

12. Object Detection and Image Segmentation

GEV6135 Deep Learning for Visual Recognition and Applications

Kibok Lee
Assistant Professor of
Applied Statistics / Statistics and Data Science
Dec 1, 2022



Assignment 6

- Due **Friday 12/2, 11:59pm KST**
- Convolutional networks
 - Modularized implementation (loss is given!)
 - BatchNorm is given, but should be plugged into DeepConvNet
- Before submitting your work, we recommend you
 - Re-download clean files
 - Copy-paste your solution to clean py
 - Re-run clean ipynb only once
- If you feel difficult, consider to take **option 2.**

Assignment 7

- Due **Wednesday 12/14, 11:59pm KST**
- PyTorch autograd and nn
 - 3 different abstraction levels
 - Residual networks
- Before submitting your work, we recommend you
 - Re-download clean files
 - Copy-paste your solution to clean py
 - Re-run clean ipynb only once
- If you feel difficult, consider to take **option 2.**

So Far: Image Classification



This image is [CC0 public domain](#)

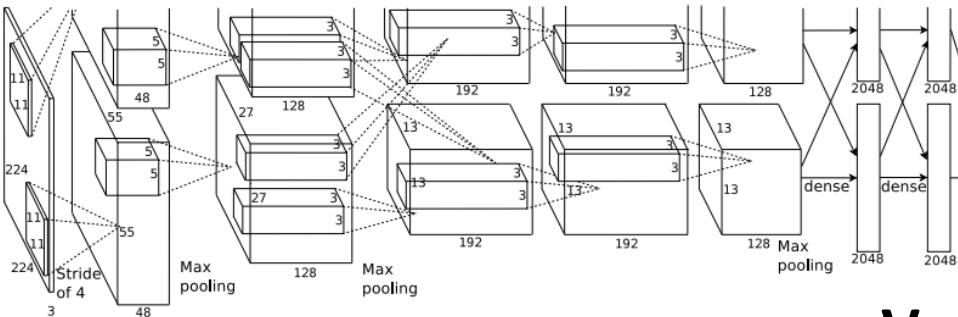


Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

Vector:
4096

Fully-Connected:
4096 to 1000

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

Computer Vision Tasks

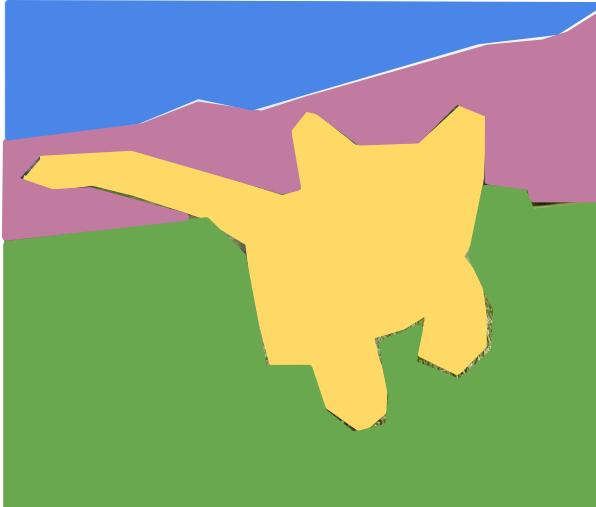
Classification



CAT

No spatial extent

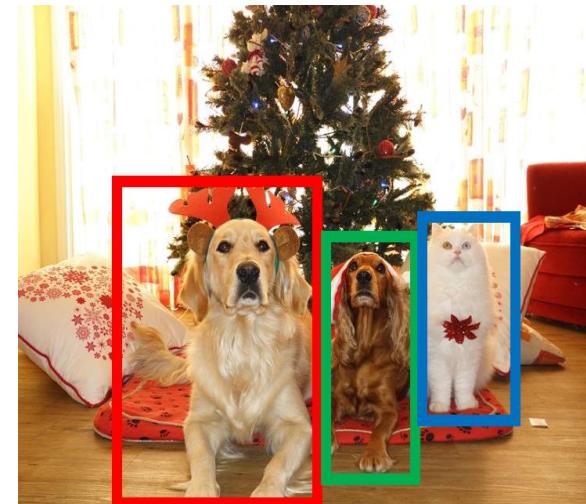
Semantic Segmentation



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](#)

Computer Vision Task: Semantic Segmentation

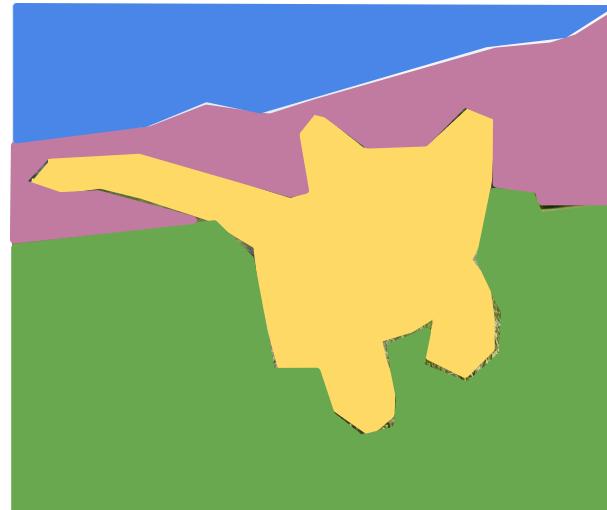
Classification



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](#)

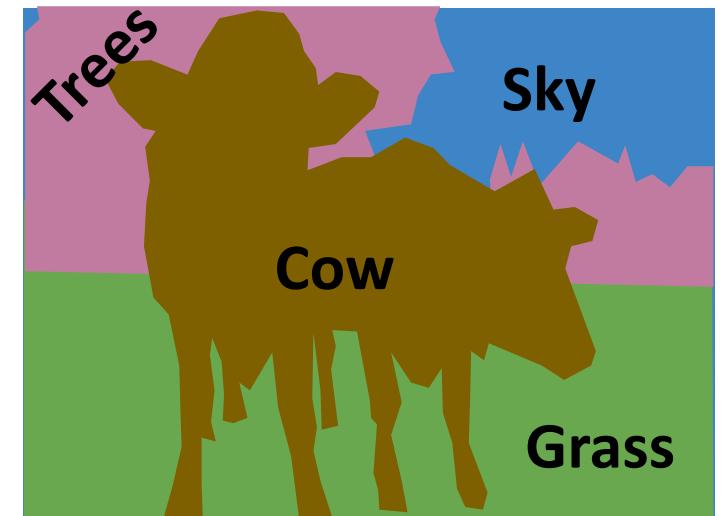
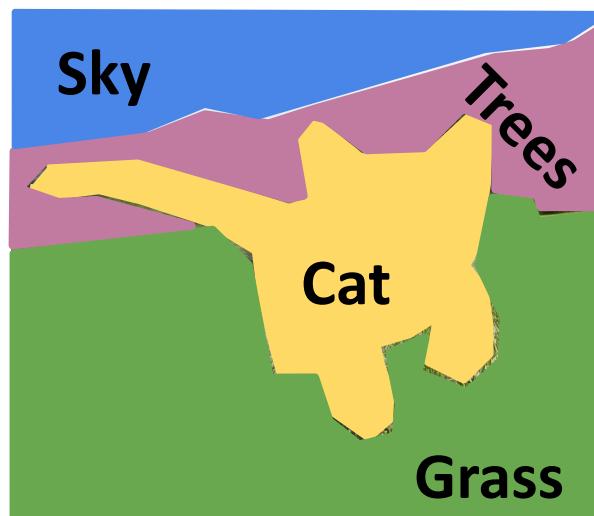
Semantic Segmentation

Label each pixel in the image with a category label

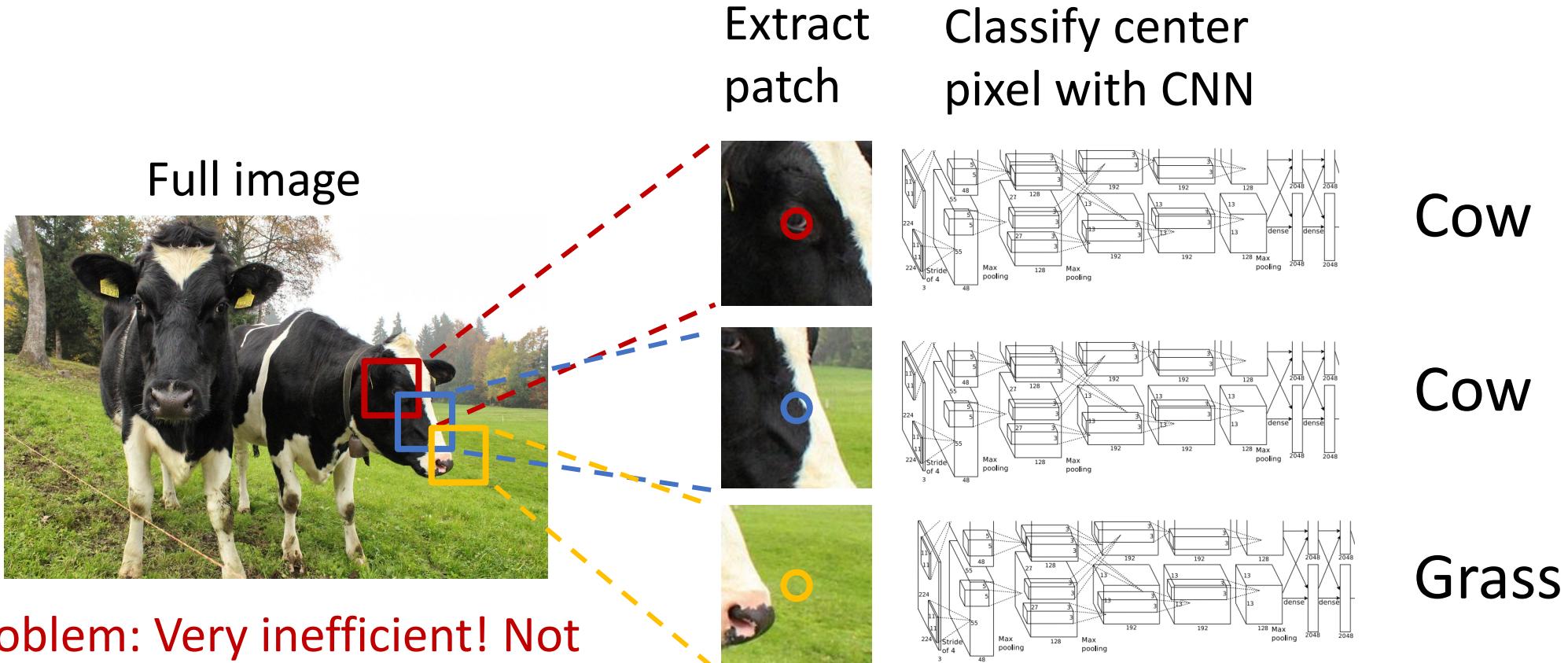
Don't differentiate instances, only care about pixels



This image is [CC0 public domain](#)



Semantic Segmentation Idea: Sliding Window

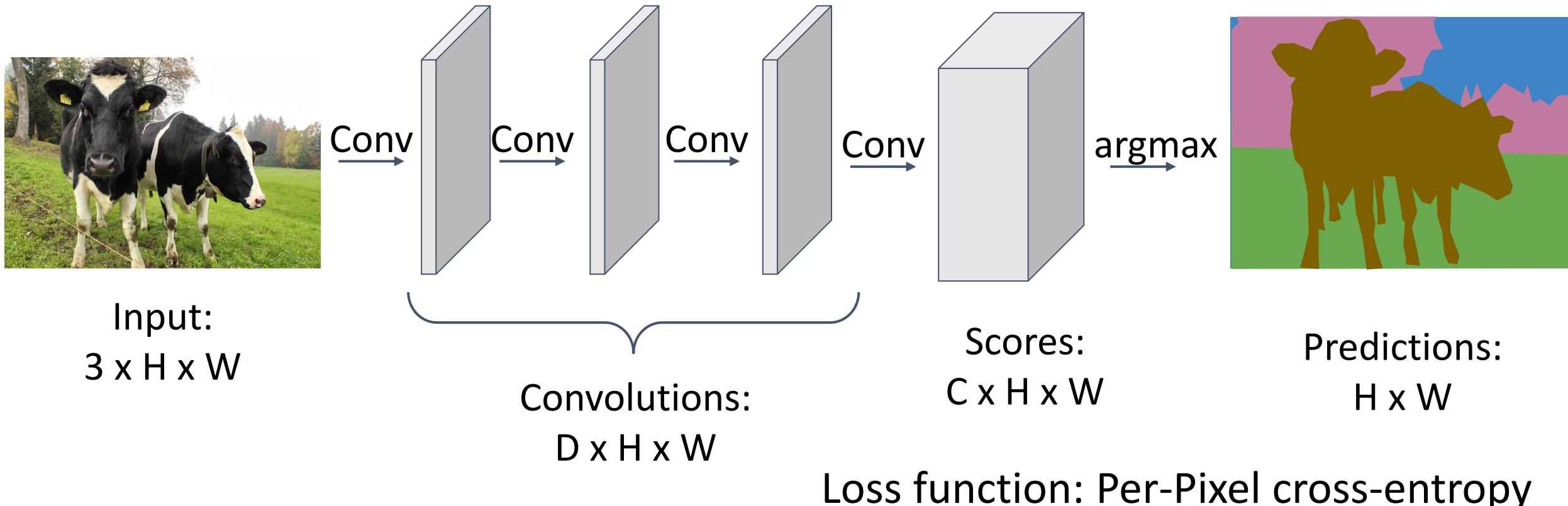


Problem: Very inefficient! Not reusing shared features between overlapping patches

Farabet et al, "Learning Hierarchical Features for Scene Labeling," TPAMI 2013
Pinheiro and Collobert, "Recurrent Convolutional Neural Networks for Scene Labeling", ICML 2014

Semantic Segmentation: Fully Convolutional Network

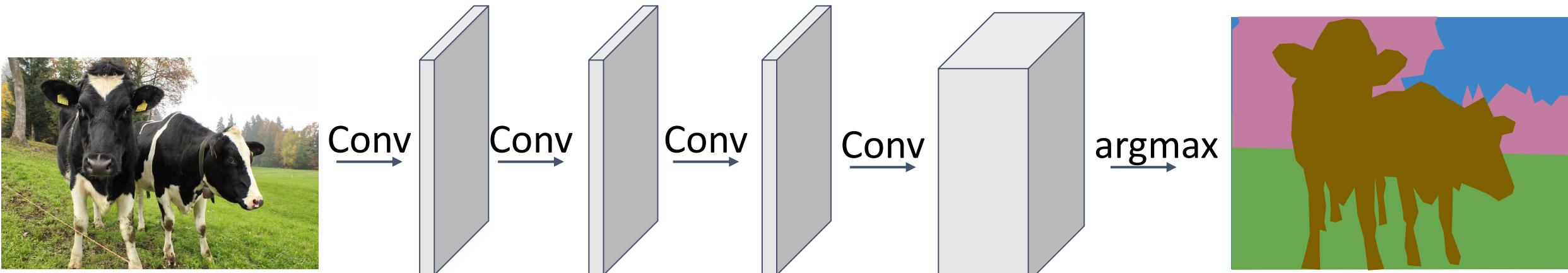
Design a network as a bunch of convolutional layers to make predictions for pixels all at once!



Long et al, "Fully convolutional networks for semantic segmentation", CVPR 2015

Semantic Segmentation: Fully Convolutional Network

Design a network as a bunch of convolutional layers to make predictions for pixels all at once!



Input: $3 \times H \times W$ **Problem #1:** Effective receptive field size is linear in number of conv layers: With L 3×3 conv layers, receptive field is $1+2L$

Problem #2: Convolution on high res images is expensive!
Recall ResNet stem aggressively downsamples

Long et al, "Fully convolutional networks for semantic segmentation", CVPR 2015

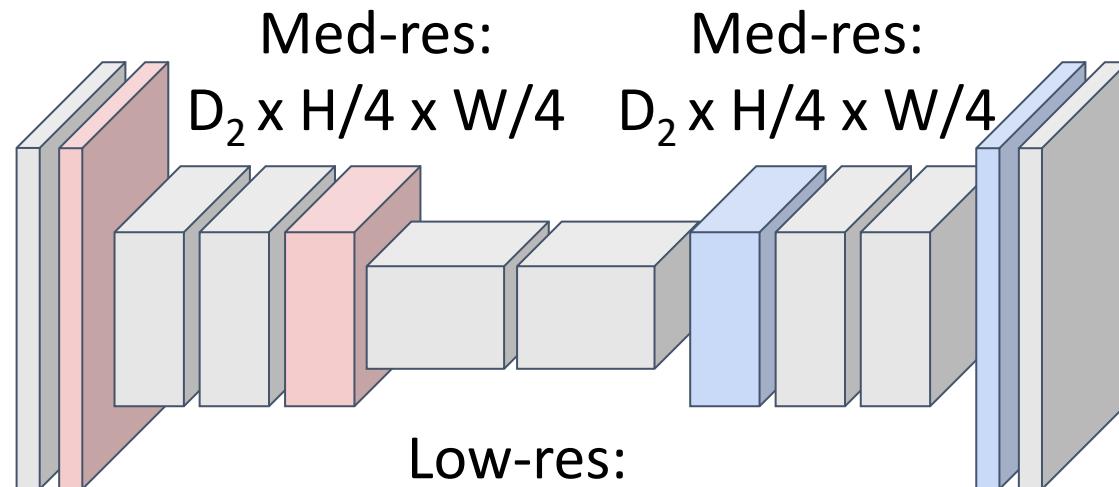
Semantic Segmentation: Fully Convolutional Network

Downsampling:
Pooling, strided
convolution



Input:
 $3 \times H \times W$

High-res:
 $D_1 \times H/2 \times W/2$



Design network as a bunch of convolutional layers, with
downsampling and **upsampling** inside the network!

Upsampling:
???



Predictions:
 $H \times W$

In-Network Upsampling: “Unpooling”

Bed of Nails

| | |
|---|---|
| 1 | 2 |
| 3 | 4 |



| | | | |
|---|---|---|---|
| 1 | 0 | 2 | 0 |
| 0 | 0 | 0 | 0 |
| 3 | 0 | 4 | 0 |
| 0 | 0 | 0 | 0 |

Nearest Neighbor

| | |
|---|---|
| 1 | 2 |
| 3 | 4 |



| | | | |
|---|---|---|---|
| 1 | 1 | 2 | 2 |
| 1 | 1 | 2 | 2 |
| 3 | 3 | 4 | 4 |
| 3 | 3 | 4 | 4 |

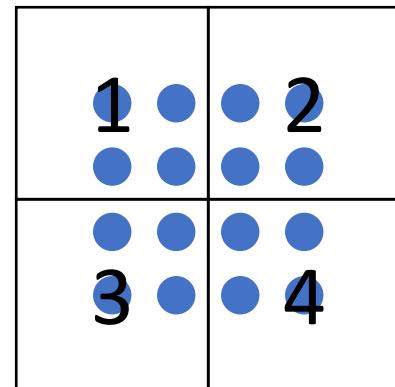
Input
 $C \times 2 \times 2$

Output
 $C \times 4 \times 4$

Input
 $C \times 2 \times 2$

Output
 $C \times 4 \times 4$

In-Network Upsampling: Bilinear Interpolation



| | | | |
|------|------|------|------|
| 1.00 | 1.25 | 1.75 | 2.00 |
| 1.50 | 1.75 | 2.25 | 2.50 |
| 2.50 | 2.75 | 3.25 | 3.50 |
| 3.00 | 3.25 | 3.75 | 4.00 |

Input: $C \times 2 \times 2$

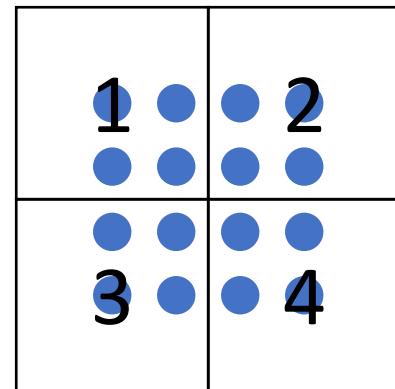
Output: $C \times 4 \times 4$

$$f_{x,y} = \sum_{i,j} f_{i,j} \max(0, 1 - |x - i|) \max(0, 1 - |y - j|) \quad i \in \{\lfloor x \rfloor - 1, \dots, \lceil x \rceil + 1\}$$

Use two closest neighbors in x and y
to construct linear approximations

$$j \in \{\lfloor y \rfloor - 1, \dots, \lceil y \rceil + 1\}$$

In-Network Upsampling: Bicubic Interpolation



| | | | |
|------|------|------|------|
| 0.68 | 1.02 | 1.56 | 1.89 |
| 1.35 | 1.68 | 2.23 | 2.56 |
| 2.44 | 2.77 | 3.32 | 3.65 |
| 3.11 | 3.44 | 3.98 | 4.32 |

Input: $C \times 2 \times 2$

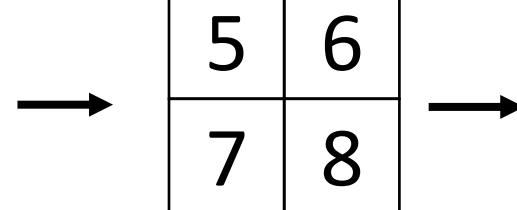
Output: $C \times 4 \times 4$

Use **three** closest neighbors in x and y to
construct **cubic** approximations
(This is how we normally resize images!)

