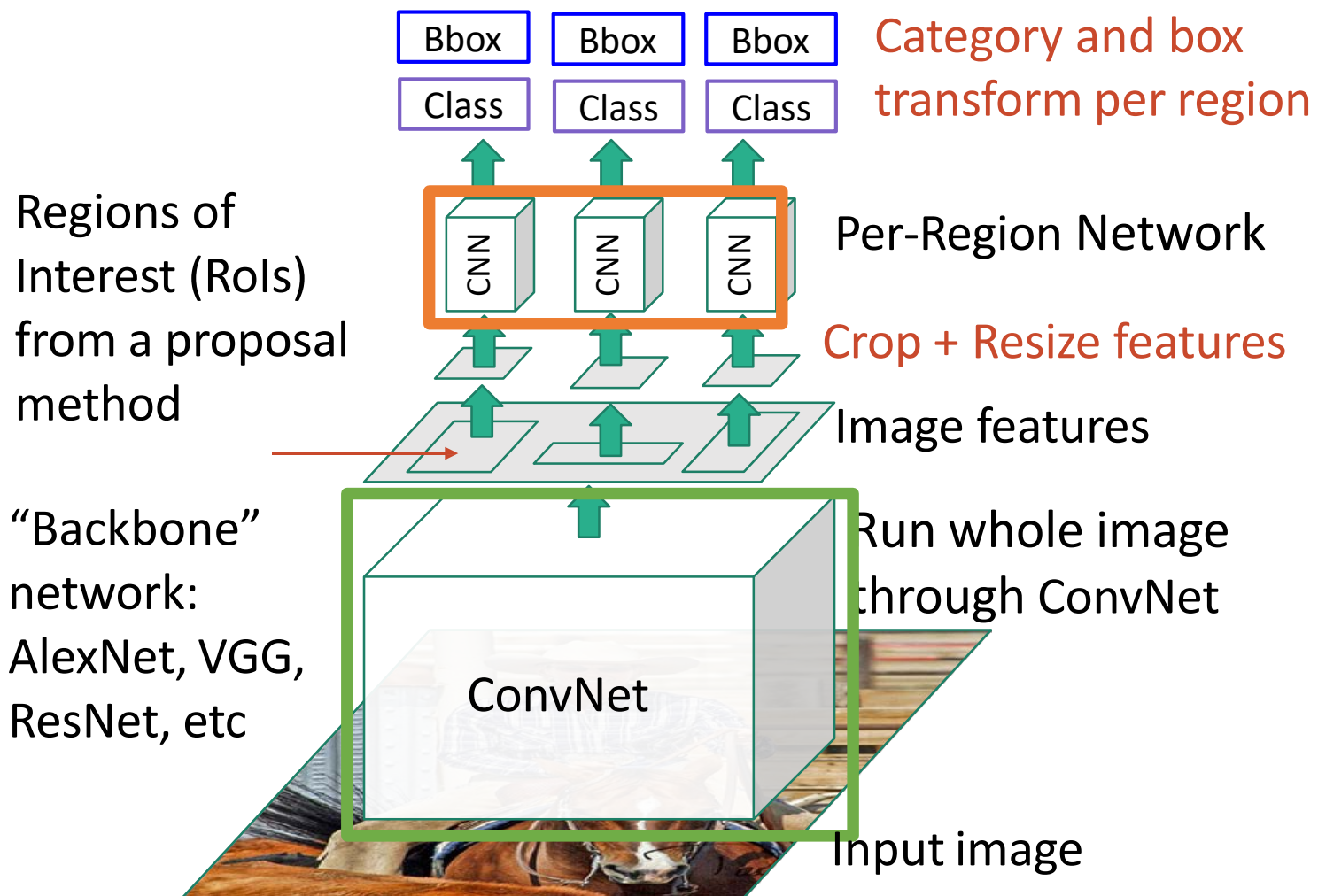


# Fast R-CNN



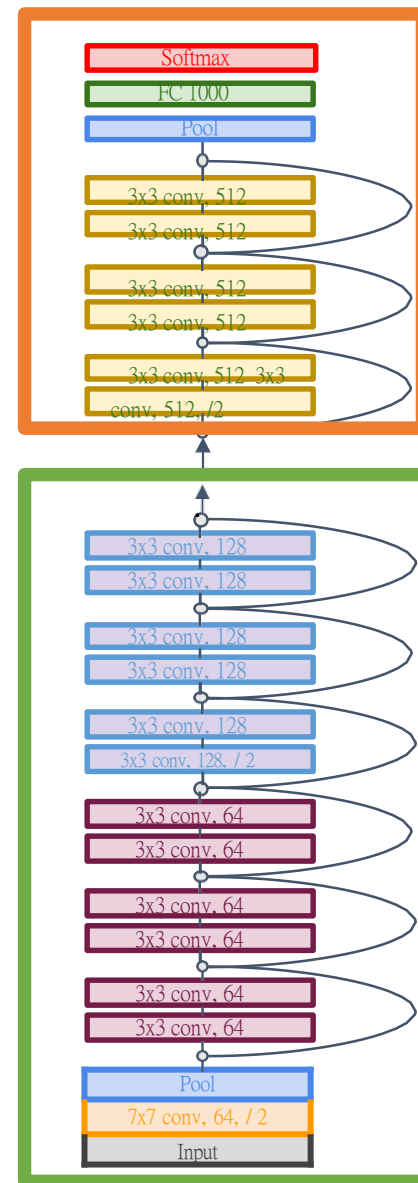
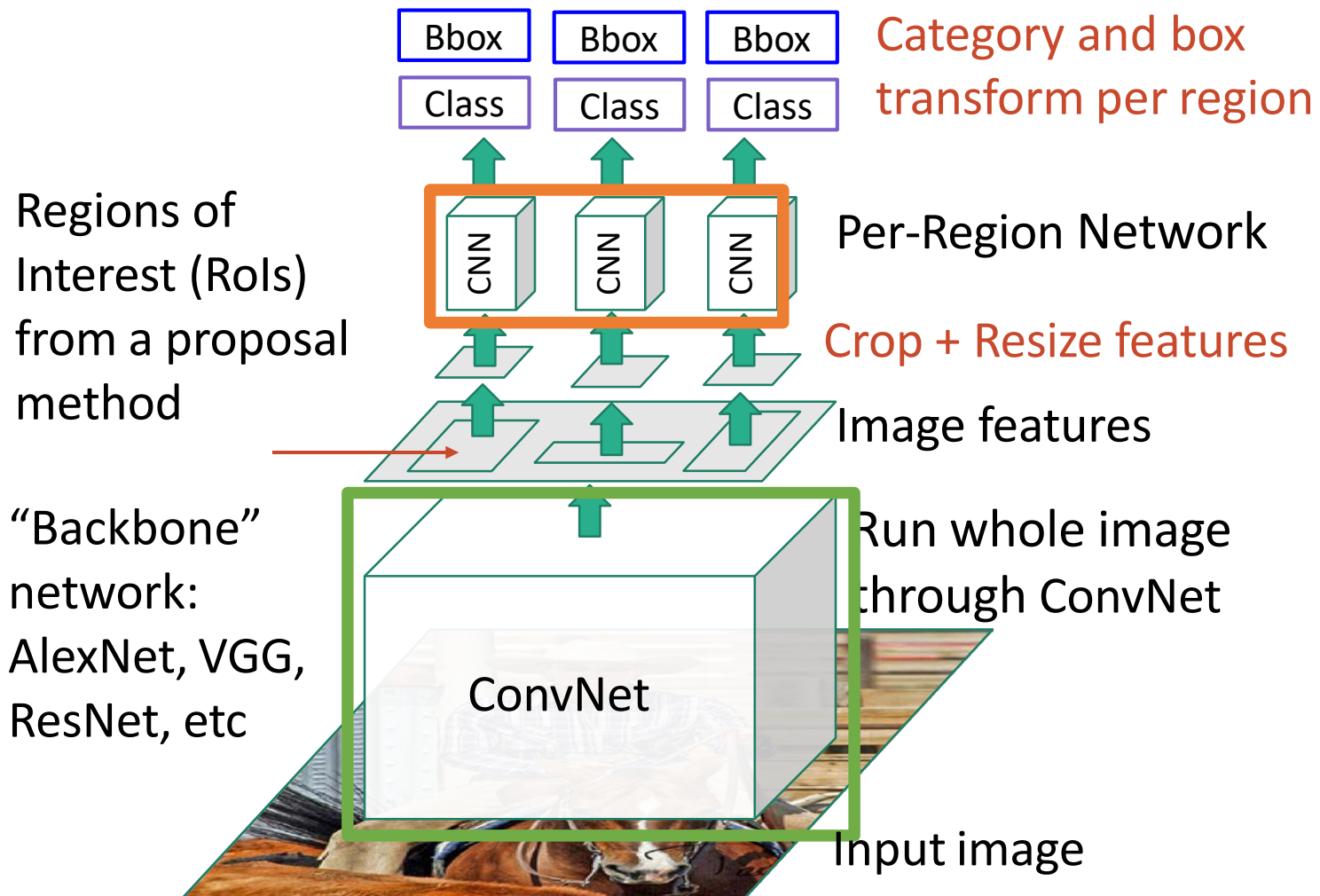


# Fast R-CNN



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

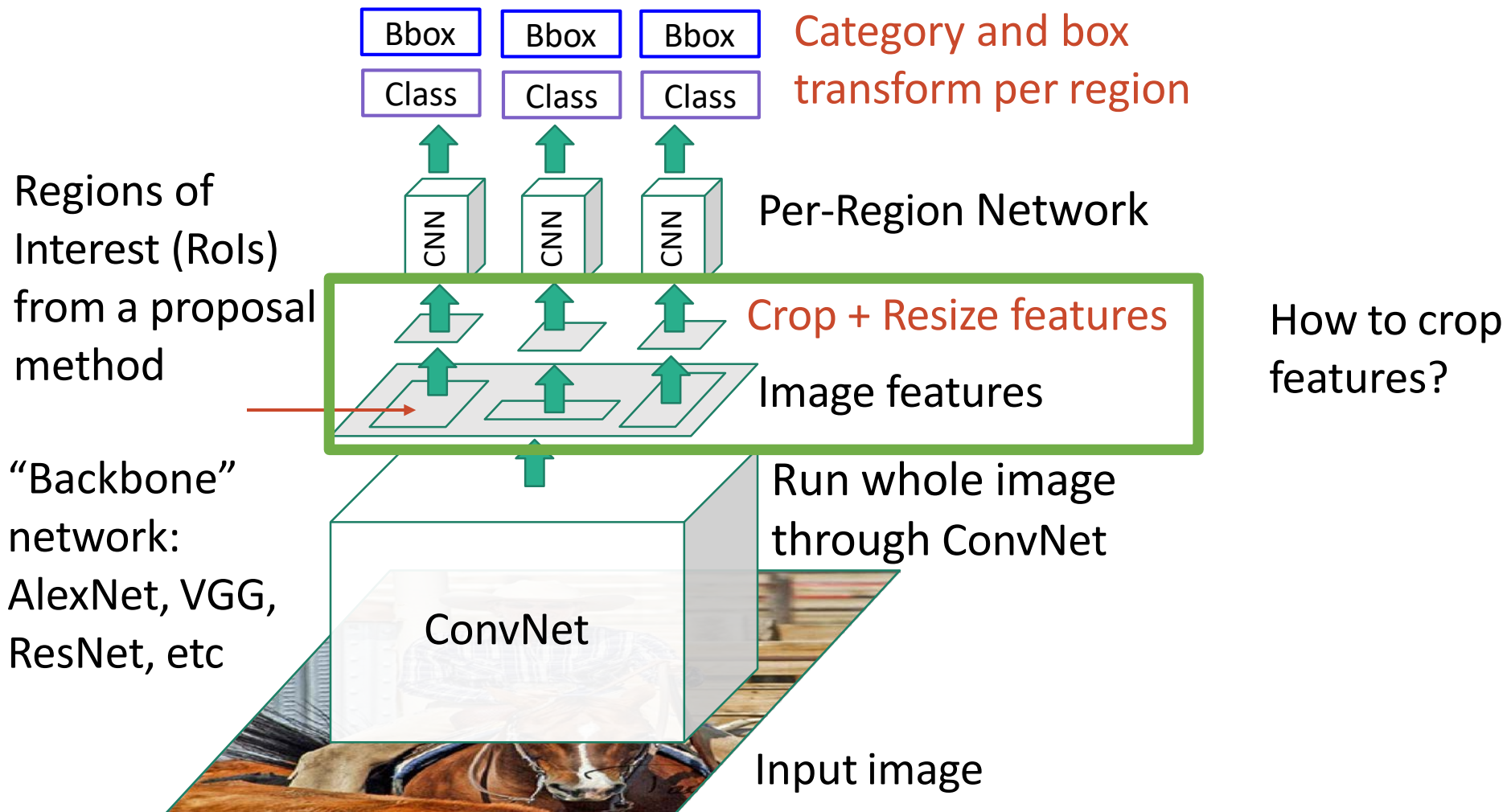
# Fast R-CNN



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



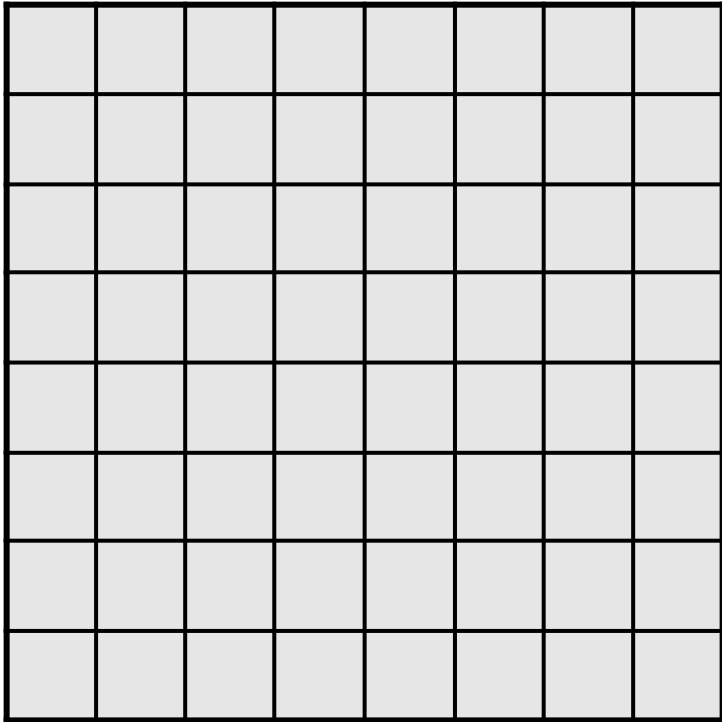
# Fast R-CNN



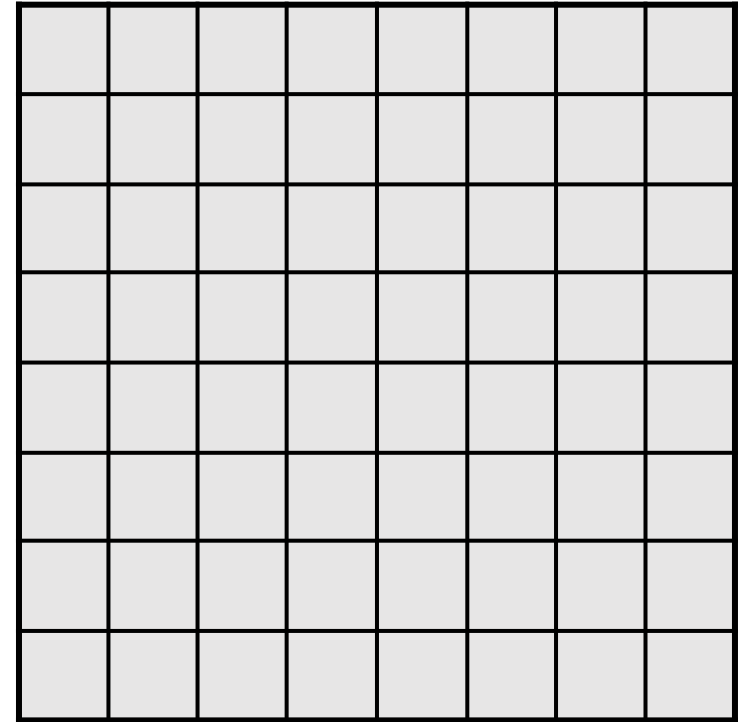
# Recall: Receptive Fields

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

3x3 Conv  
Stride 1, pad 1



Input Image: 8 x 8

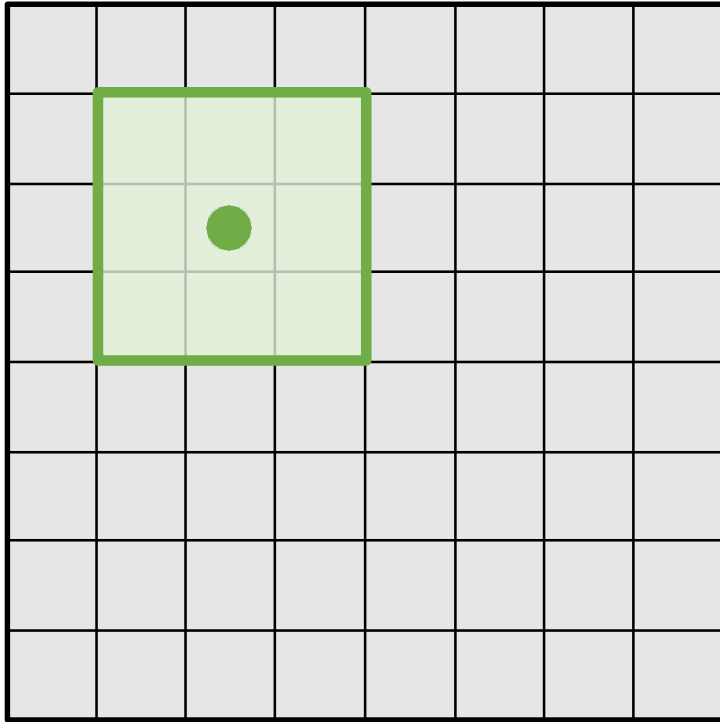


Output Image: 8 x 8

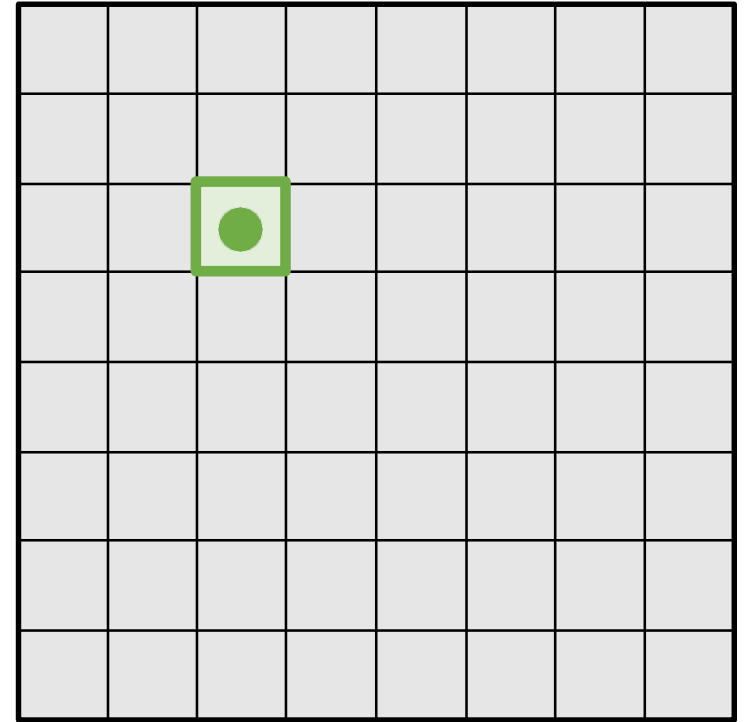
# Recall: Receptive Fields

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

3x3 Conv  
Stride 1, pad 1



Input Image: 8 x 8



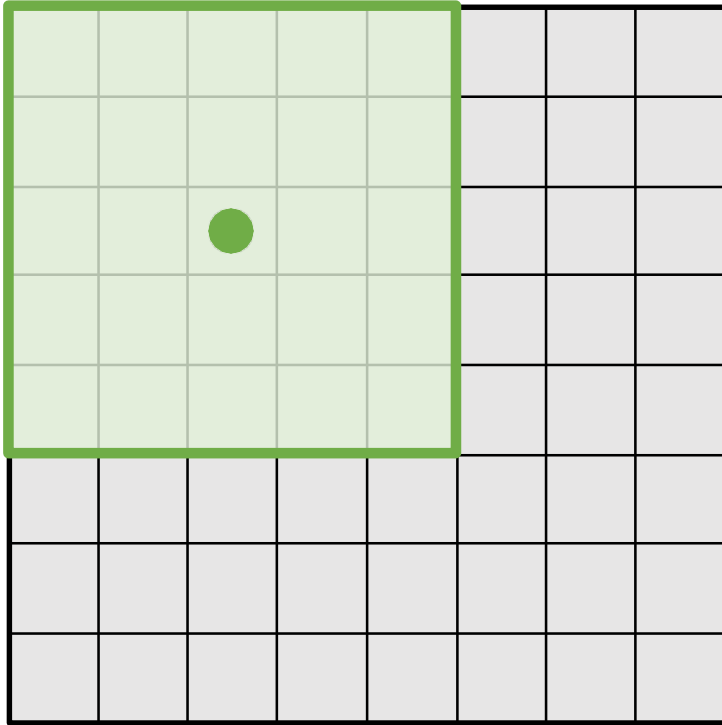
Output Image: 8 x 8

# Recall: Receptive Fields

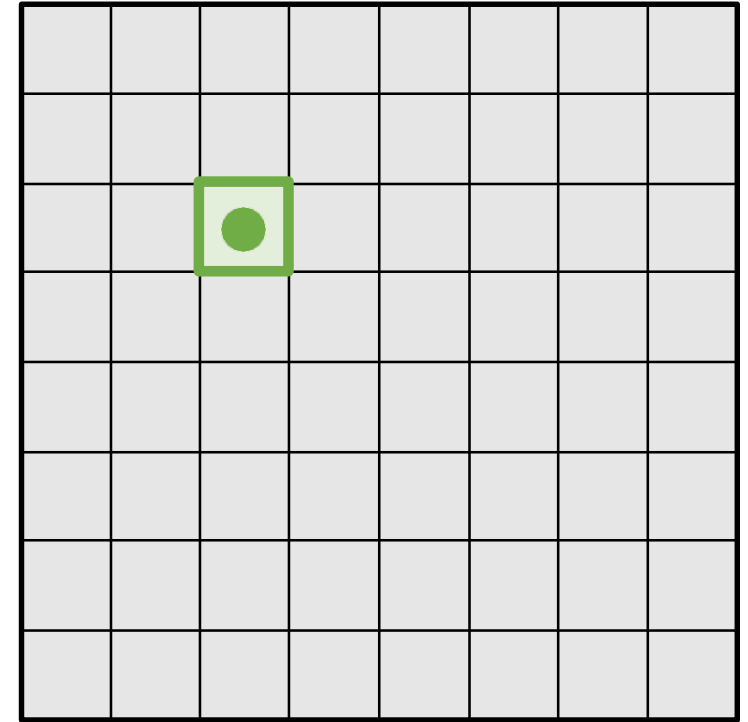
Every position in the output feature map depends on a 5x5 receptive field in the input