In-Network Upsampling: “Max Unpooling”

Max Pooling: Remember which position had the max

| | | | |
|---|---|---|---|
| 1 | 2 | 6 | 3 |
| 3 | 5 | 2 | 1 |
| 1 | 2 | 2 | 1 |
| 7 | 3 | 4 | 8 |

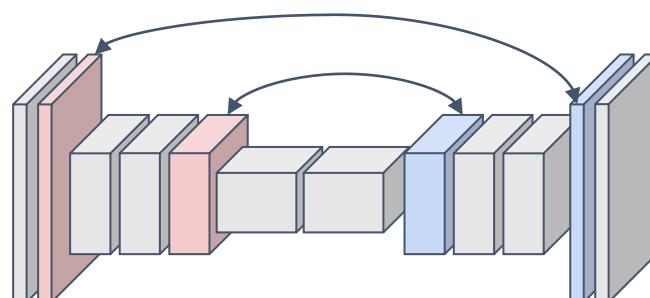


Rest
of
net

| | |
|---|---|
| 1 | 2 |
| 3 | 4 |



| | | | |
|---|---|---|---|
| 0 | 0 | 2 | 0 |
| 0 | 1 | 0 | 0 |
| 0 | 0 | 0 | 0 |
| 3 | 0 | 0 | 4 |

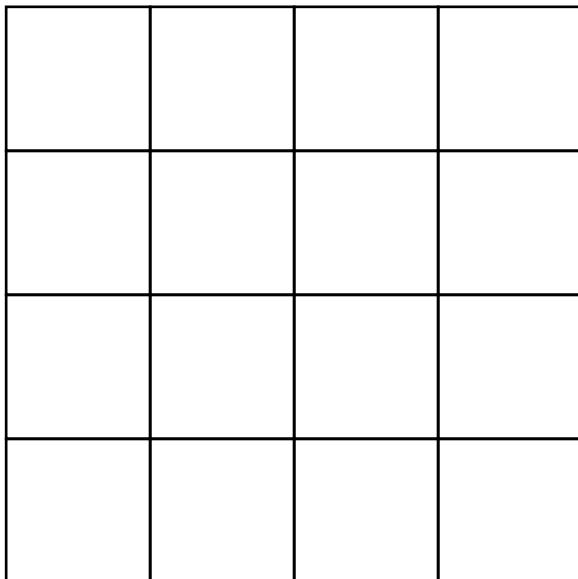


Pair each downsampling layer with an upsampling layer

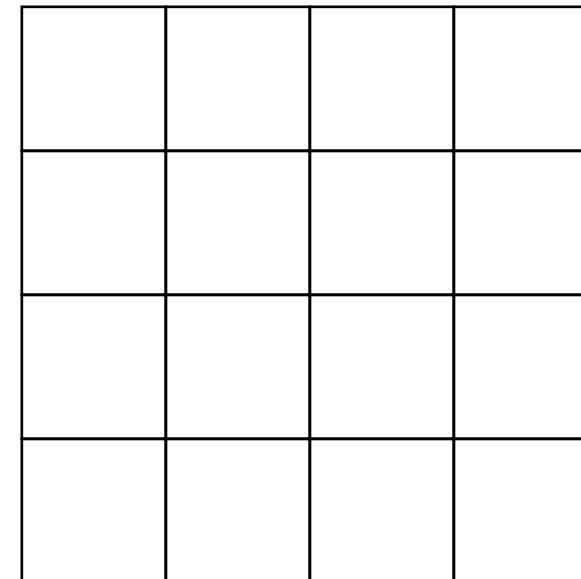
Noh et al, “Learning Deconvolution Network for Semantic Segmentation”, ICCV 2015

Learnable Upsampling: Transposed Convolution

Recall: Normal 3×3 convolution, stride 1, pad 1



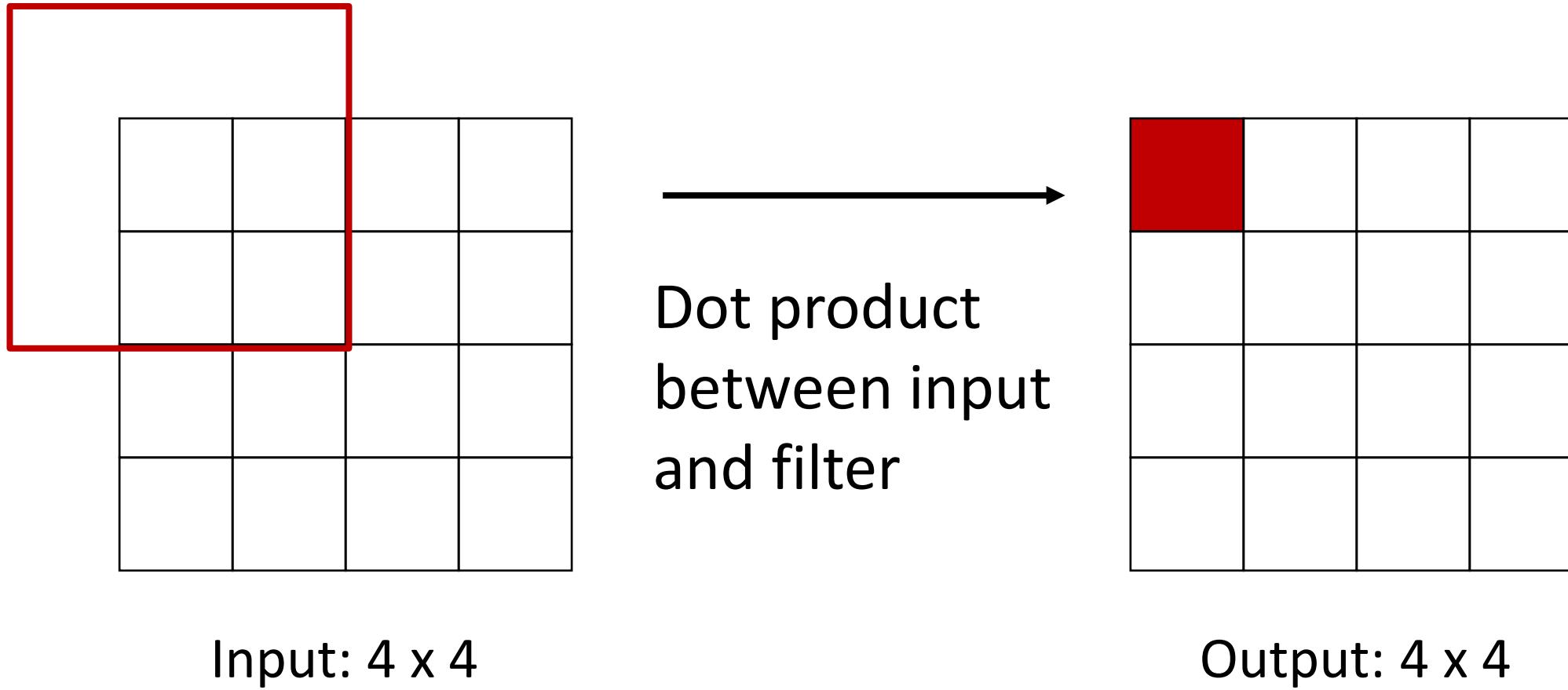
Input: 4×4



Output: 4×4

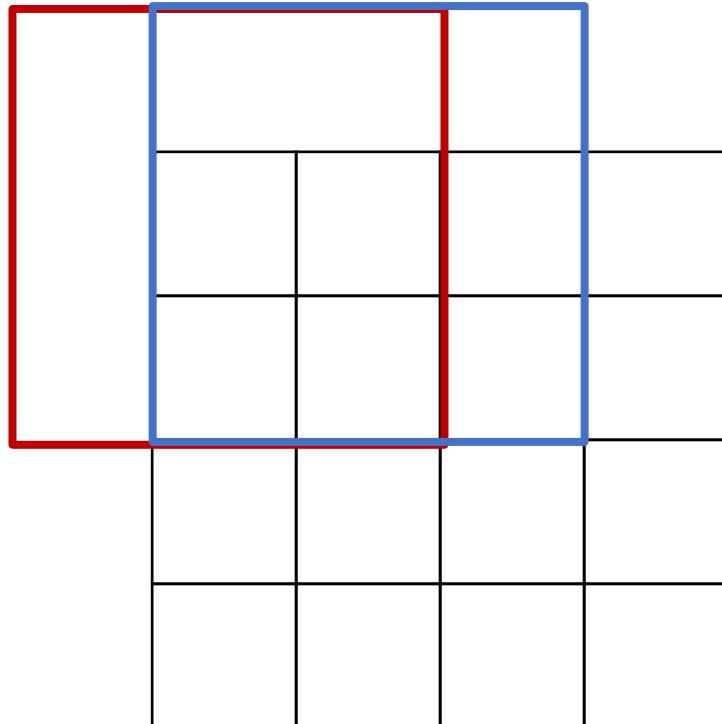
Learnable Upsampling: Transposed Convolution

Recall: Normal 3×3 convolution, stride 1, pad 1



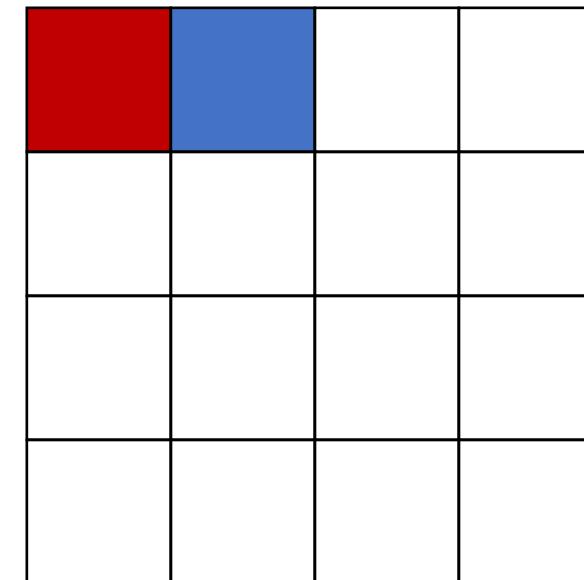
Learnable Upsampling: Transposed Convolution

Recall: Normal 3×3 convolution, stride 1, pad 1



Input: 4×4

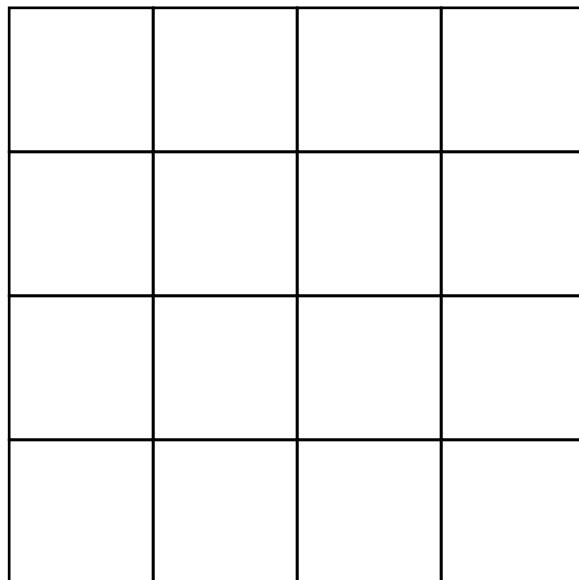
Dot product
between input
and filter



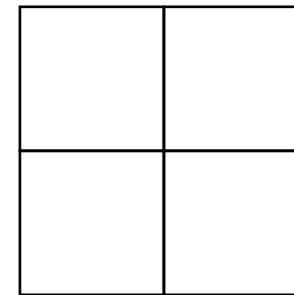
Output: 4×4

Learnable Upsampling: Transposed Convolution

Recall: Normal 3×3 convolution, stride 2, pad 1



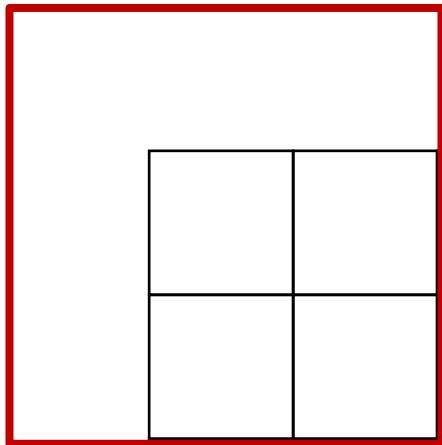
Input: 4×4



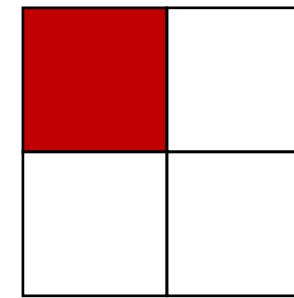
Output: 2×2

Learnable Upsampling: Transposed Convolution

Recall: Normal 3×3 convolution, stride 2, pad 1



Dot product
between input
and filter

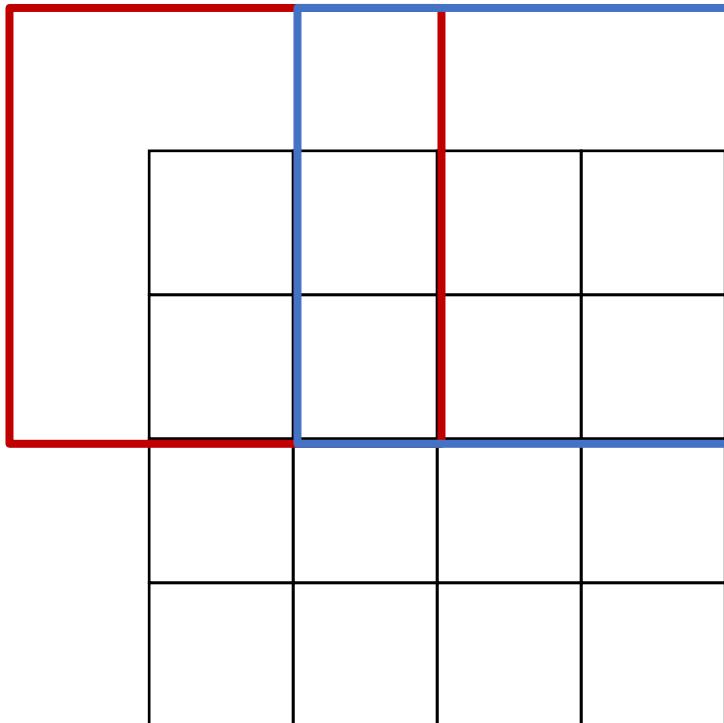


Input: 4×4

Output: 2×2

Learnable Upsampling: Transposed Convolution

Recall: Normal 3×3 convolution, stride 2, pad 1

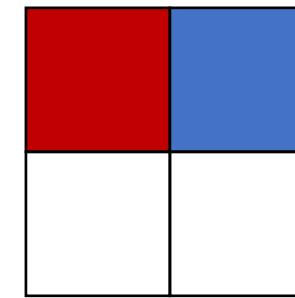


Input: 4×4

Convolution with stride > 1 is “Learnable Downsampling”
Can we use stride < 1 for “Learnable Upsampling”?



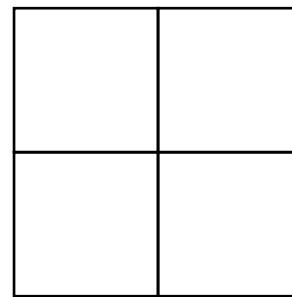
Dot product
between input
and filter



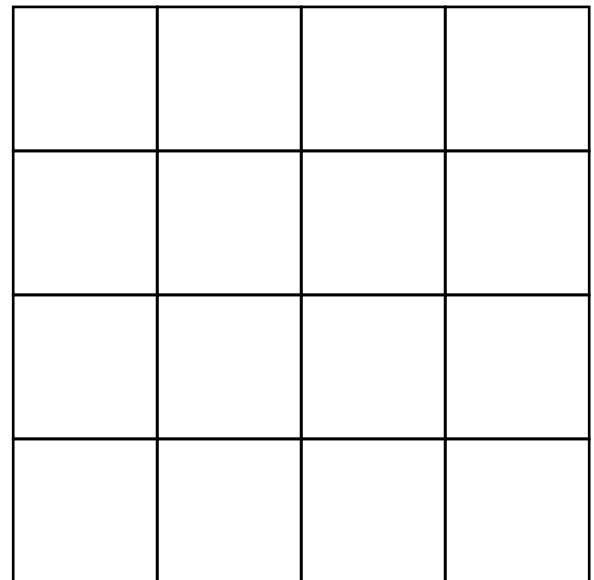
Output: 2×2

Learnable Upsampling: Transposed Convolution

3 x 3 convolution transpose, stride 2



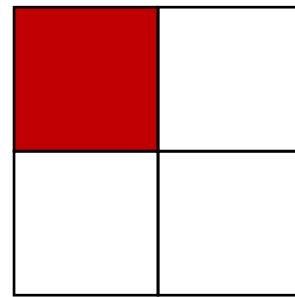
Input: 2 x 2



Output: 4 x 4

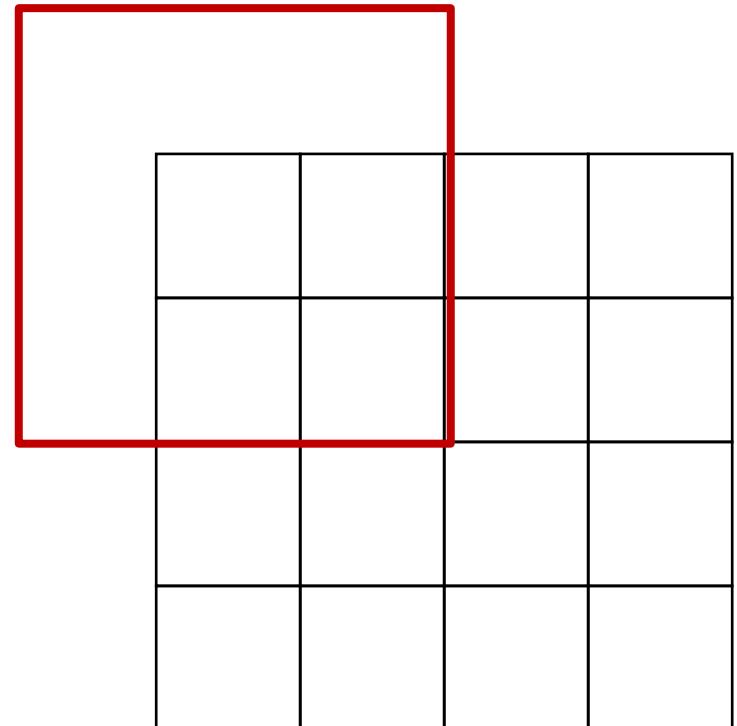
Learnable Upsampling: Transposed Convolution