3x3 Conv  
Stride 1, pad 1

3x3 Conv  
Stride 1, pad 1

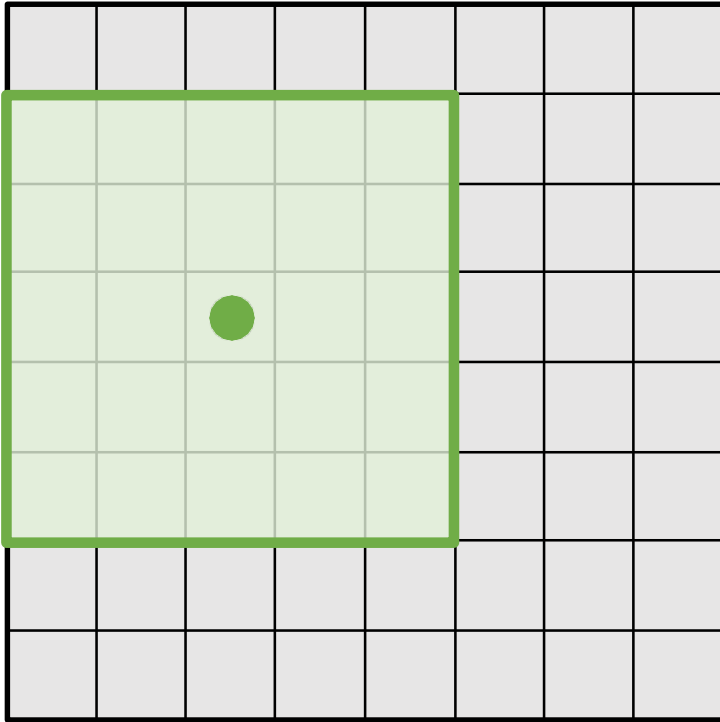


Input Image: 8 x 8



Output Image: 8 x 8

# Recall: Receptive Fields

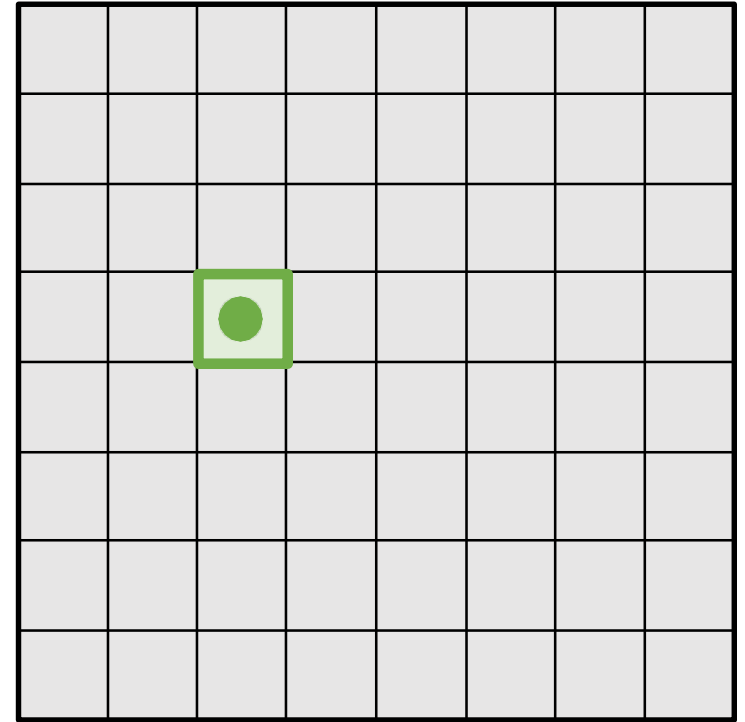


Input Image: 8 x 8

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

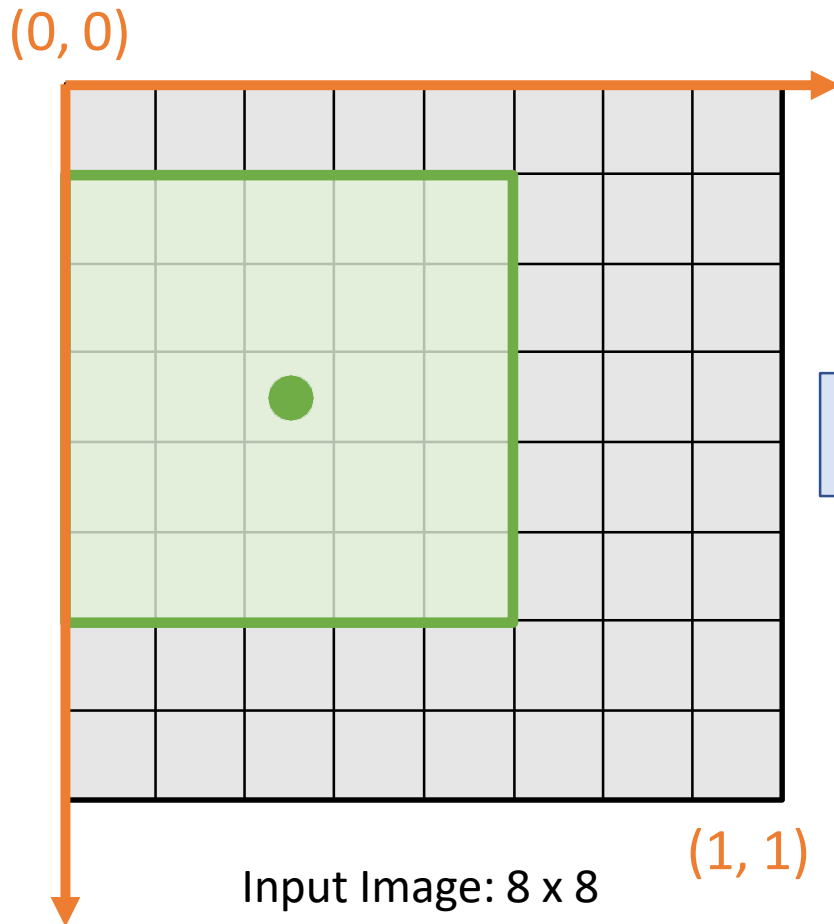
3x3 Conv  
Stride 1, pad 1

3x3 Conv  
Stride 1, pad 1



Output Image: 8 x 8

# Recall: Receptive Fields

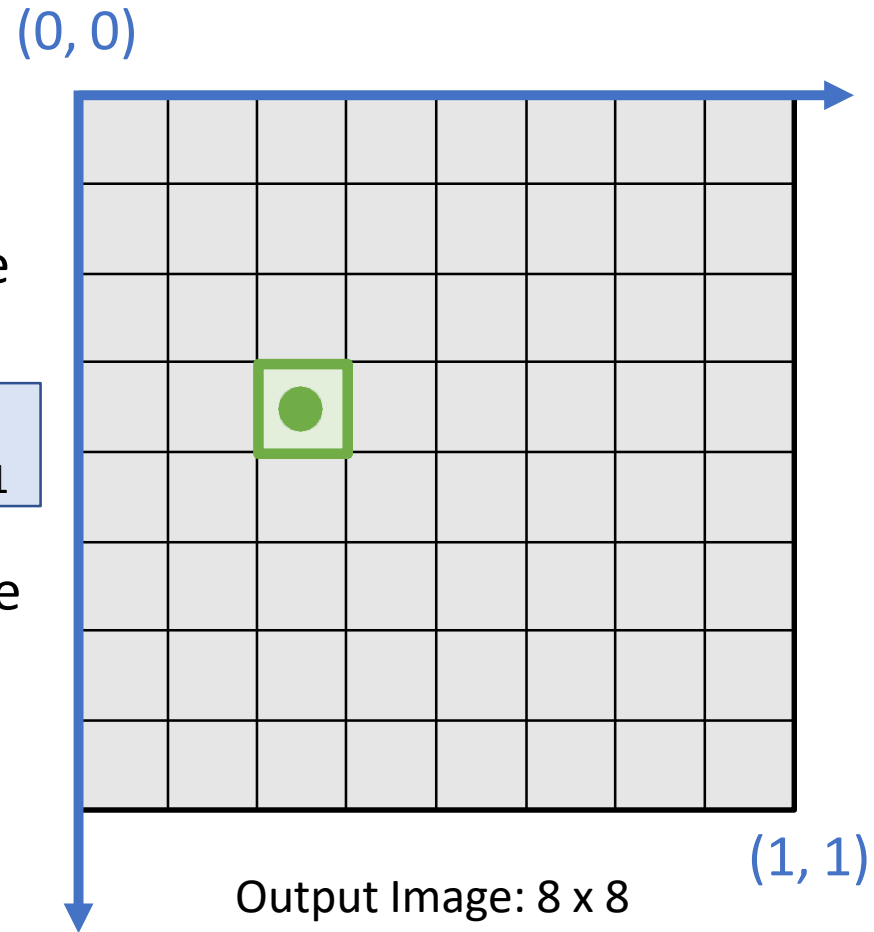


Moving one unit in the output space also moves the receptive field by one

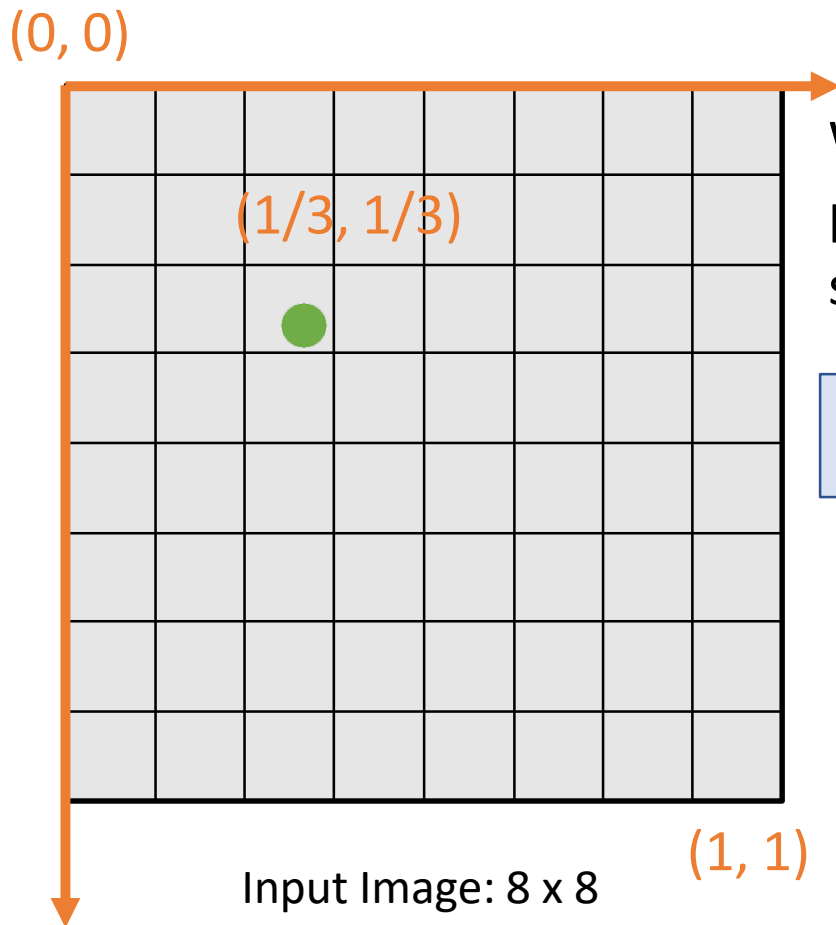
3x3 Conv  
Stride 1, pad 1

3x3 Conv  
Stride 1, pad 1

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



# Projecting Points

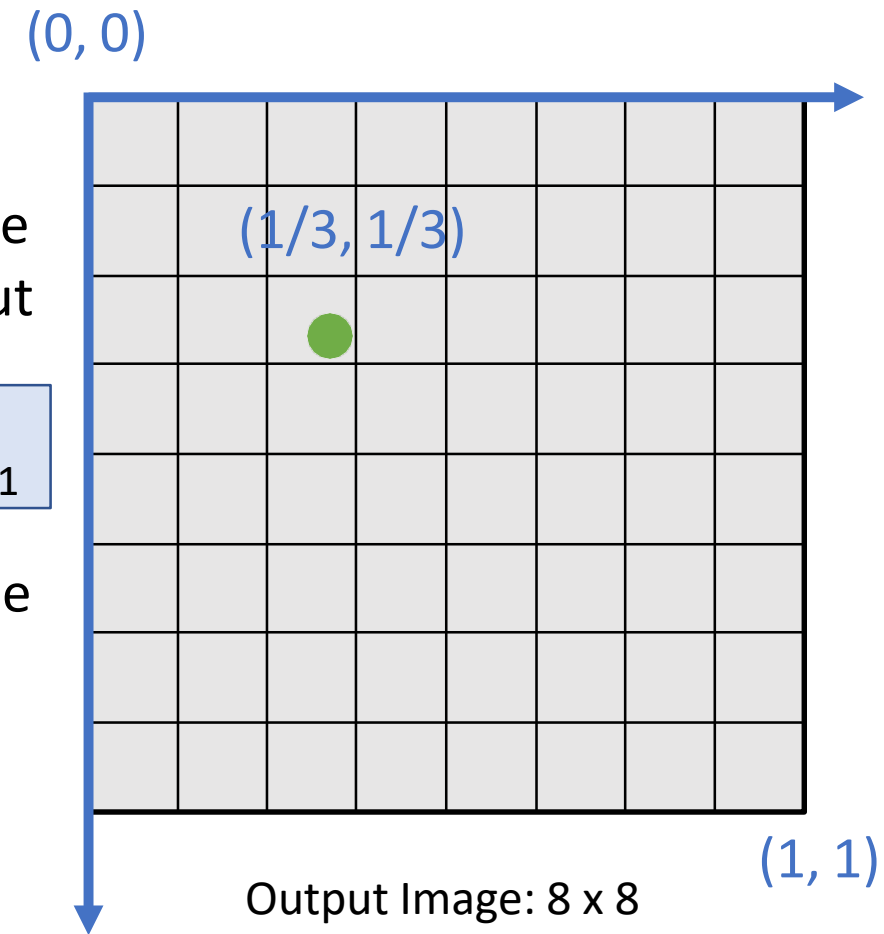


We can align arbitrary points between coordinate system of input and output

3x3 Conv  
Stride 1, pad 1

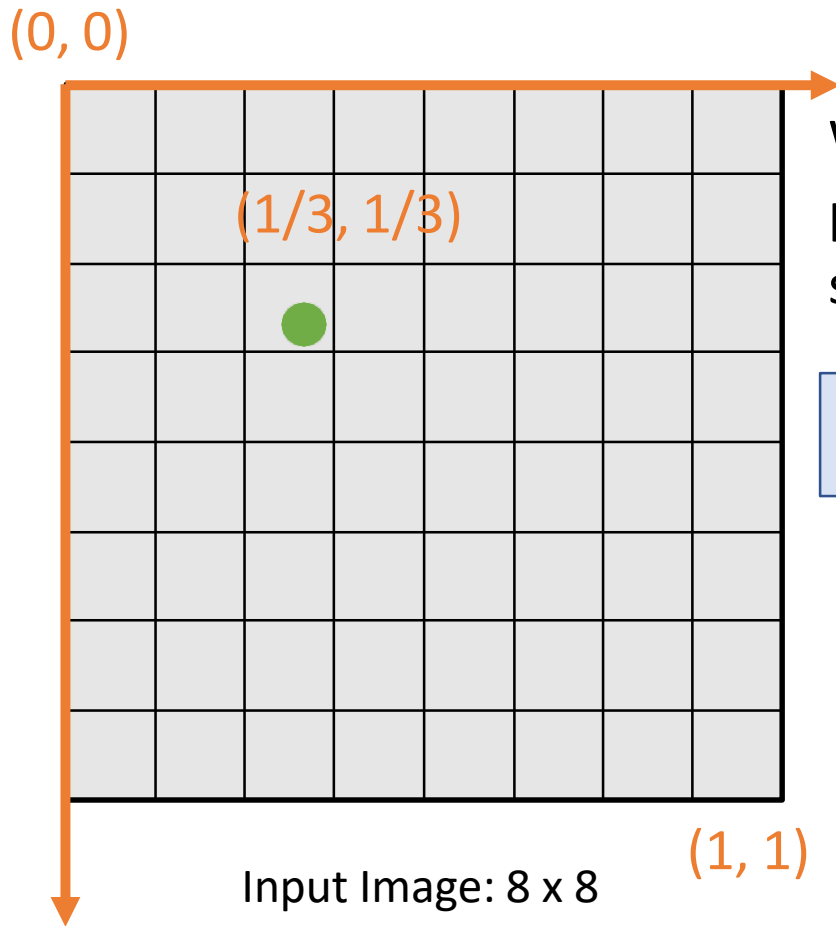
3x3 Conv  
Stride 1, pad 1

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



# Projecting Points

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

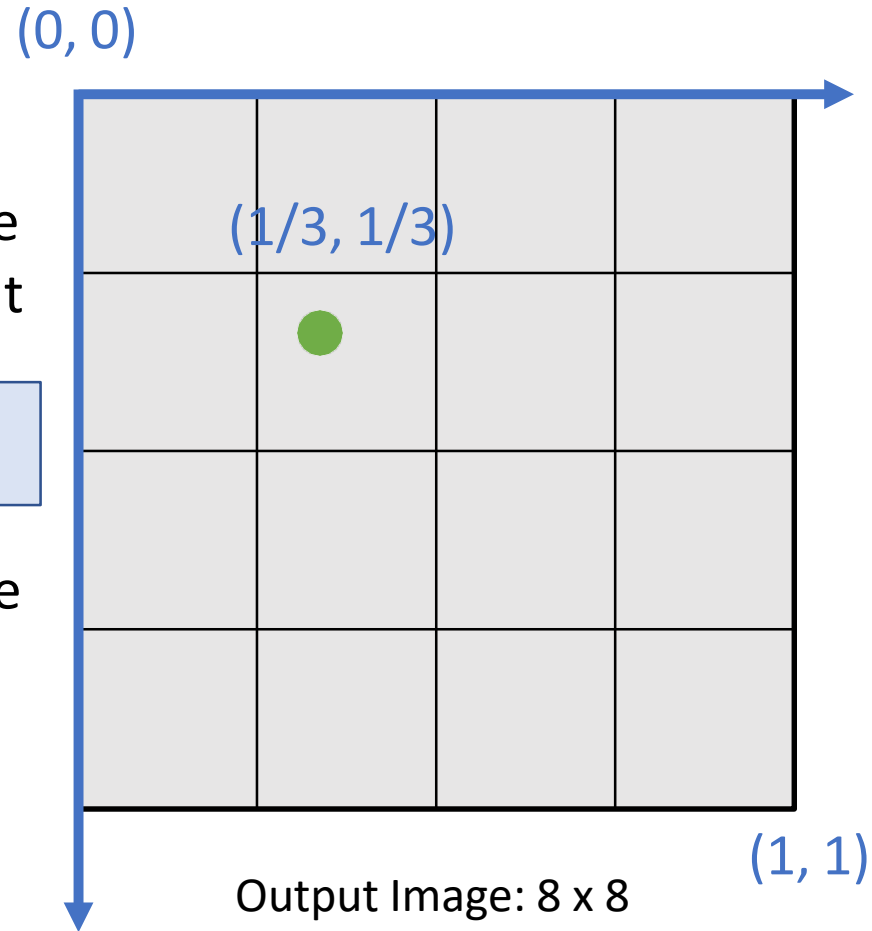


We can align arbitrary points between coordinate system of input and output

3x3 Conv  
Stride 1, pad 1

2x2 MaxPool  
Stride 2

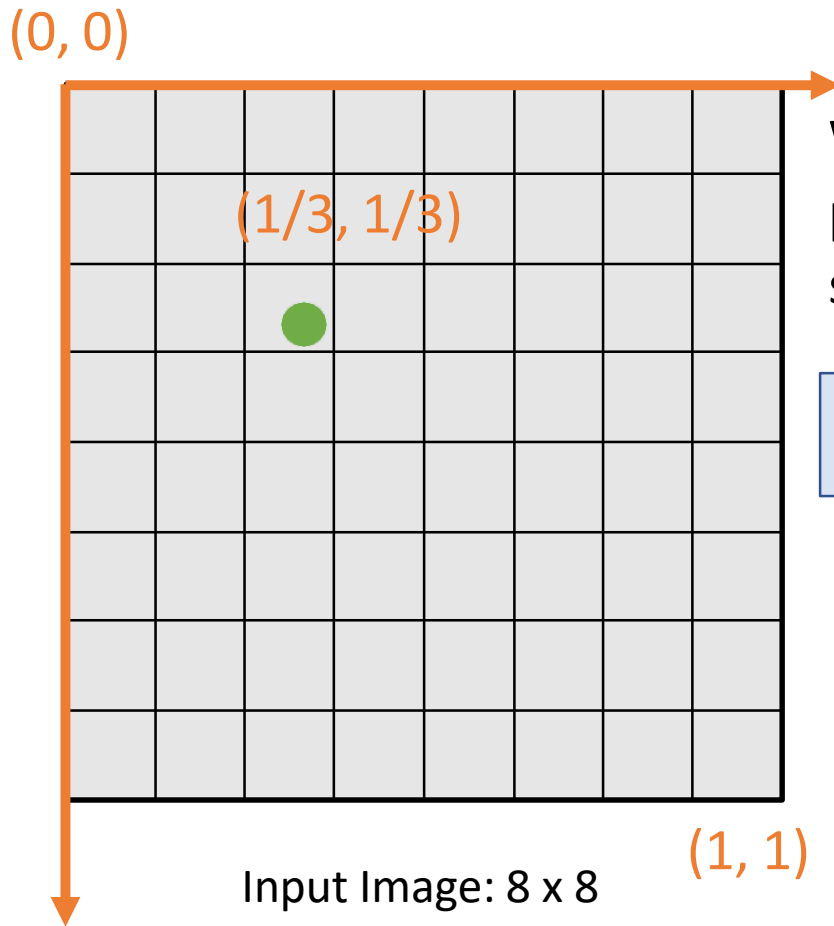
There is a correspondence between the **coordinate system of the input** and the **coordinate system of the output**





# Projecting Points

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

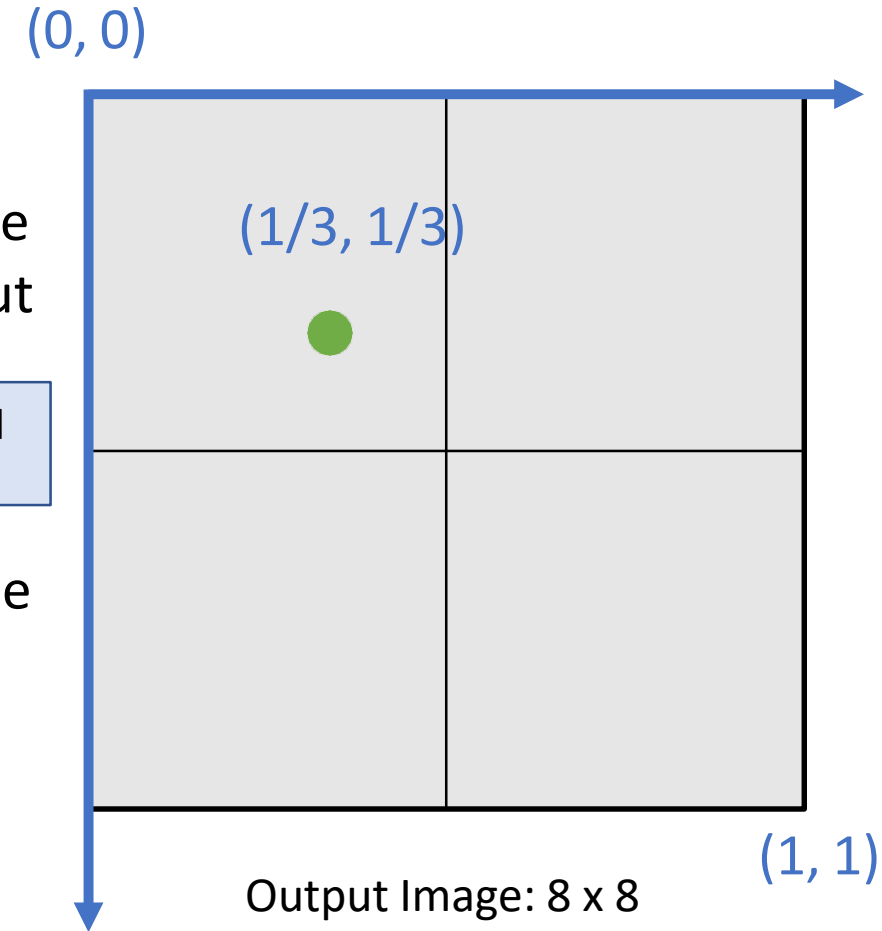


We can align arbitrary points between coordinate system of input and output

3x3 Conv  
Stride 1, pad 1

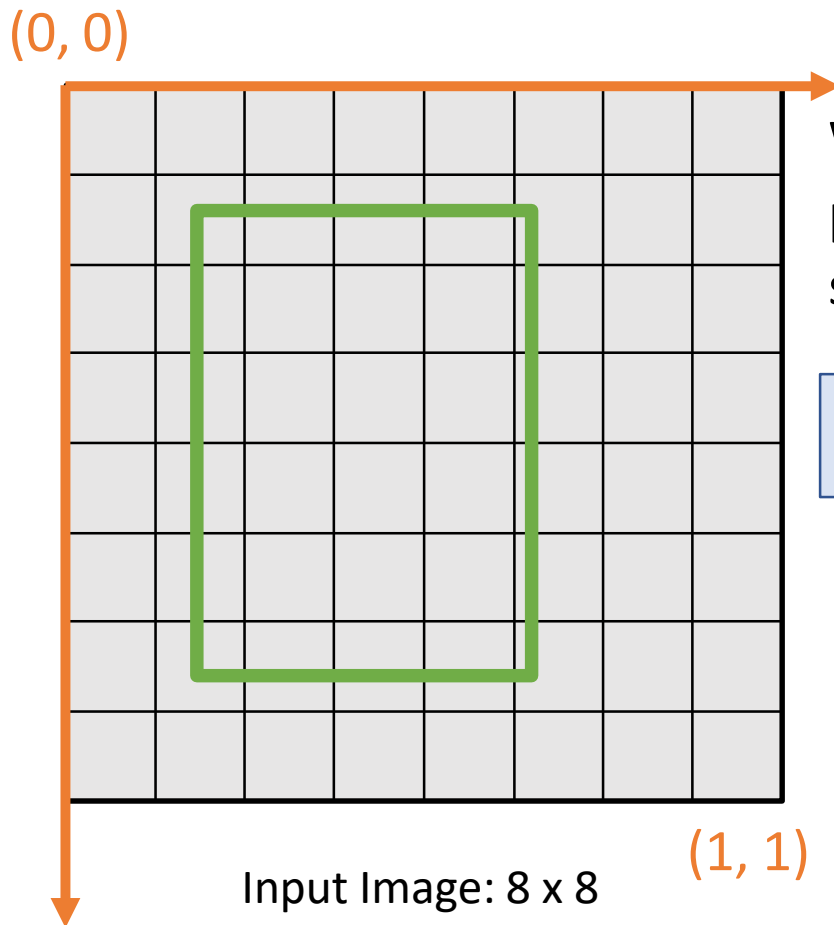
4x4 MaxPool  
Stride 4

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



# Projecting Boxes

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

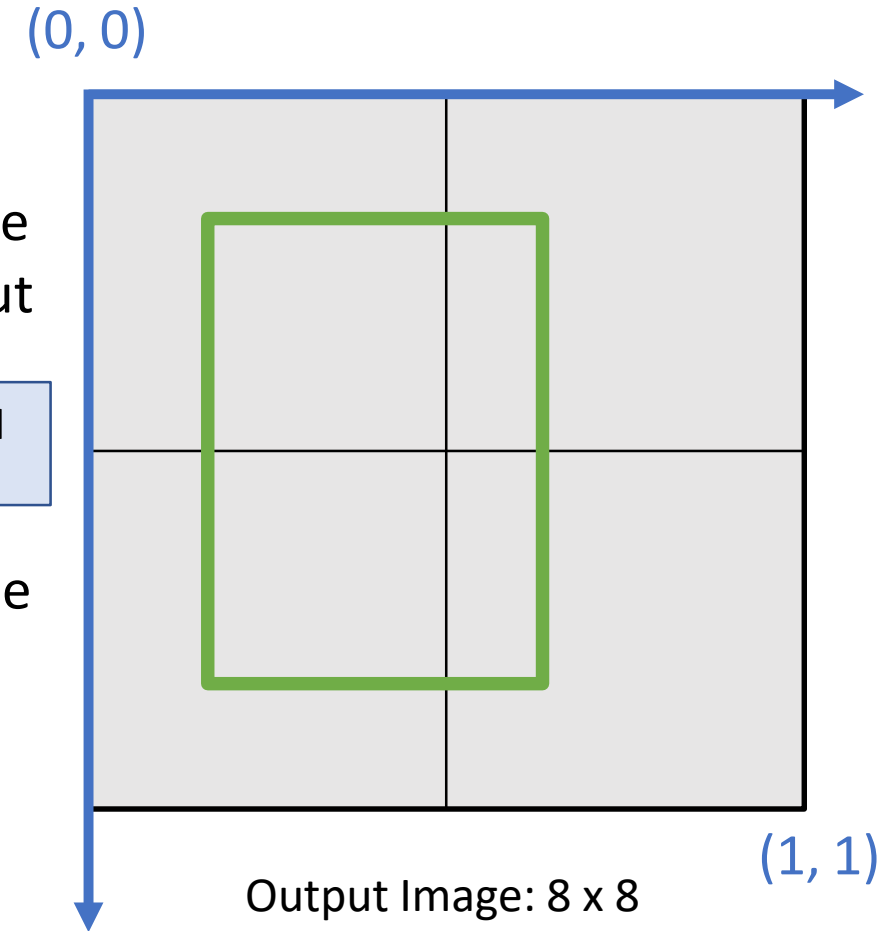


We can align arbitrary points between coordinate system of input and output

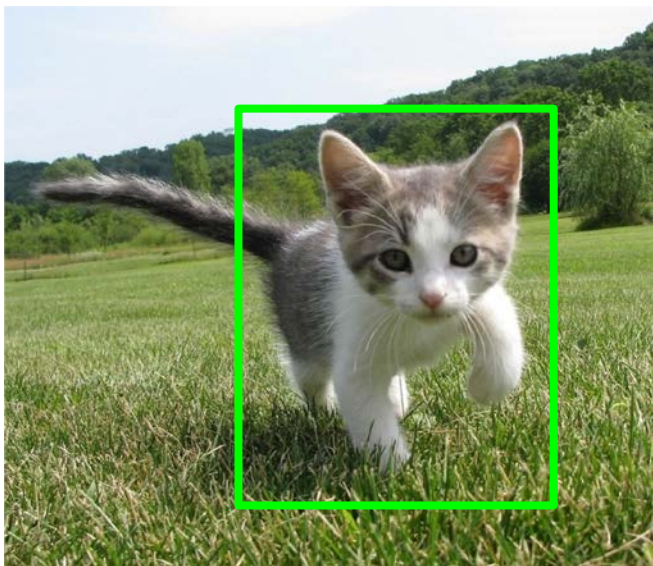
3x3 Conv  
Stride 1, pad 1

4x4 MaxPool  
Stride 4

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

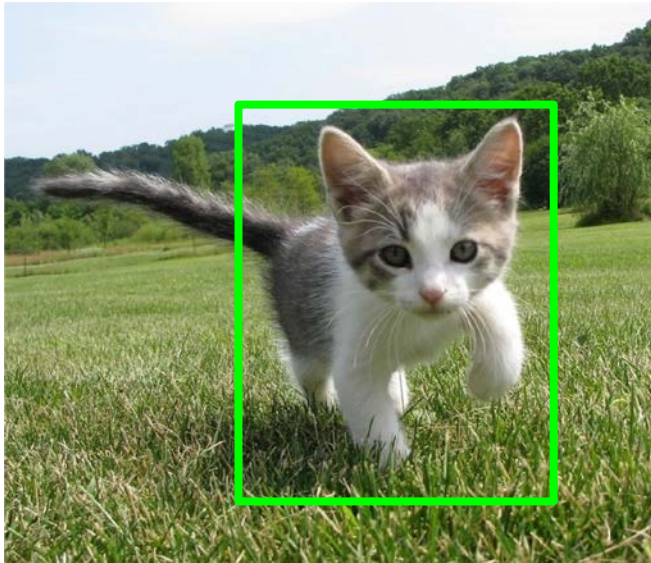


# Cropping Features: RoI Pool



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

# Cropping Features: RoI Pool



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

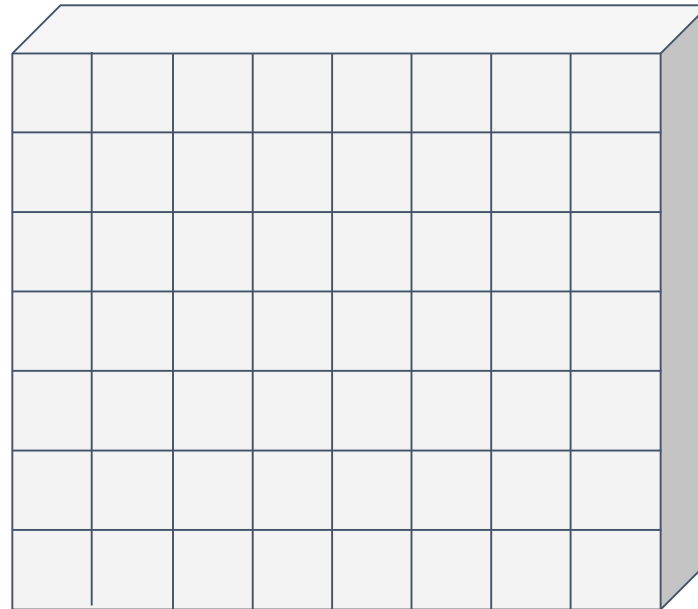
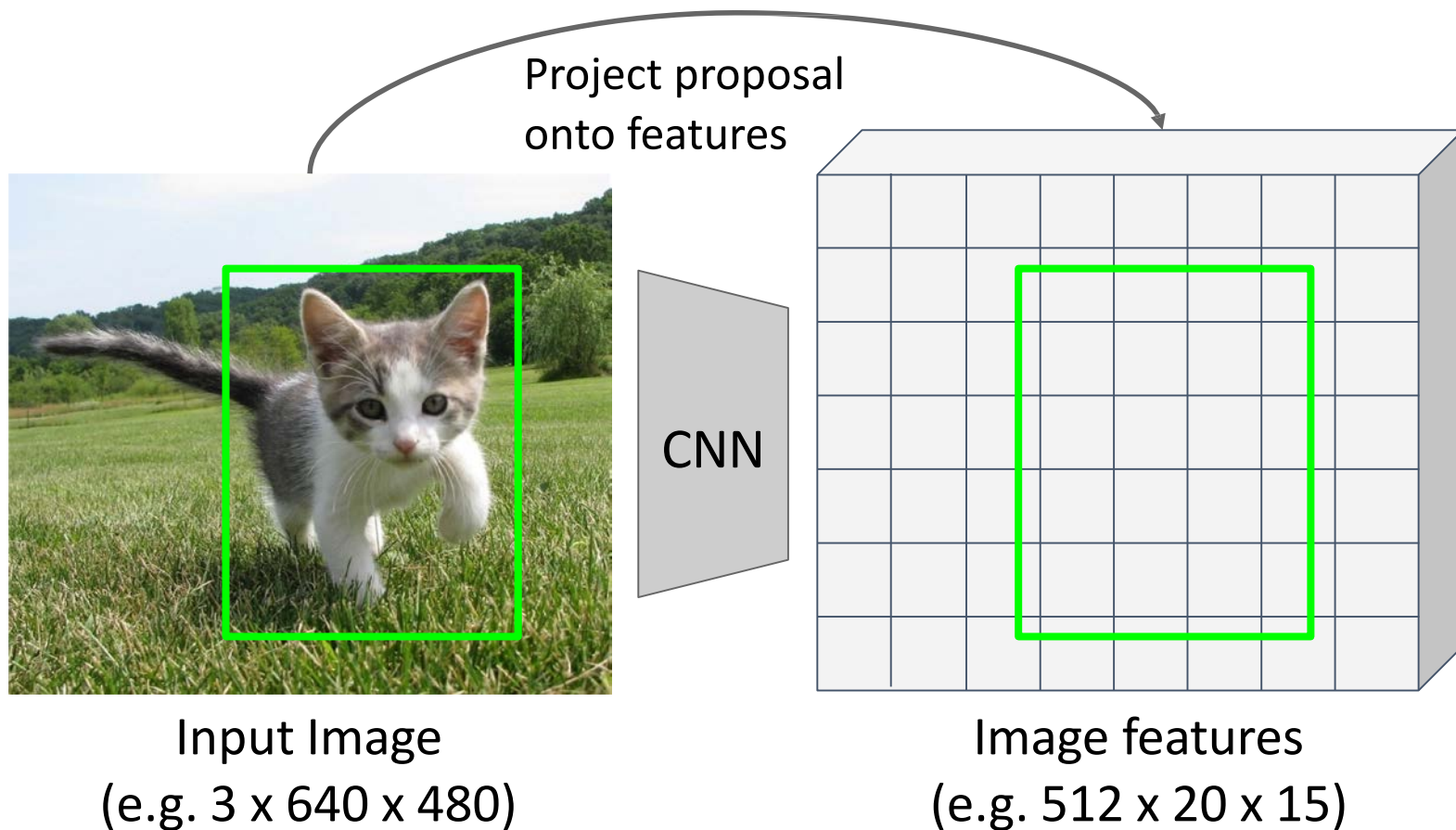


Image features  
(e.g. 512 x 20 x 15)

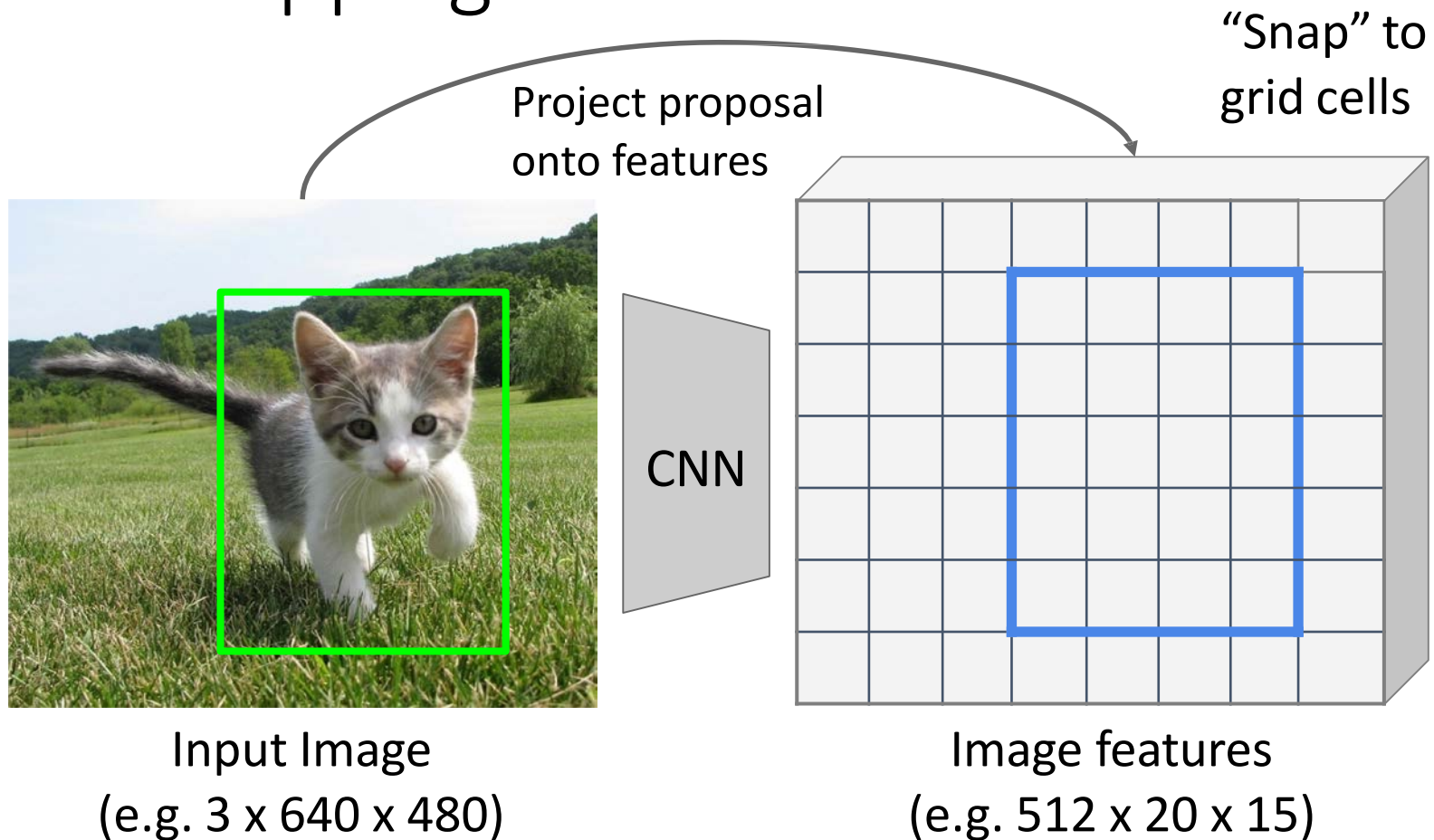
Want features for the  
box of a fixed size  
(2x2 in this example,  
7x7 or 14x14 in practice)

# Cropping Features: RoI Pool



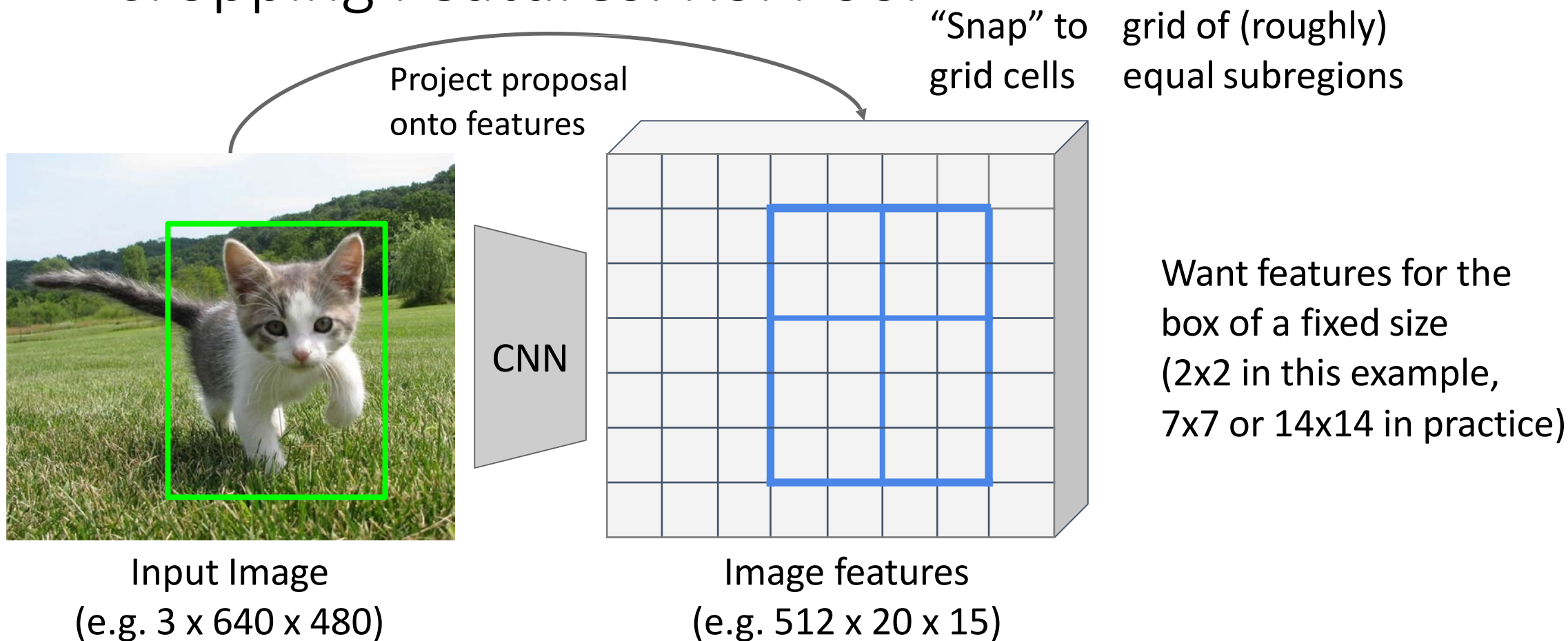
Want features for the box of a fixed size (2x2 in this example, 7x7 or 14x14 in practice)

# Cropping Features: RoI Pool



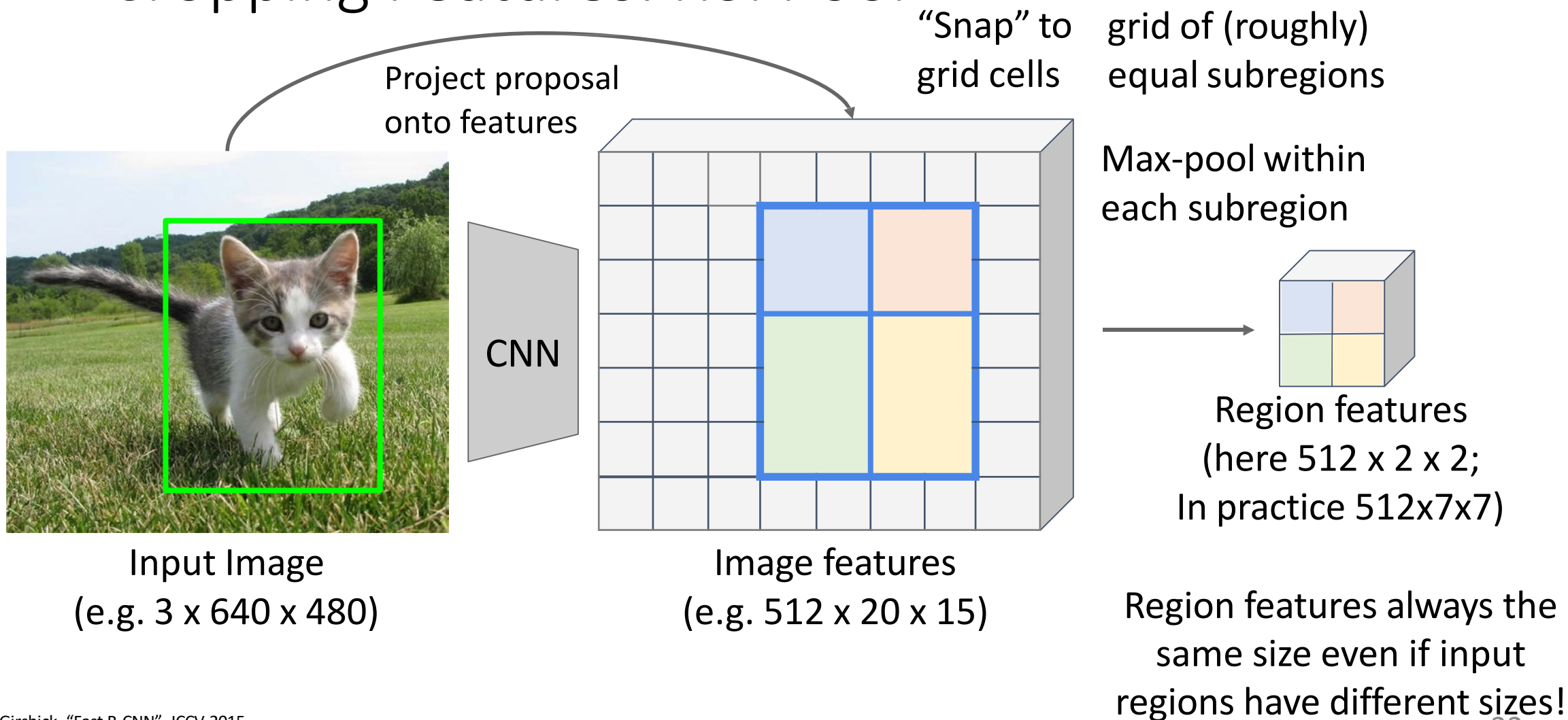
Want features for the box of a fixed size (2x2 in this example, 7x7 or 14x14 in practice)

# Cropping Features: RoI Pool



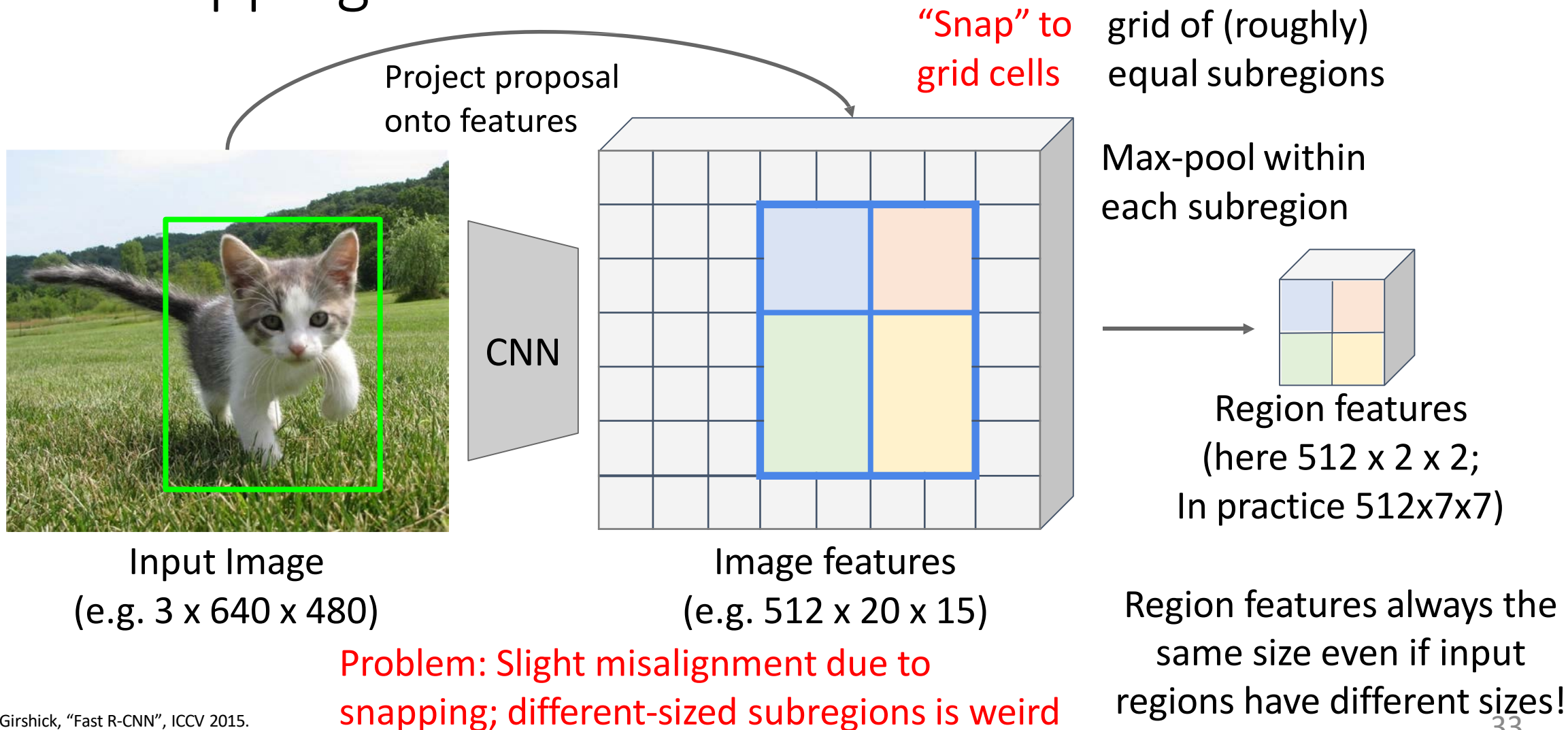


# Cropping Features: RoI Pool



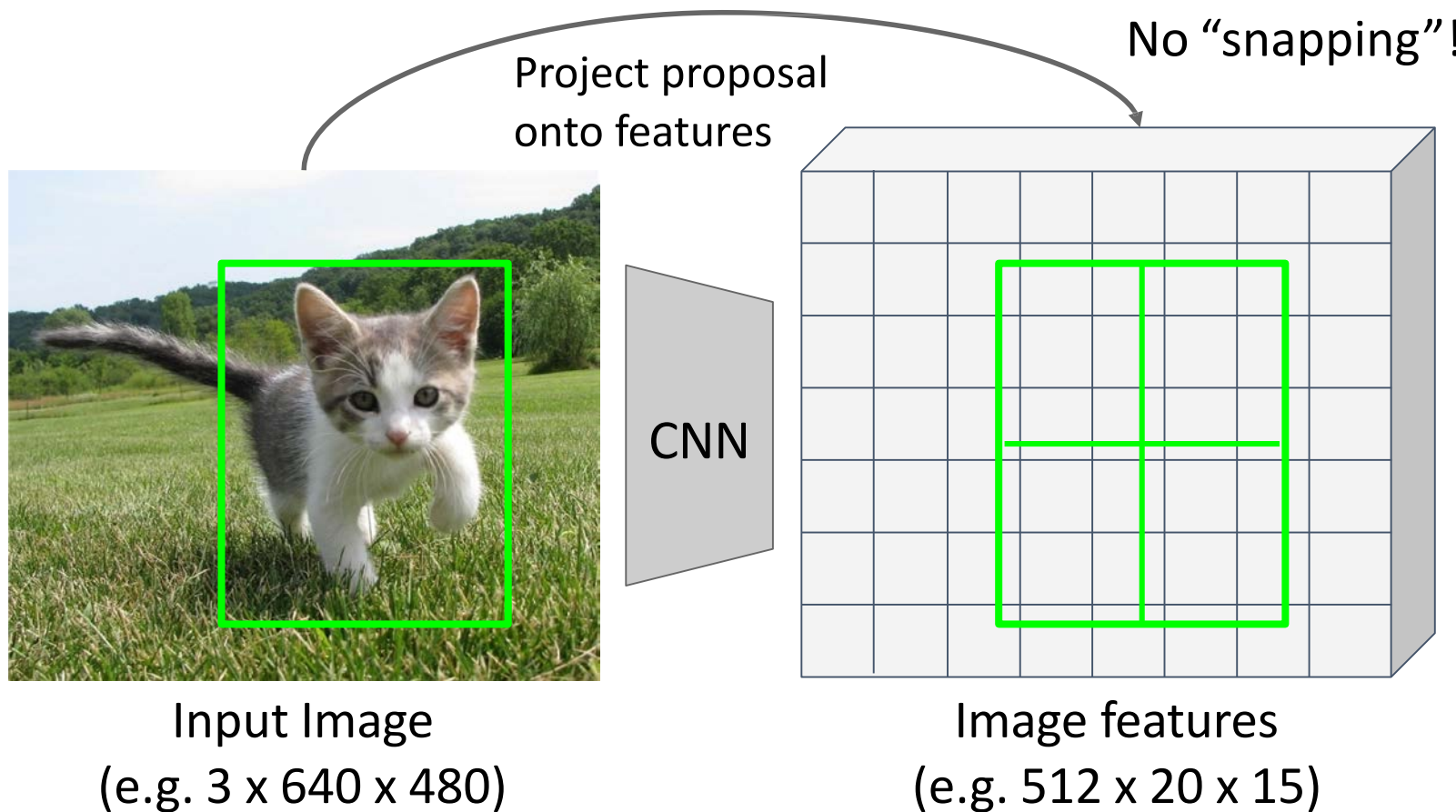


# Cropping Features: RoI Pool



# Cropping Features: RoI Align

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



Want features for the  
box of a fixed size  
(2x2 in this example,  
7x7 or 14x14 in practice)

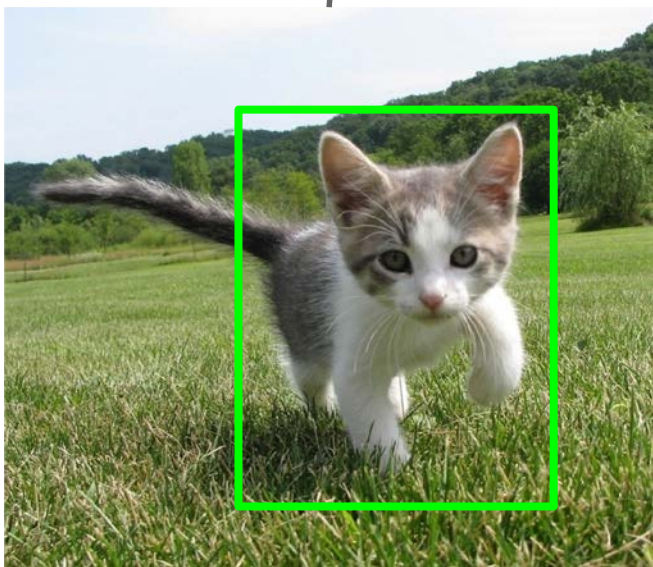
# Cropping Features: RoI Align

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

No “snapping”!

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

Project proposal  
onto features



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

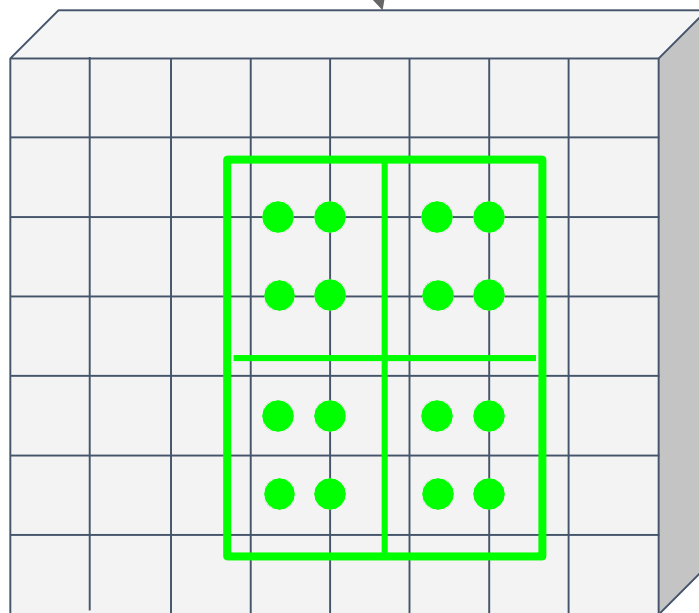
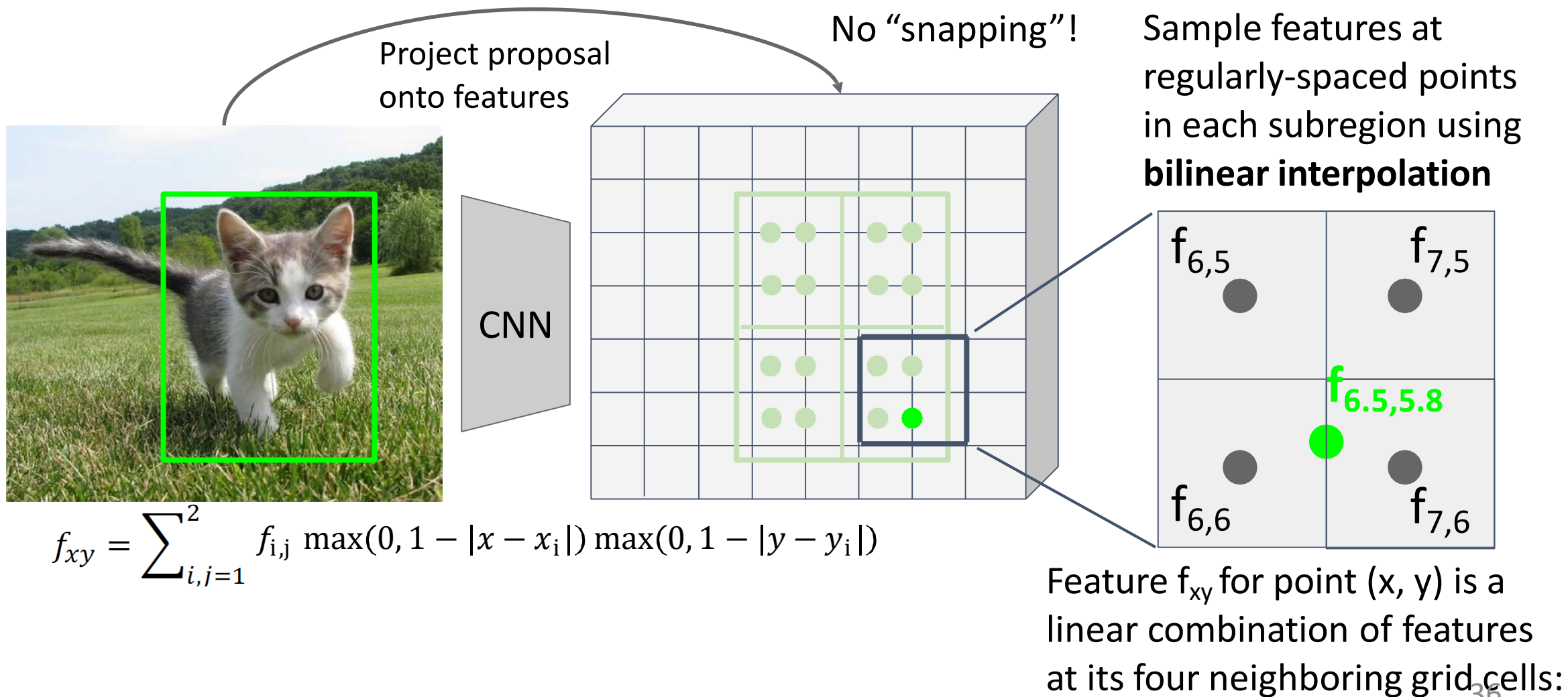


Image features  
(e.g. 512 x 20 x 15)

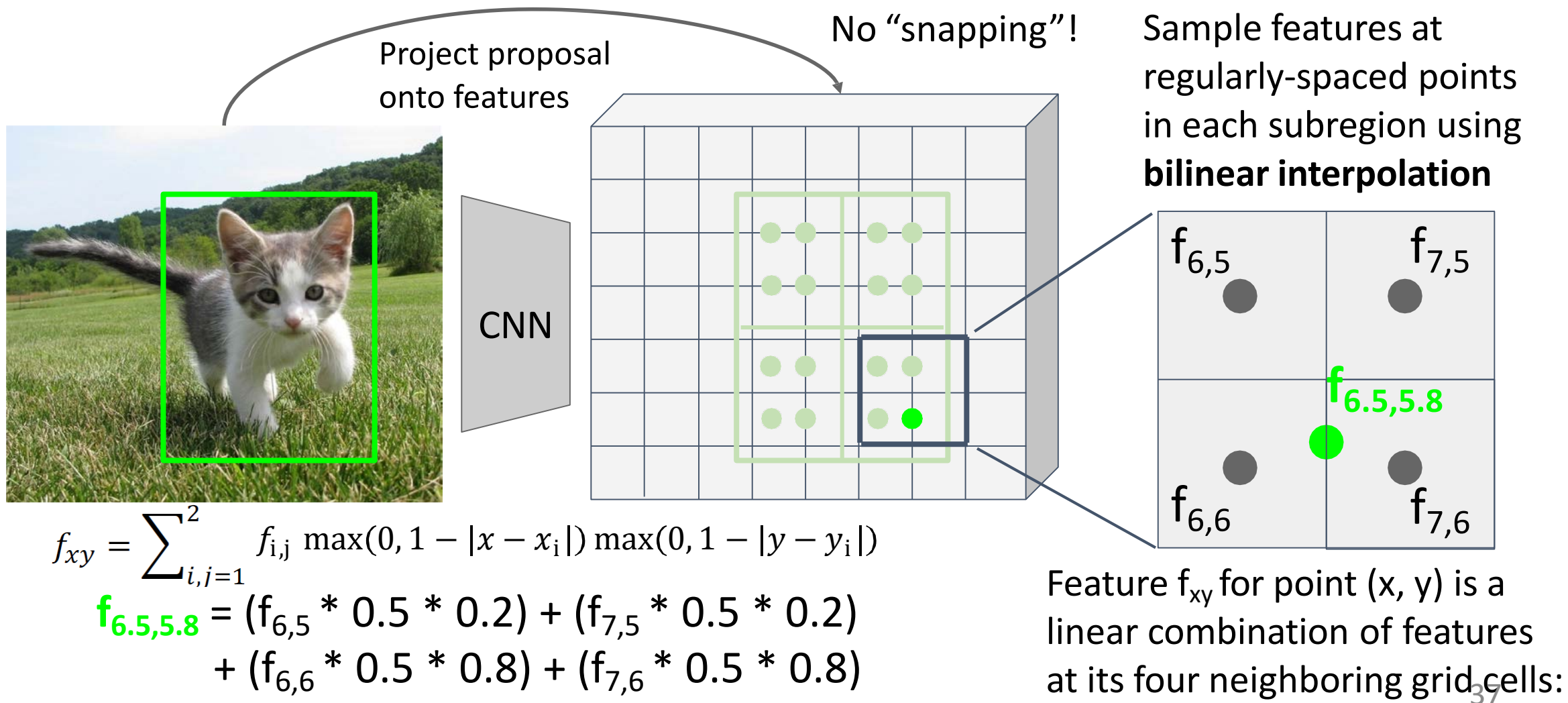
# Cropping Features: RoI Align

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



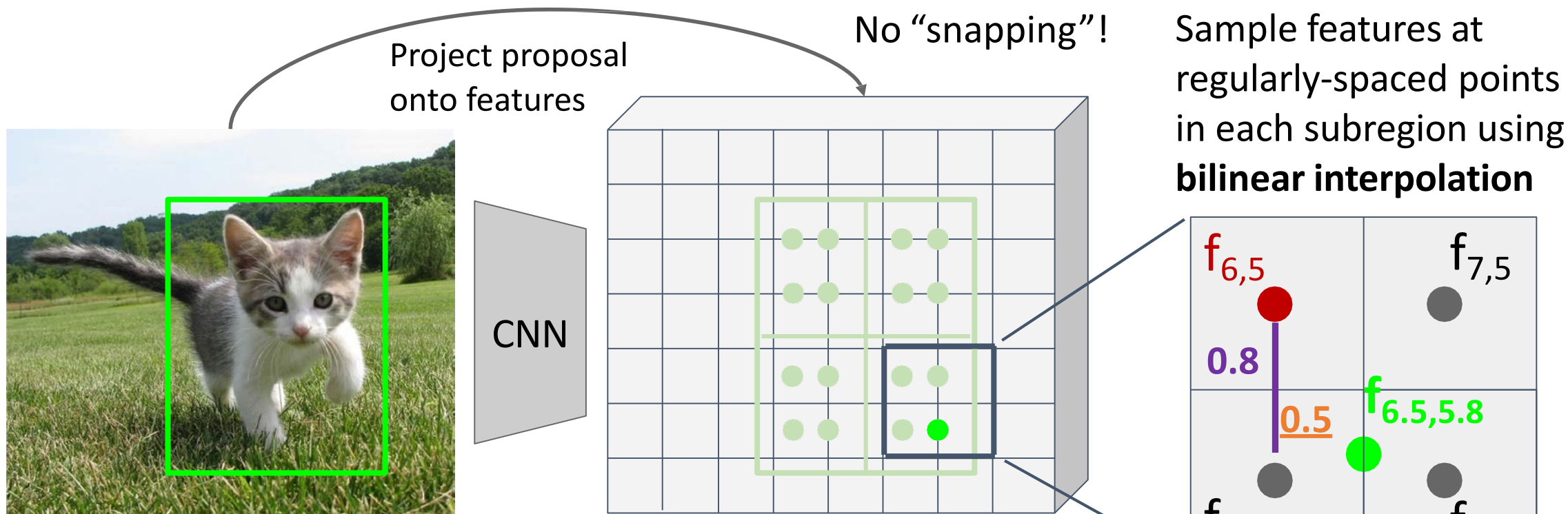
# Cropping Features: RoI Align

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





# Cropping Features: RoI Align

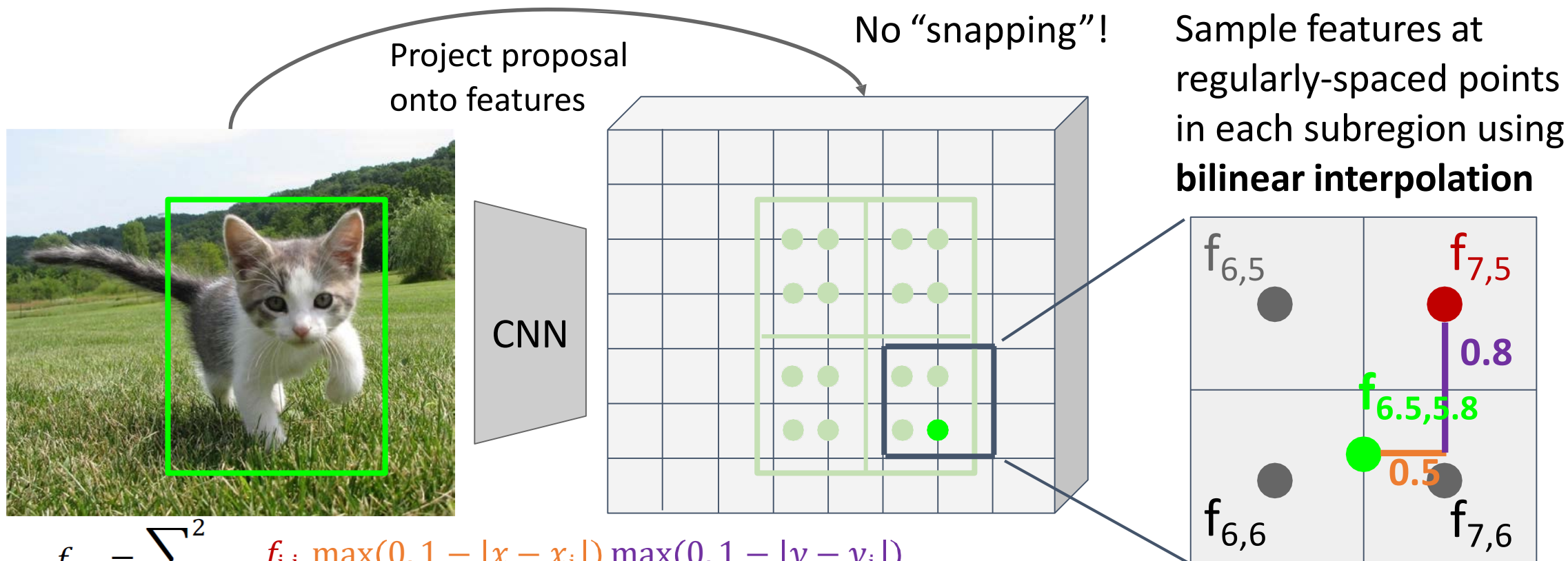


$$f_{xy} = \sum_{i,j=1}^2 f_{i,j} \max(0, 1 - |x - x_i|) \max(0, 1 - |y - y_i|)$$

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

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

# Cropping Features: RoI Align

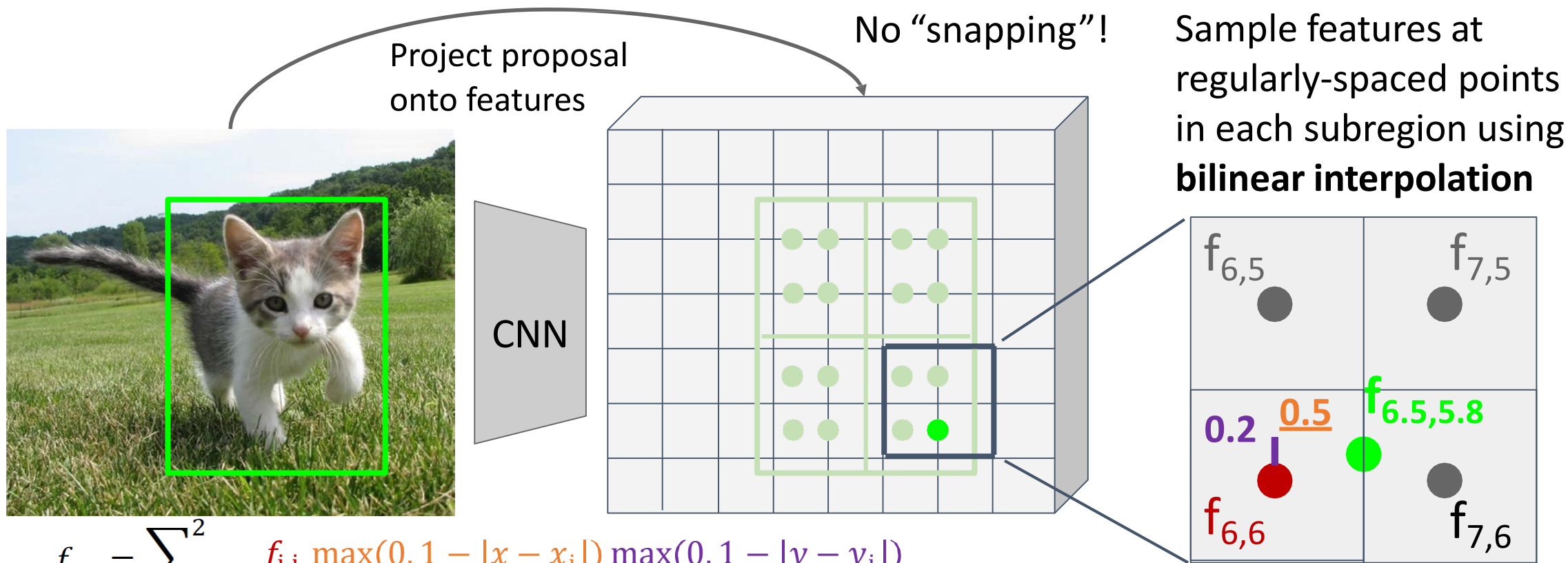


$$f_{xy} = \sum_{i,j=1}^2 f_{i,j} \max(0, 1 - |x - x_i|) \max(0, 1 - |y - y_i|)$$

$$f_{6.5,5.8} = (f_{6,5} * 0.5 * 0.2) + (f_{7,5} * 0.5 * 0.2) + (f_{6,6} * 0.5 * 0.8) + (f_{7,6} * 0.5 * 0.8)$$

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

# Cropping Features: RoI Align



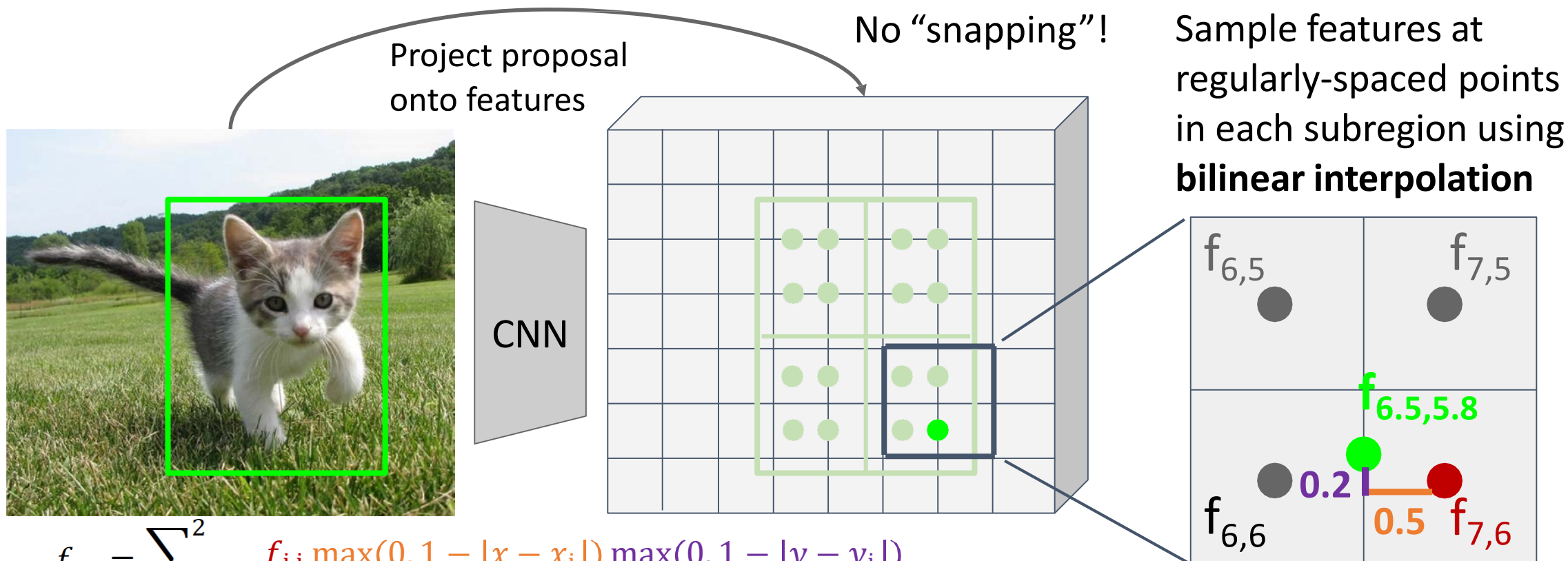
$$f_{xy} = \sum_{i,j=1}^2 f_{i,j} \max(0, 1 - |x - x_i|) \max(0, 1 - |y - y_i|)$$