3 x 3 convolution transpose, stride 2



Input: 2 x 2

→
Weight filter by
input value and
copy to output

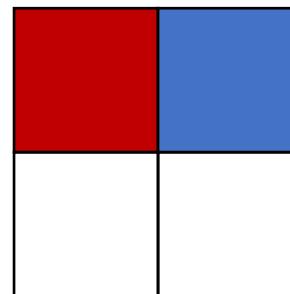


Output: 4 x 4

Learnable Upsampling: Transposed Convolution

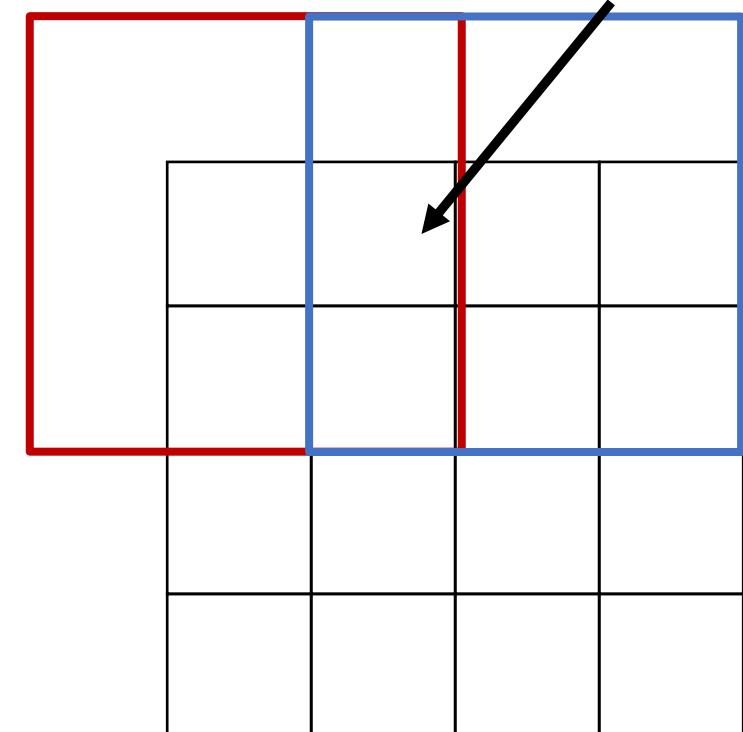
3 x 3 convolution transpose, stride 2

Filter moves 2 pixels in output
for every 1 pixel in input



Input: 2 x 2

Weight filter by
input value and
copy to output

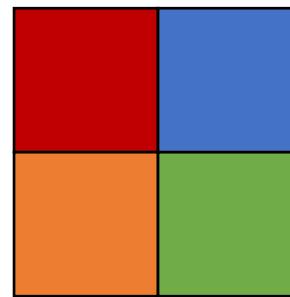


Output: 4 x 4

Learnable Upsampling: Transposed Convolution

3 x 3 convolution transpose, stride 2

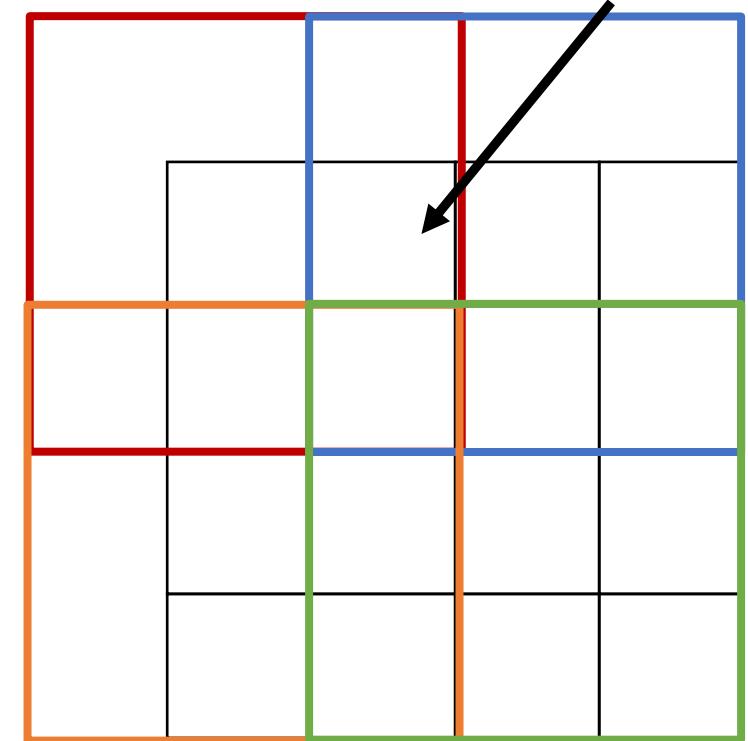
This gives 5x5 output – need to trim one pixel from top and left to give 4x4 output



Input: 2 x 2

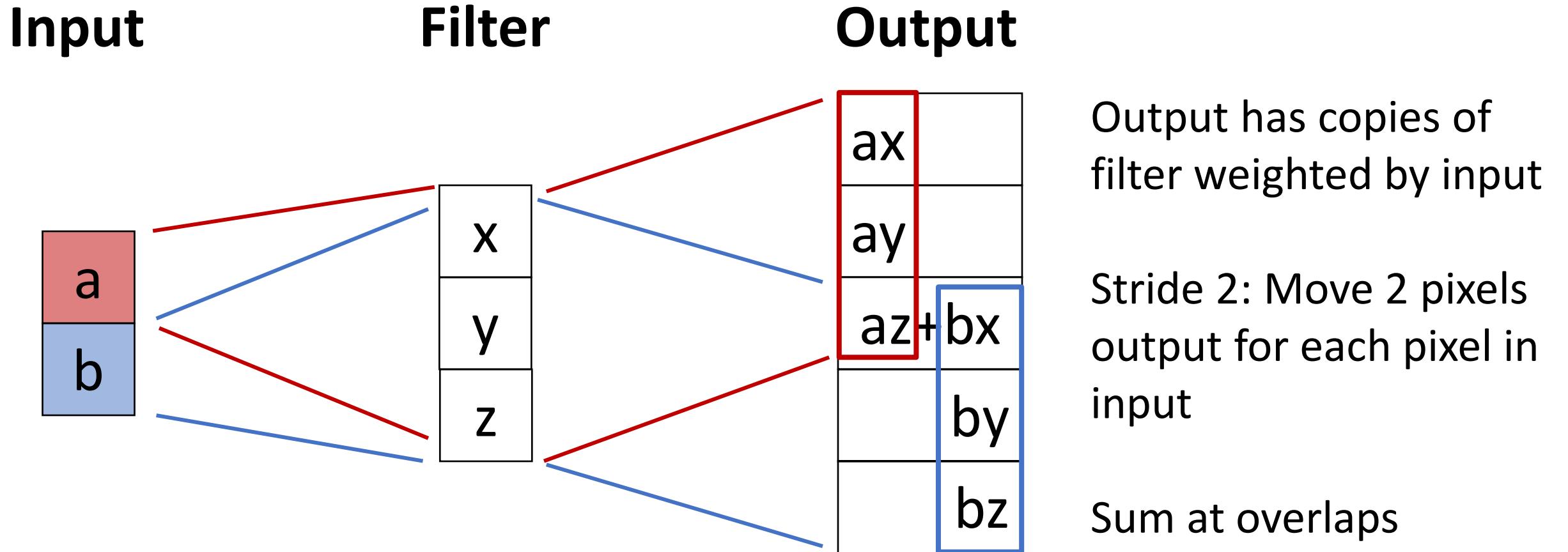
Weight filter by
input value and
copy to output

Sum where
output overlaps



Output: 4 x 4

Transposed Convolution: 1D example



Transposed Convolution: 1D example

Input

$$\begin{matrix} a \\ b \end{matrix}$$

Filter

$$\begin{matrix} x \\ y \\ z \end{matrix}$$

Output

$$\begin{matrix} ax \\ ay \\ az+bx \\ by \\ bz \end{matrix}$$

This has many names:

- Deconvolution (bad)!
- Upconvolution
- Fractionally strided convolution
- Backward strided convolution
- Transposed Convolution (best name)

Convolution as Matrix Multiplication (1D Example)

We can express convolution in terms of a matrix multiplication

$$x * a = Xa$$

$$\begin{bmatrix} x & y & z & 0 & 0 & 0 \\ 0 & x & y & z & 0 & 0 \\ 0 & 0 & x & y & z & 0 \\ 0 & 0 & 0 & x & y & z \end{bmatrix} \begin{bmatrix} 0 \\ a \\ b \\ c \\ d \\ 0 \end{bmatrix} = \begin{bmatrix} ay + bz \\ ax + by + cz \\ bx + cy + dz \\ cx + dy \end{bmatrix}$$

Transposed convolution multiplies by the transpose of the same matrix:

$$x *^T a = X^T a$$

$$\begin{bmatrix} x & 0 & 0 & 0 \\ y & x & 0 & 0 \\ z & y & x & 0 \\ 0 & z & y & x \\ 0 & 0 & z & y \\ 0 & 0 & 0 & z \end{bmatrix} \begin{bmatrix} a \\ b \\ c \\ d \end{bmatrix} = \begin{bmatrix} ax \\ ay + bx \\ az + by + cx \\ bz + cy + dx \\ cz + dy \\ dz \end{bmatrix}$$

Example: 1D conv, kernel size=3, stride=1, padding=1

When stride=1, transposed conv is just a regular conv (with different padding rules)

Convolution as Matrix Multiplication (1D Example)

We can express convolution in terms of a matrix multiplication

$$x * a = Xa$$

$$\begin{bmatrix} x & y & z & 0 & 0 & 0 \\ 0 & 0 & x & y & z & 0 \end{bmatrix} \begin{bmatrix} 0 \\ a \\ b \\ c \\ d \\ 0 \end{bmatrix} = \begin{bmatrix} ay + bz \\ bx + cy + dz \end{bmatrix}$$

Example: 1D conv, kernel size=3, stride=2, padding=1

Transposed convolution multiplies by the transpose of the same matrix:

$$x *^T a = X^T a$$

$$\begin{bmatrix} x & 0 \\ y & 0 \\ z & x \\ 0 & y \\ 0 & z \\ 0 & 0 \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} = \begin{bmatrix} ax \\ ay \\ az + bx \\ by \\ bz \\ 0 \end{bmatrix}$$

When stride>1, transposed convolution cannot be expressed as normal conv

Recall: PyTorch Convolution Layer

```
CLASS torch.nn.Conv2d(in_channels, out_channels, kernel_size, stride=1, padding=0, dilation=1,  
groups=1, bias=True, padding_mode='zeros', device=None, dtype=None) [SOURCE]
```

Applies a 2D convolution over an input signal composed of several input planes.

In the simplest case, the output value of the layer with input size (N, C_{in}, H, W) and output $(N, C_{\text{out}}, H_{\text{out}}, W_{\text{out}})$ can be precisely described as:

$$\text{out}(N_i, C_{\text{out}_j}) = \text{bias}(C_{\text{out}_j}) + \sum_{k=0}^{C_{\text{in}}-1} \text{weight}(C_{\text{out}_j}, k) \star \text{input}(N_i, k)$$

where \star is the valid 2D cross-correlation operator, N is a batch size, C denotes a number of channels, H is a height of input planes in pixels, and W is width in pixels.

Recall: PyTorch Convolution Layers

```
CLASS torch.nn.Conv2d(in_channels, out_channels, kernel_size, stride=1, padding=0, dilation=1,  
    groups=1, bias=True, padding_mode='zeros', device=None, dtype=None) [SOURCE]
```

```
CLASS torch.nn.Conv1d(in_channels, out_channels, kernel_size, stride=1, padding=0, dilation=1,  
    groups=1, bias=True, padding_mode='zeros', device=None, dtype=None) [SOURCE]
```

```
CLASS torch.nn.Conv3d(in_channels, out_channels, kernel_size, stride=1, padding=0, dilation=1,  
    groups=1, bias=True, padding_mode='zeros', device=None, dtype=None) [SOURCE]
```

PyTorch Transposed Convolution Layers

CLASS `torch.nn.ConvTranspose2d(in_channels, out_channels, kernel_size, stride=1, padding=0, output_padding=0, groups=1, bias=True, dilation=1, padding_mode='zeros', device=None, dtype=None)` [\[SOURCE\]](#)

CLASS `torch.nn.ConvTranspose1d(in_channels, out_channels, kernel_size, stride=1, padding=0, output_padding=0, groups=1, bias=True, dilation=1, padding_mode='zeros', device=None, dtype=None)` [\[SOURCE\]](#)

CLASS `torch.nn.ConvTranspose3d(in_channels, out_channels, kernel_size, stride=1, padding=0, output_padding=0, groups=1, bias=True, dilation=1, padding_mode='zeros', device=None, dtype=None)` [\[SOURCE\]](#)

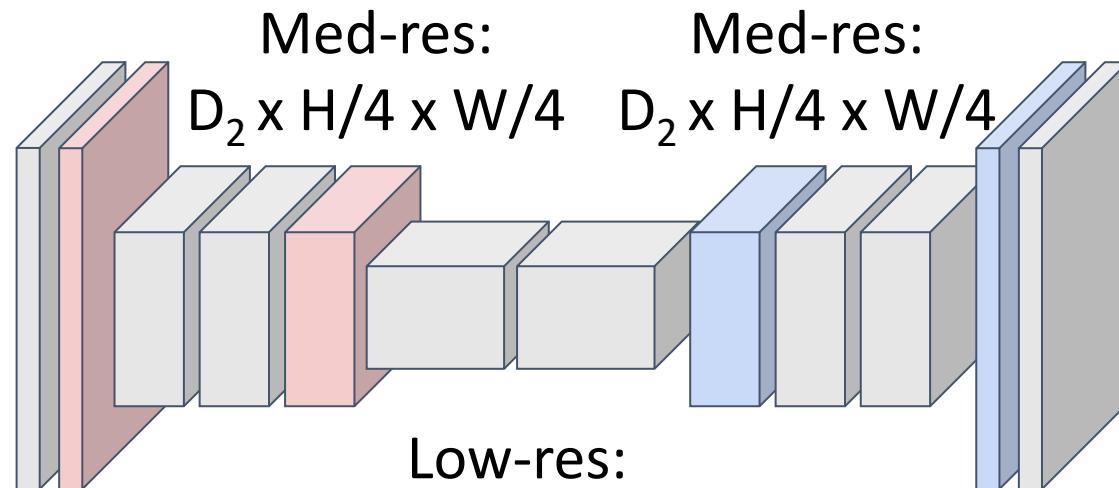
Semantic Segmentation: Fully Convolutional Network

Downsampling:
Pooling, strided
convolution



Input:
 $3 \times H \times W$

High-res:
 $D_1 \times H/2 \times W/2$



Design network as a bunch of convolutional layers, with
downsampling and **upsampling** inside the network!

Upsampling:
interpolation,
transposed conv



Predictions:
 $H \times W$

Loss function: Per-Pixel cross-entropy

Computer Vision Task: Semantic Segmentation

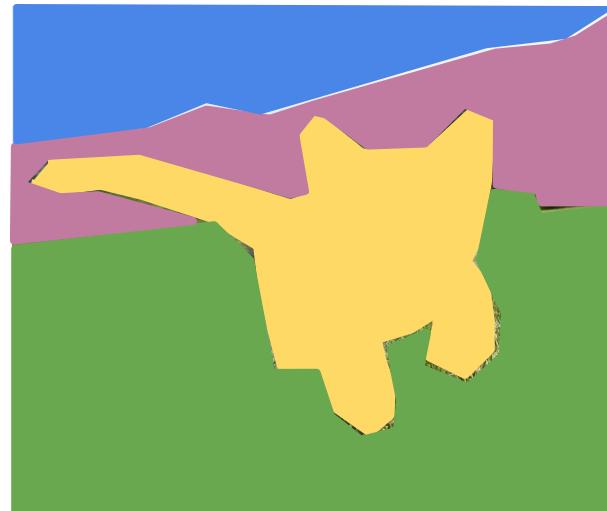
Classification



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](#)

Computer Vision Task: Object Detection

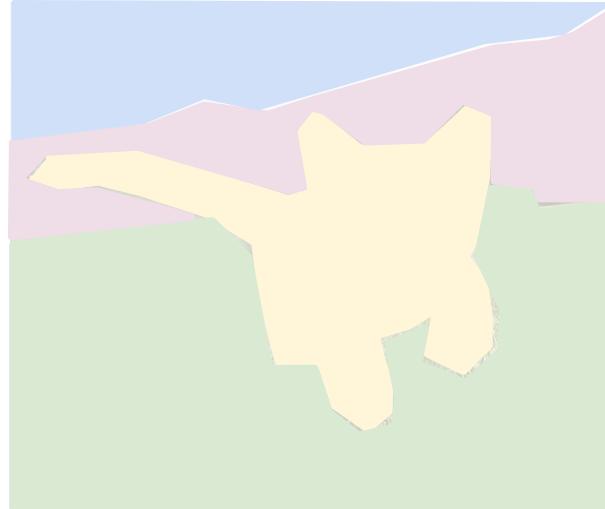
Classification



CAT

No spatial extent

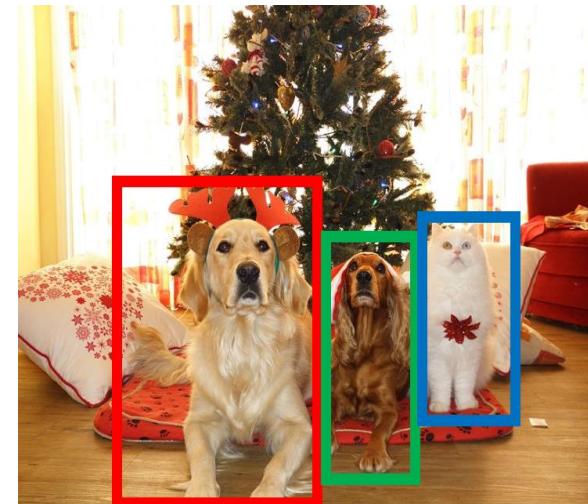
Semantic
Segmentation



GRASS, CAT, TREE,
SKY

No objects, just pixels

Object
Detection



DOG, DOG, CAT

Multiple Objects

Instance
Segmentation



DOG, DOG, CAT

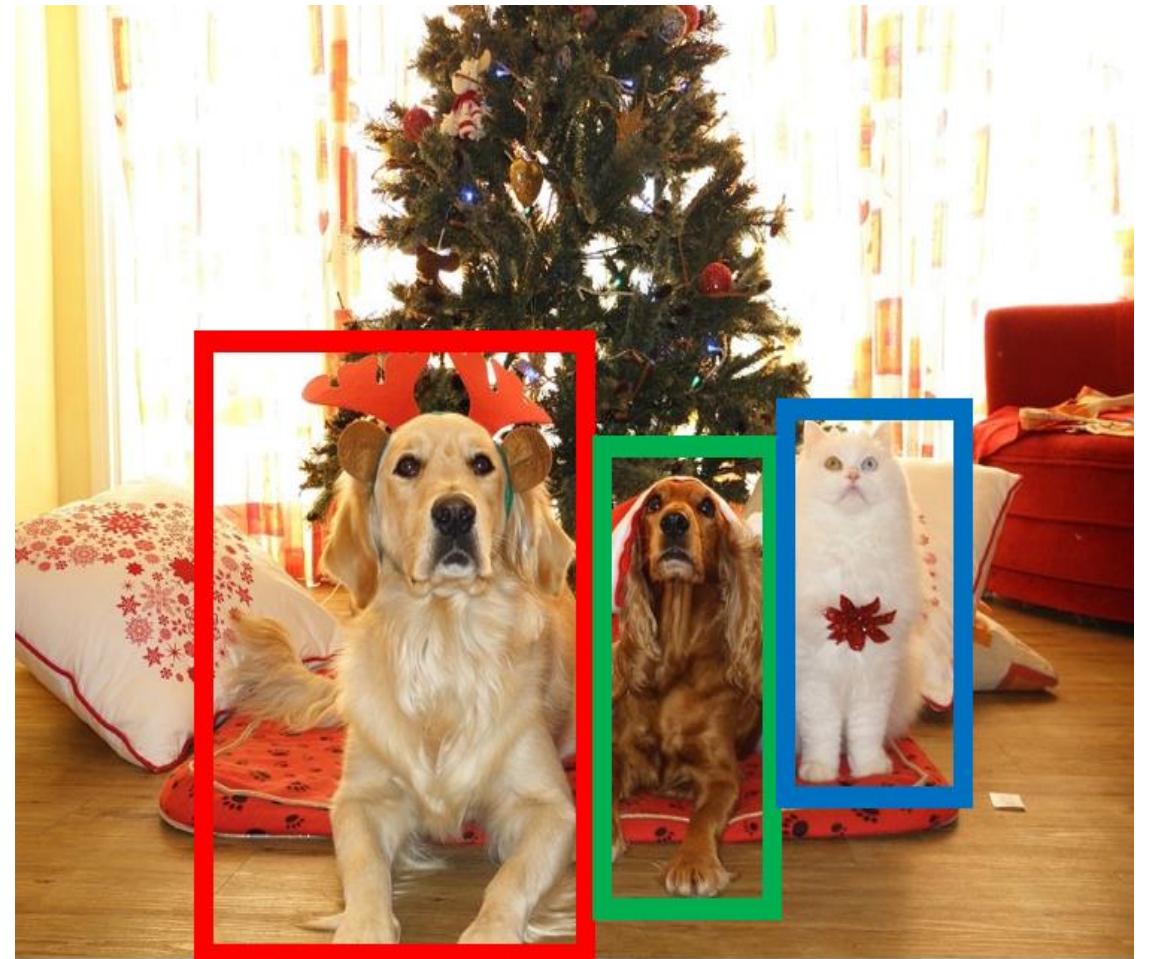
[This image is CC0 public domain](#)

Object Detection: Task Definition

Input: Single RGB Image

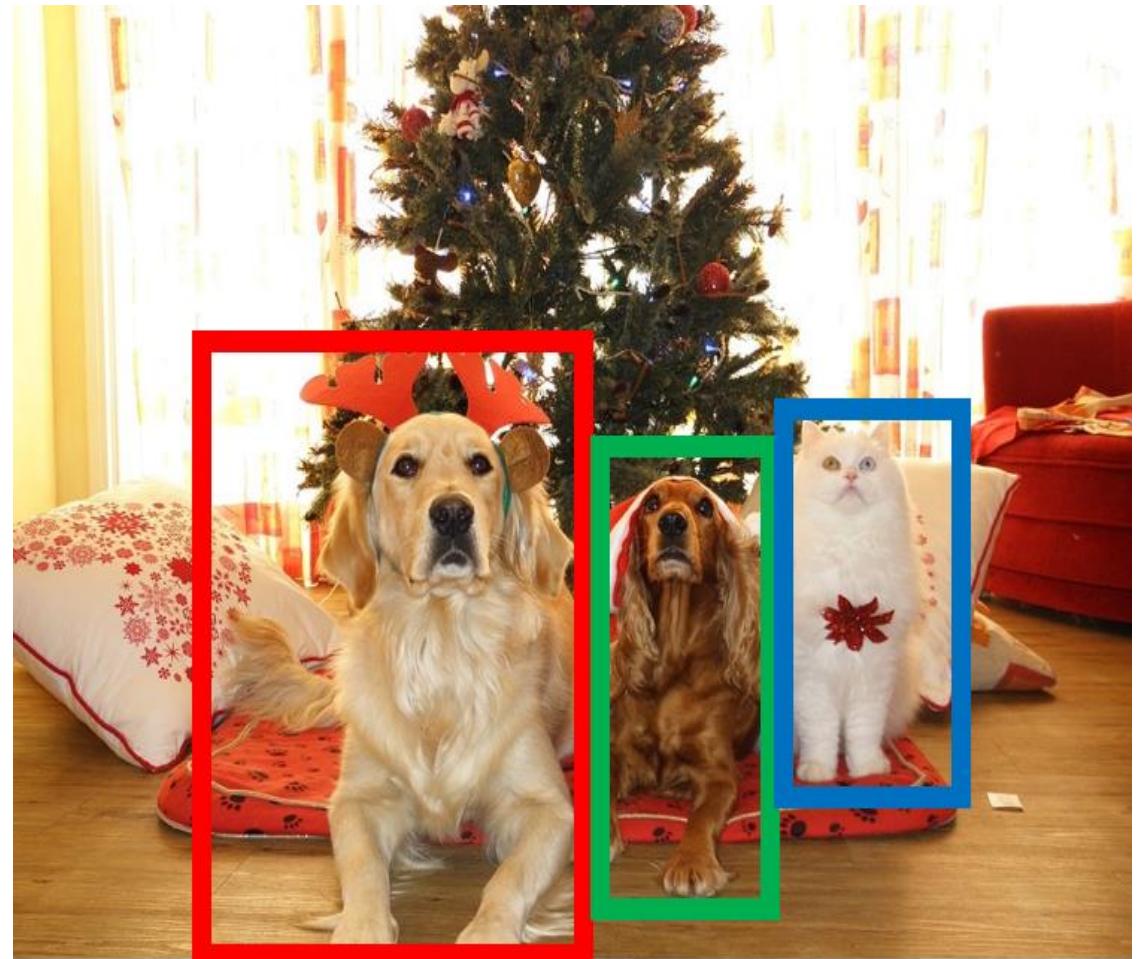
Output: A set of detected objects;
For each object, predict:

1. Category label (from fixed, known set of categories)
2. Bounding box (four numbers: x, y, width, height)



Object Detection: Challenges

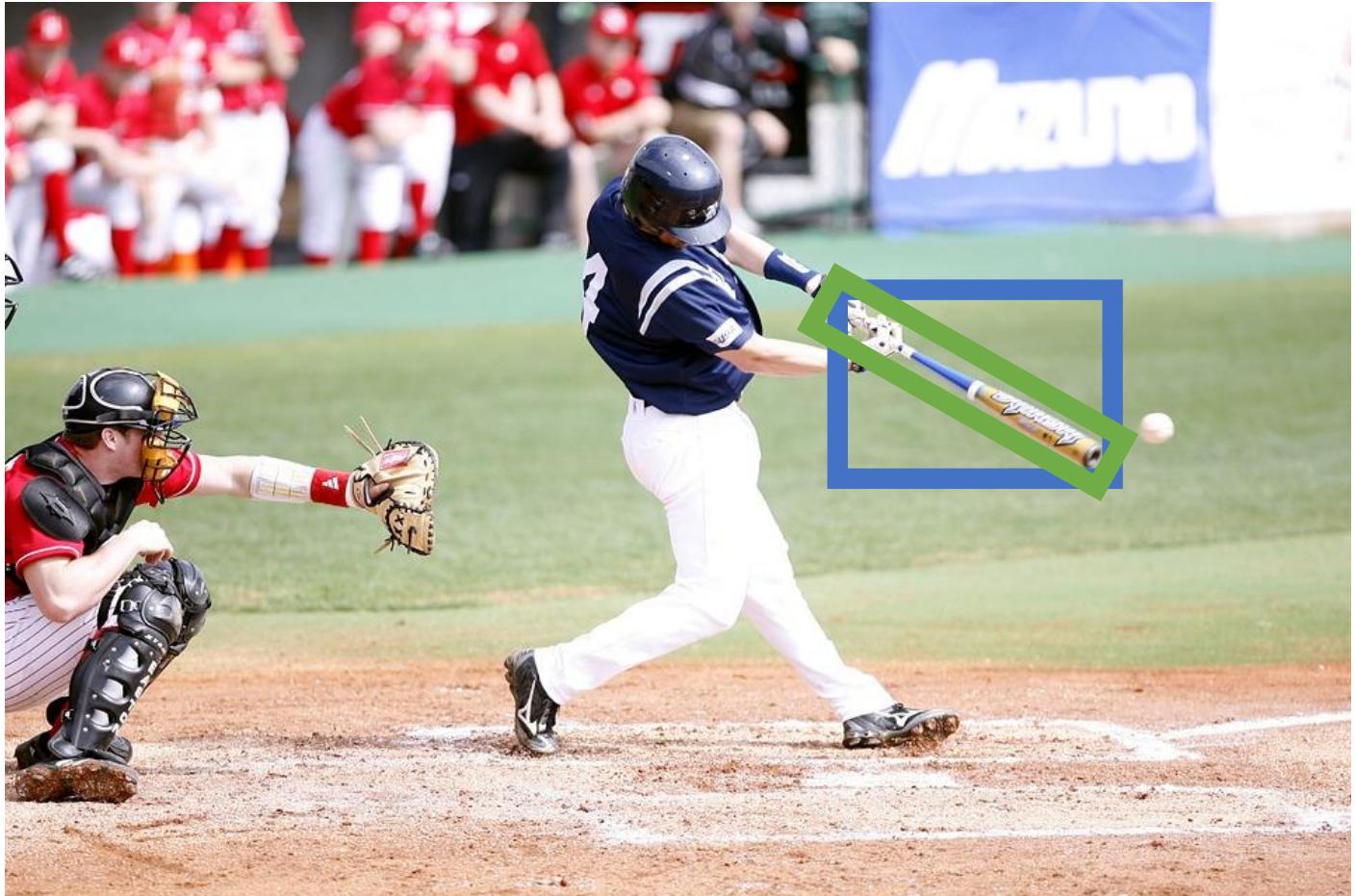
- **Multiple outputs:** Need to output variable numbers of objects per image
- **Multiple types of output:** Need to predict “what” (category label) as well as “where” (bounding box)
- **Large images:** Classification works at 224x224; need higher resolution for detection, often ~800x600



Bounding Boxes

Bounding boxes are typically *axis-aligned*

Oriented boxes are much less common

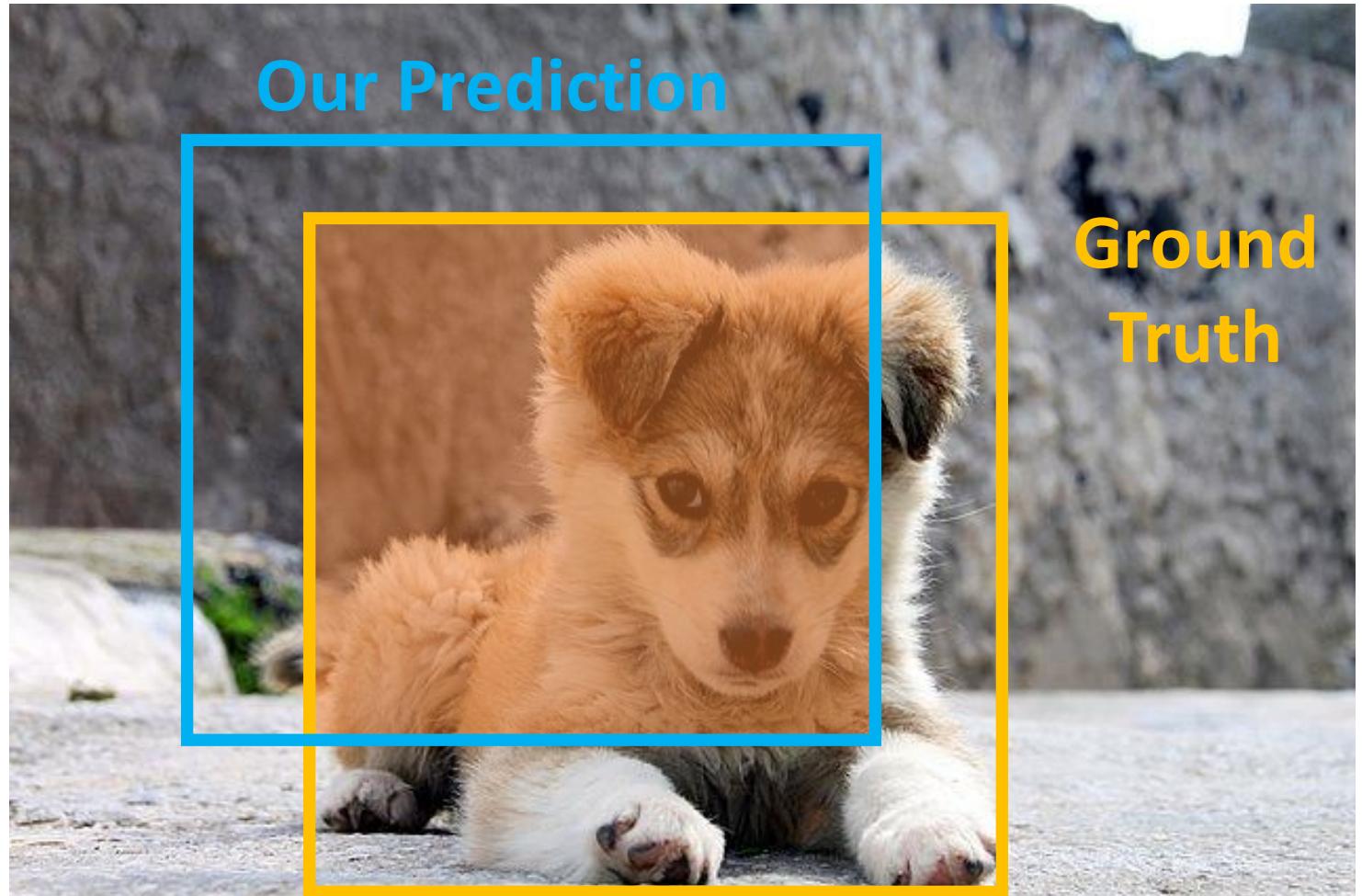


Comparing Boxes: Intersection over Union (IoU)

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

Intersection over Union (IoU)
(Also called “Jaccard similarity” or
“Jaccard index”):

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



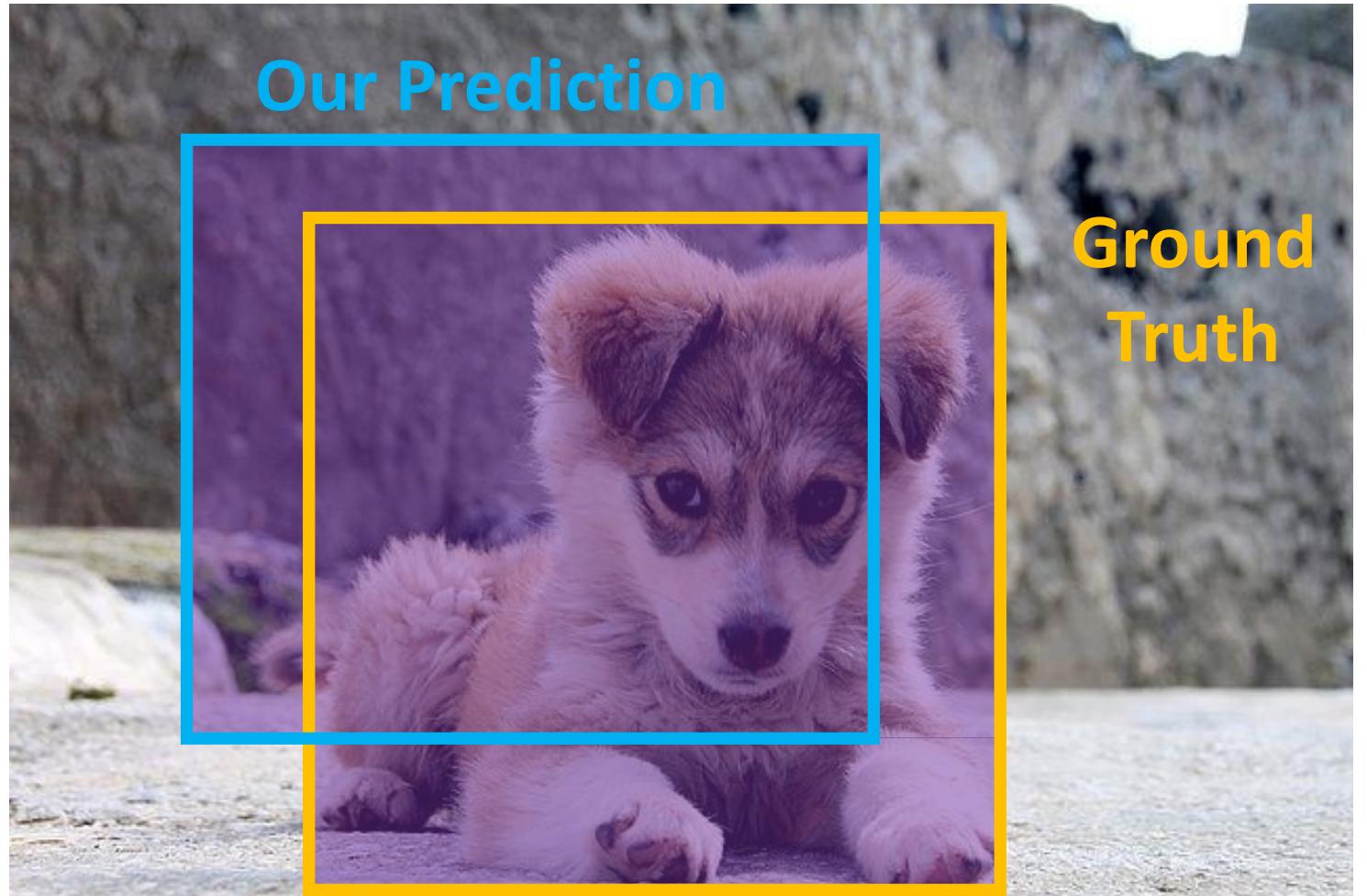
[Puppy image](#) is licensed under [CC-A 2.0 Generic license](#). Bounding boxes and text added by Justin Johnson.

Comparing Boxes: Intersection over Union (IoU)

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

Intersection over Union (IoU)
(Also called “Jaccard similarity” or
“Jaccard index”):

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



[Puppy image](#) is licensed under [CC-A 2.0 Generic license](#). Bounding boxes and text added by Justin Johnson.

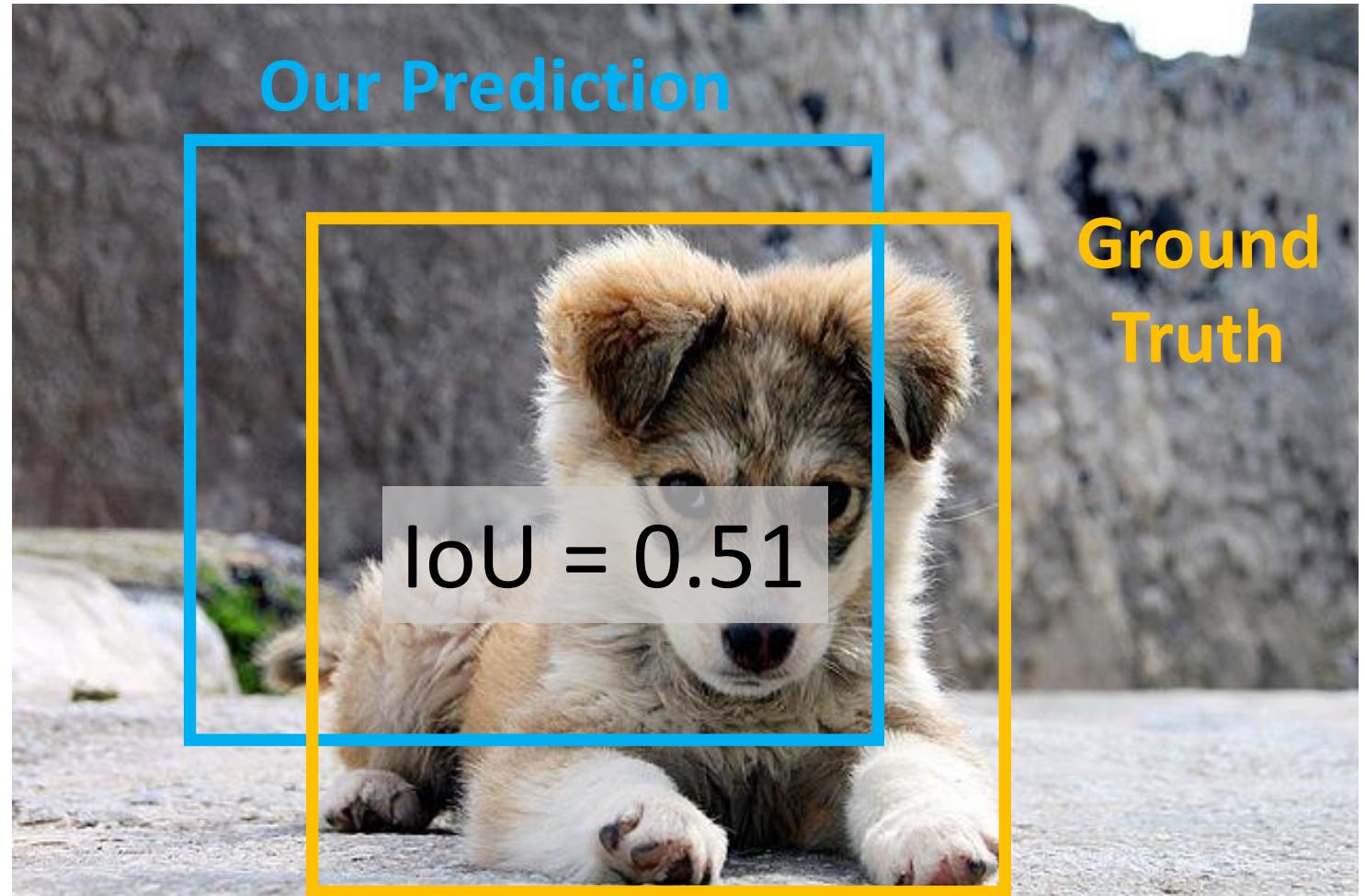
Comparing Boxes: Intersection over Union (IoU)

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

Intersection over Union (IoU)
(Also called “Jaccard similarity” or
“Jaccard index”):

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

IoU > 0.5 is “decent”



[Puppy image](#) is licensed under [CC-A 2.0 Generic license](#). Bounding boxes and text added by Justin Johnson.