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

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



# Cropping Features: RoI Align

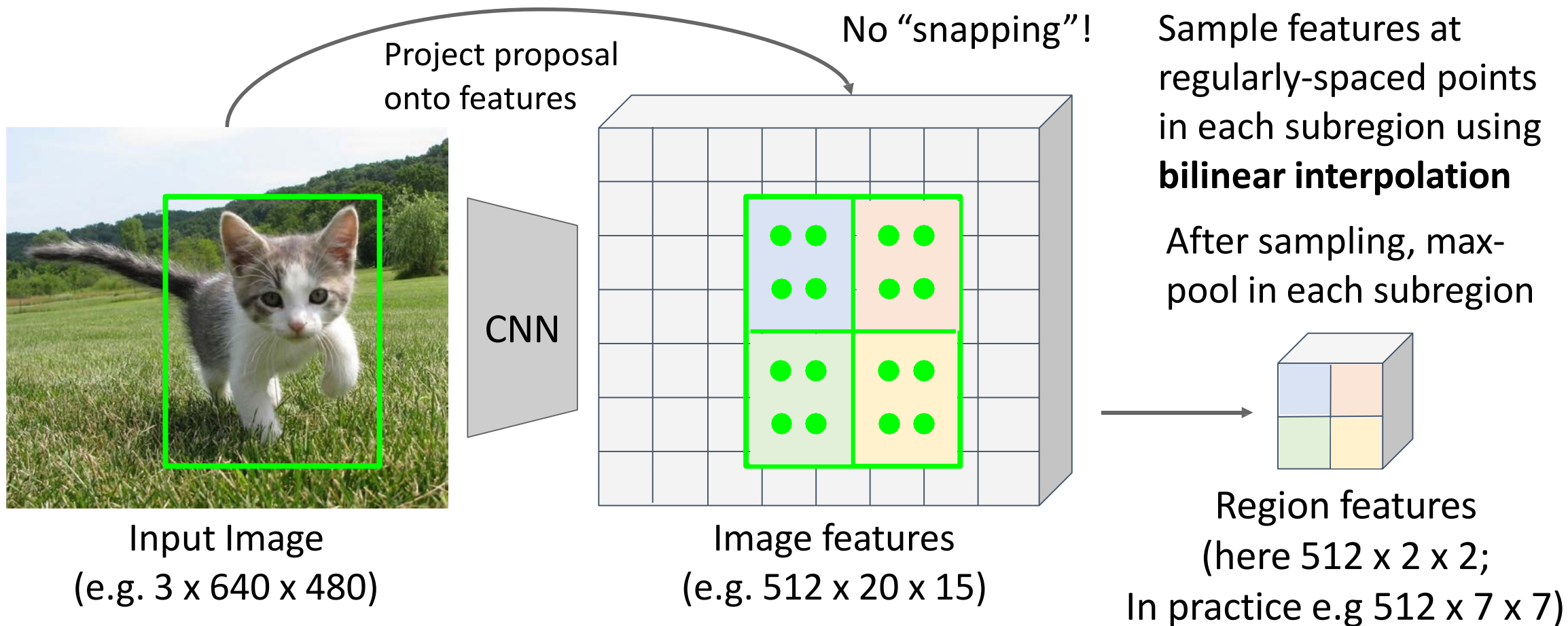


$$f_{xy} = \sum_{i,j=1}^2 f_{i,j} \max(0, 1 - |x - x_i|) \max(0, 1 - |y - y_i|)$$

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

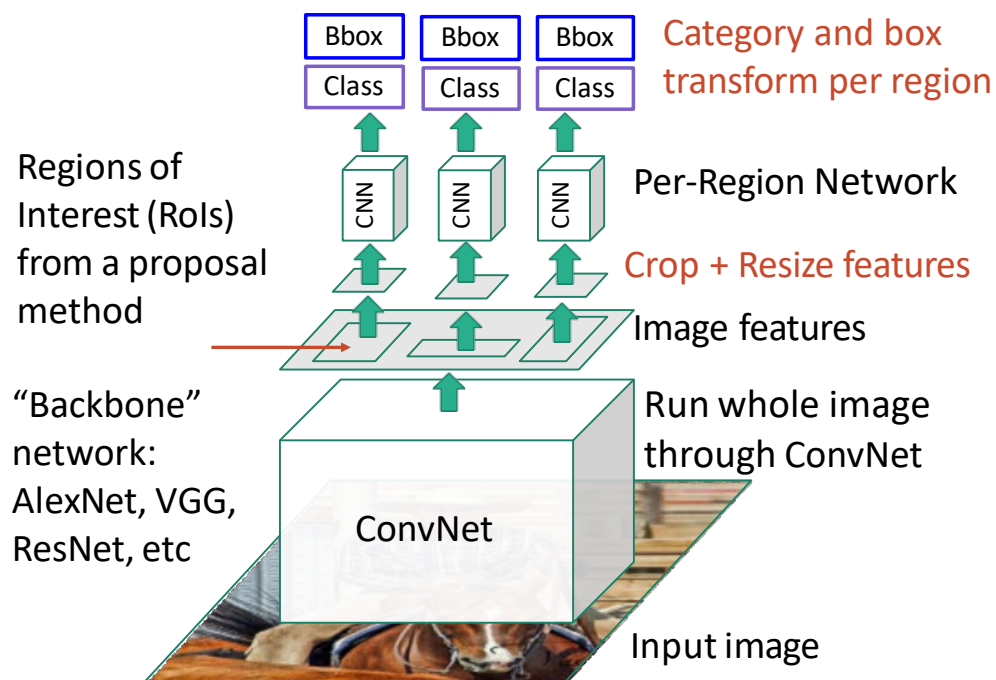
Feature  $f_{xy}$  for point  $(x, y)$  is a linear combination of features at its four neighboring grid cells:

# Cropping Features: RoI Align

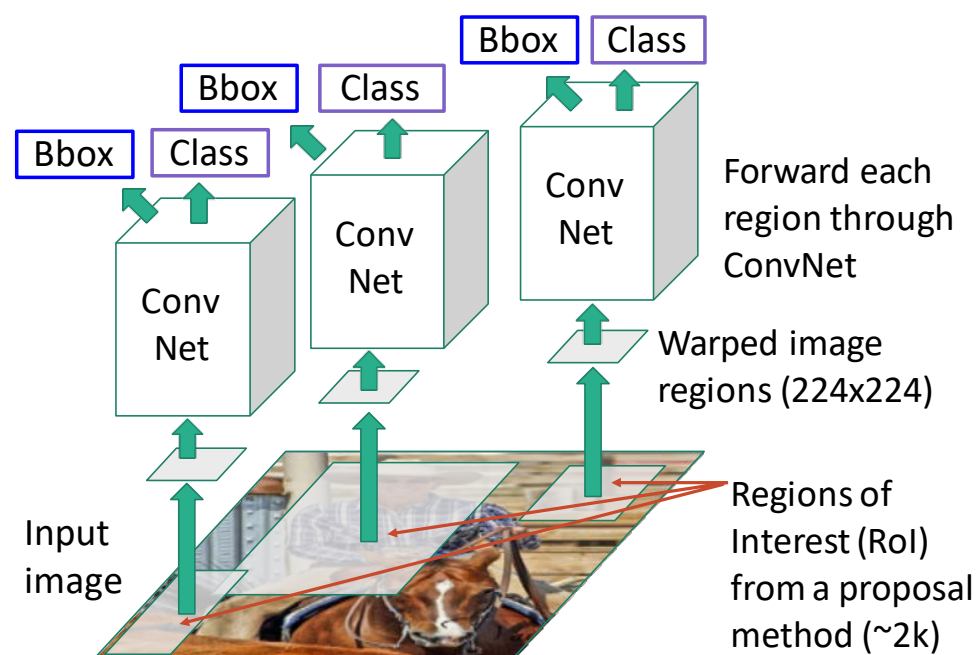


# Fast R-CNN vs “Slow” R-CNN

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

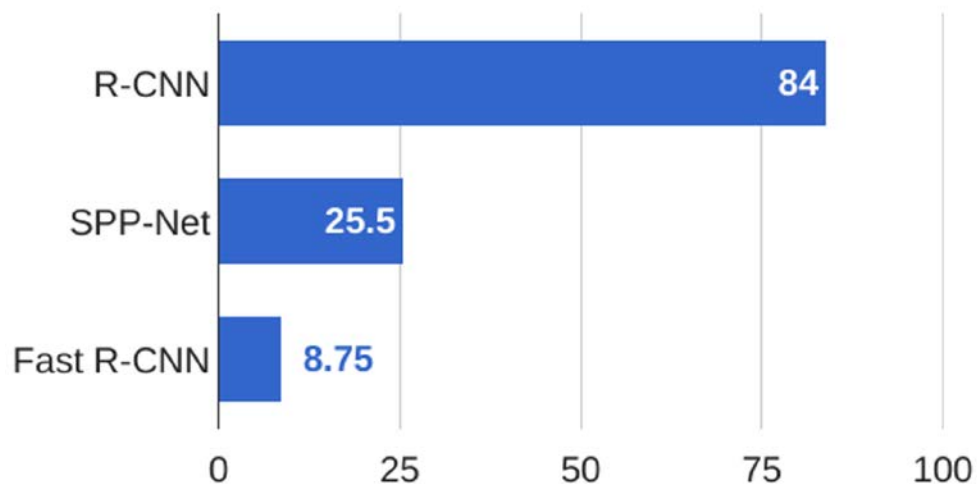


**“Slow” R-CNN:** Apply differentiable cropping to shared image features

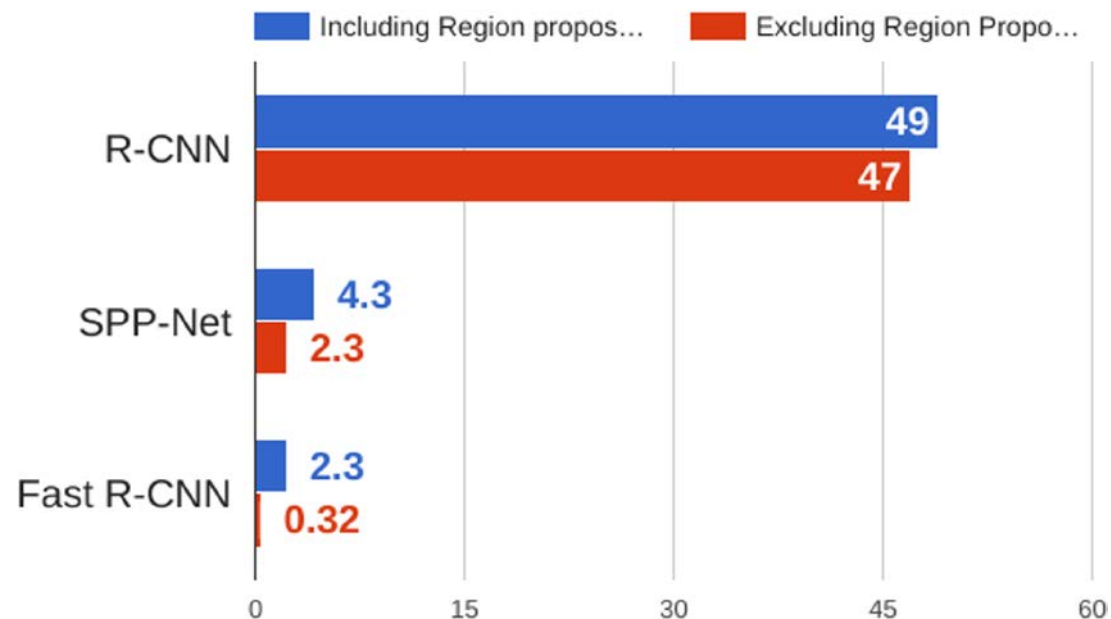


# Fast R-CNN vs “Slow” R-CNN

## Training time (Hours)



## Test time (seconds)



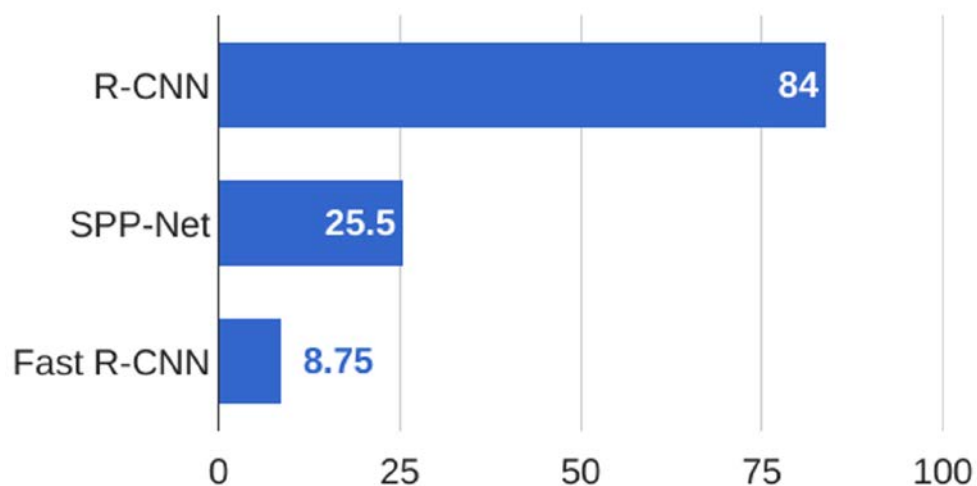
Girshick et al, “Rich feature hierarchies for accurate object detection and semantic segmentation”, CVPR 2014.

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

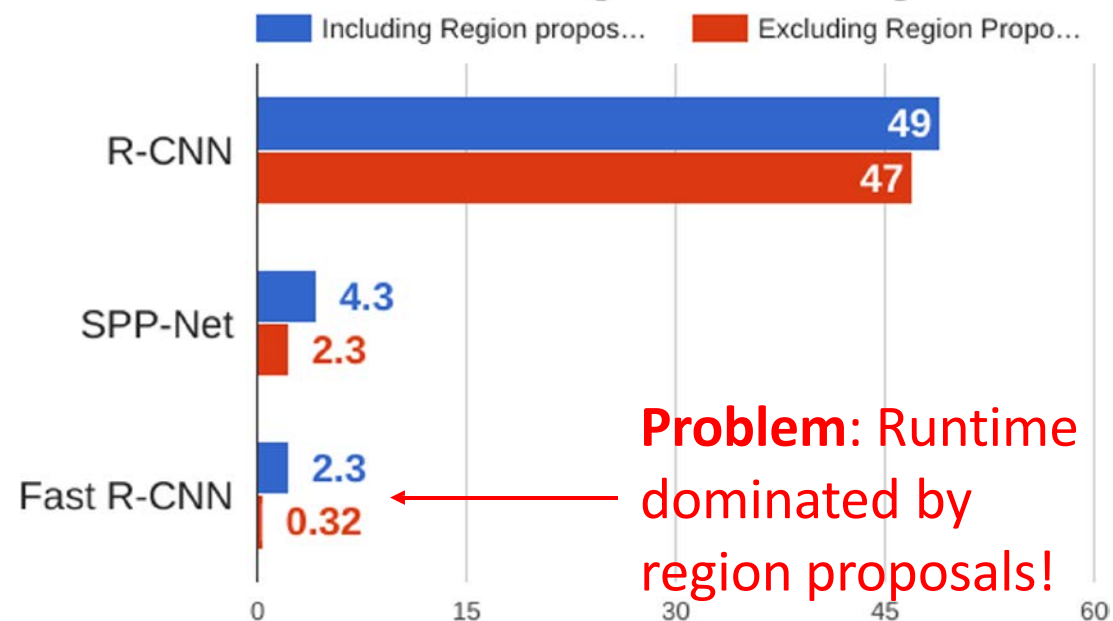
Girshick, “Fast R-CNN”, ICCV 2015

# Fast R-CNN vs “Slow” R-CNN

## Training time (Hours)



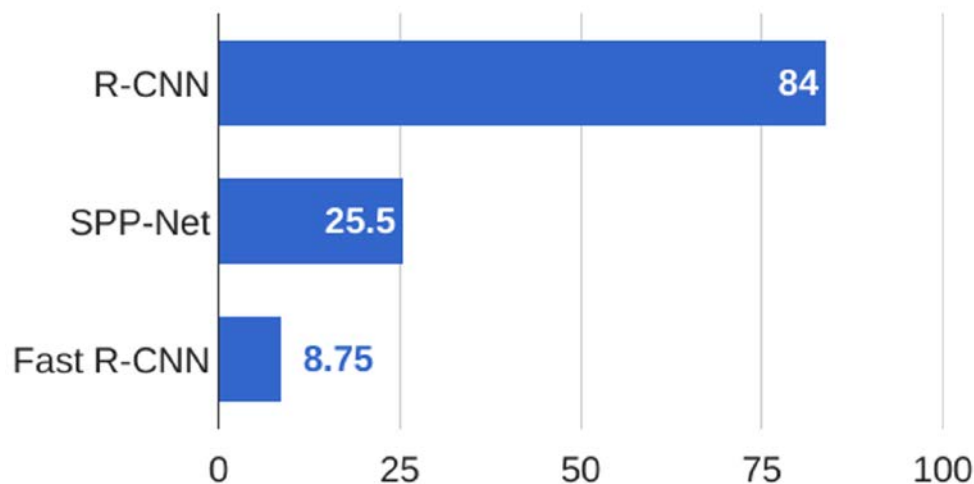
## Test time (seconds)



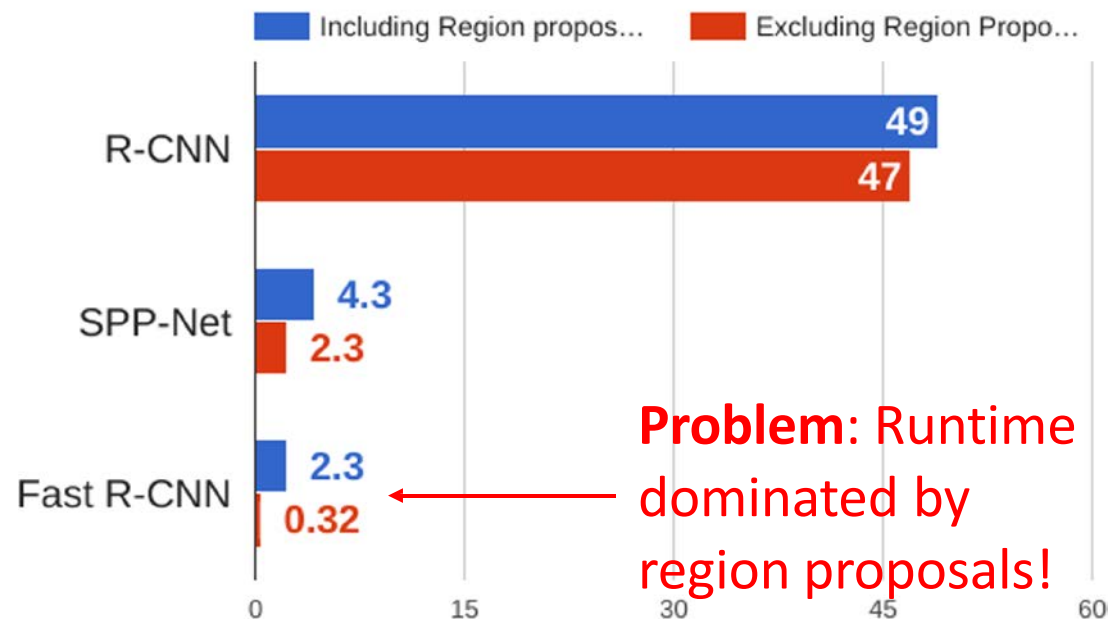
**Problem: Runtime dominated by region proposals!**

# Fast R-CNN vs “Slow” R-CNN

## Training time (Hours)



## Test time (seconds)



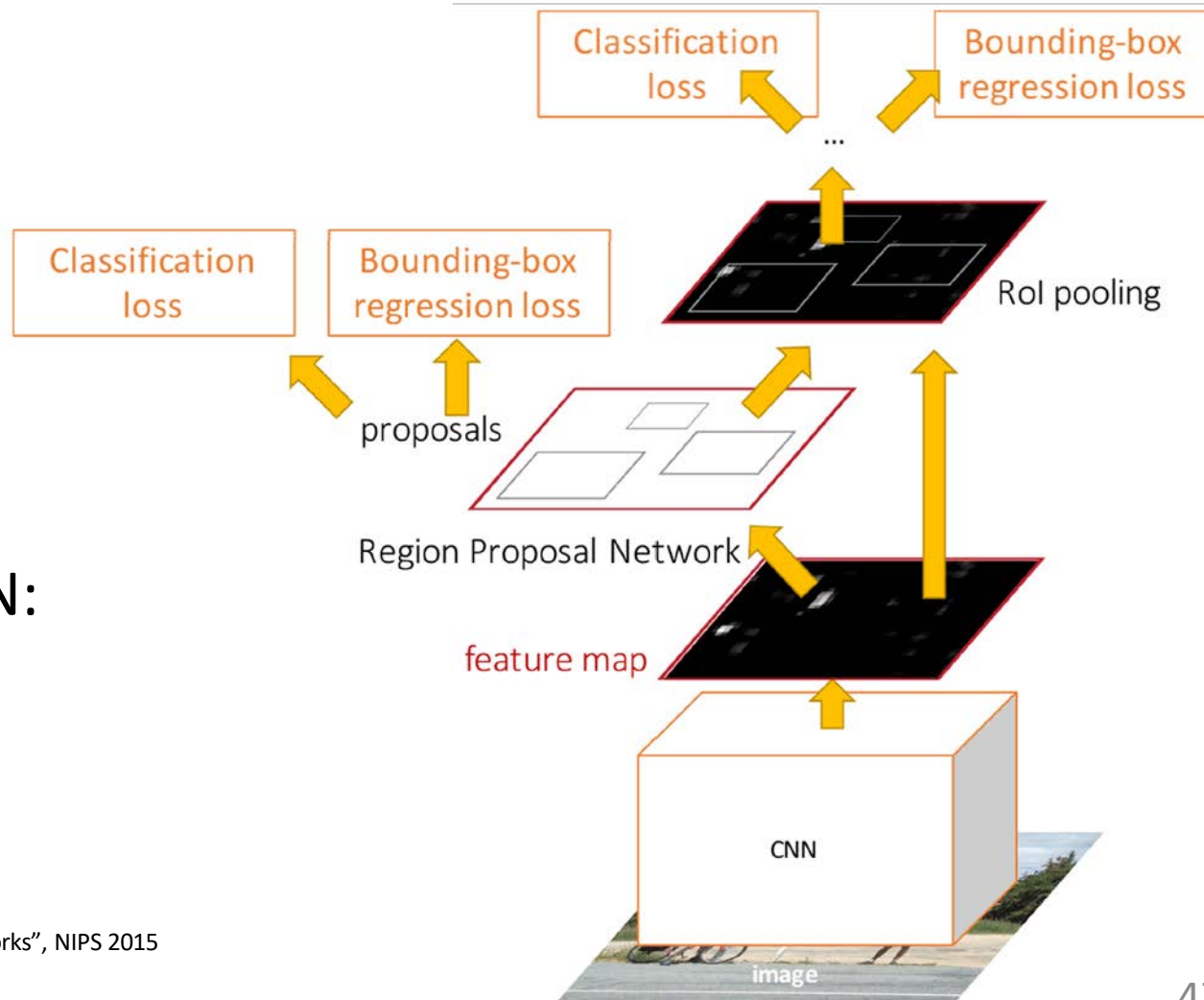
**Problem:** Runtime dominated by region proposals!

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

# 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





# Region Proposal Network (RPN)

Run backbone CNN to get  
features aligned to input image



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

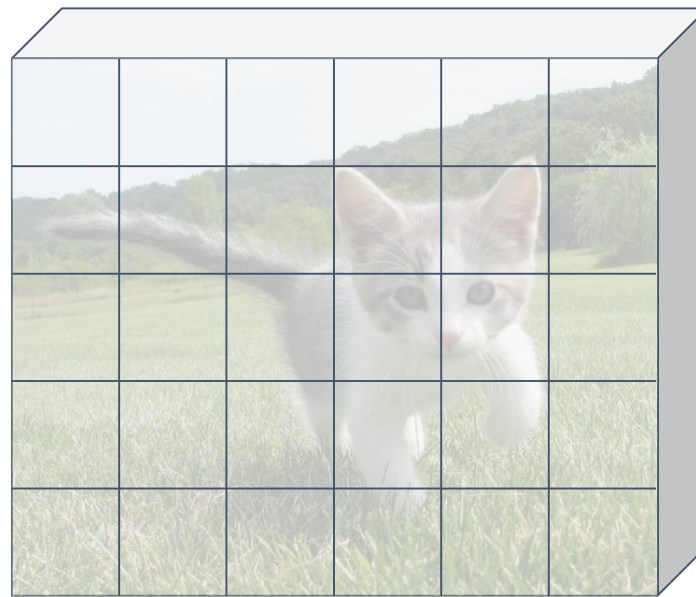
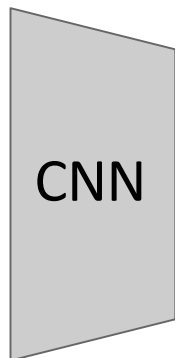


Image features  
(e.g. 512 x 5 x 6)



# Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



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

Each feature corresponds to a point in the input

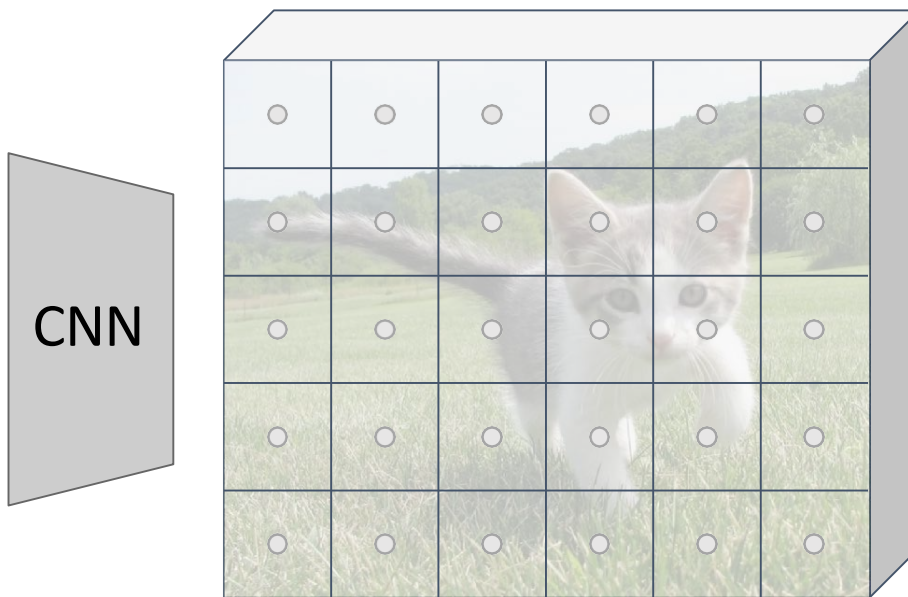
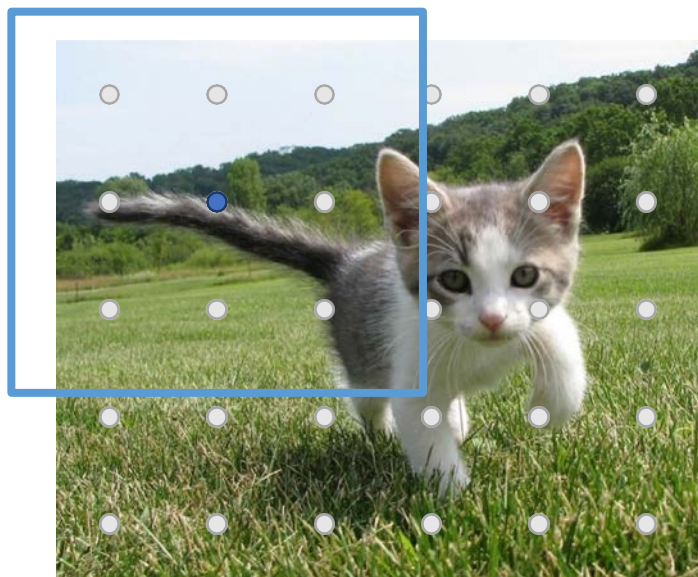


Image features  
(e.g. 512 x 5 x 6)

# Region Proposal Network (RPN)

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

Run backbone CNN to get features aligned to input image



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

CNN

Each feature corresponds to a point in the input

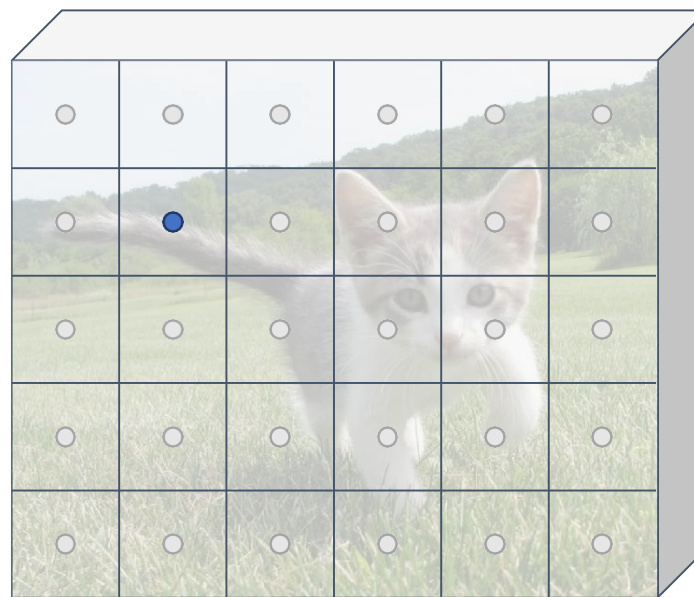
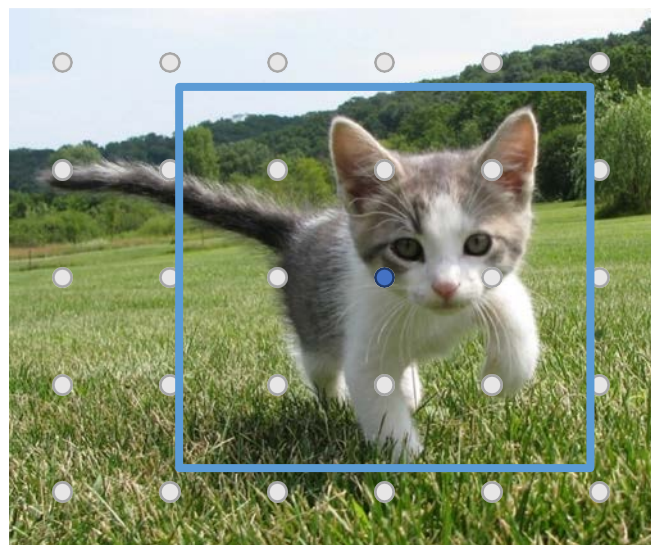


Image features  
(e.g. 512 x 5 x 6)

# Region Proposal Network (RPN)

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

Run backbone CNN to get features aligned to input image



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



Each feature corresponds to a point in the input

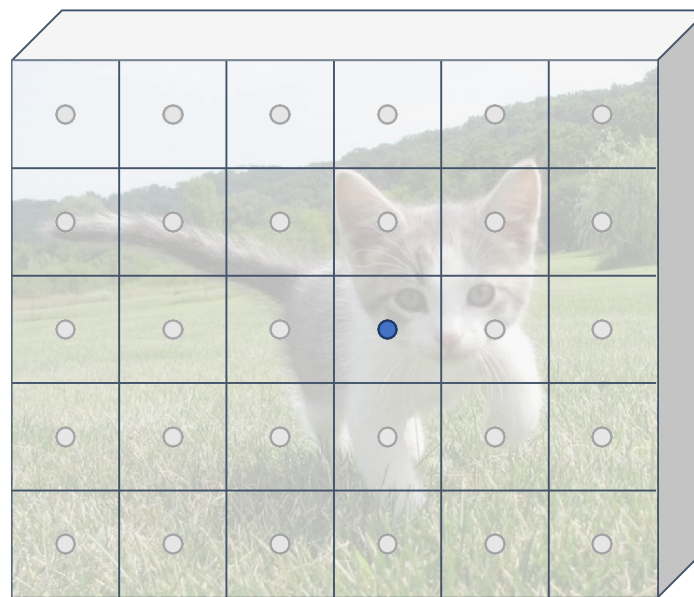
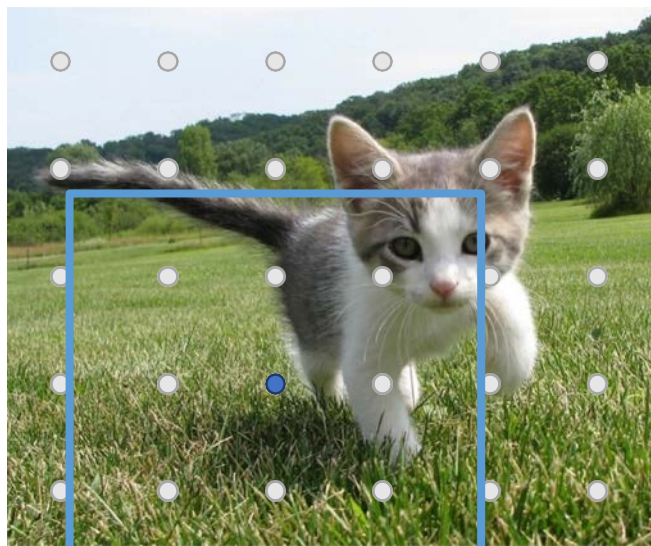


Image features  
(e.g. 512 x 5 x 6)



# Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



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



Each feature corresponds to a point in the input

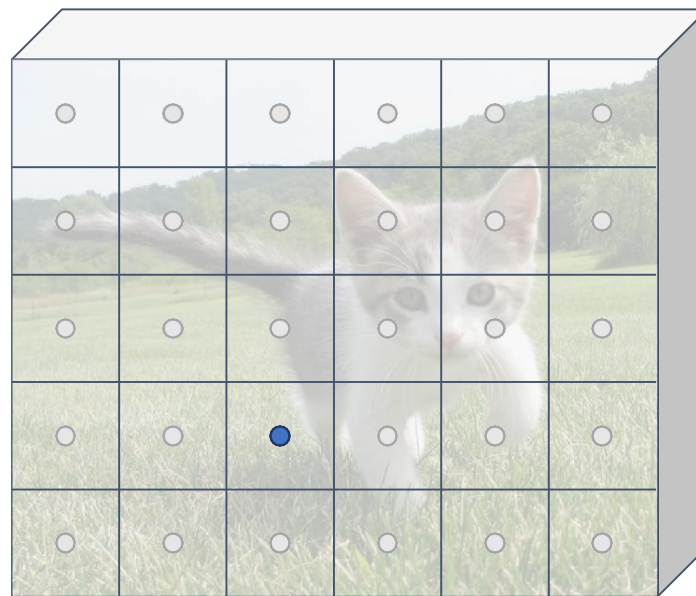
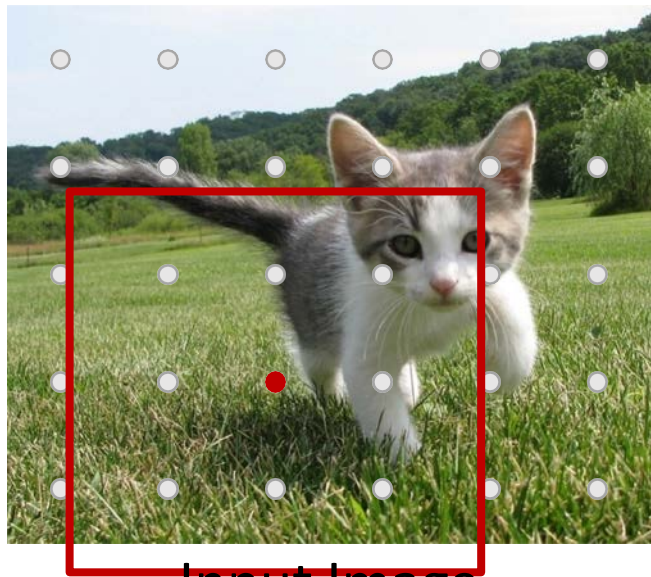


Image features  
(e.g. 512 x 5 x 6)

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 x 640 x 480)



Each feature corresponds to a point in the input

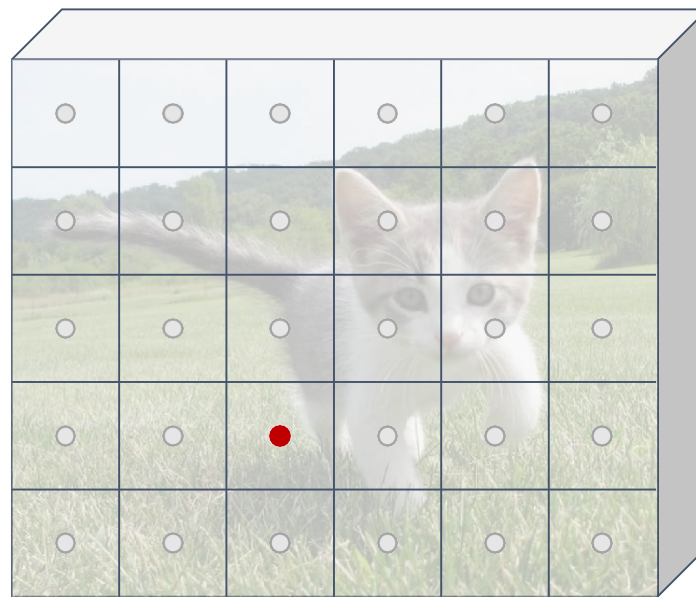


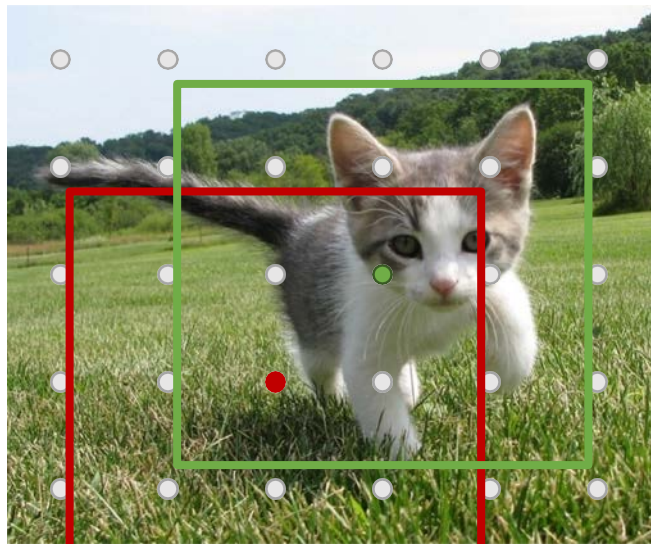
Image features  
(e.g. 512 x 5 x 6)

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

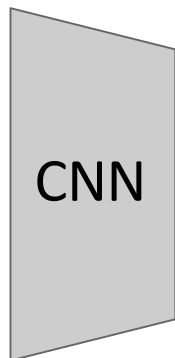
Classify each anchor as  
**positive (object)** or  
**negative (no object)**

# Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



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



Each feature corresponds to a point in the input

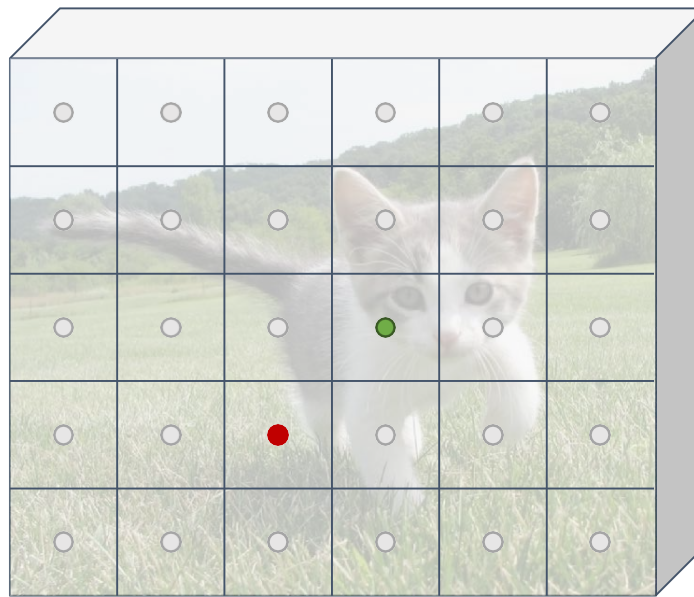


Image features  
(e.g. 512 x 5 x 6)

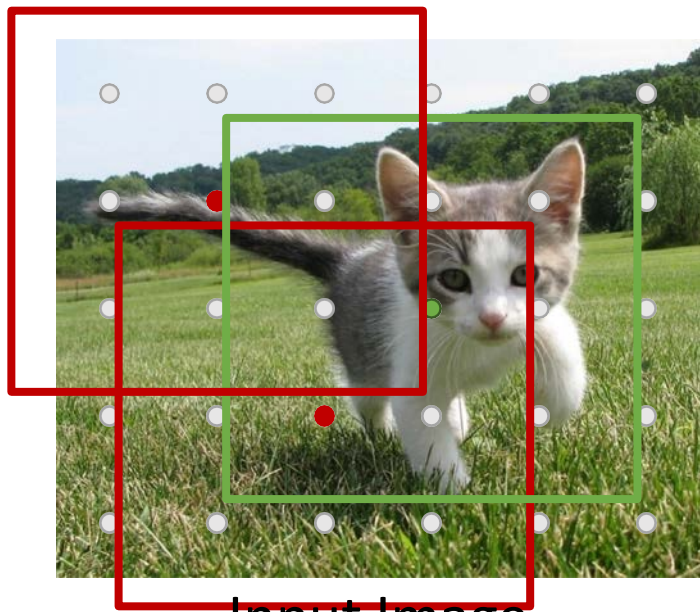
Imagine an **anchor box** of fixed size at each point in the feature map

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 x 640 x 480)



Each feature corresponds to a point in the input

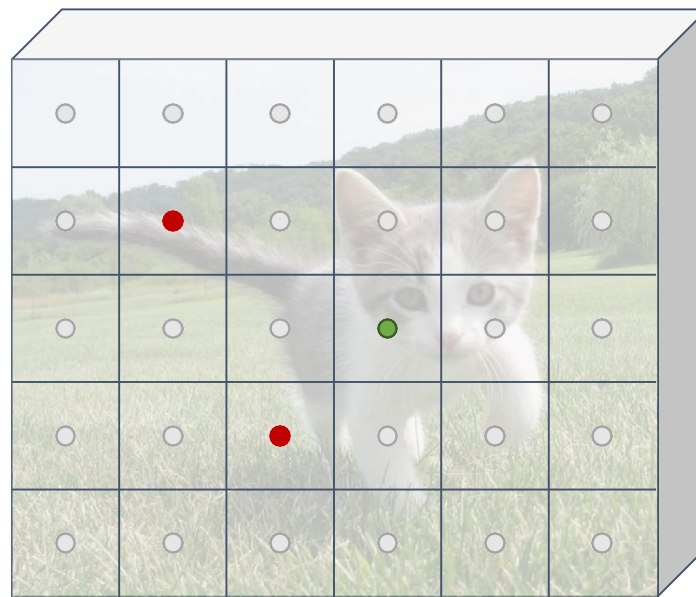


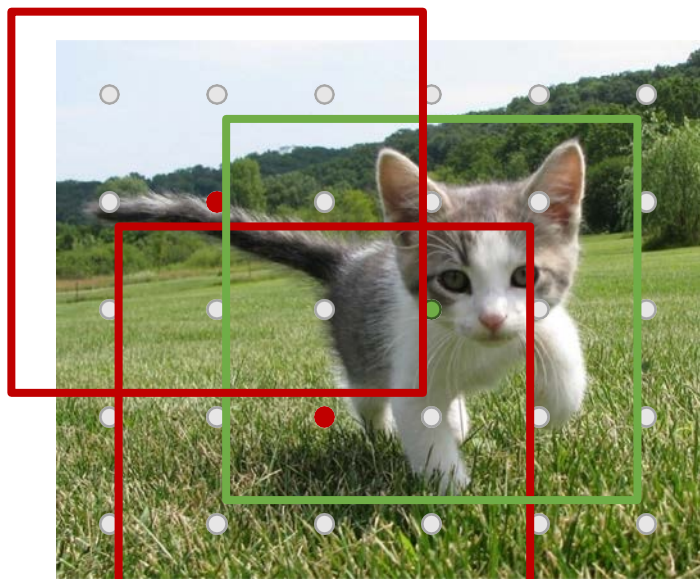
Image features  
(e.g. 512 x 5 x 6)

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 x 640 x 480)



Each feature corresponds to a point in the input

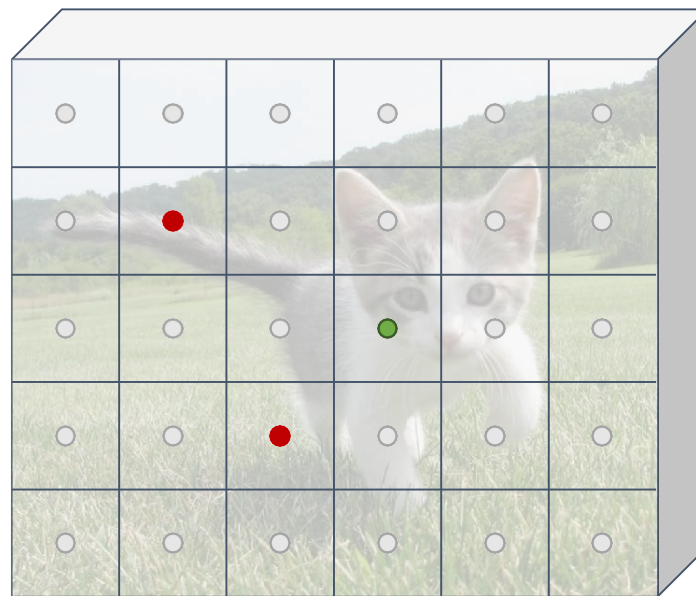


Image features  
(e.g. 512 x 5 x 6)

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



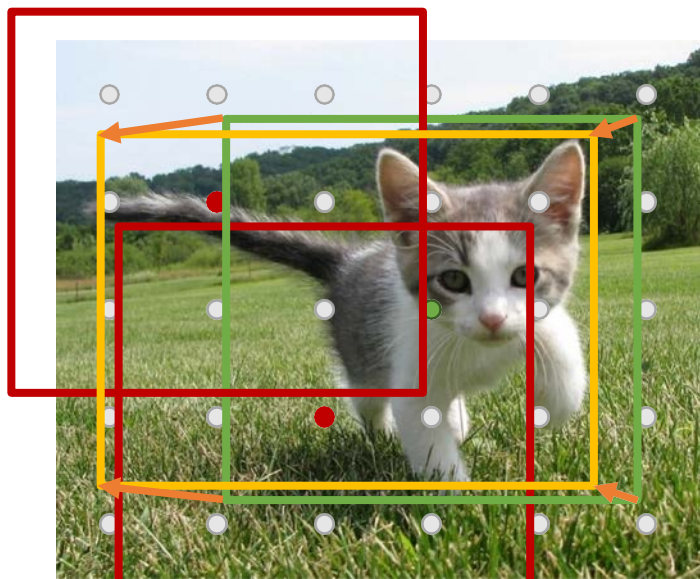
Anchor is  
object?  
2 x 5 x 6

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 x 640 x 480)



Each feature corresponds to a point in the input

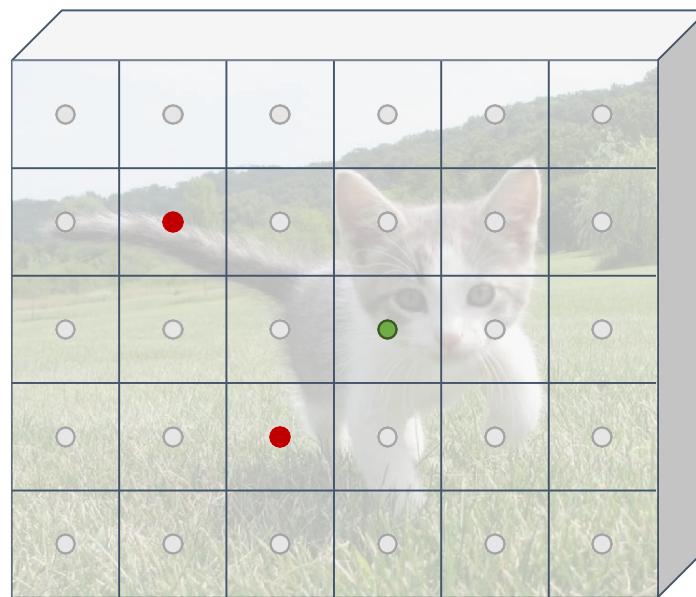


Image features  
(e.g. 512 x 5 x 6)

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

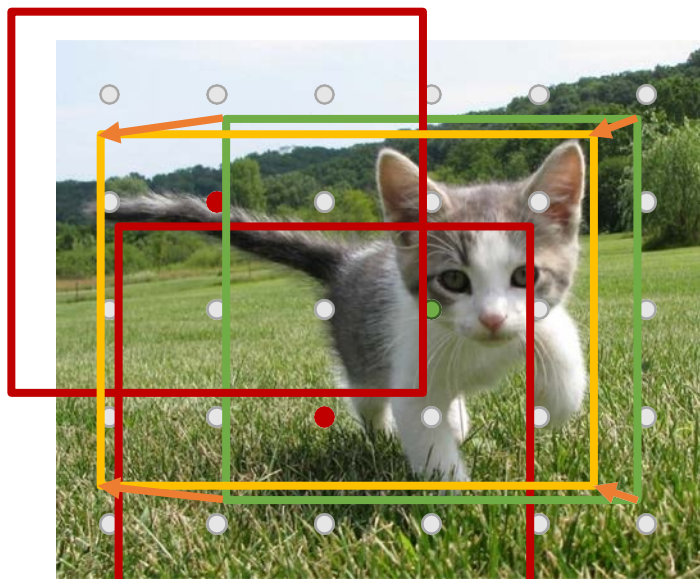


Anchor is  
object?  
2 x 5 x 6

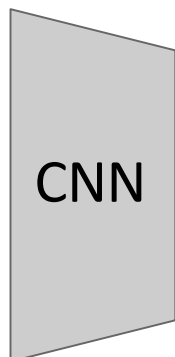
Classify each anchor as  
**positive (object)** or  
**negative (no object)**

## Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



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



Each feature corresponds to a point in the input

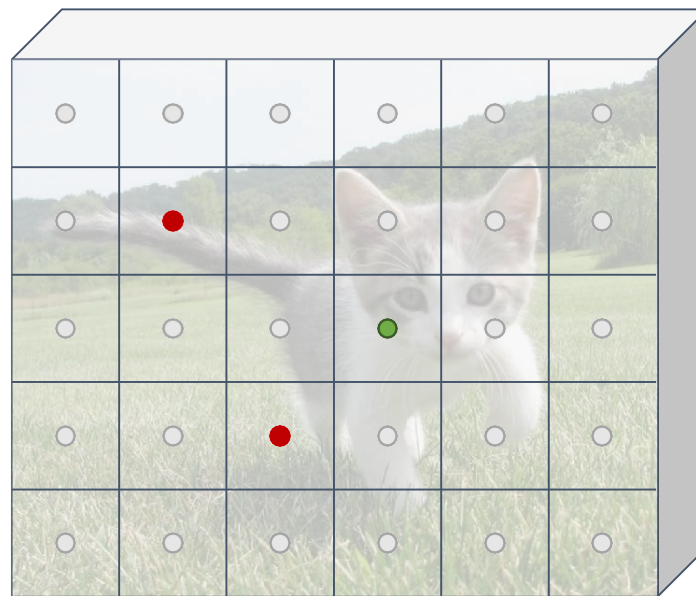


Image features  
(e.g. 512 x 5 x 6)

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



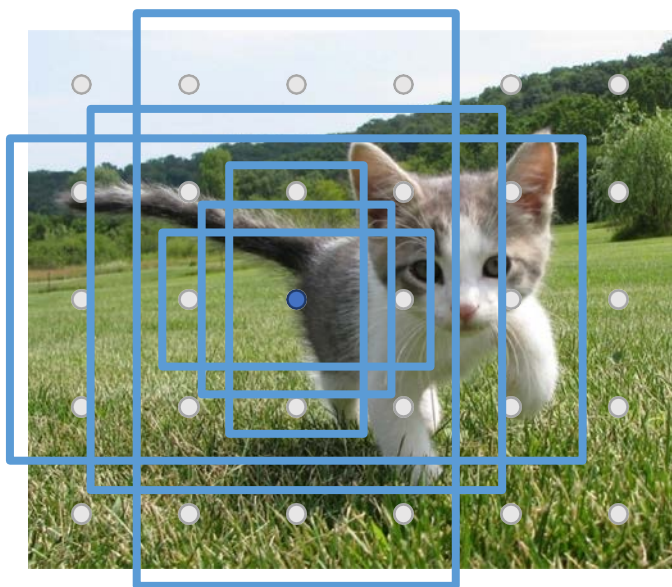
Anchor is  
object?  
2 x 5 x 6

Anchor  
transforms  
4 x 5 x 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 x 640 x 480)



Each feature corresponds to a point in the input

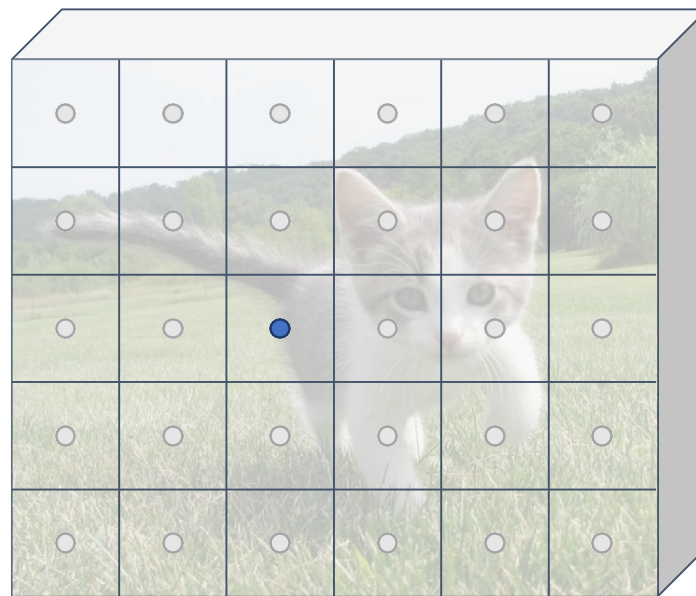


Image features  
(e.g. 512 x 5 x 6)

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



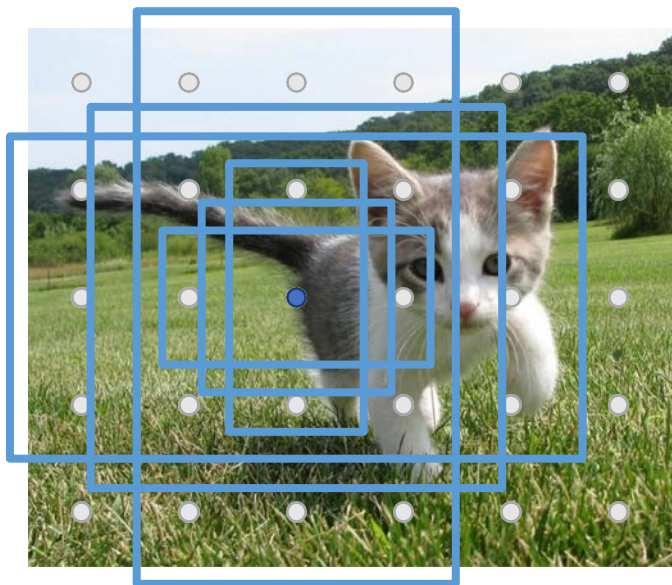
Anchor is  
object?  
 $2K \times 5 \times 6$

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



## Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



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



Each feature corresponds to a point in the input

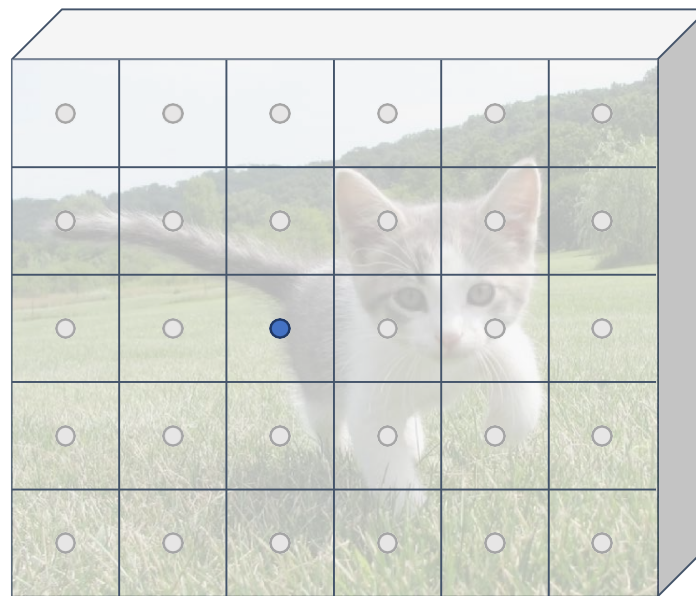


Image features  
(e.g.  $512 \times 5 \times 6$ )

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



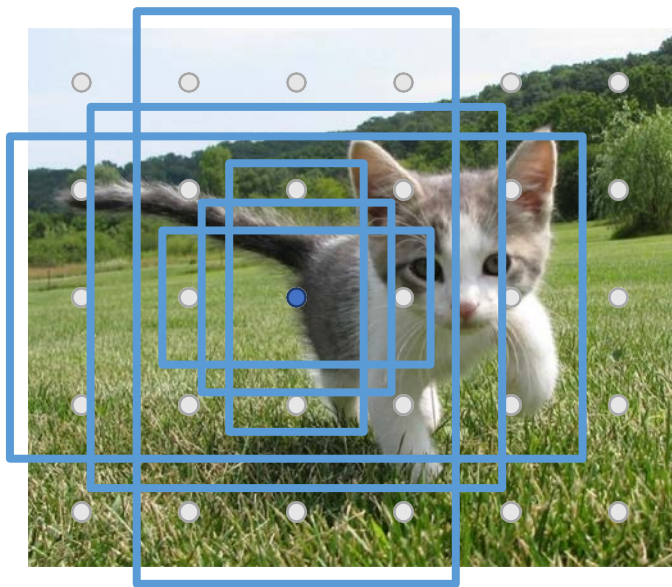
Anchor is object?  
 $2K \times 5 \times 6$

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



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



Each feature corresponds to a point in the input

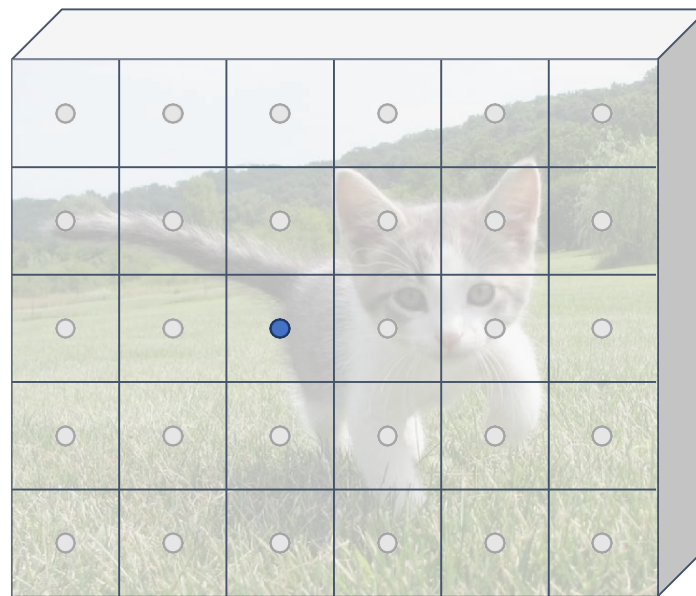


Image features  
(e.g. 512 x 5 x 6)

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



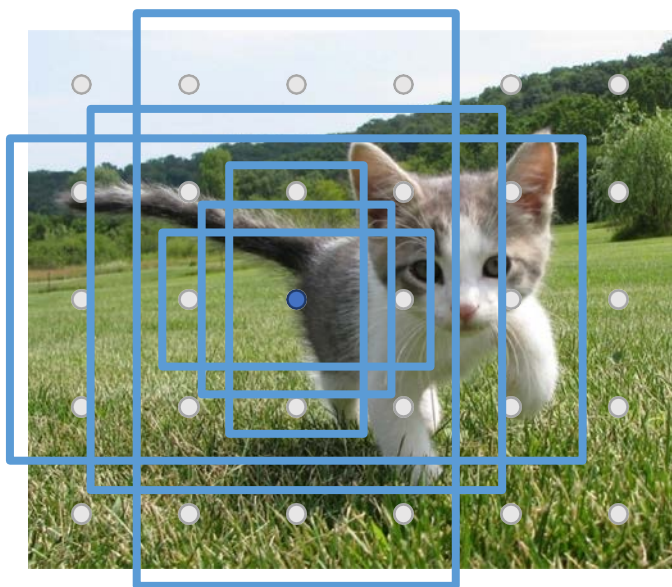
Anchor is  
object?  
 $2K \times 5 \times 6$

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

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



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



Each feature corresponds to a point in the input

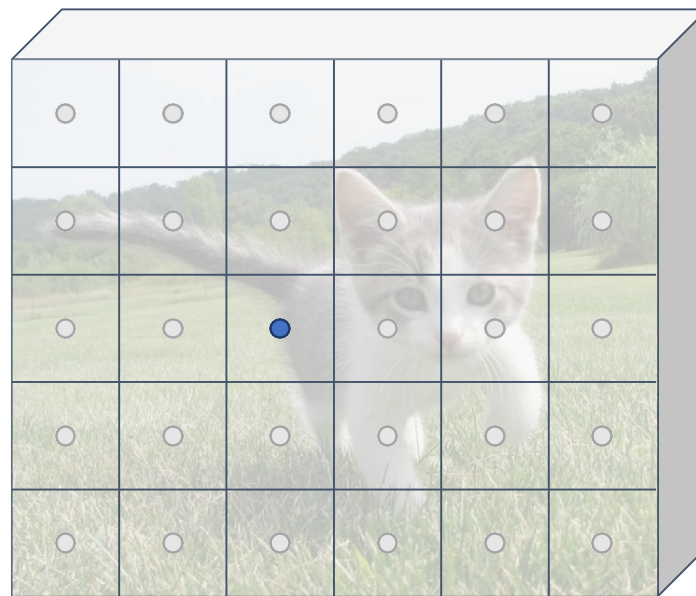


Image features  
(e.g. 512 x 5 x 6)

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?  
 $2K \times 5 \times 6$

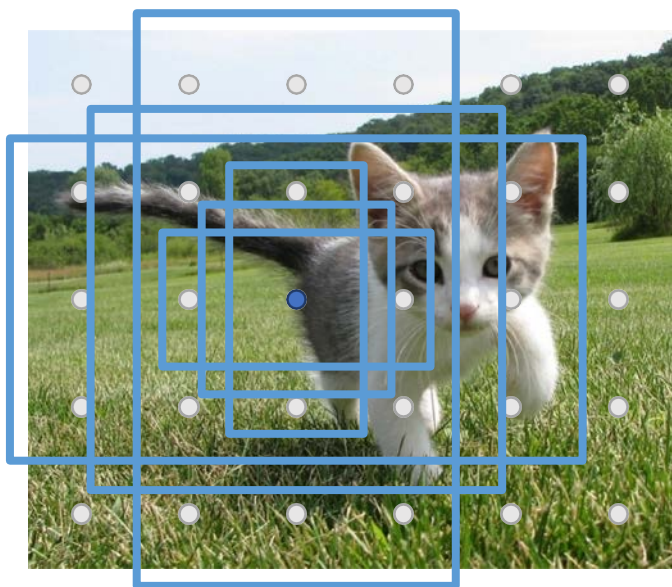
Anchor transforms  
 $4K \times 5 \times 6$

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



## Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



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



Each feature corresponds to a point in the input

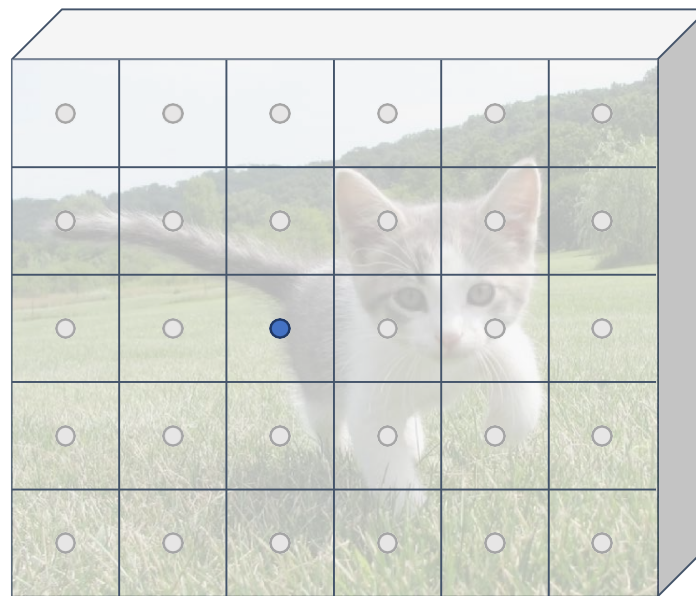


Image features  
(e.g. 512 x 5 x 6)