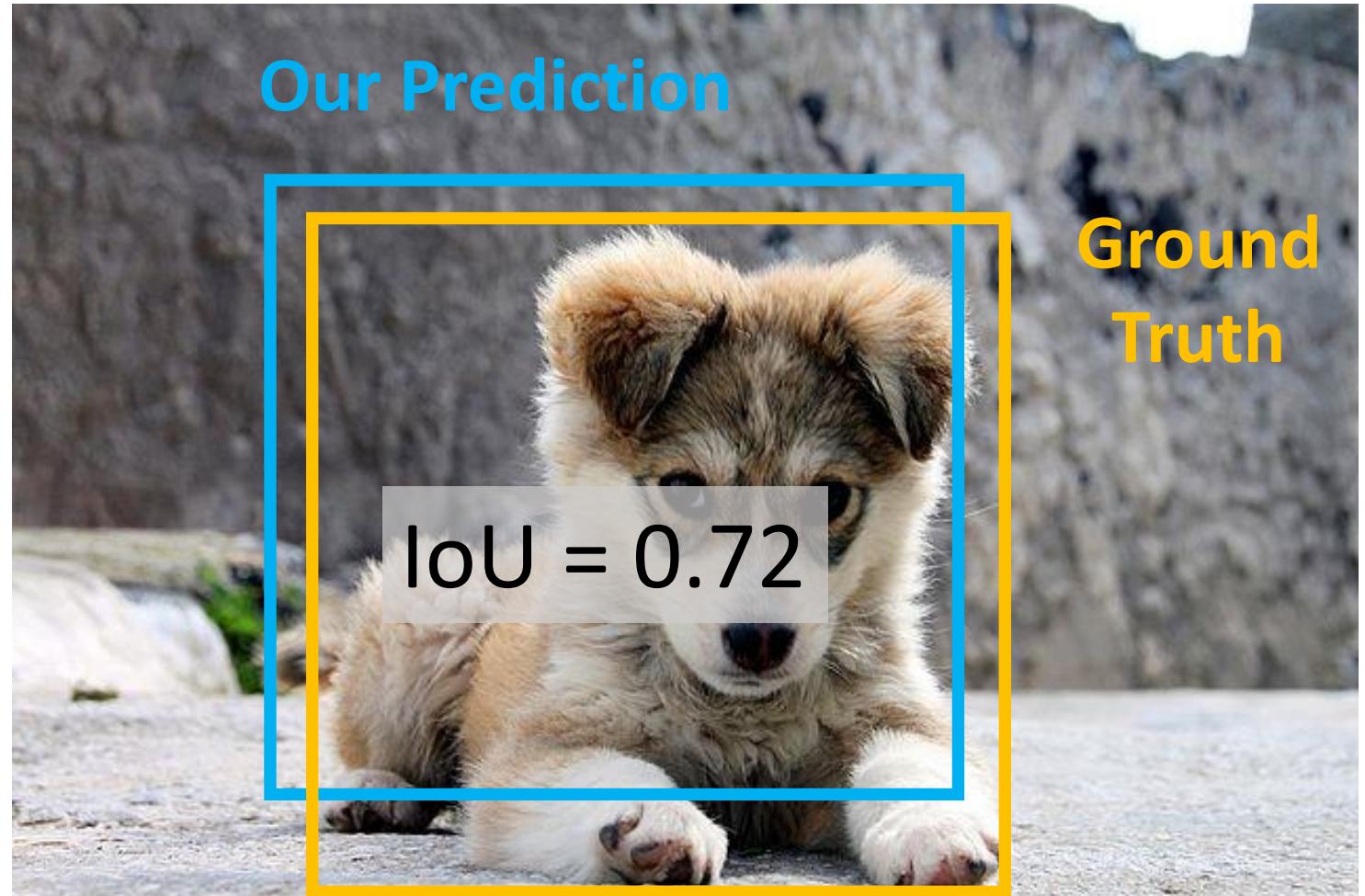
Comparing Boxes: Intersection over Union (IoU)

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

Intersection over Union (IoU)
(Also called “Jaccard similarity” or
“Jaccard index”):

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

IoU > 0.5 is “decent”,
IoU > 0.7 is “pretty good”,



[Puppy image](#) is licensed under [CC-A 2.0 Generic license](#). Bounding boxes and text added by Justin Johnson.

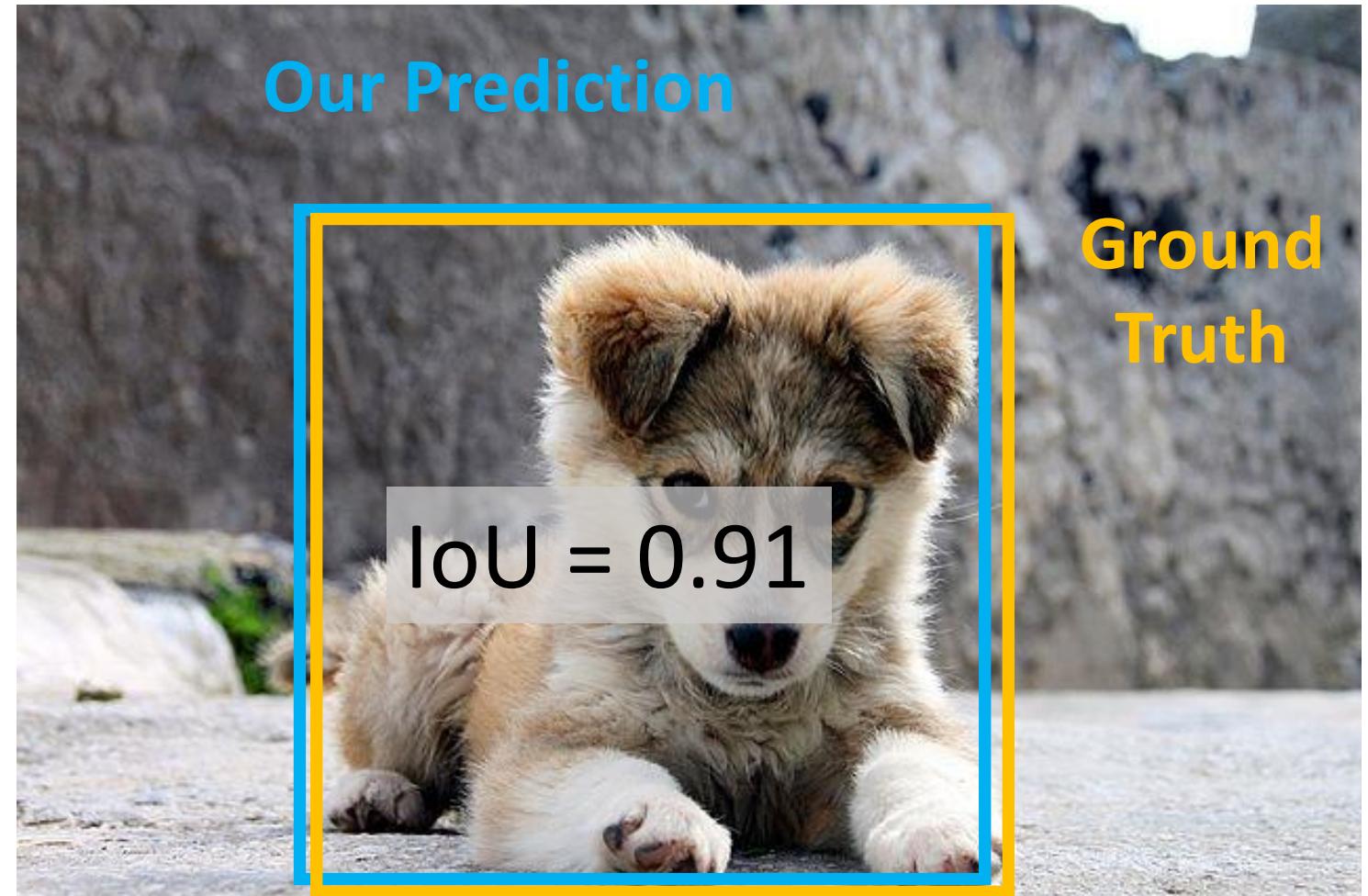
Comparing Boxes: Intersection over Union (IoU)

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

Intersection over Union (IoU)
(Also called “Jaccard similarity” or
“Jaccard index”):

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

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

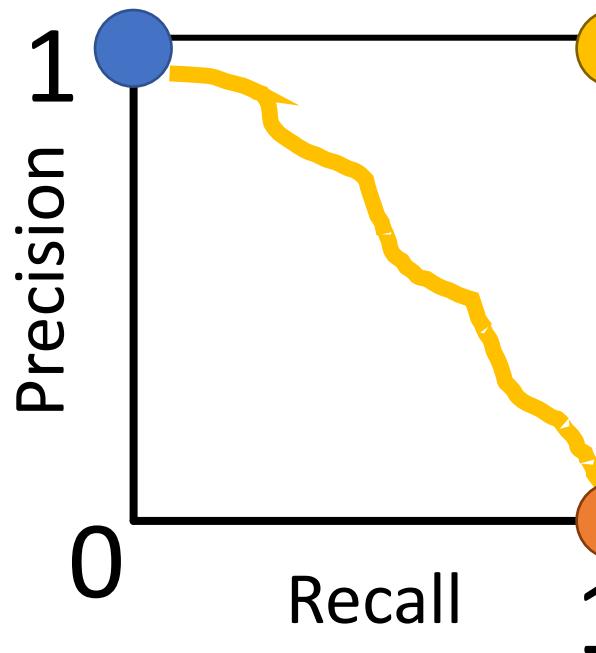


[Puppy image](#) is licensed under [CC-A 2.0 Generic license](#). Bounding boxes and text added by Justin Johnson.

Evaluating Detectors

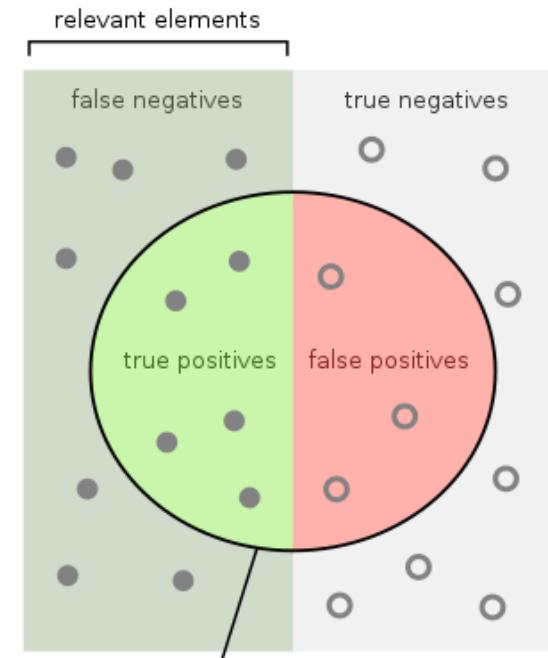
- True detection: high IoU
- Precision: # true detections / # all detections
- Recall: # true detections / # ground truths

Reject everything: no mistakes



**Summarize by area under curve
(avg. precision)**

**Accept everything:
Miss nothing**



$$\text{Precision} = \frac{\text{How many relevant items are retrieved?}}{\text{How many retrieved items are relevant?}}$$
$$\text{Recall} = \frac{\text{How many relevant items are retrieved?}}{\text{How many relevant items are there?}}$$

Image source: [wikipedia](#)

Object Detection Progress

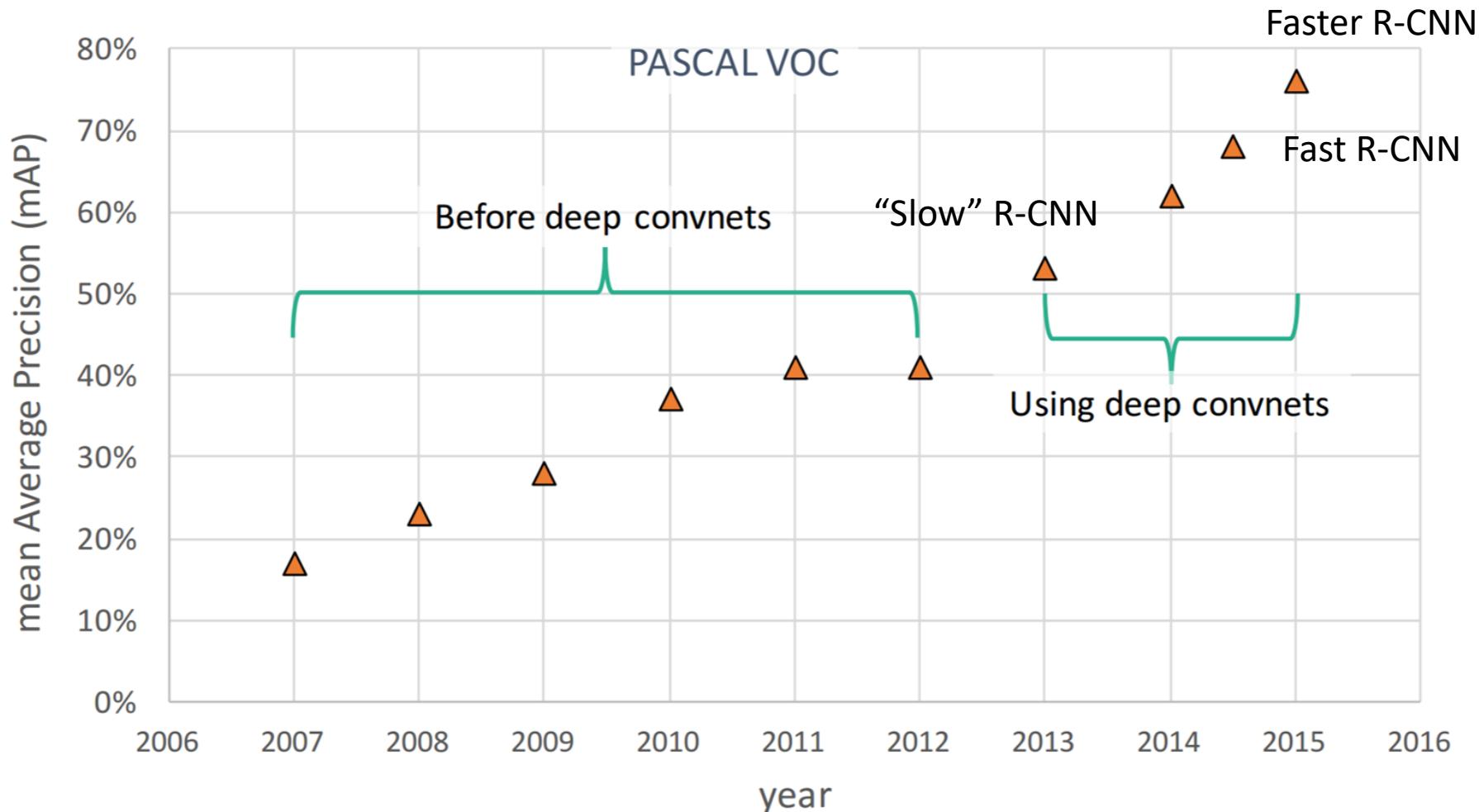
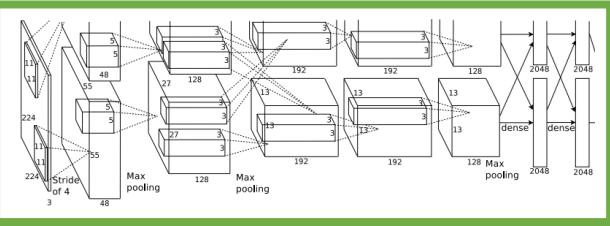


Figure copyright Ross Girshick, 2015.
Reproduced with permission.

Detecting a single object

Often pretrained
on ImageNet
(Transfer learning)



Vector:

4096

Treat localization as a
regression problem!

This image is CC0 public domain

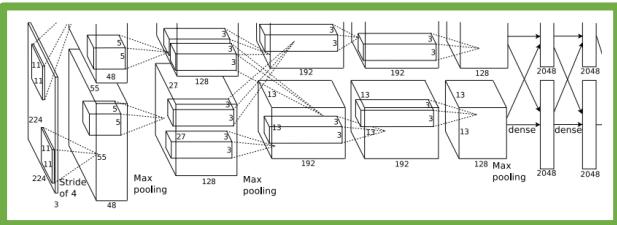
Detecting a single object “What”



This image is CC0 public domain

Treat localization as a
regression problem!

Often pretrained
on ImageNet
(Transfer learning)



Vector:

4096

Fully
Connected:
4096 to 1000

Class Scores

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

Correct label:
Cat

Softmax
Loss

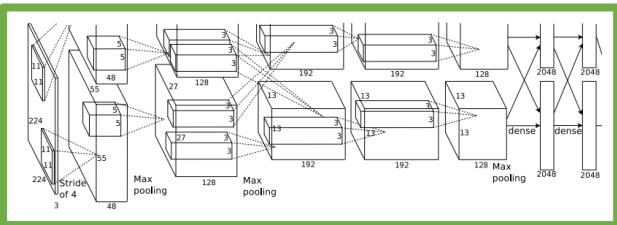
Detecting a single object “What”



This image is CC0 public domain

Treat localization as a
regression problem!

Often pretrained
on ImageNet
(Transfer learning)



Vector:
4096

Fully
Connected:
4096 to 1000

Class Scores

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

Correct label:
Cat

Softmax
Loss

“Where”

Fully
Connected:
4096 to 4

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

L2 Loss

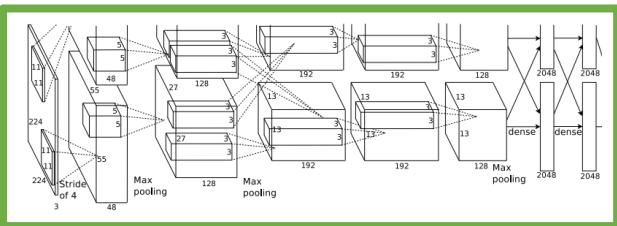
Correct box:
(x', y', w', h')

Detecting a single object “What”



This image is CC0 public domain

Often pretrained
on ImageNet
(Transfer learning)



Vector:
4096

Fully
Connected:
4096 to 1000

Class Scores

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

Correct label:
Cat

Softmax
Loss

Multitask
Loss

Weighted
Sum

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

Loss

L2 Loss

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

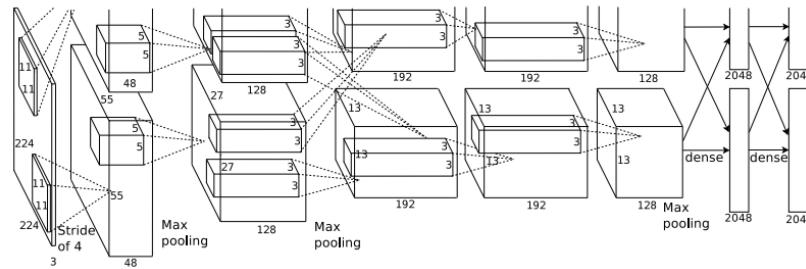
Treat localization as a
regression problem!

Problem: Images can have
more than one object!

“Where”

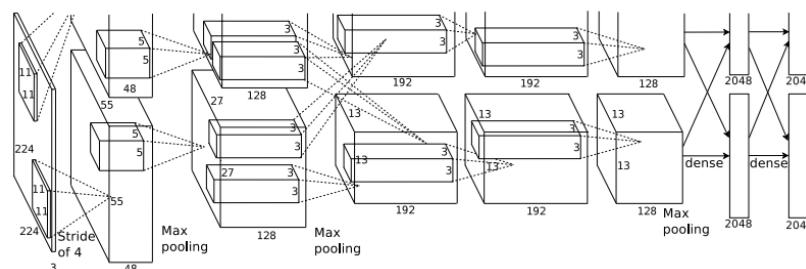
Detecting Multiple Objects

Need different numbers
of outputs per image



CAT: (x, y, w, h)

4 numbers

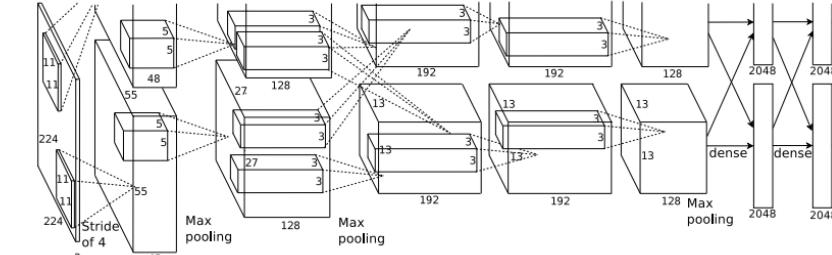


DOG: (x, y, w, h)

12 numbers

DOG: (x, y, w, h)

CAT: (x, y, w, h)



DUCK: (x, y, w, h)

Many
numbers!

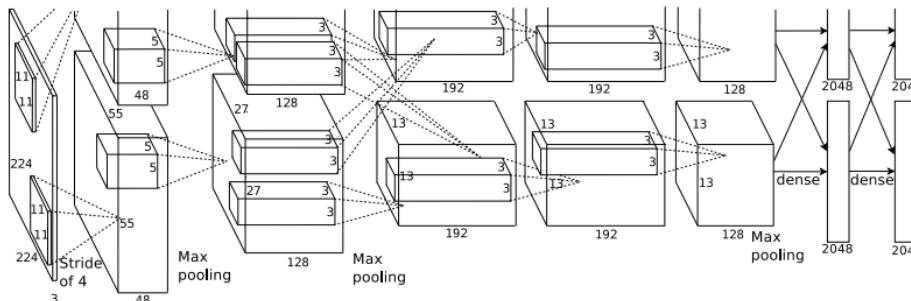
....

[Duck image](#) is free to use under the [Pixabay license](#)

Detecting Multiple Objects: Sliding Window



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

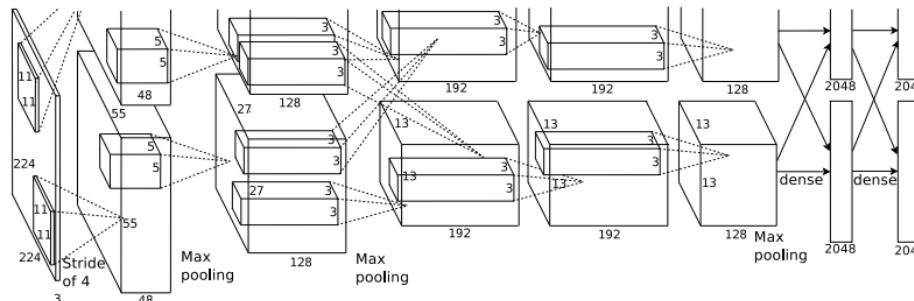


Dog? NO
Cat? NO
Background? YES

Detecting Multiple Objects: Sliding Window



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

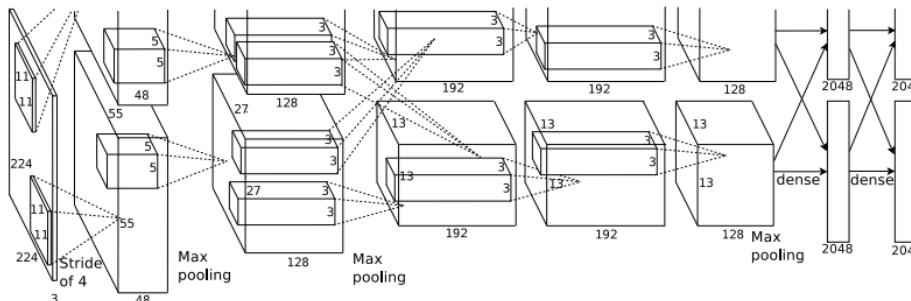


Dog? YES
Cat? NO
Background? NO

Detecting Multiple Objects: Sliding Window



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

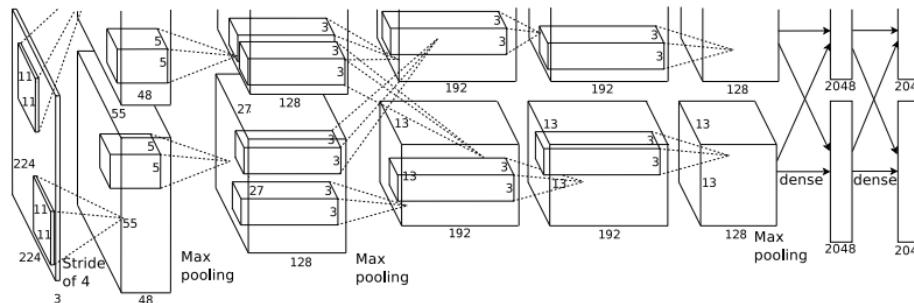


Dog? YES
Cat? NO
Background? NO

Detecting Multiple Objects: Sliding Window



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



Dog? NO
Cat? YES
Background? NO

Detecting Multiple Objects: Sliding Window



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

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

Consider a box of size $h \times w$:

Possible x positions: $W - w + 1$

Possible y positions: $H - h + 1$

Possible positions:

$(W - w + 1) * (H - h + 1)$

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

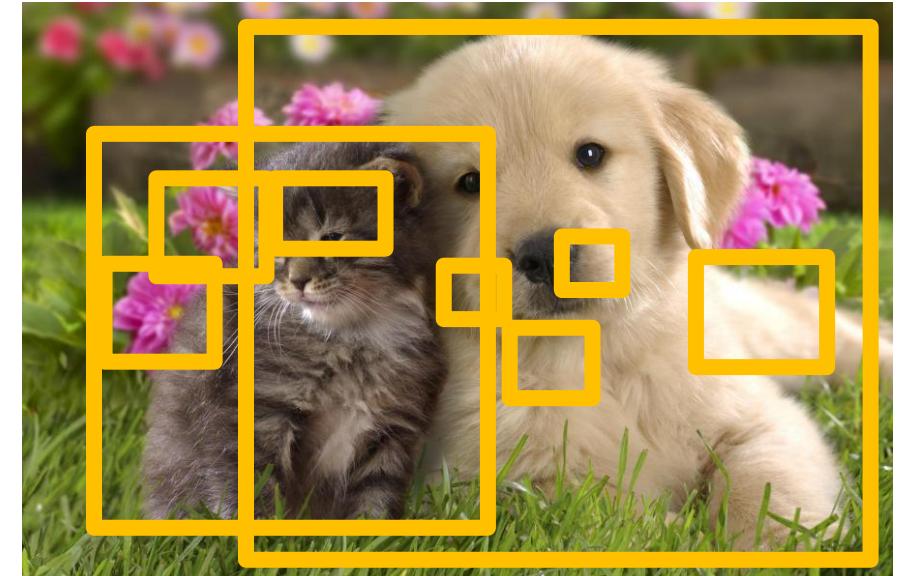
Total possible boxes:

$$\sum_{h=1}^H \sum_{w=1}^W (W - w + 1)(H - h + 1)$$

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

Region Proposals

- Find a small set of boxes that are likely to cover all objects
- Often based on heuristics: e.g., look for “blob-like” image regions
- Relatively fast to run; e.g., Selective Search gives 2000 region proposals in a few seconds on CPU



Alexe et al, “Measuring the objectness of image windows”, TPAMI 2012

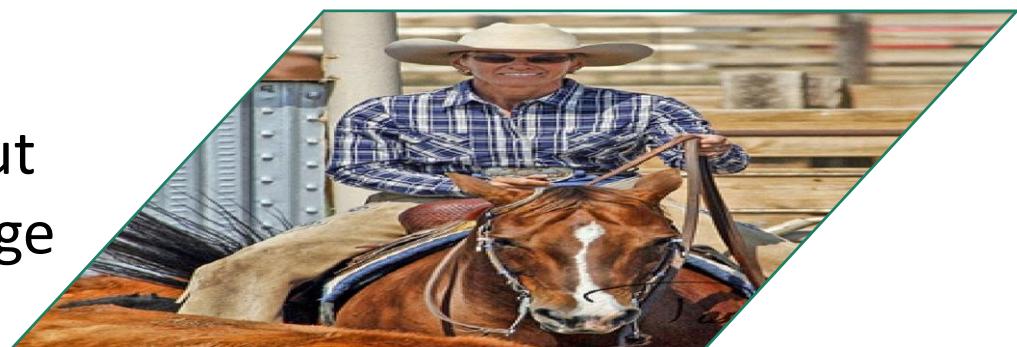
Uijlings et al, “Selective Search for Object Recognition”, IJCV 2013

Cheng et al, “BING: Binarized normed gradients for objectness estimation at 300fps”, CVPR 2014

Zitnick and Dollar, “Edge boxes: Locating object proposals from edges”, ECCV 2014

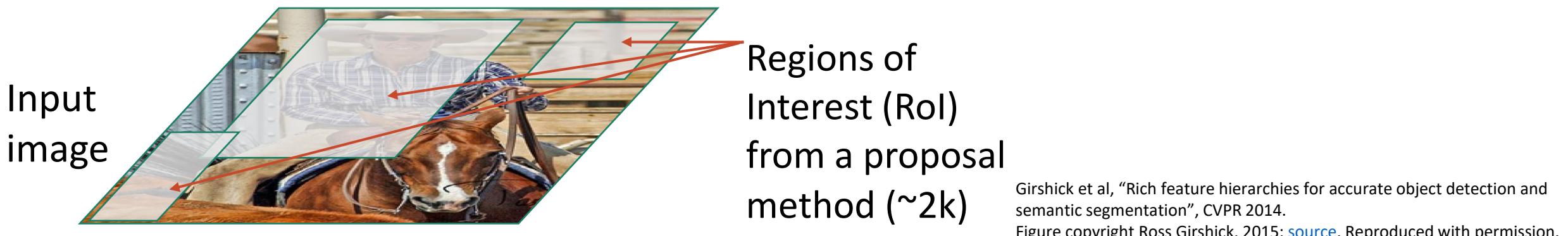
R-CNN: Region-Based CNN

Input
image

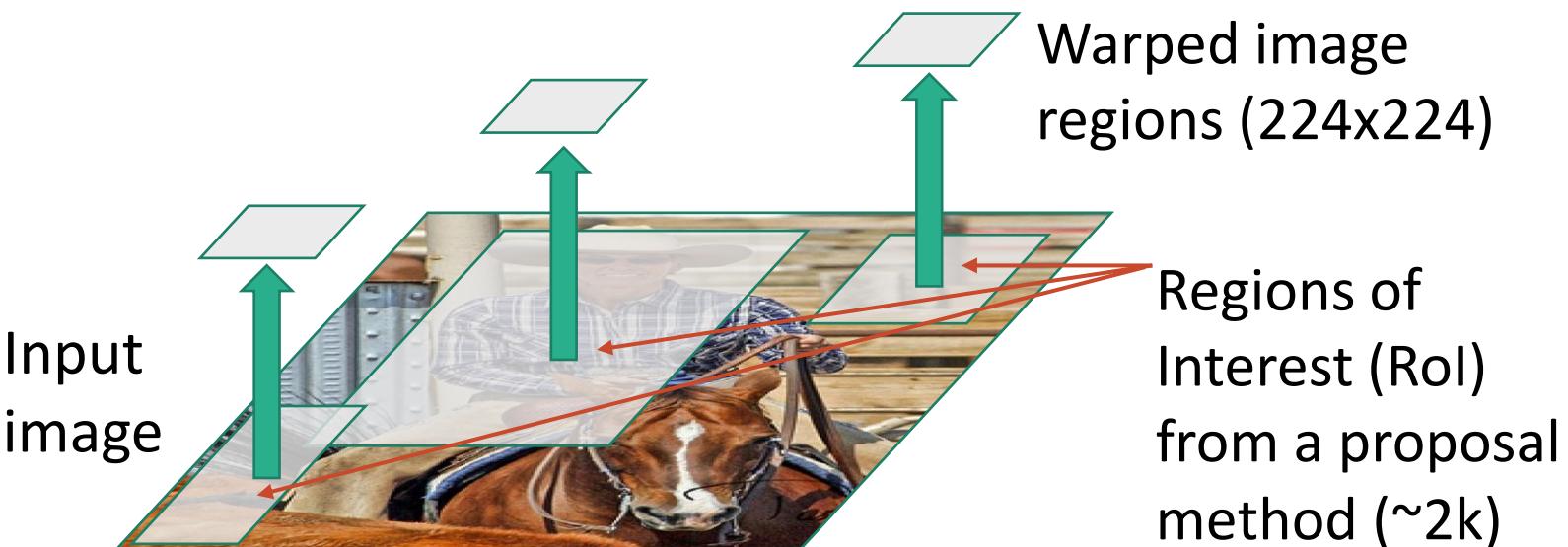


Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.
Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.

R-CNN: Region-Based CNN

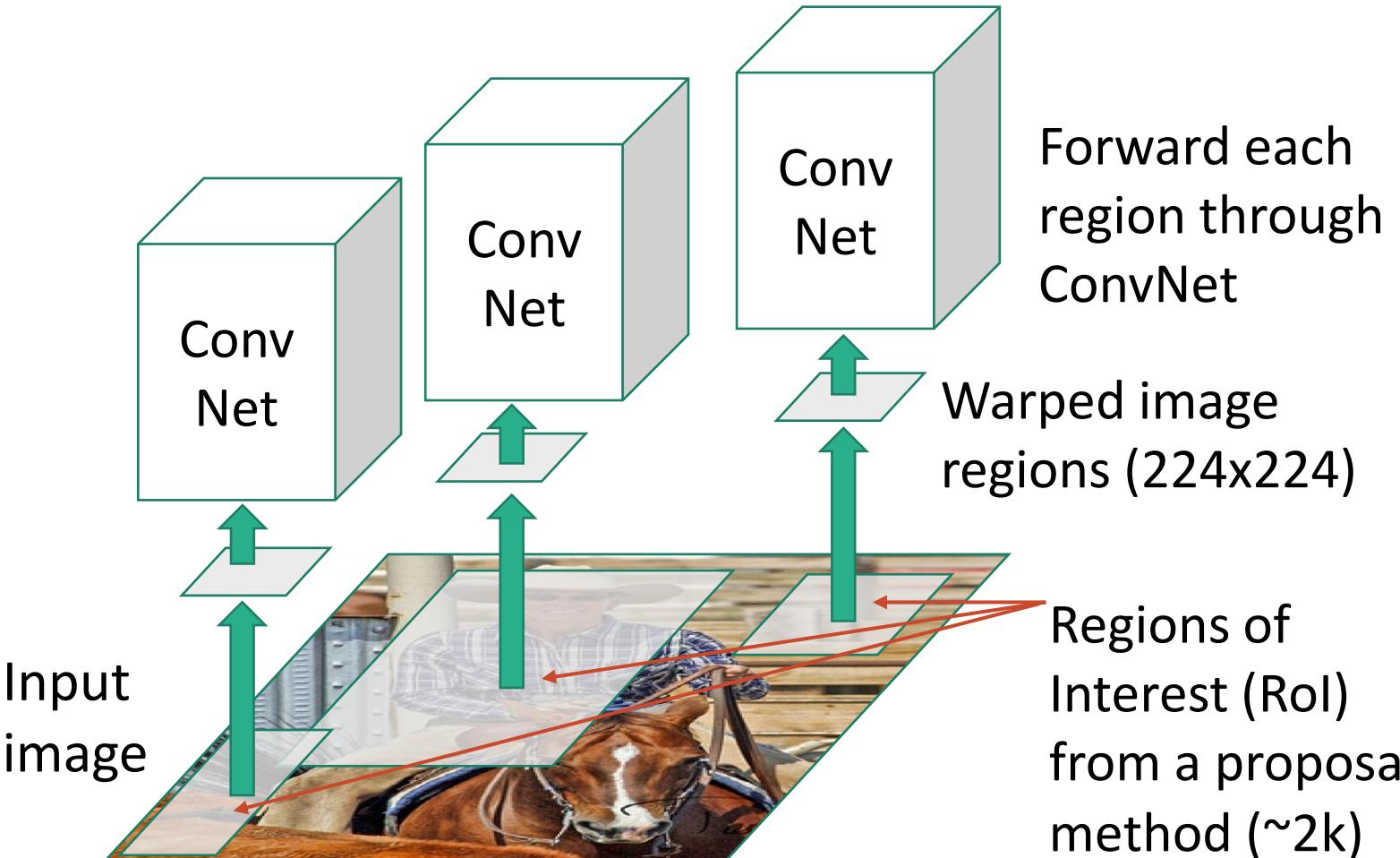


R-CNN: Region-Based CNN



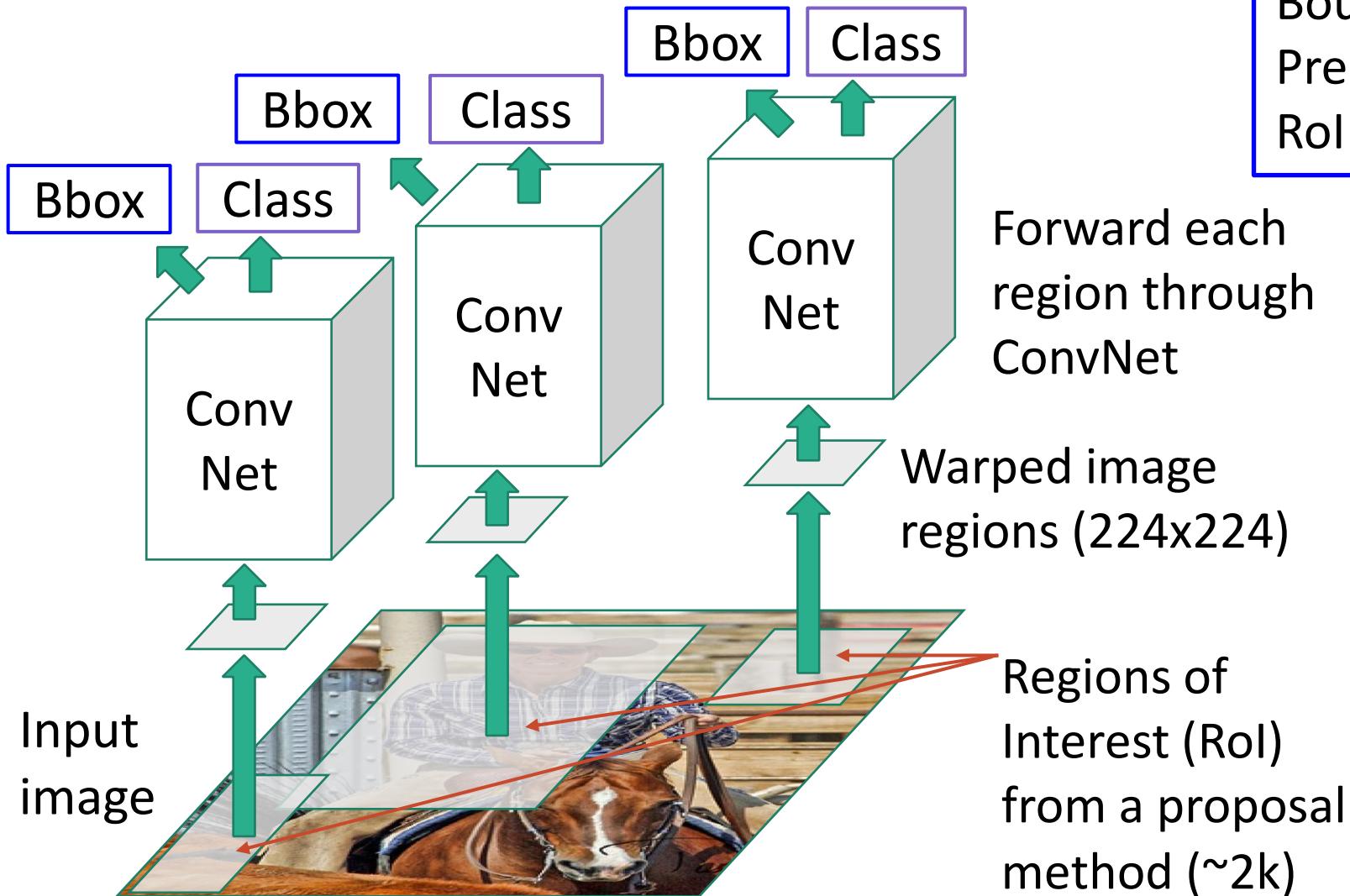
Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.
Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.

R-CNN: Region-Based CNN



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

R-CNN: Region-Based CNN

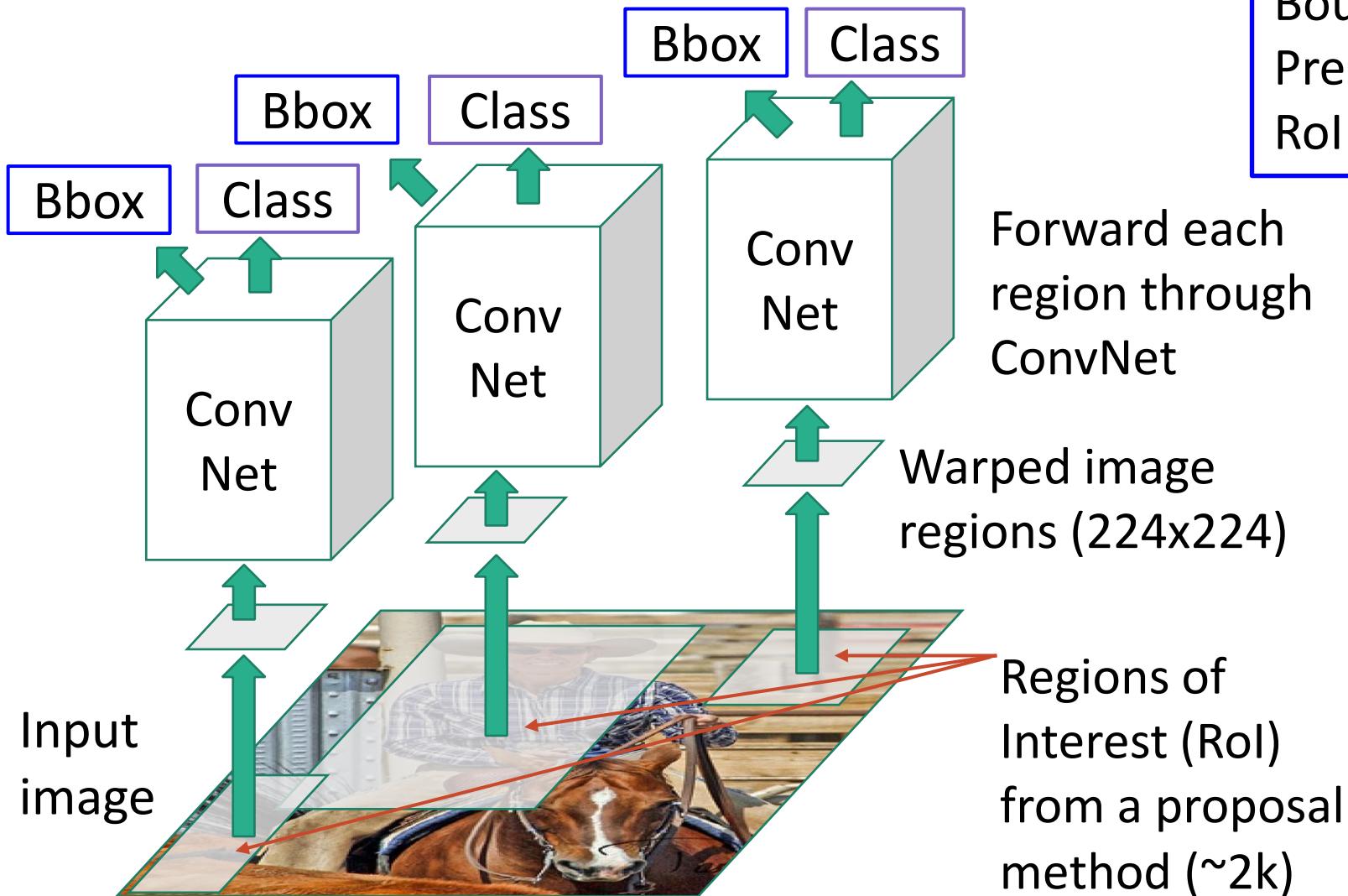


Classify each region

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

Girshick et al, “Rich feature hierarchies for accurate object detection and semantic segmentation”, CVPR 2014.
Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.