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



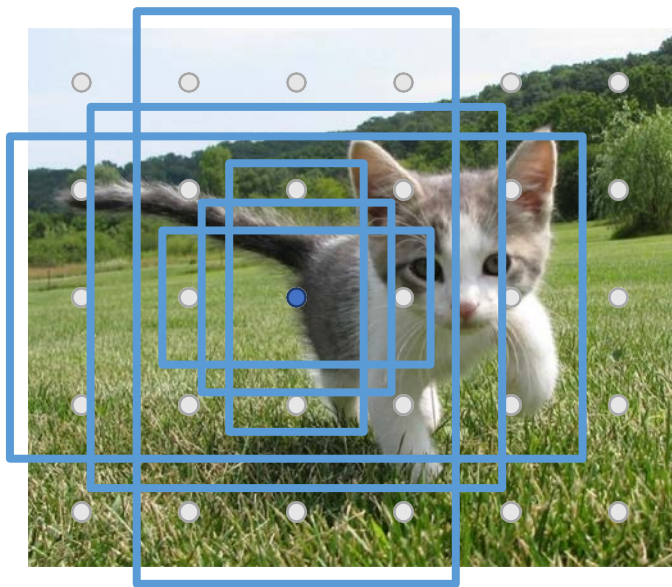
Anchor is  
object?  
 $2K \times 5 \times 6$

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

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



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



Each feature corresponds to a point in the input

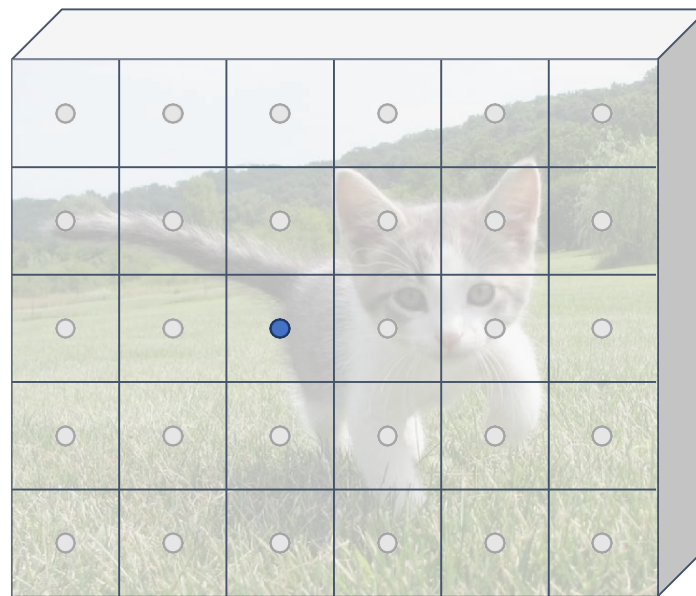


Image features  
(e.g. 512 x 5 x 6)

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  
 $4K \times 5 \times 6$

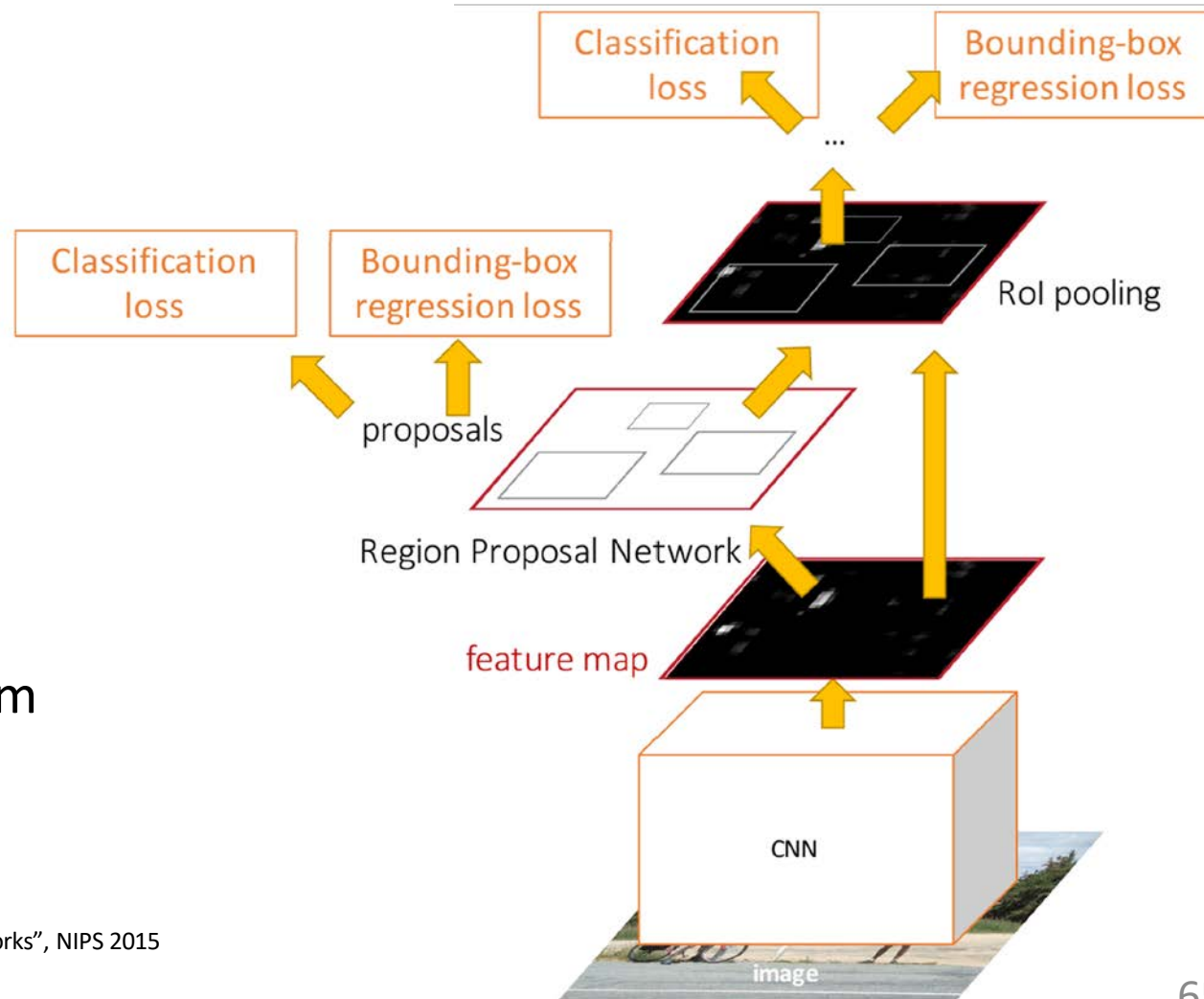
At test-time, sort all  $K \times 5 \times 6$  boxes by their positive score, take top 300 as our region proposals



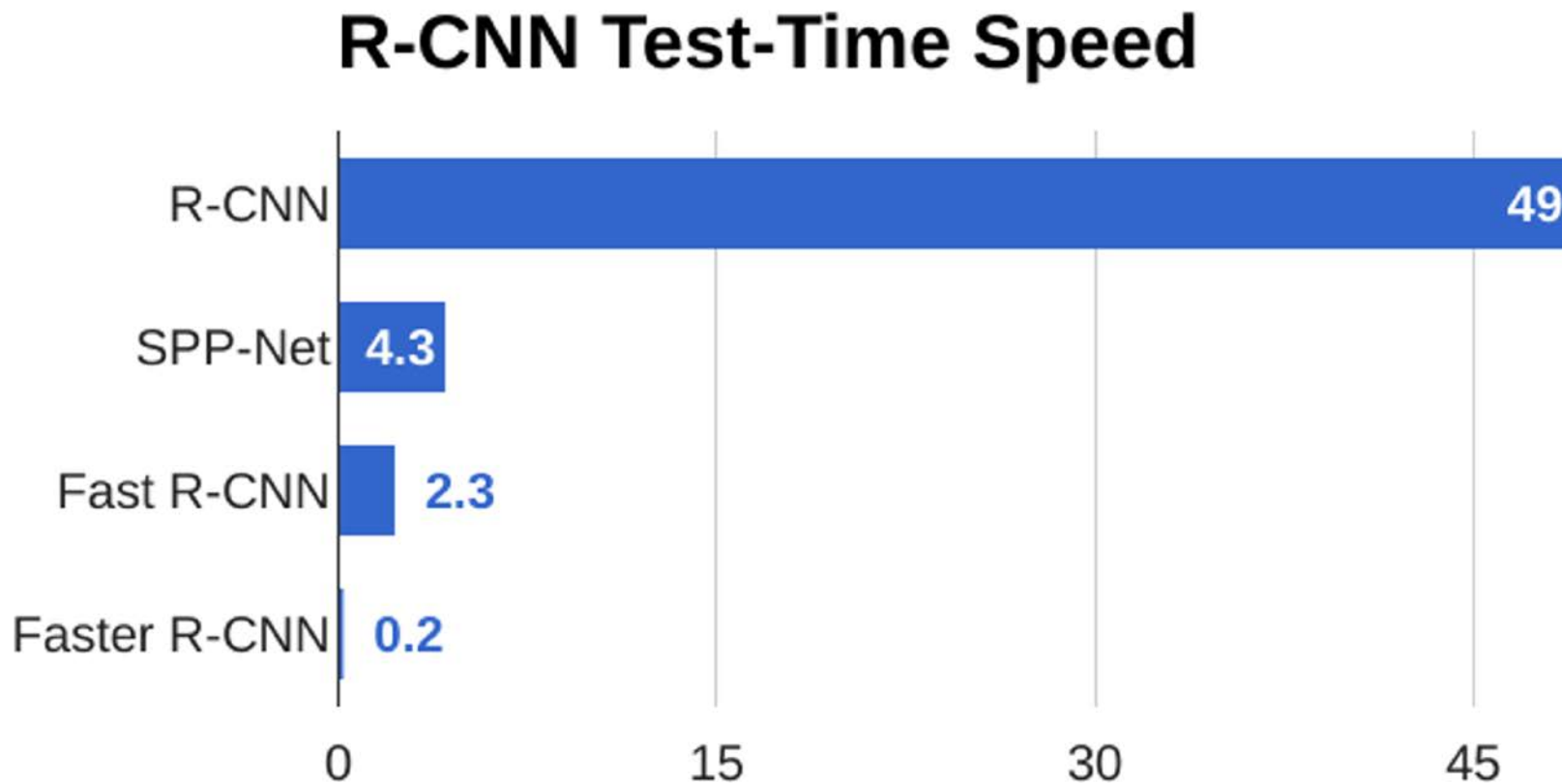
# Faster R-CNN: Learnable Region Proposals

Jointly train with 4 losses:

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



# Faster R-CNN: Learnable Region Proposals



# Faster R-CNN: Learnable Region Proposals

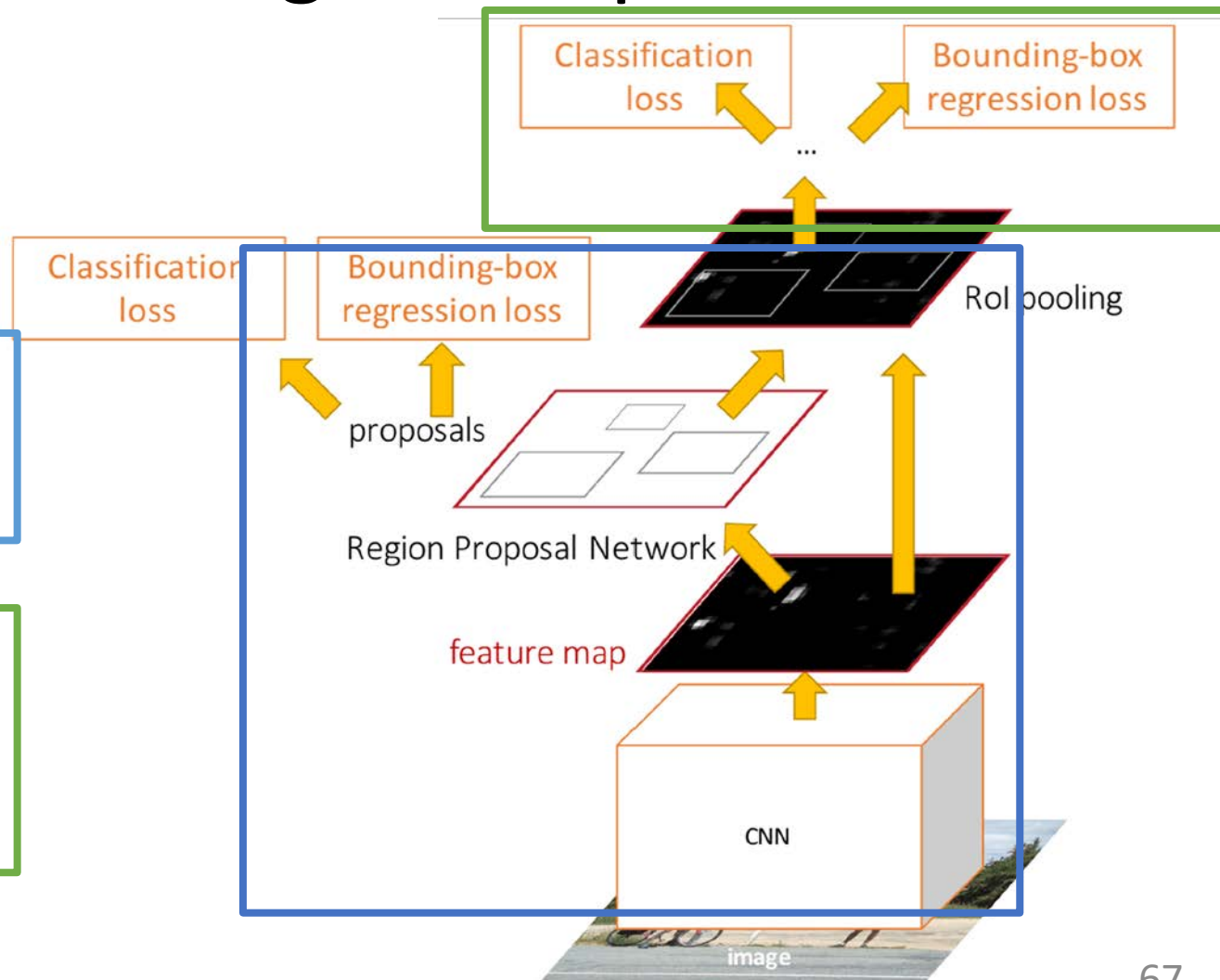
Faster R-CNN is a  
**Two-stage object detector**

First stage: Run once per image

- Backbone network
- Region proposal network

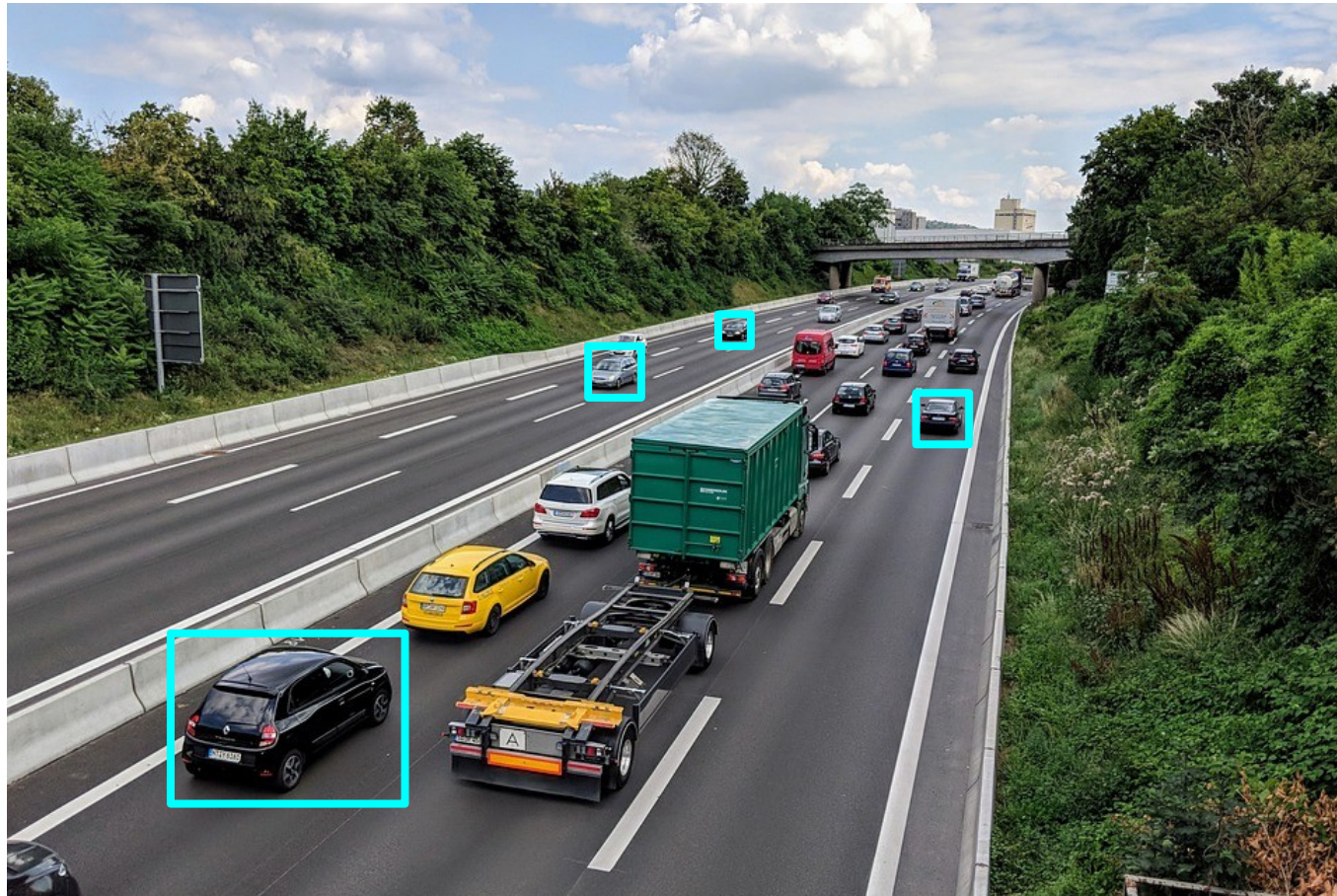
Second stage: Run once per region

- 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**



Object  
Detector



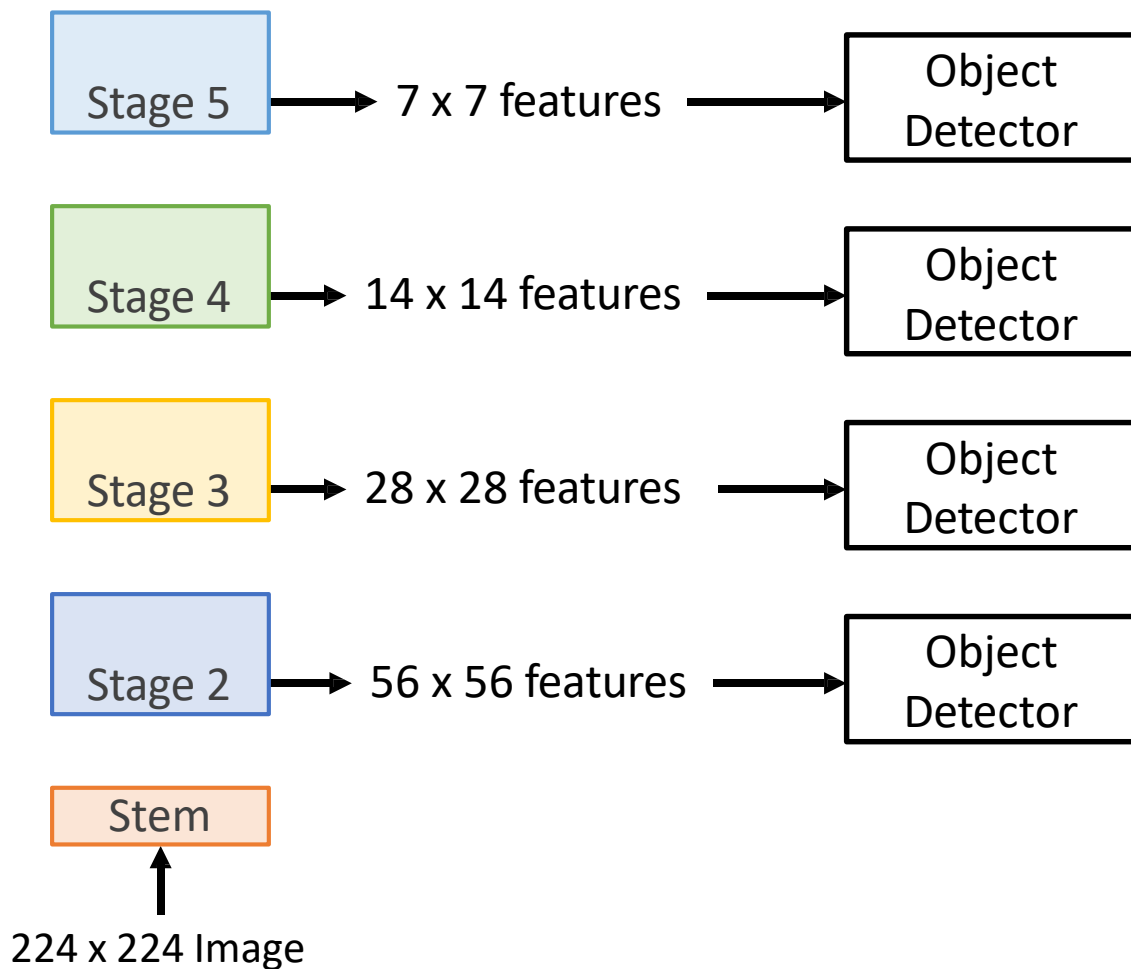
Object  
Detector



Object  
Detector

# Dealing with Scale: Multiscale Features

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

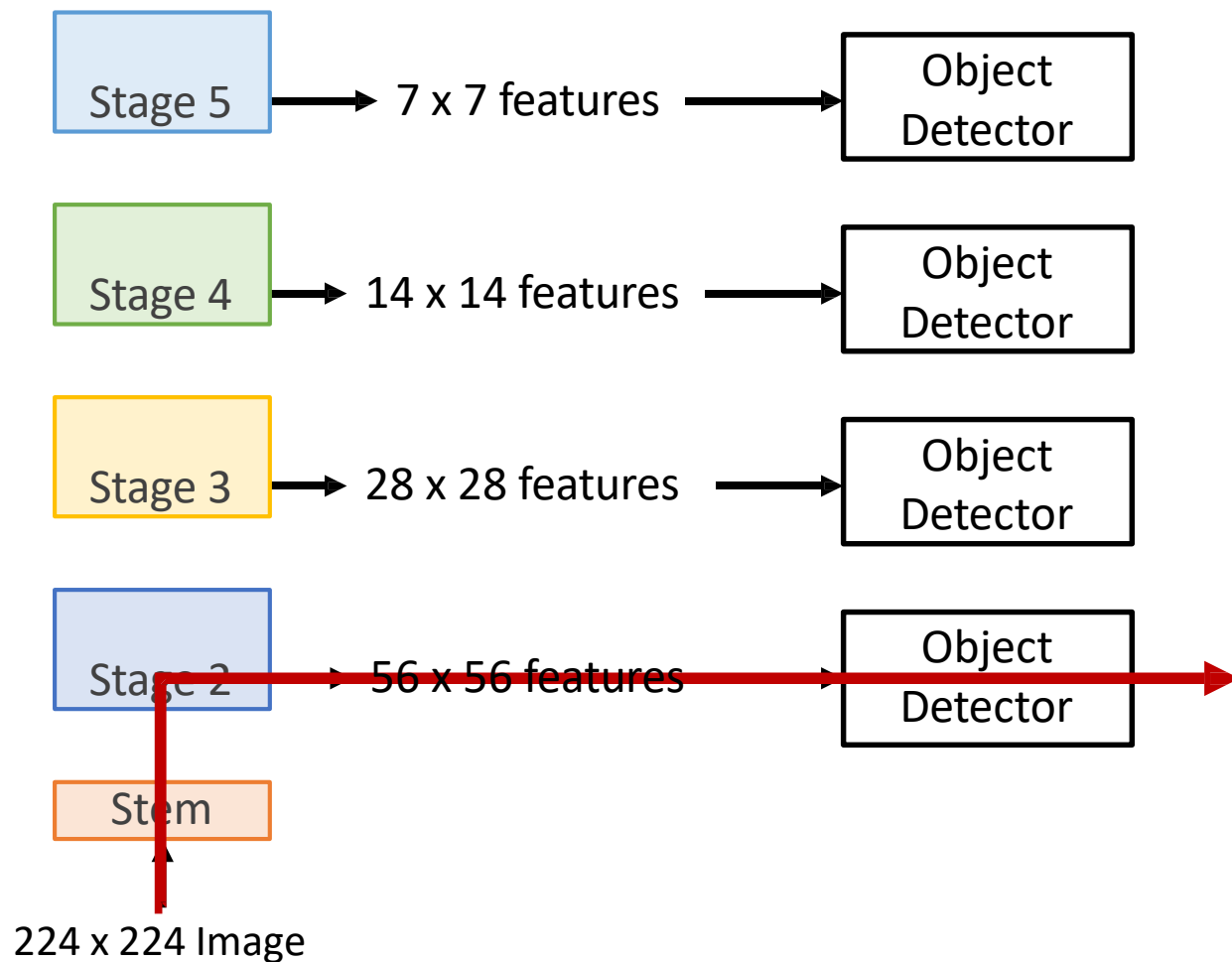




# Dealing with Scale: Multiscale Features

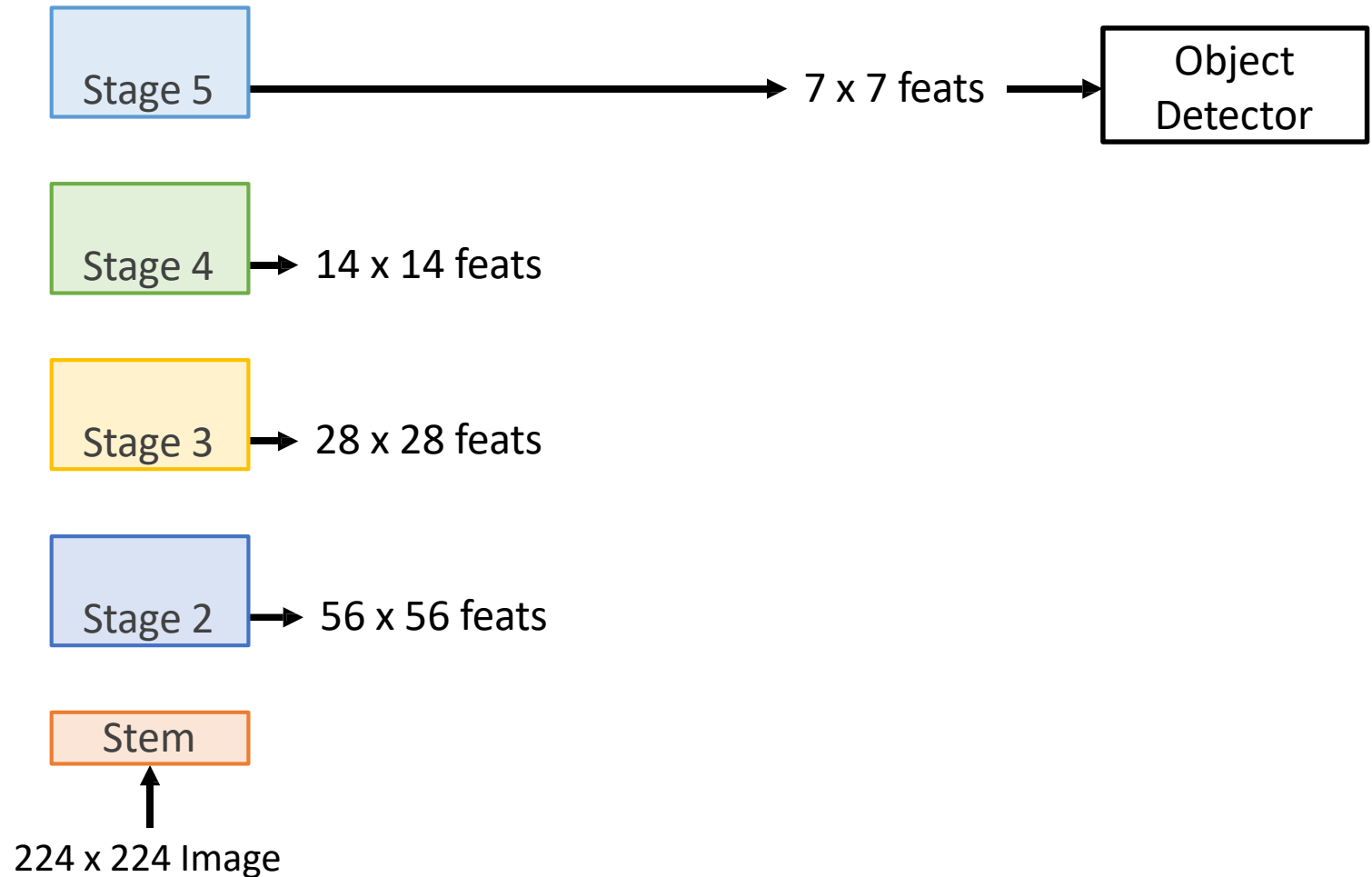
CNNs have multiple *stages* that operate at different resolutions. Attach an independent detector to the features at each level

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



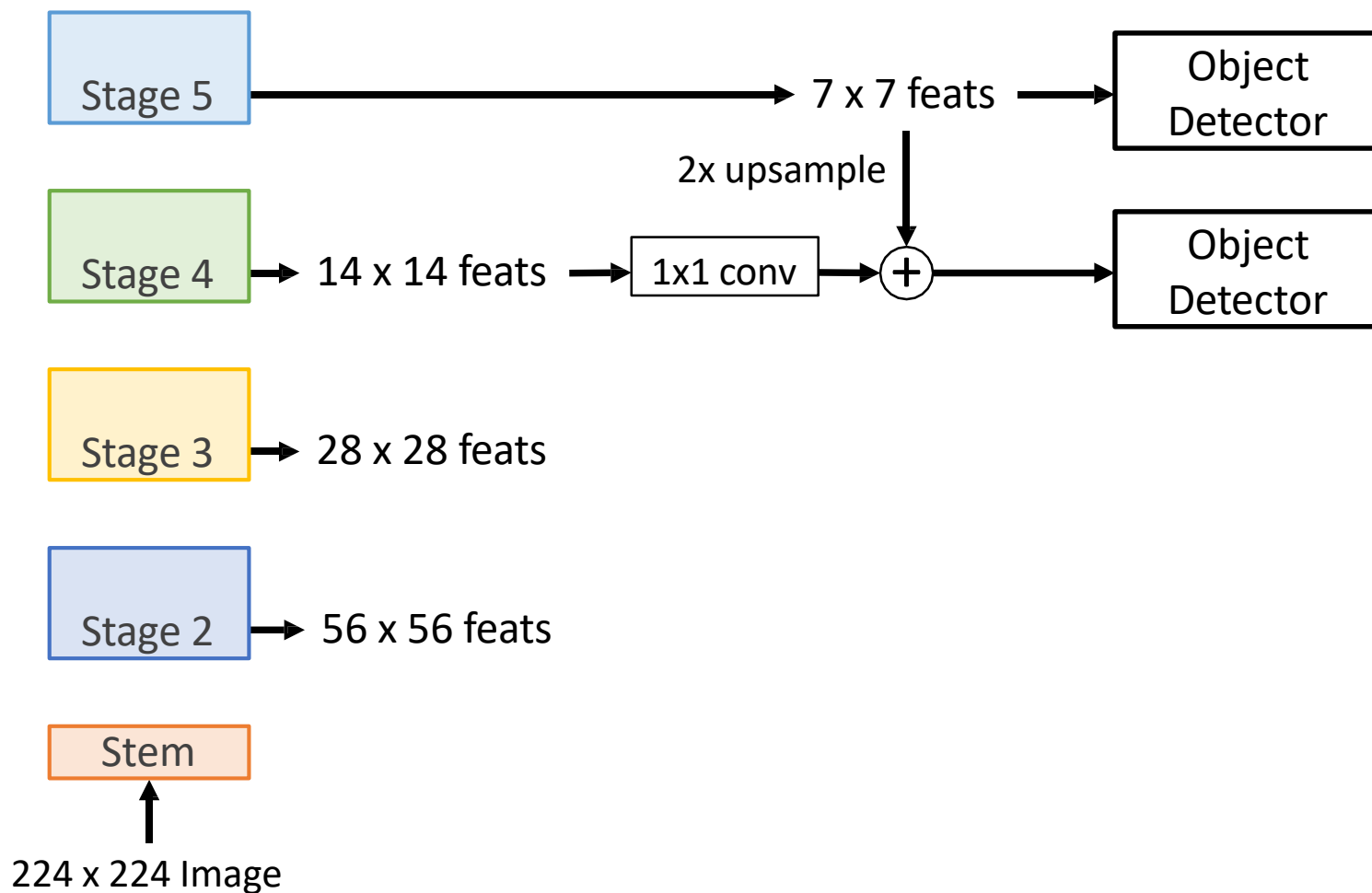
# Dealing with Scale: Feature Pyramid Network

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



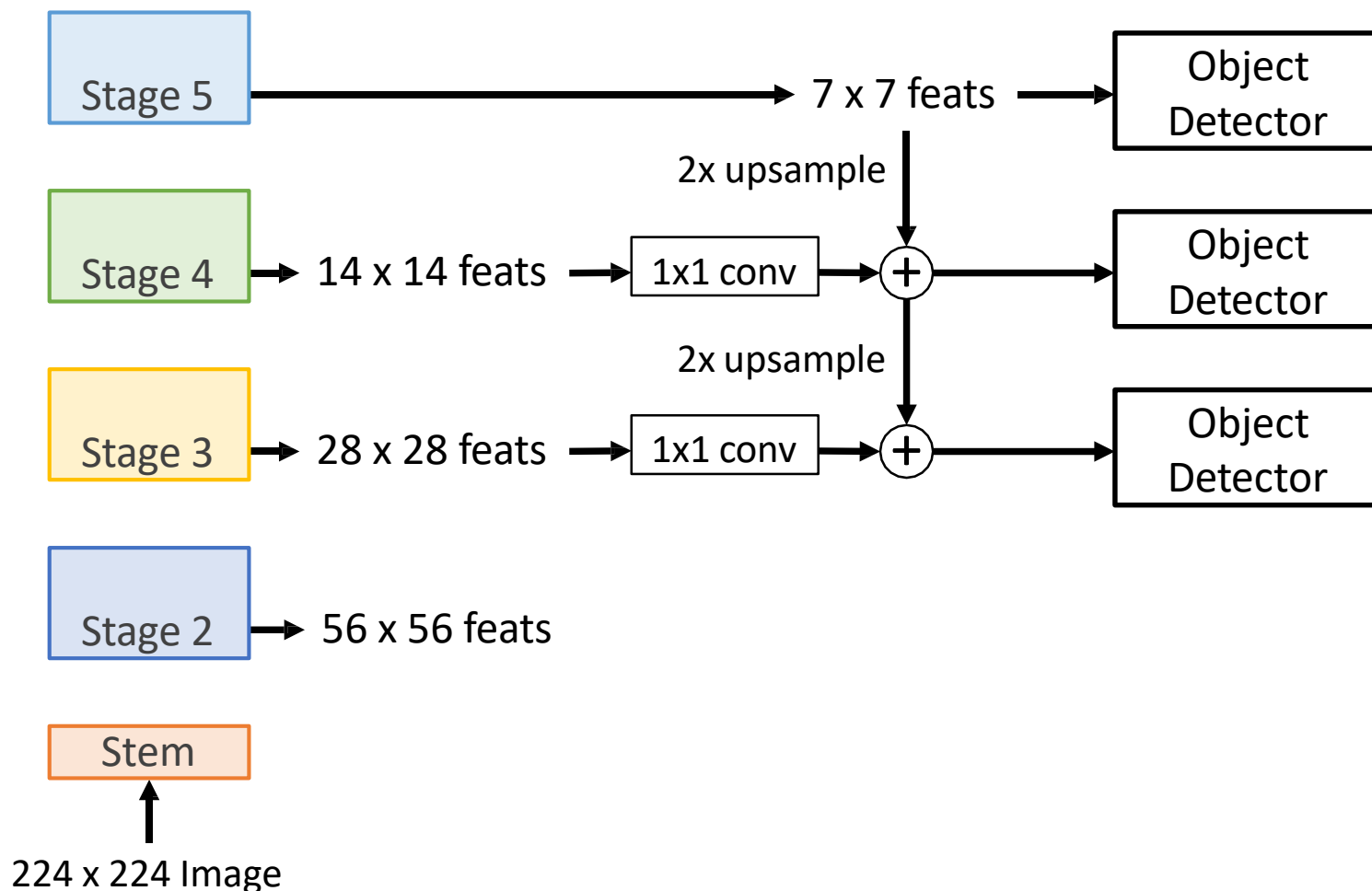
# Dealing with Scale: Feature Pyramid Network

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



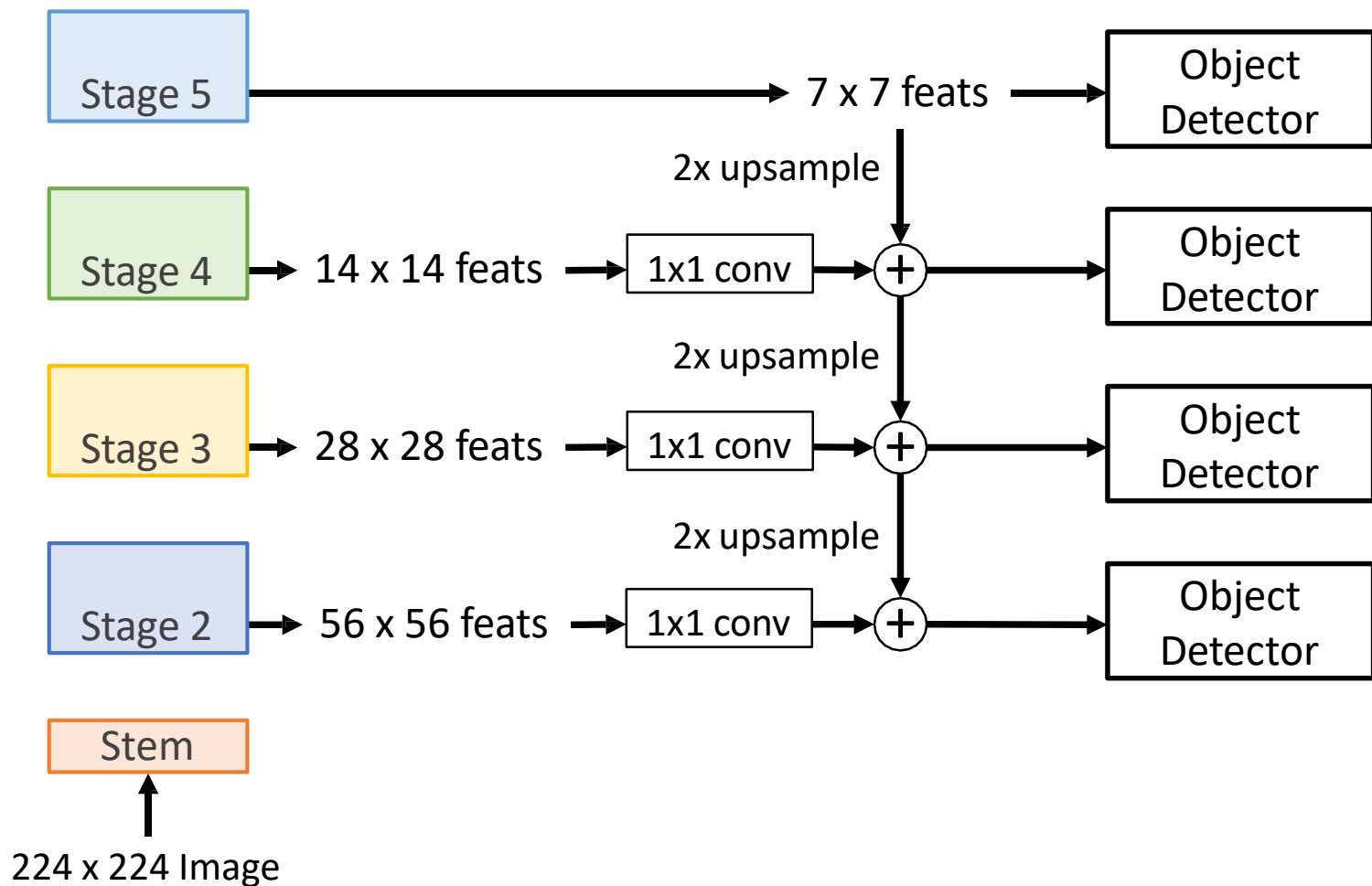
# Dealing with Scale: Feature Pyramid Network

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



# Dealing with Scale: Feature Pyramid Network

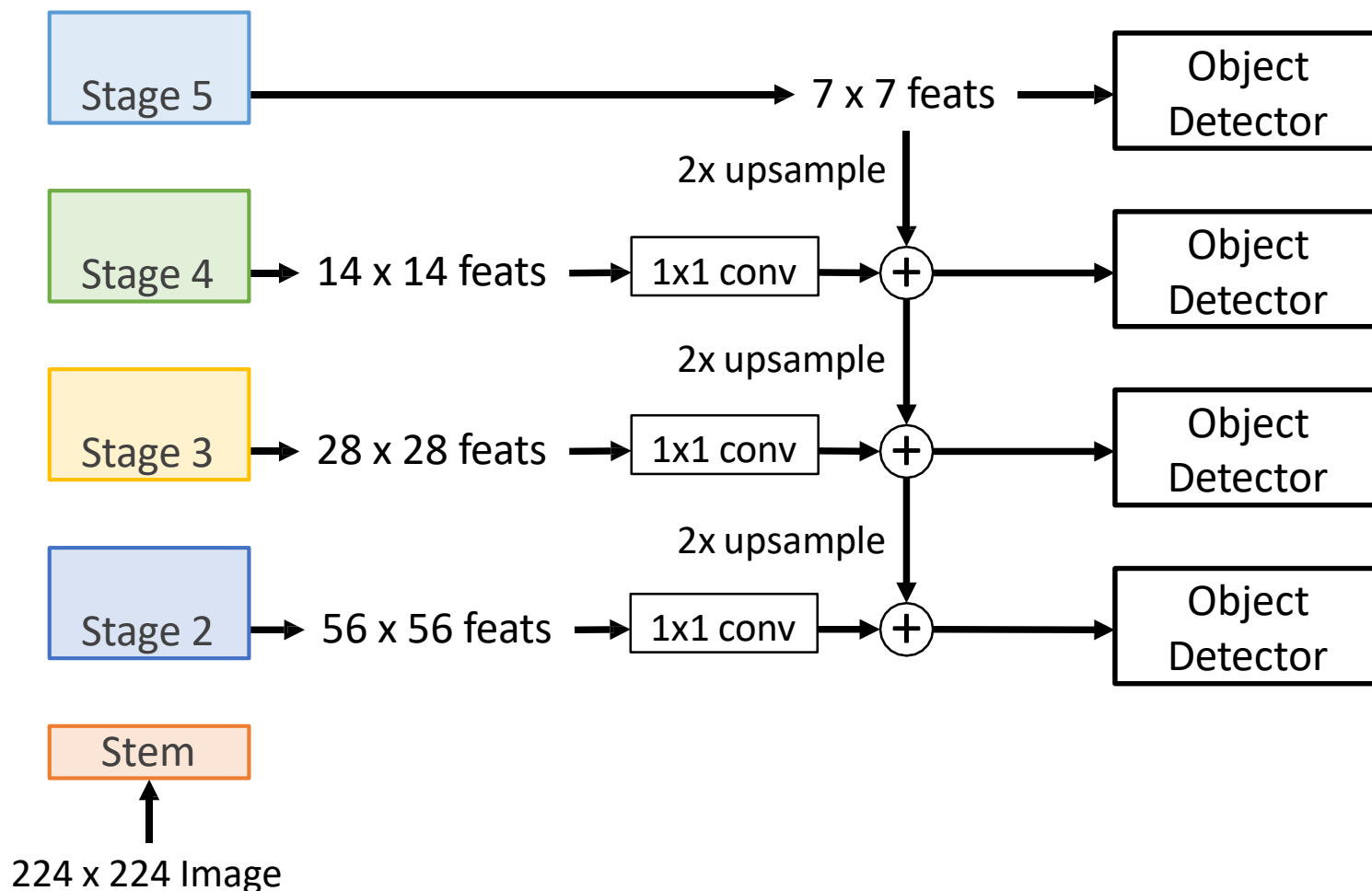
Add *top down* connections that feed information from high level features back down to lower level features



# Dealing with Scale: Feature Pyramid Network

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

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

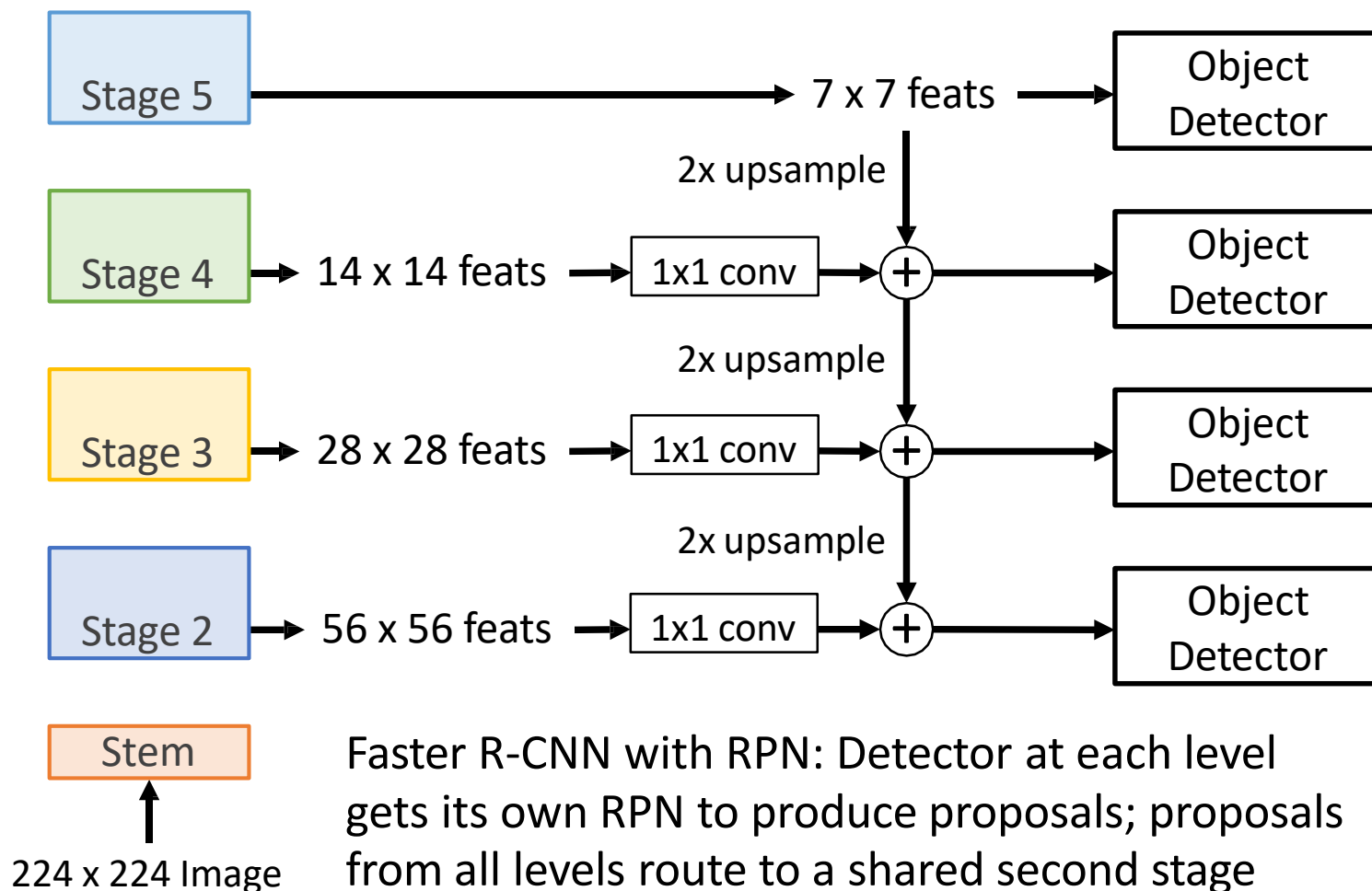




# Dealing with Scale: Feature Pyramid Network

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



# Fast<sub>er</sub> R-CNN: Learnable Region Proposals

Question: Do we really need the second stage?

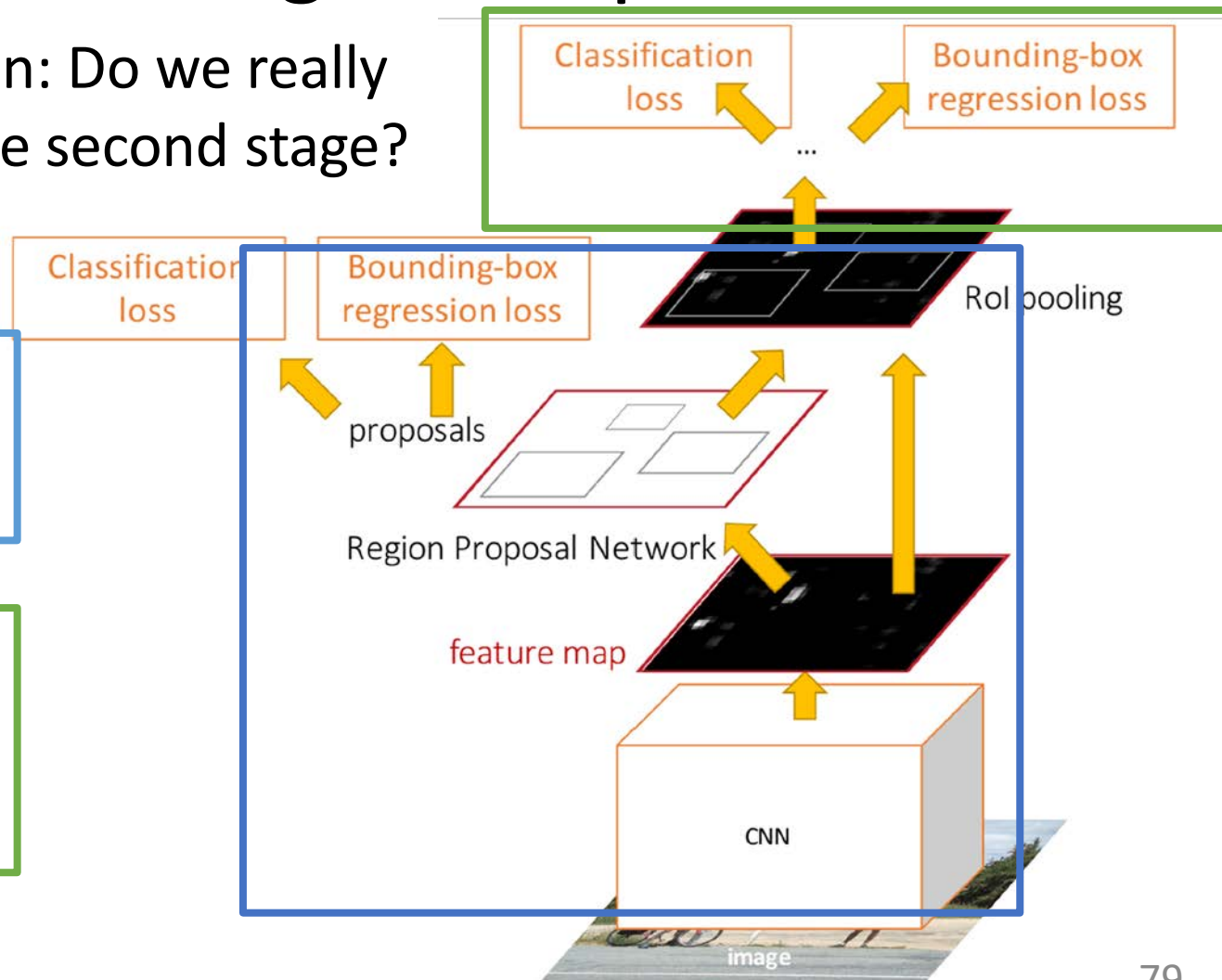
Faster R-CNN is a  
**Two-stage object detector**

First stage: Run once per image

- Backbone network
- Region proposal network

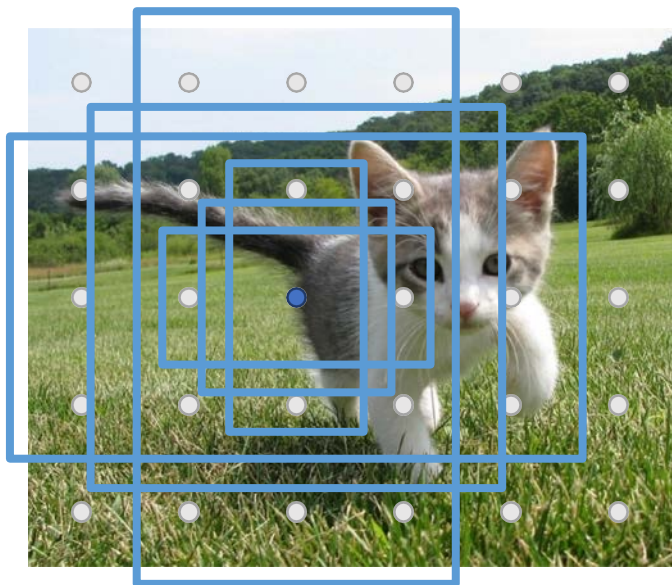
Second stage: Run once per region

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



## Single-Stage Detectors: RetinaNet

Run backbone CNN to get features aligned to input image



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



Each feature corresponds to a point in the input

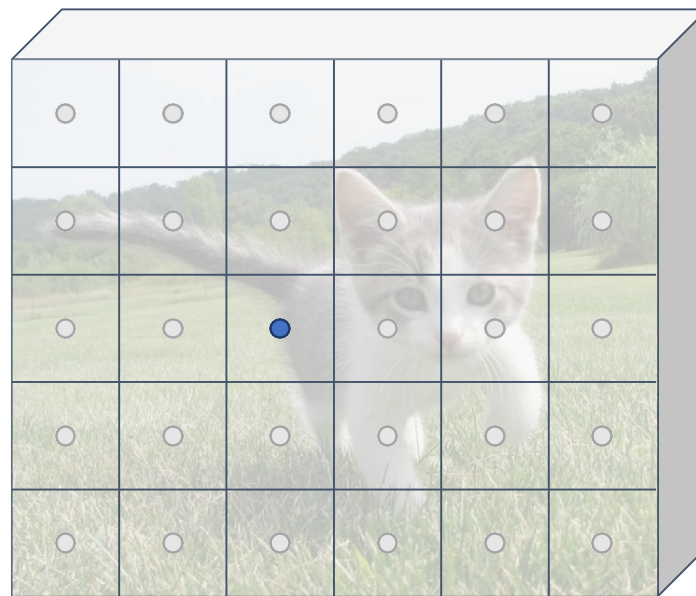


Image features  
(e.g. 512 x 5 x 6)

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

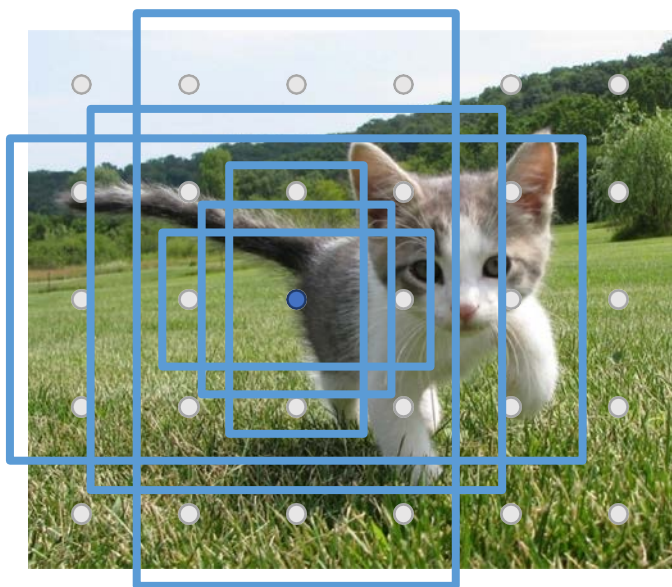


Anchor  
→ classification  
 $2K \times (C+1) \times 5 \times 6$

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

## Single-Stage Detectors: RetinaNet

Run backbone CNN to get features aligned to input image



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

Each feature corresponds to a point in the input

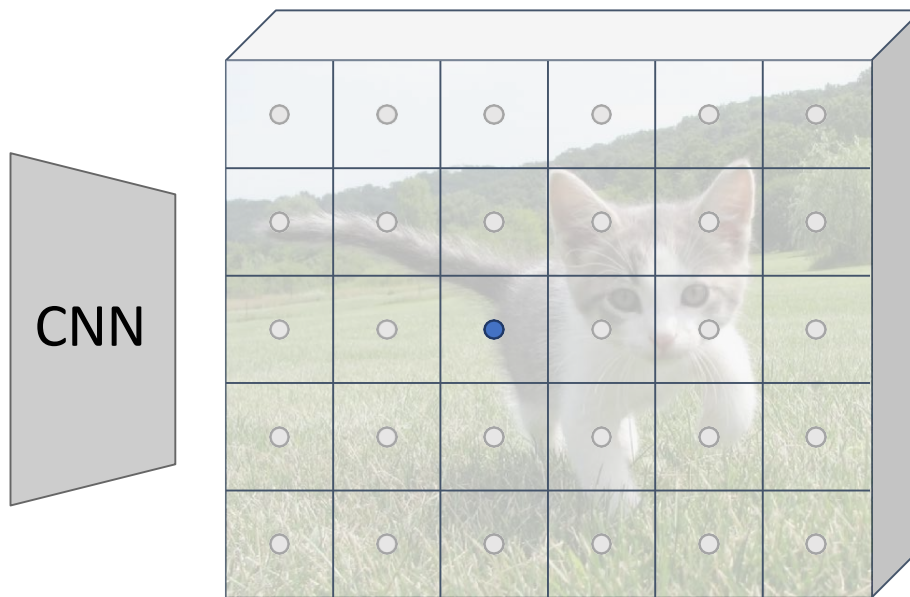
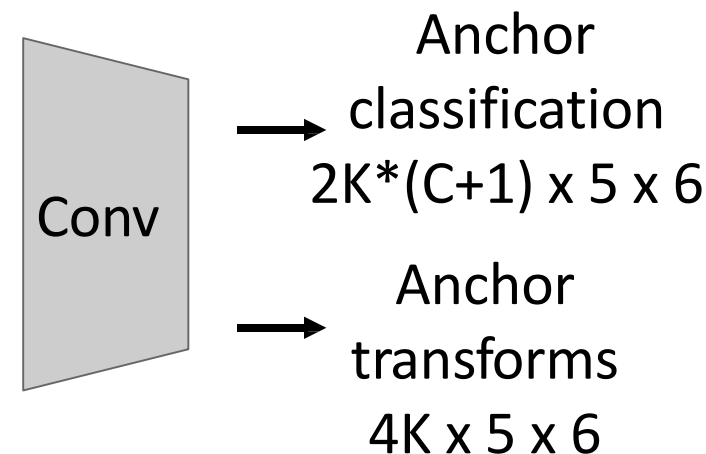


Image features  
(e.g. 512 x 5 x 6)

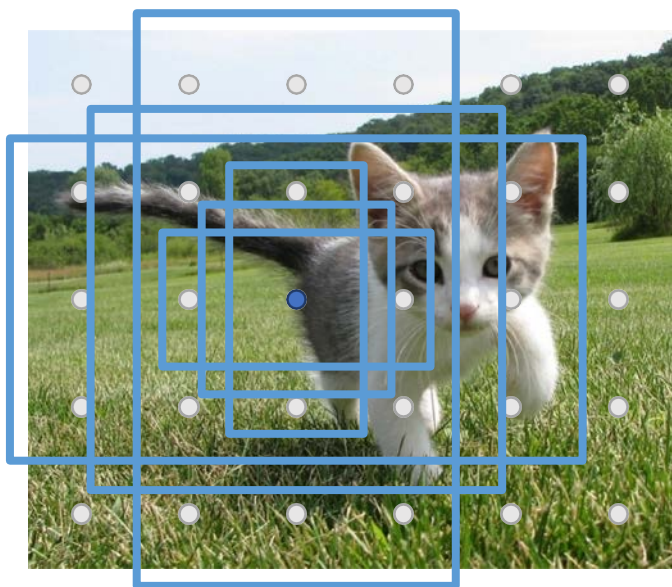
Problem: class imbalance – many more background anchors vs non-background



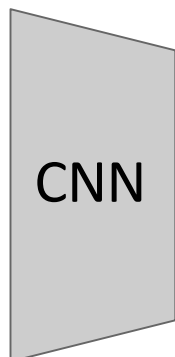


## Single-Stage Detectors: RetinaNet

Run backbone CNN to get features aligned to input image



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



Each feature corresponds to a point in the input

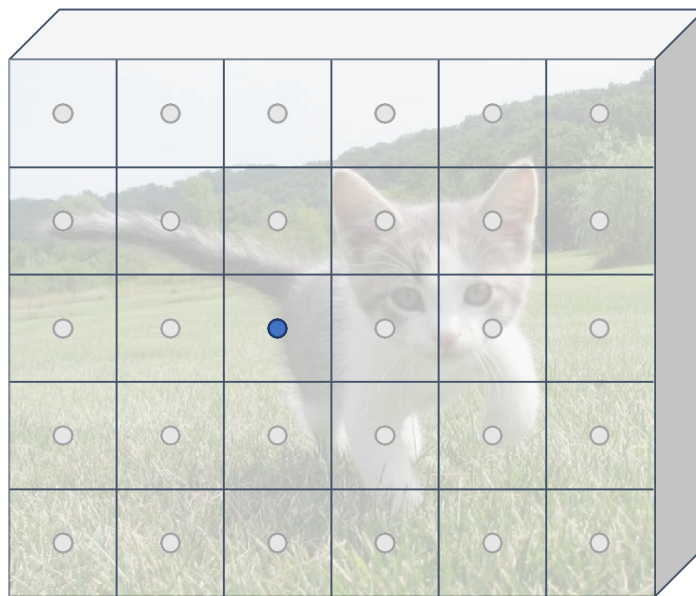


Image features  
(e.g. 512 x 5 x 6)