R-CNN: Region-Based CNN



Classify each region

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

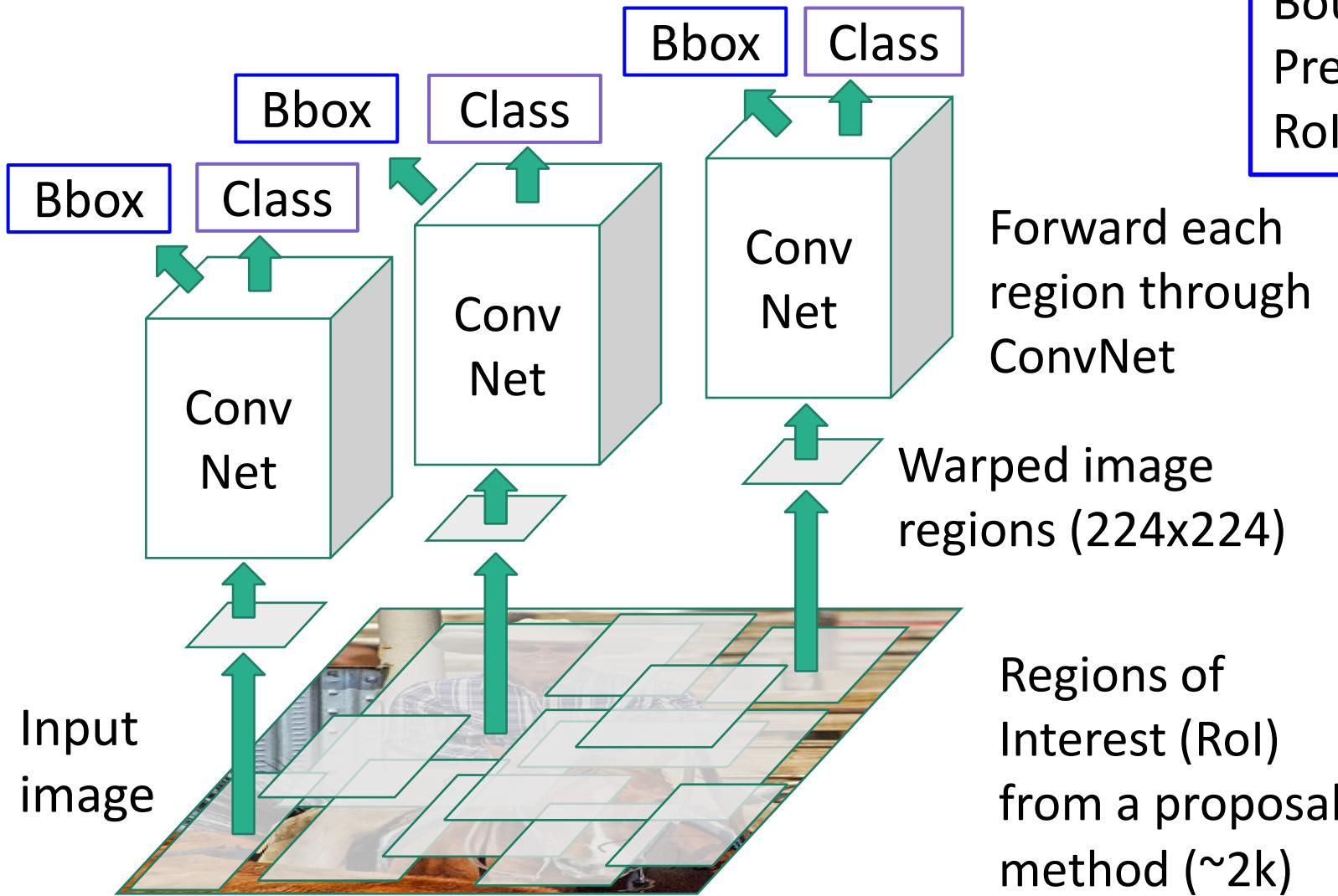
Region proposal: (p_x, p_y, p_h, p_w)
Transform: (t_x, t_y, t_h, t_w)
Output box: (b_x, b_y, b_h, b_w)

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

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

Girshick et al, “Rich feature hierarchies for accurate object detection and semantic segmentation”, CVPR 2014.
Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.

R-CNN: Region-Based CNN



Classify each region

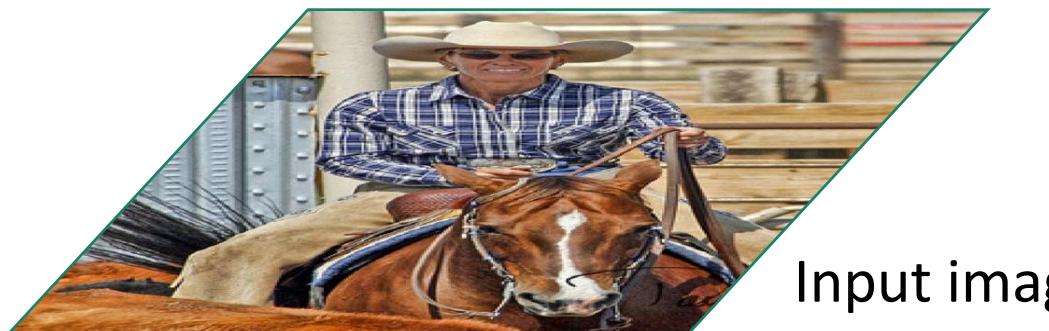
Bounding box regression:
Predict “transform” to correct the
Roi: 4 numbers (t_x, t_y, t_h, t_w)

Problem: Very slow! Need to do $\sim 2k$ forward passes through CNN per image.

Idea: Pass the image through ConvNet before cropping!

Girshick et al, “Rich feature hierarchies for accurate object detection and semantic segmentation”, CVPR 2014.
Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.

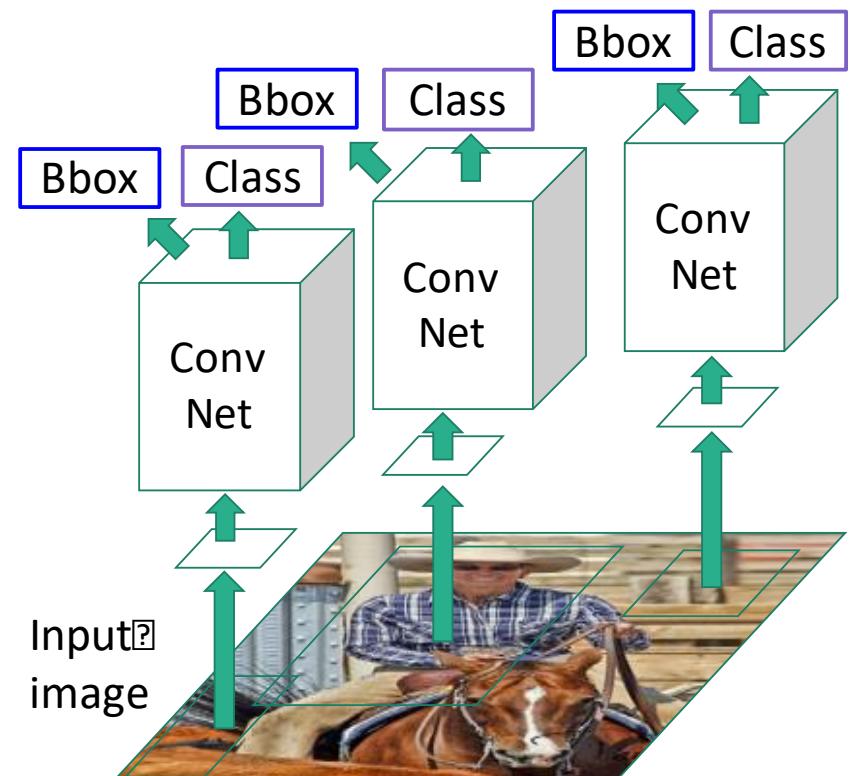
Fast R-CNN



Input image

“Slow” R-CNN

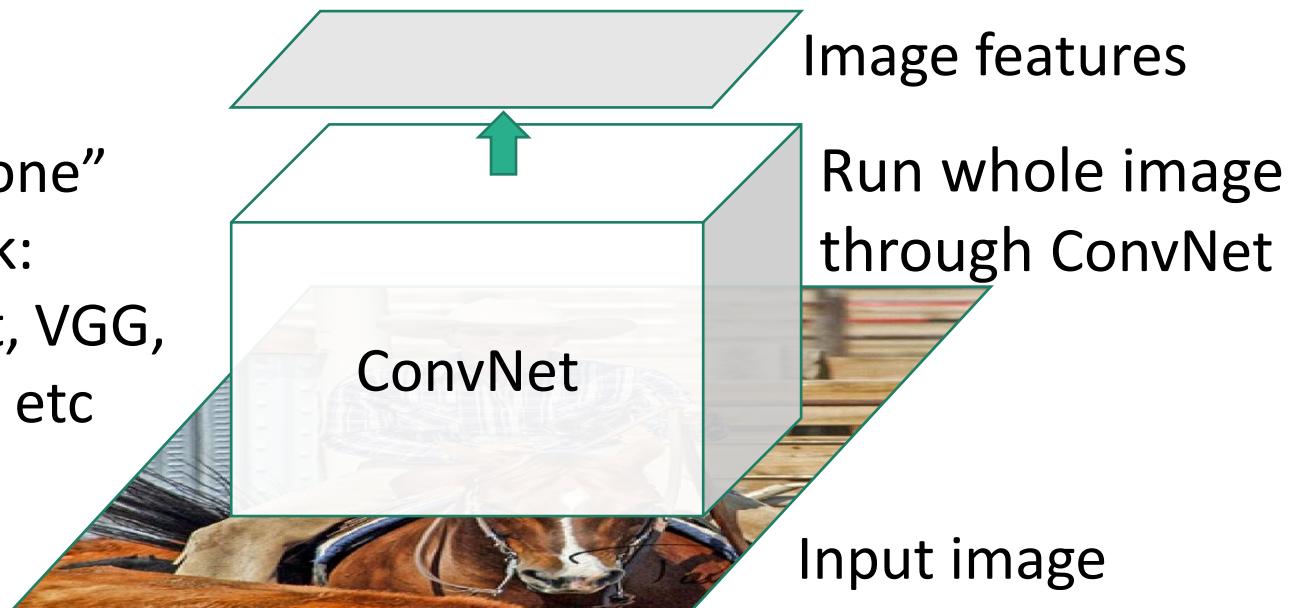
Process each region
independently



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

Fast R-CNN

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

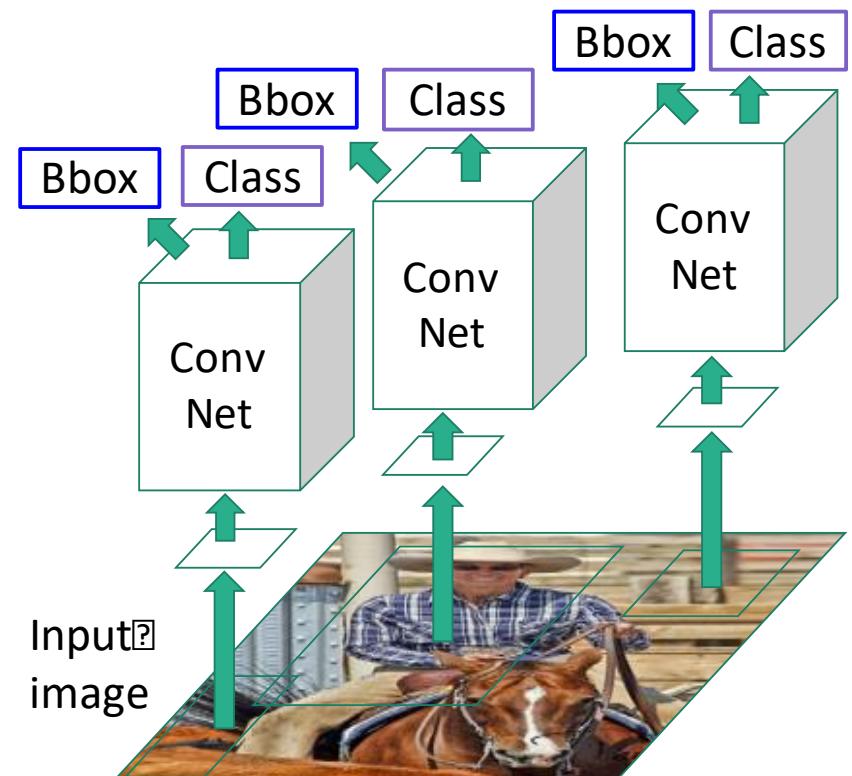


.