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

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



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

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

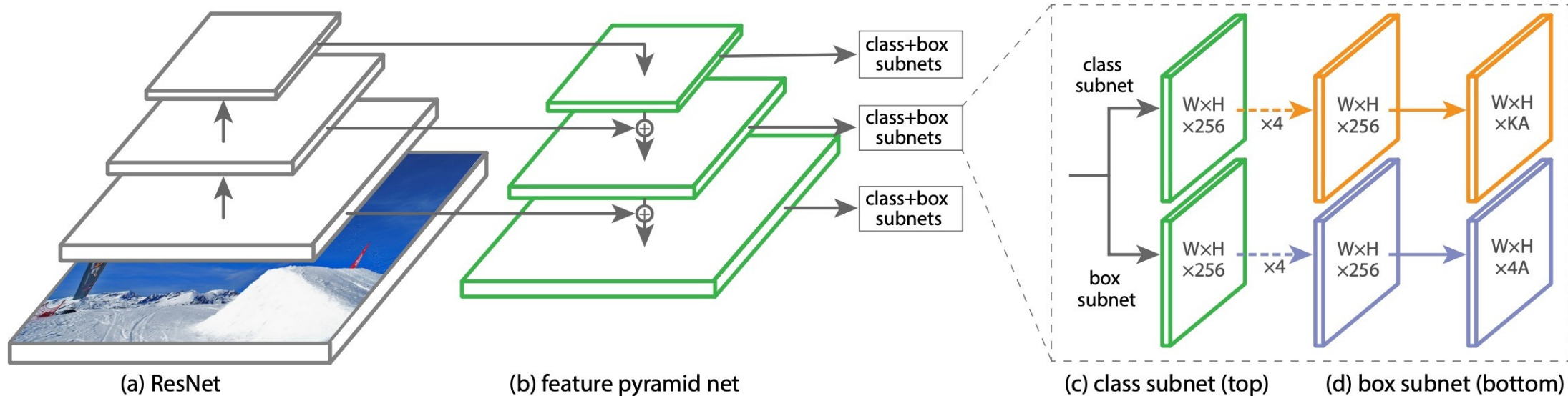
$$\text{CE}(p_t) = -\log(p_t)$$

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



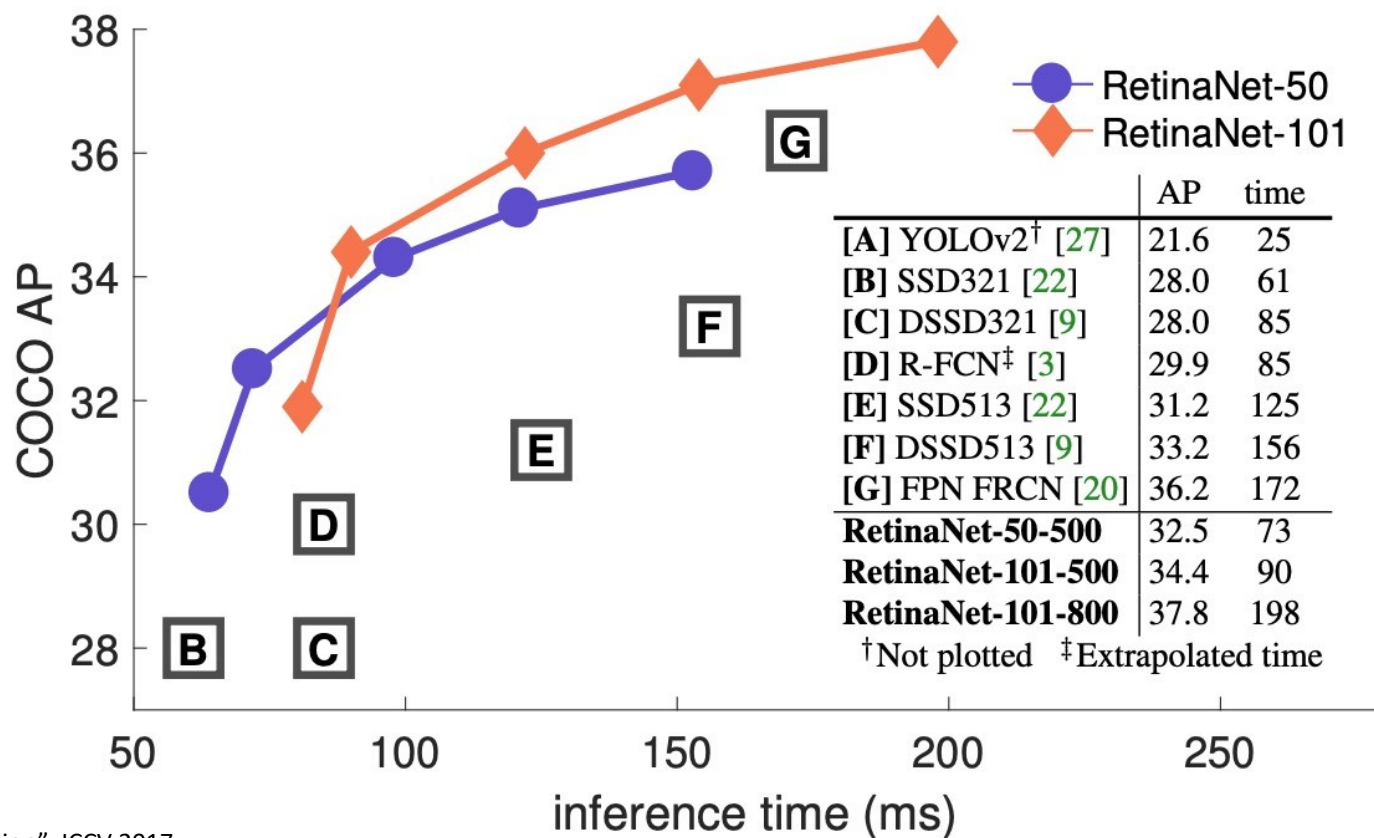
# Single-Stage Detectors: RetinaNet

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



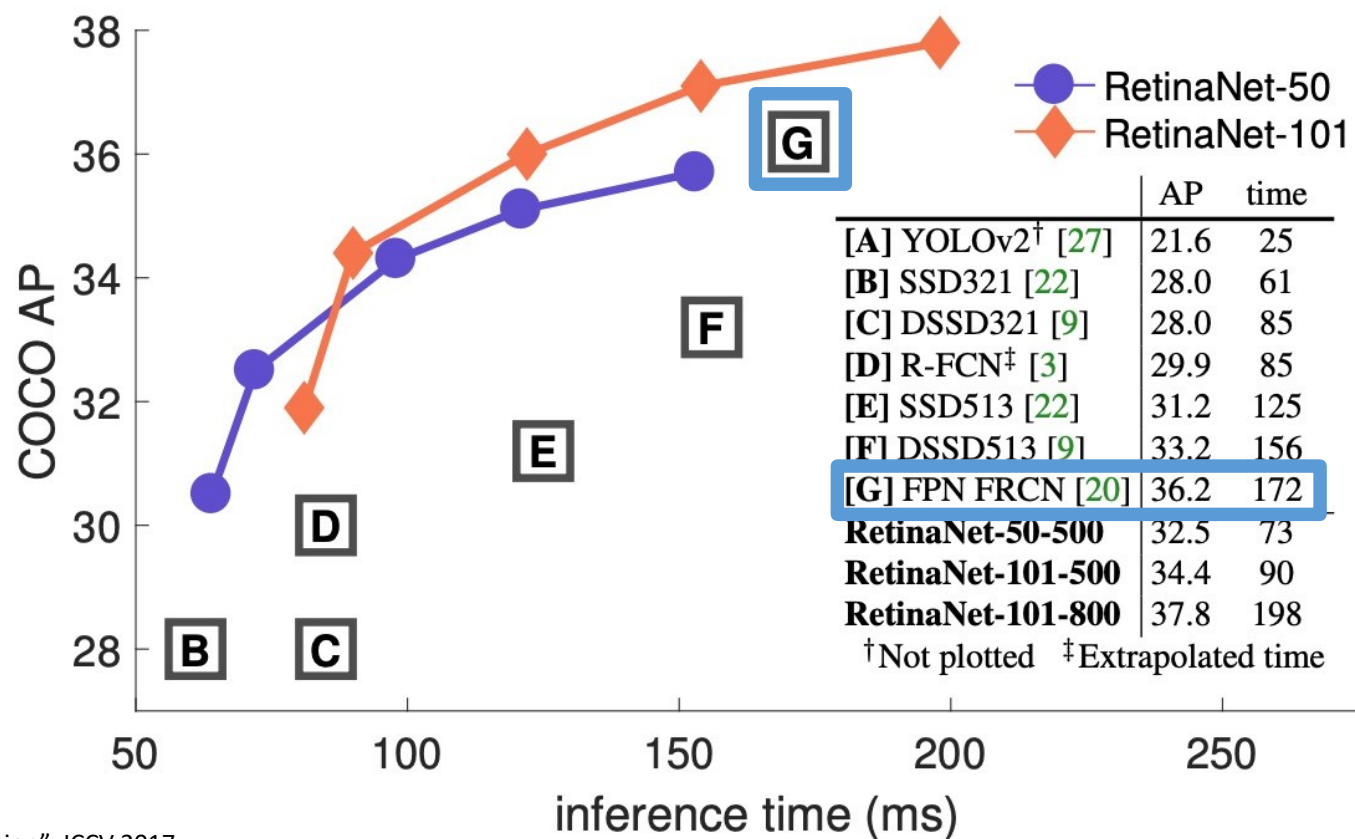
# Single-Stage Detectors: RetinaNet

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

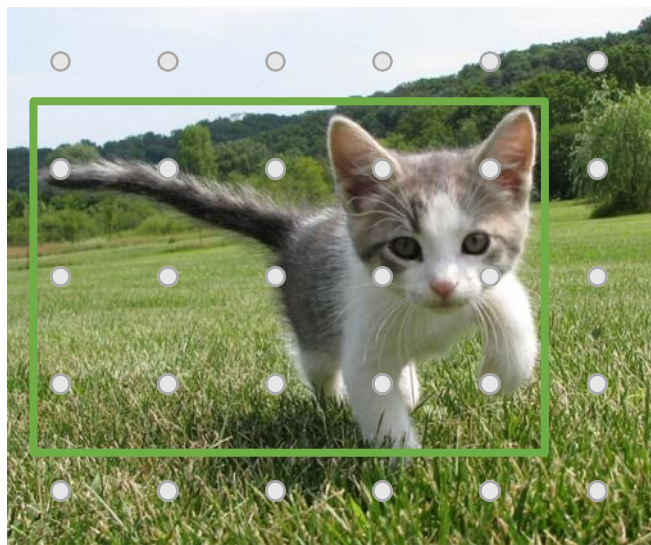


Faster R-CNN  
with Feature  
Pyramid  
Network

# Single-Stage Detectors: FCOS

“Anchor-free” detector

Run backbone CNN to get features aligned to input image



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

Each feature corresponds to a point in the input

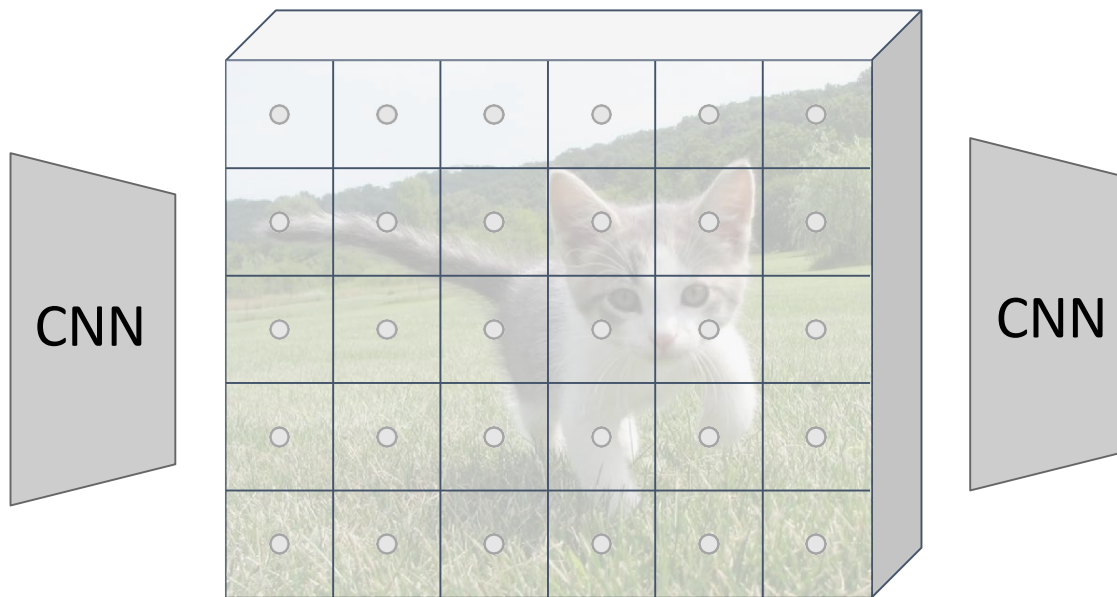
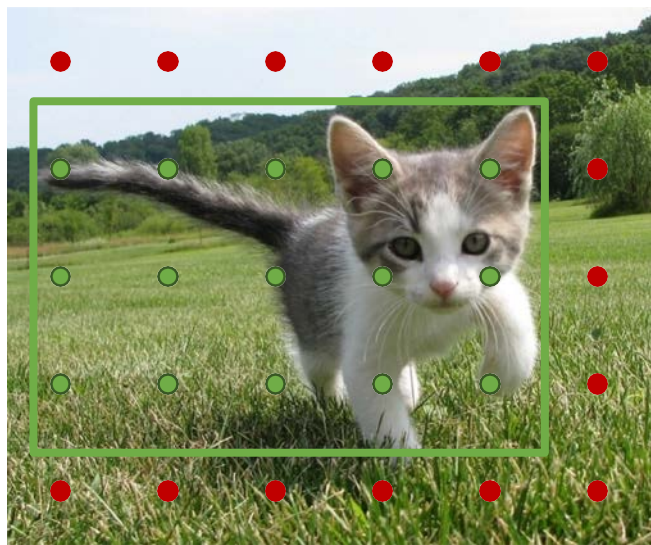


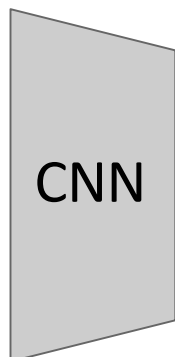
Image features  
(e.g. 512 x 5 x 6)

# Single-Stage Detectors: FCOS

Run backbone CNN to get features aligned to input image



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



Each feature corresponds to a point in the input

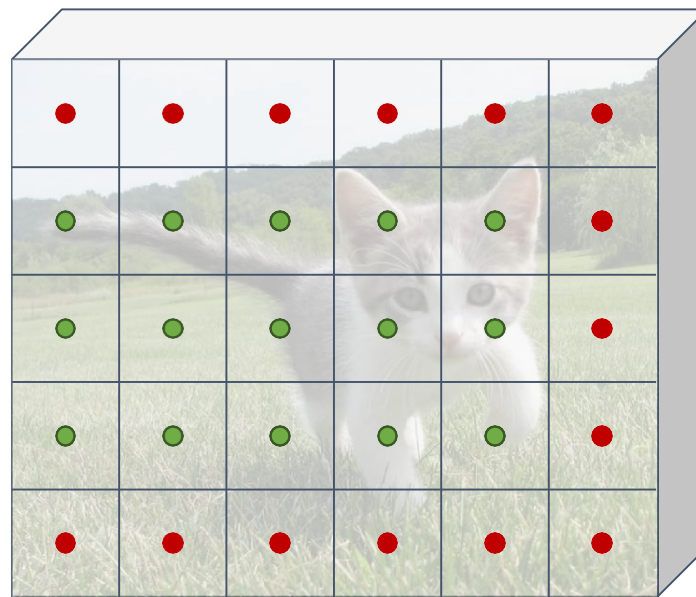


Image features  
(e.g. 512 x 5 x 6)

“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

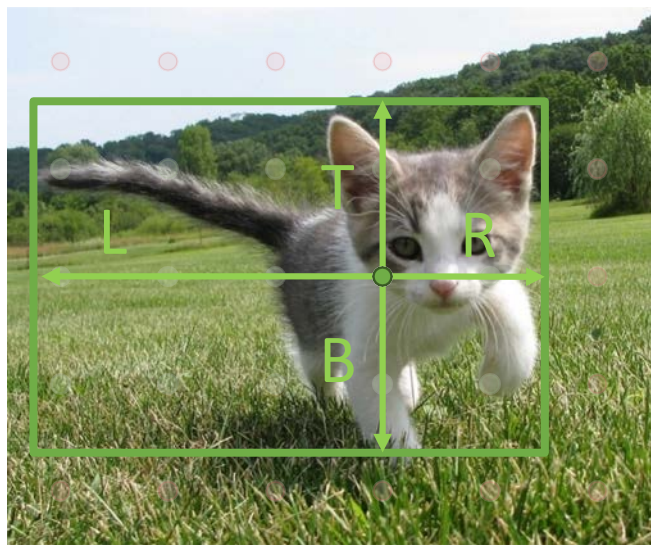


→ Class scores  
 $C \times 5 \times 6$



# Single-Stage Detectors: FCOS

Run backbone CNN to get features aligned to input image



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

Each feature corresponds to a point in the input

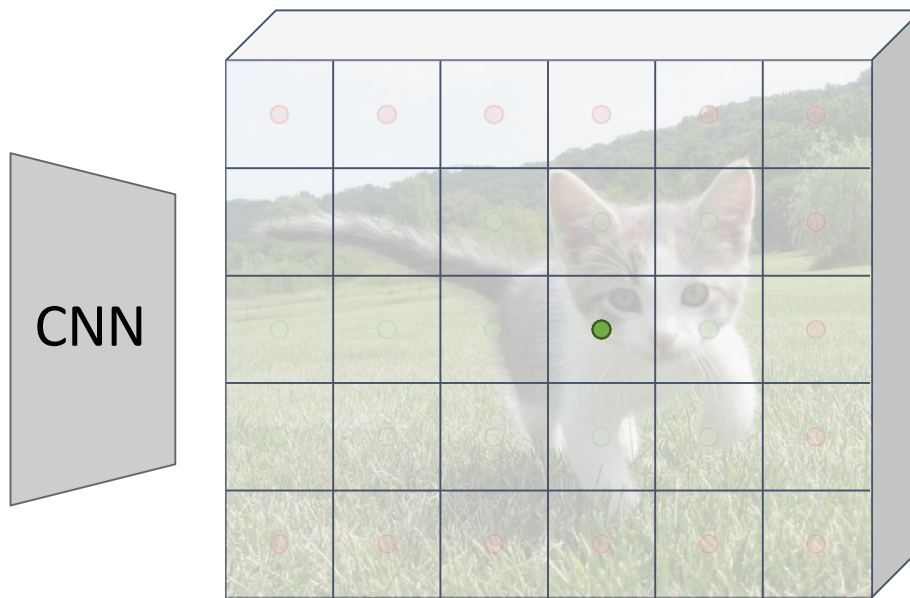


Image features  
(e.g. 512 x 5 x 6)

“Anchor-free” detector

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



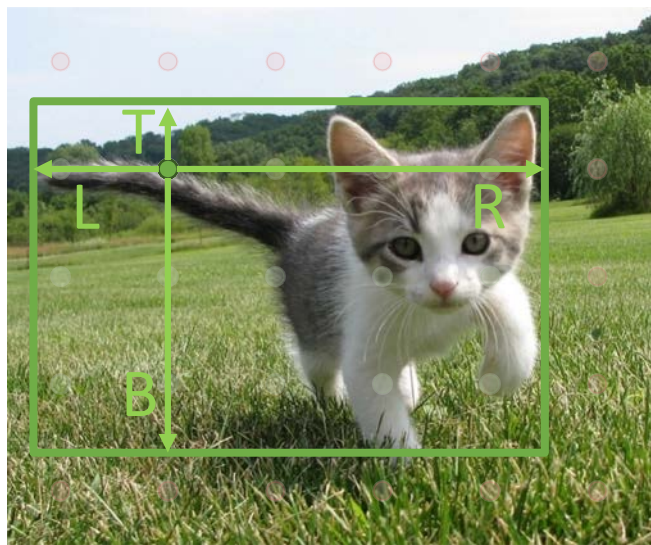
Class scores  
→ C x 5 x 6

Box edges  
→ 4 x 5 x 6



# Single-Stage Detectors: FCOS

Run backbone CNN to get features aligned to input image



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

Each feature corresponds to a point in the input

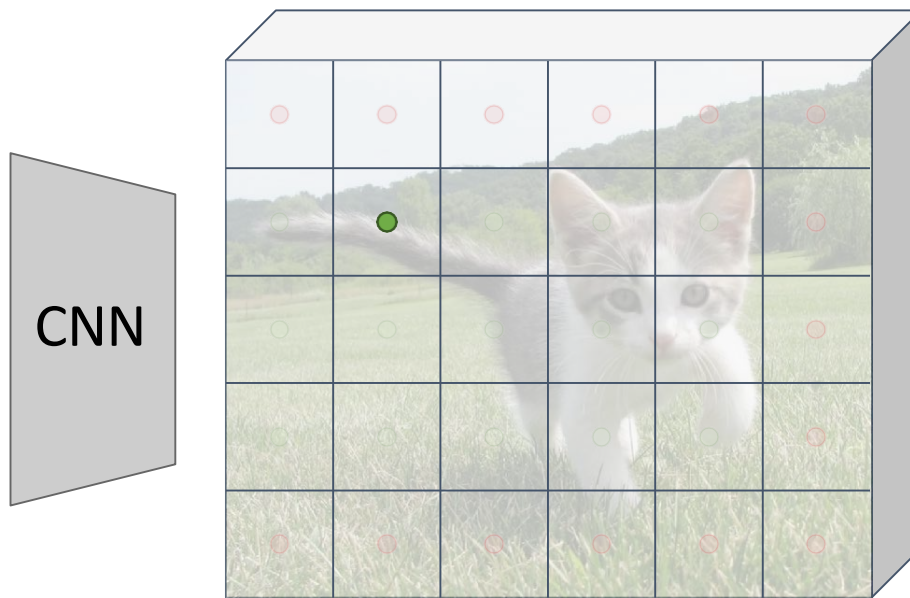


Image features  
(e.g. 512 x 5 x 6)

“Anchor-free” detector

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

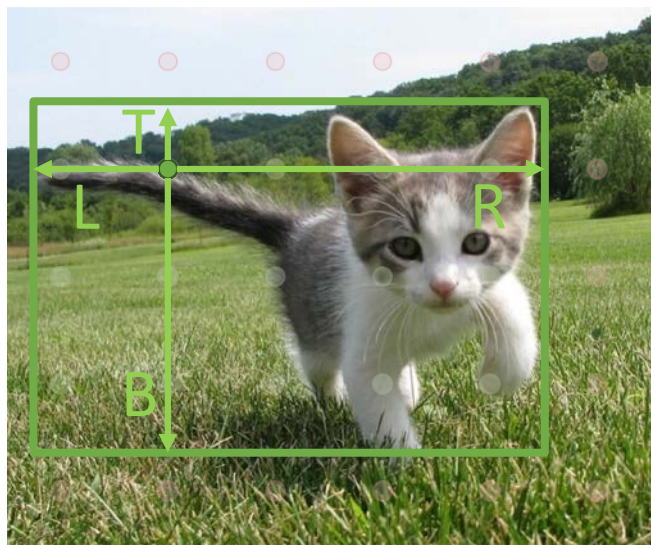


→ Class scores  
 $C \times 5 \times 6$

→ Box edges  
 $4 \times 5 \times 6$

# Single-Stage Detectors: FCOS

Run backbone CNN to get features aligned to input image



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

Each feature corresponds to a point in the input

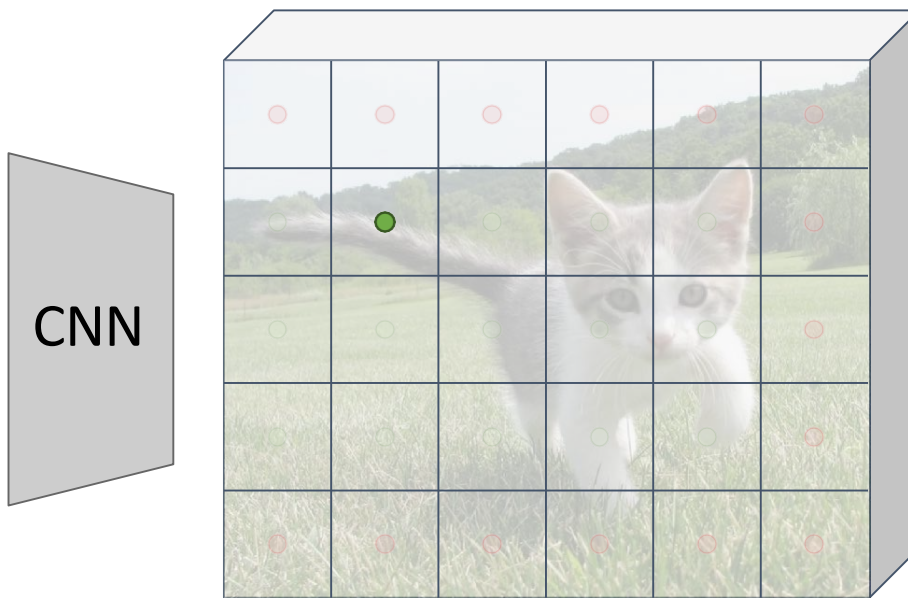
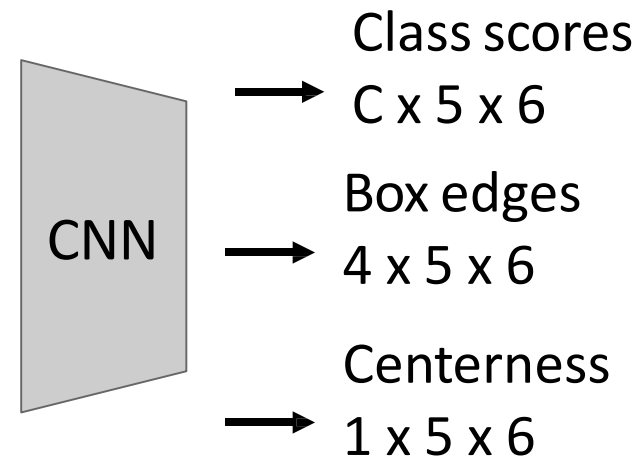


Image features  
(e.g. 512 x 5 x 6)

“Anchor-free” detector

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

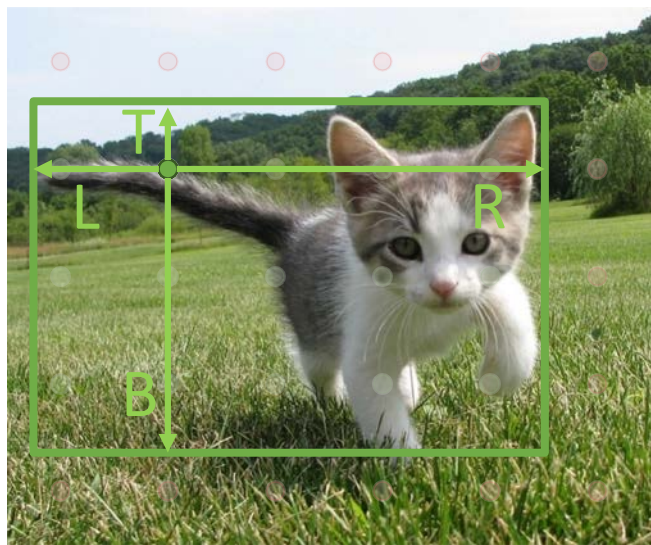


$$centerness = \sqrt{\frac{\min(L, R)}{\max(L, R)} \cdot \frac{\min(T, B)}{\max(T, B)}}$$

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 x 640 x 480)

Each feature corresponds to a point in the input

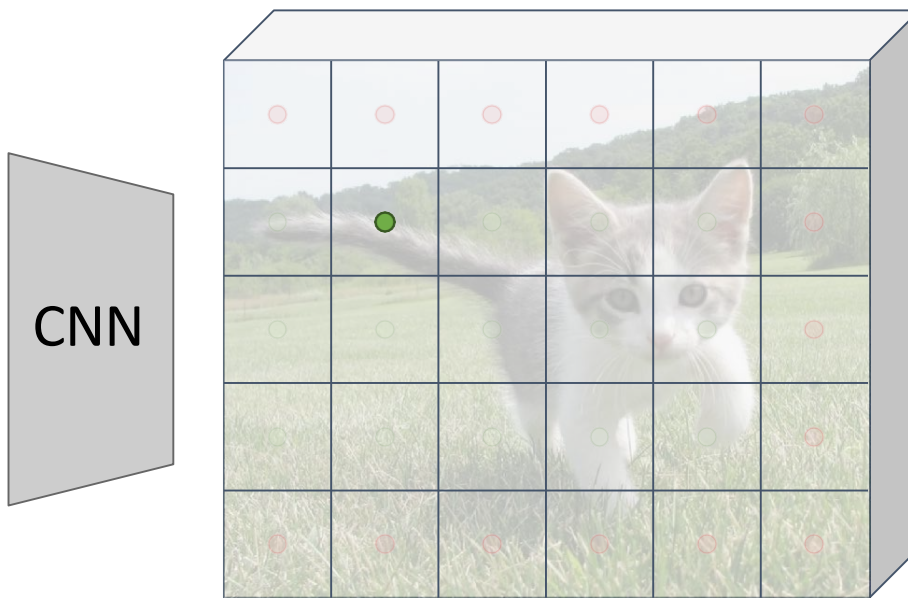


Image features  
(e.g. 512 x 5 x 6)

“Anchor-free” detector

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

Class scores  
→ C x 5 x 6  
Box edges  
→ 4 x 5 x 6  
Centerness  
→ 1 x 5 x 6

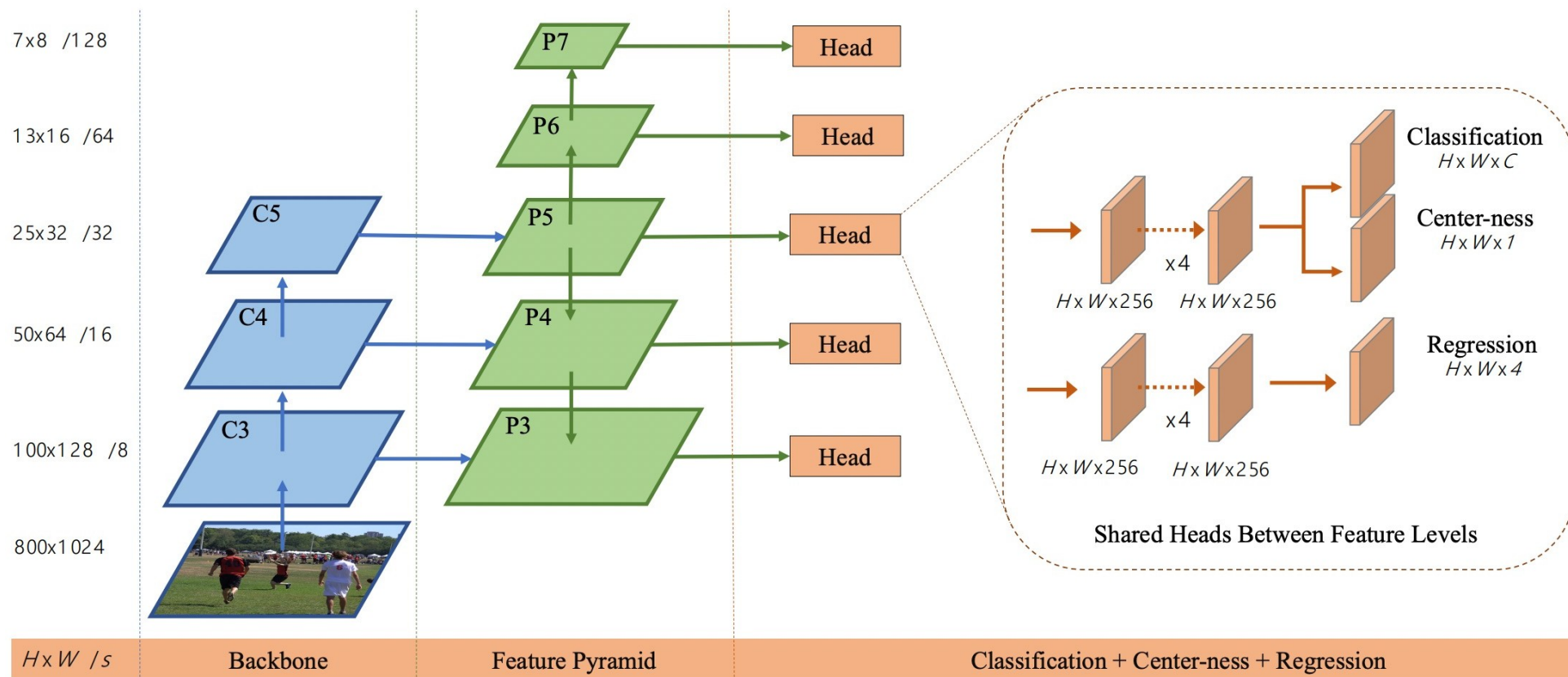
$$centerness = \sqrt{\frac{\min(L, R)}{\max(L, R)} \cdot \frac{\min(T, B)}{\max(T, B)}}$$

Ranges from 1 at box center to 0 at box edge



# Single-Stage Detectors: FCOS

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



# Summary: Beyond Image Classification

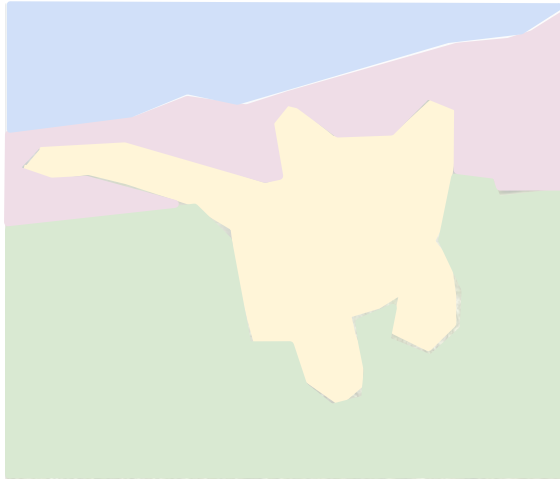
## Classification



**CAT**

No spatial extent

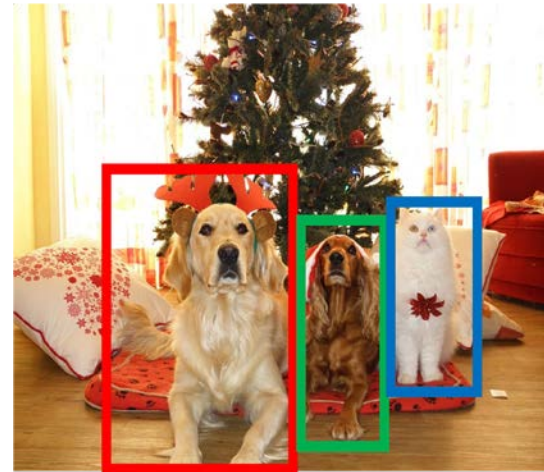
## Semantic Segmentation



GRASS, CAT, TREE,  
SKY

No objects, just pixels

## Object Detection



DOG, DOG, CAT

Multiple Objects

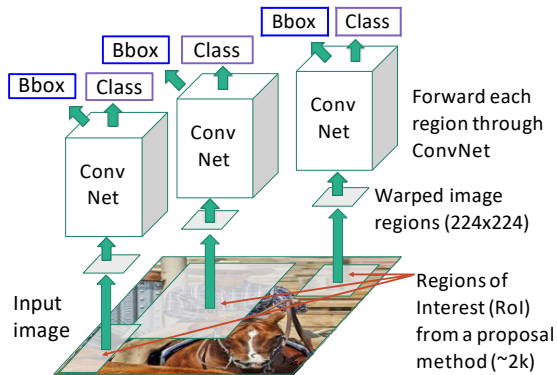
## Instance Segmentation



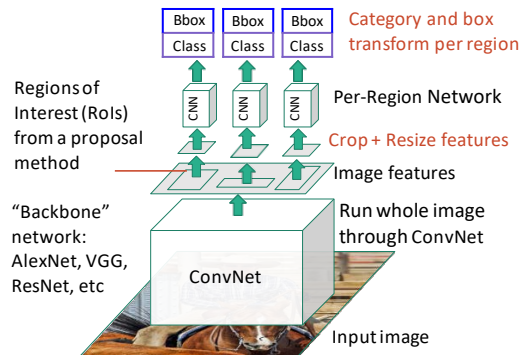
DOG, DOG, CAT

# Summary

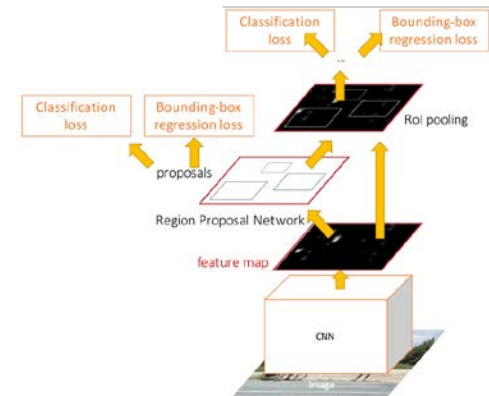
**“Slow” R-CNN:** Run CNN independently for each region



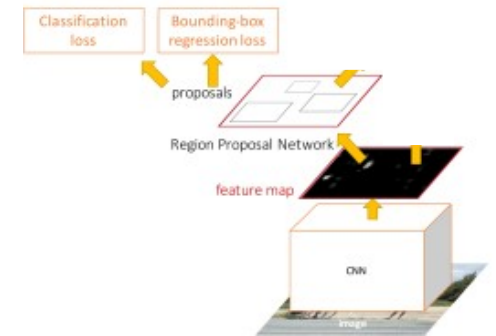
**Fast R-CNN:** Apply differentiable cropping to shared image features



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



**Single-Stage:** Fully convolutional detector



With anchors: RetinaNet  
Anchor-Free: FCOS



Next time:  
Image and Instance Segmentation