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

“Slow” R-CNN

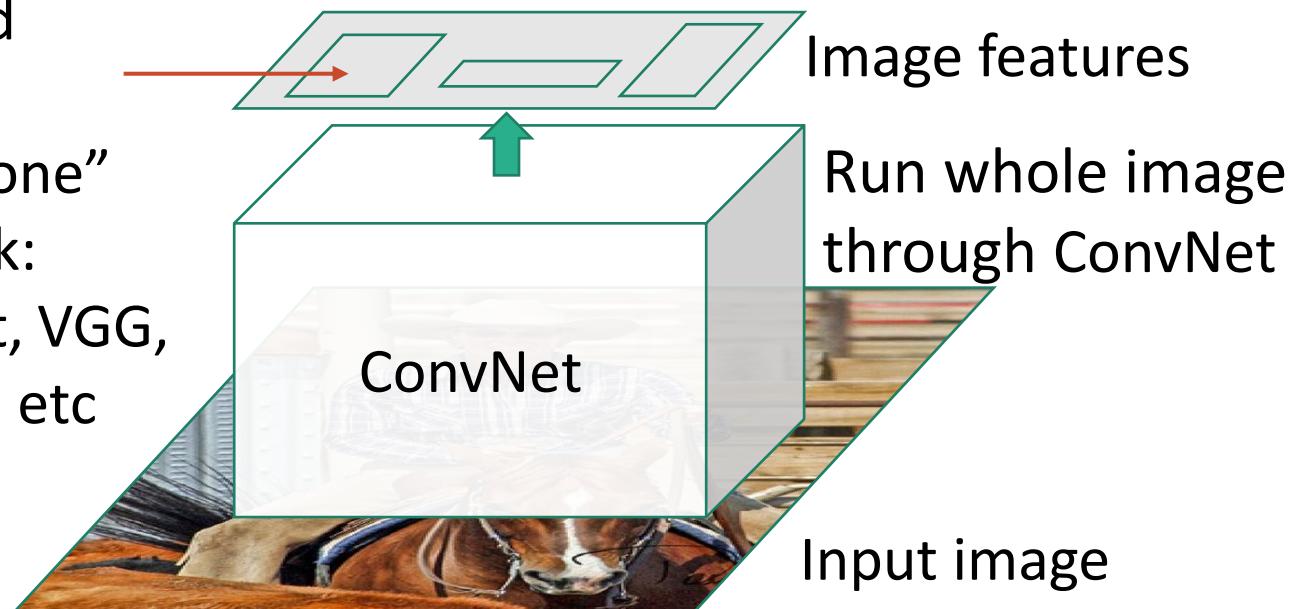
Process each region independently



Fast R-CNN

Regions of Interest (RoIs)
from a proposal
method

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

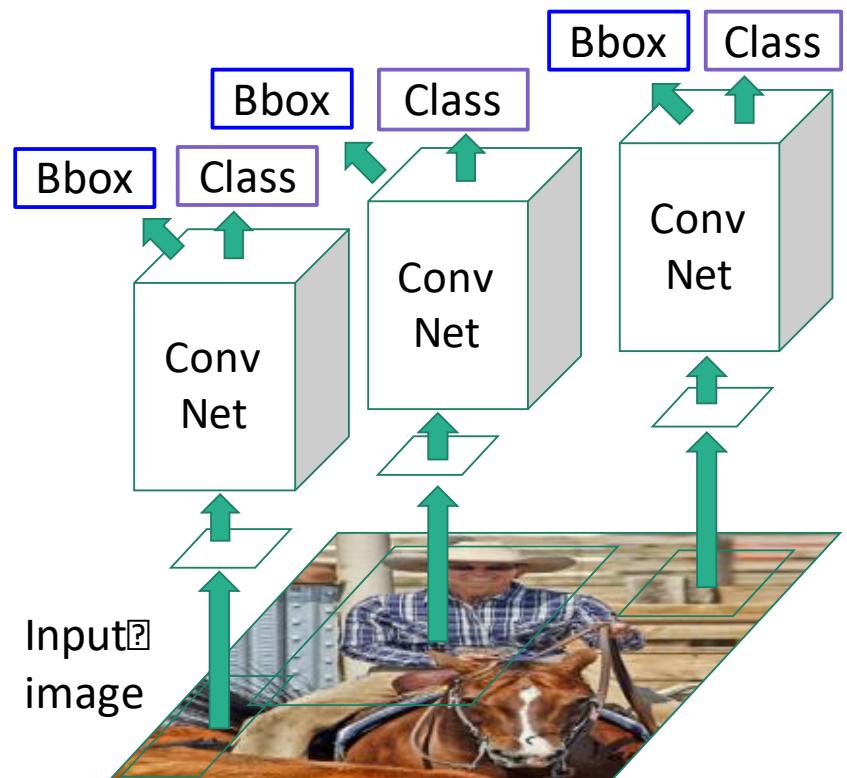


Run whole image
through ConvNet

Input image

“Slow” R-CNN

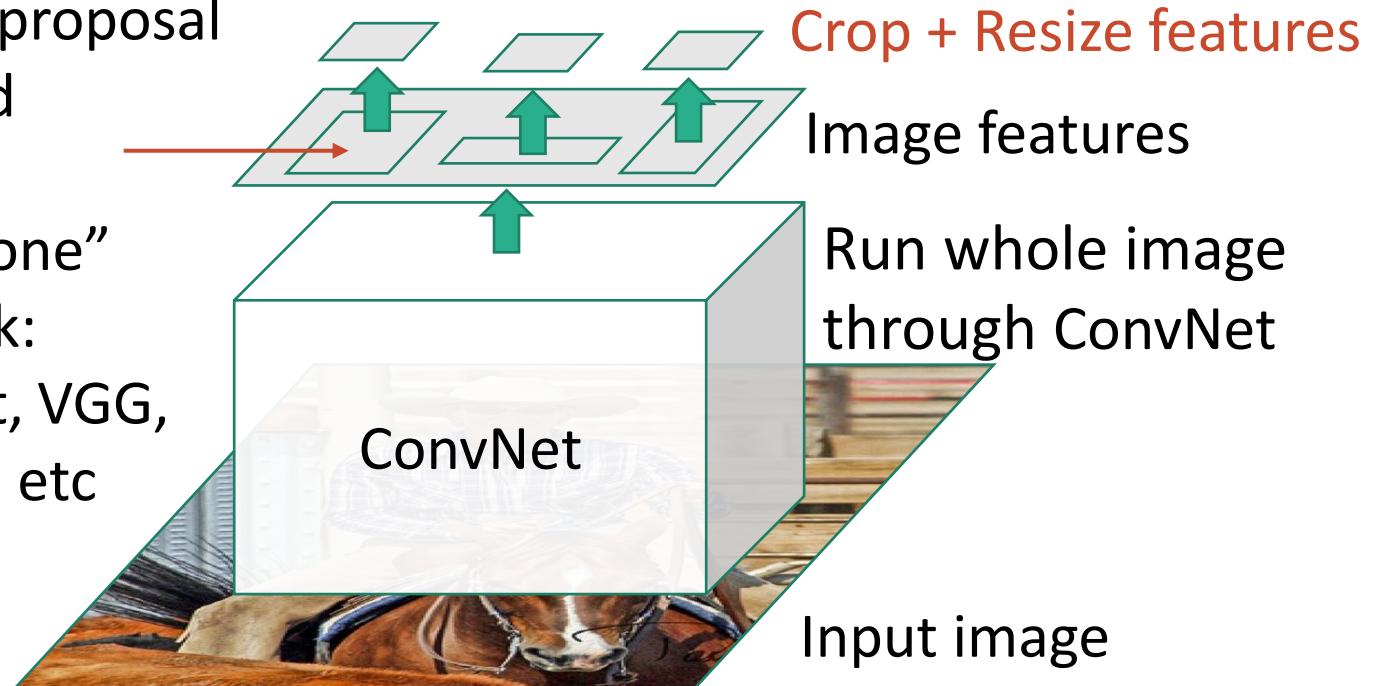
Process each region
independently



Fast R-CNN

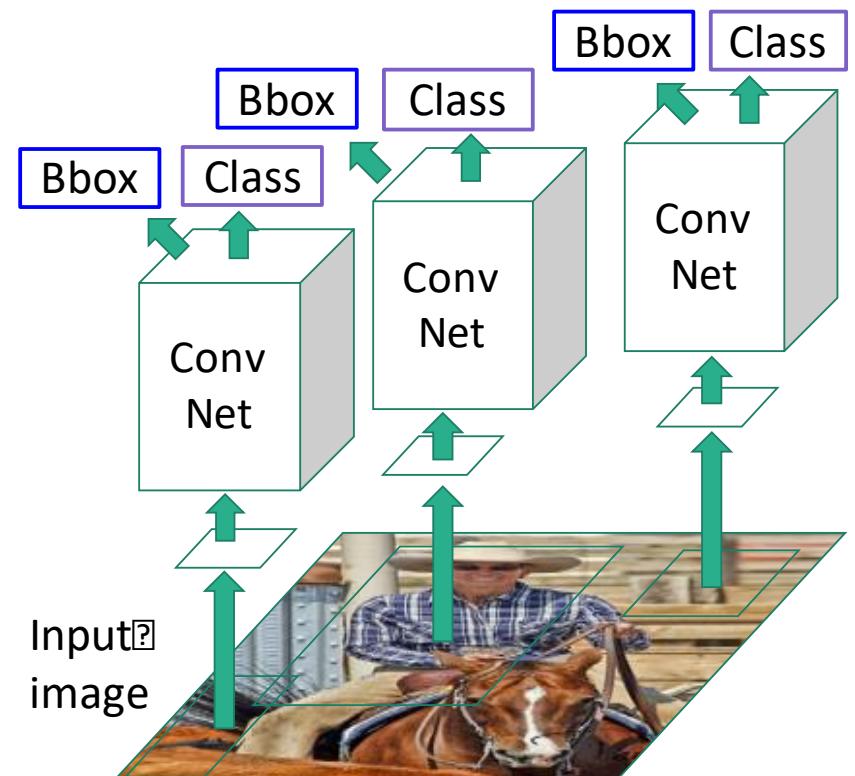
Regions of Interest (RoIs)
from a proposal
method

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



“Slow” R-CNN

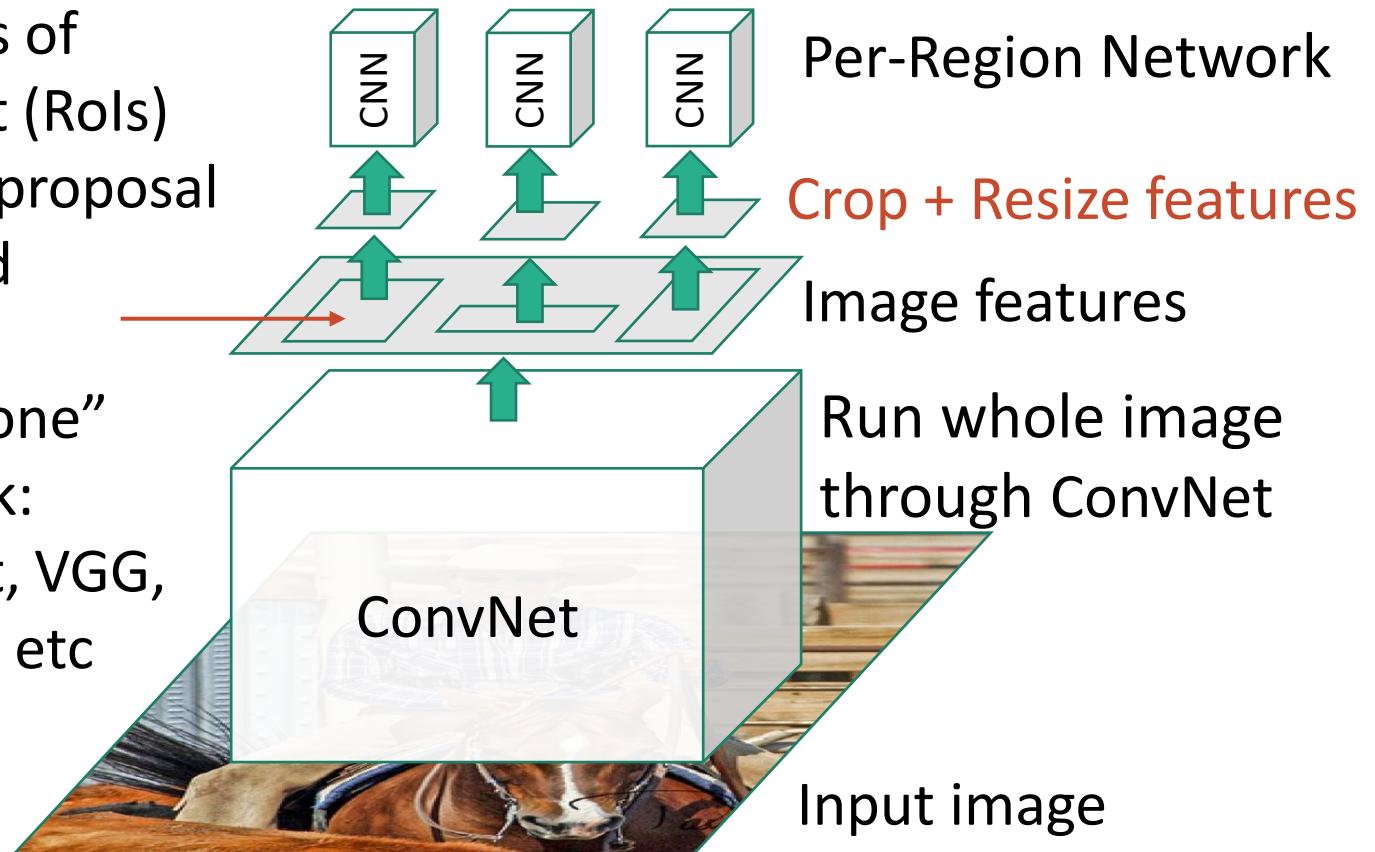
Process each region
independently



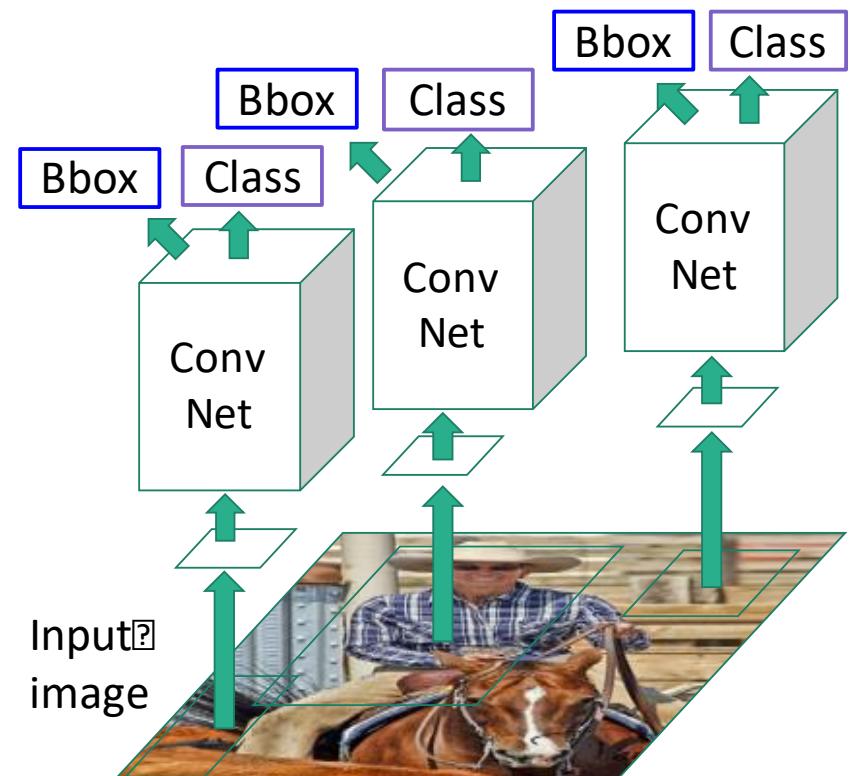
Girshick, “Fast R-CNN”, ICCV 2015. Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.

Fast R-CNN

Regions of Interest (RoIs) from a proposal method
“Backbone” network: AlexNet, VGG, ResNet, etc

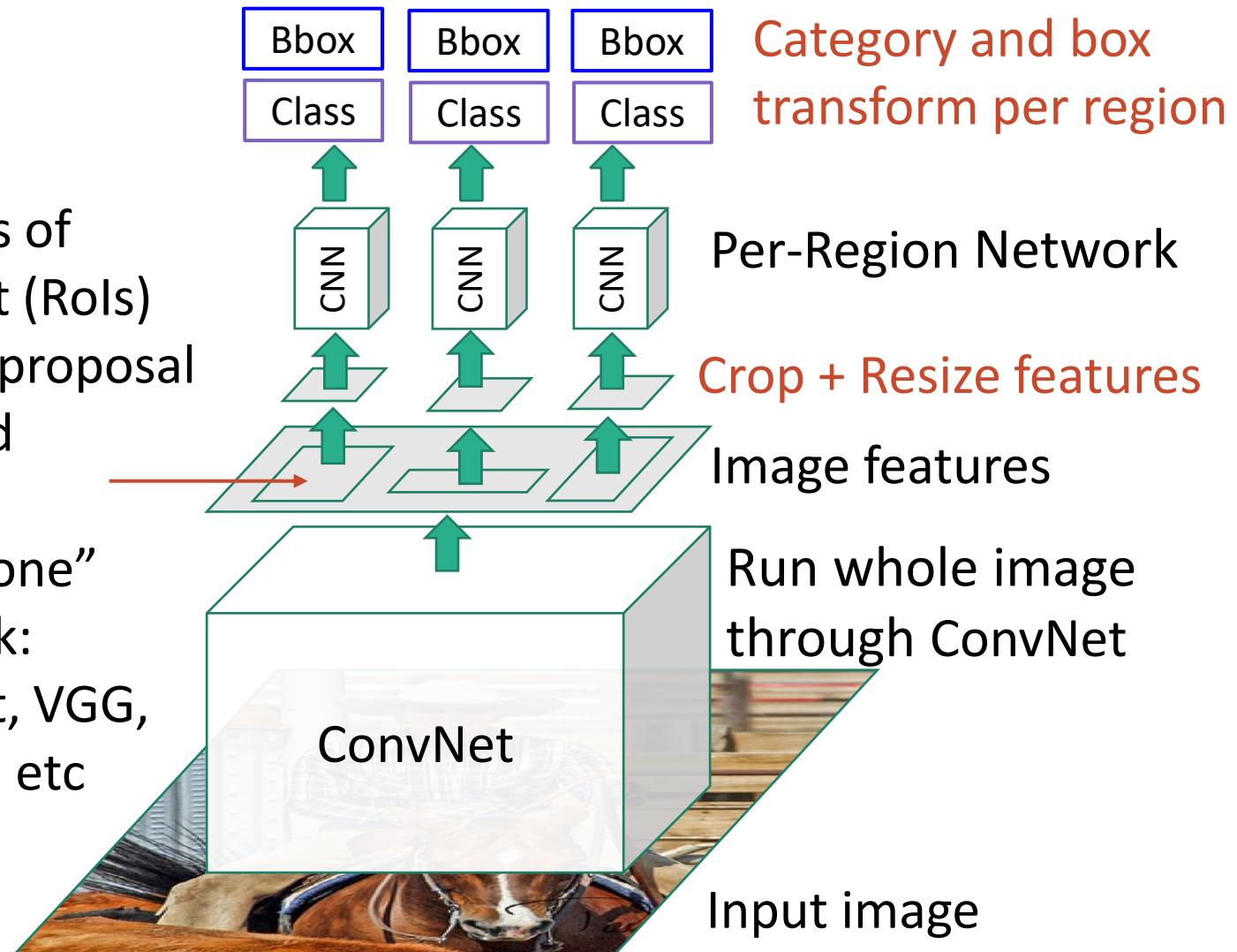


“Slow” R-CNN
Process each region independently

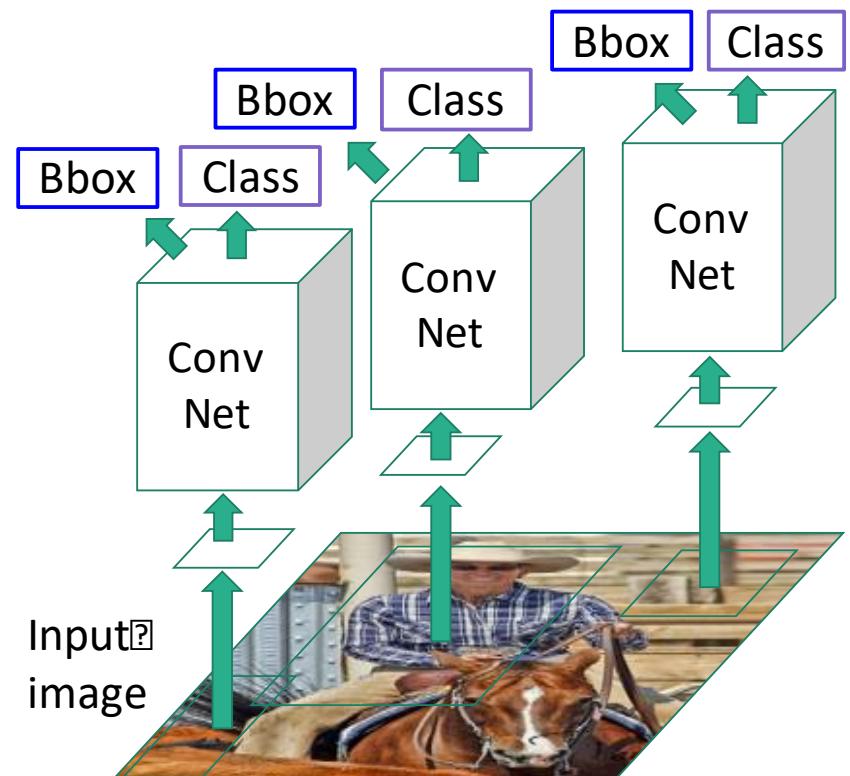


Fast R-CNN

Regions of Interest (Rois) from a proposal method
“Backbone” network: AlexNet, VGG, ResNet, etc



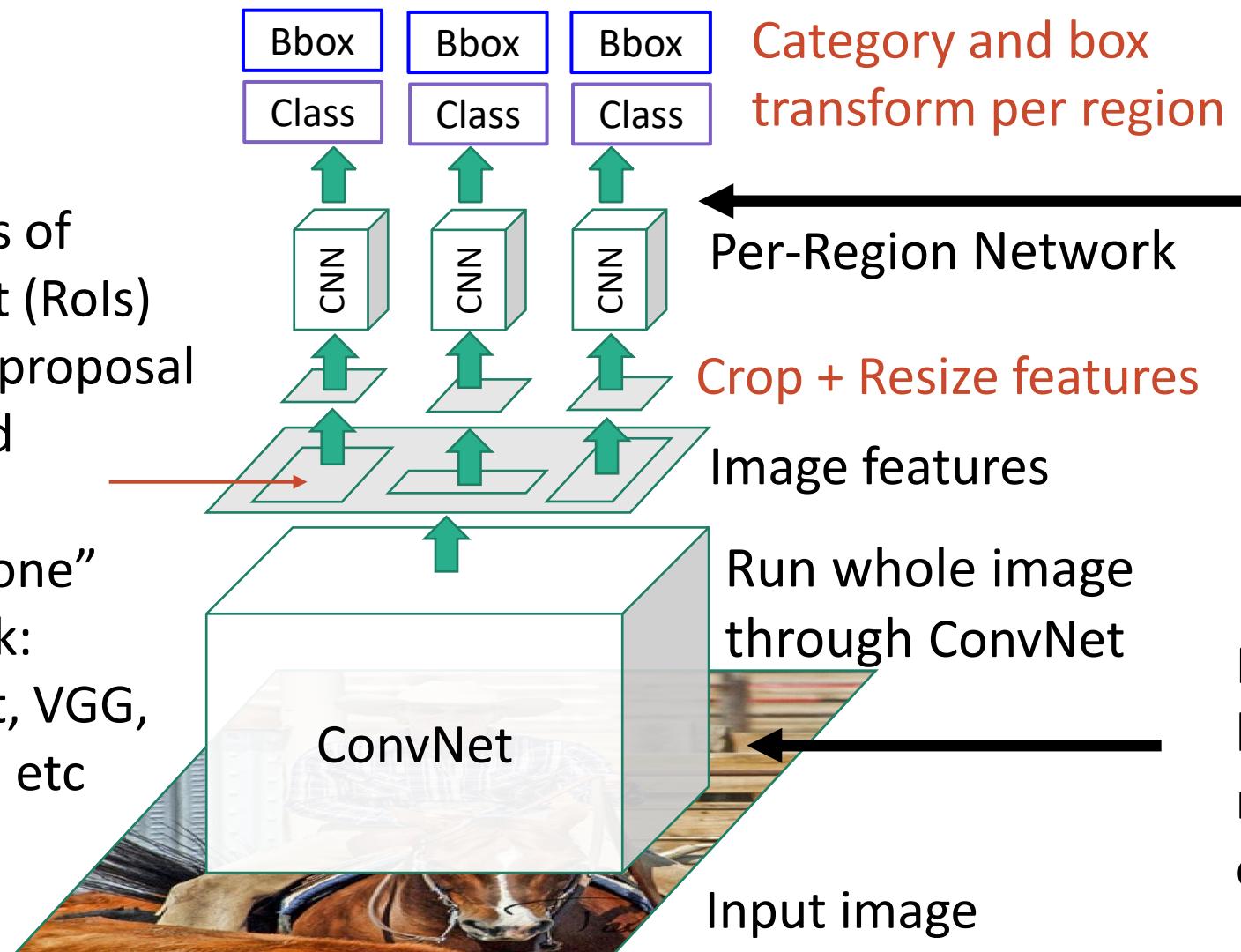
“Slow” R-CNN
Process each region independently



Fast R-CNN

Regions of Interest (RoIs)
from a proposal
method

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



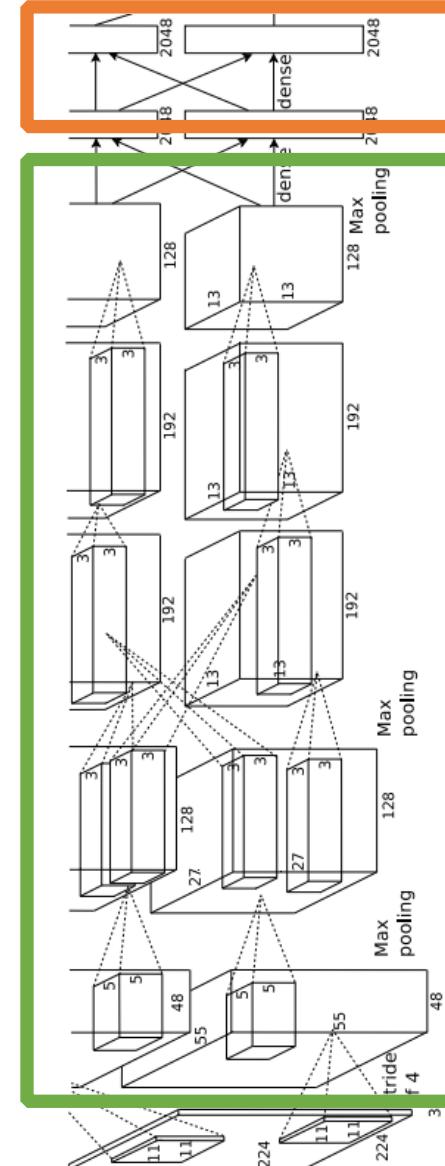
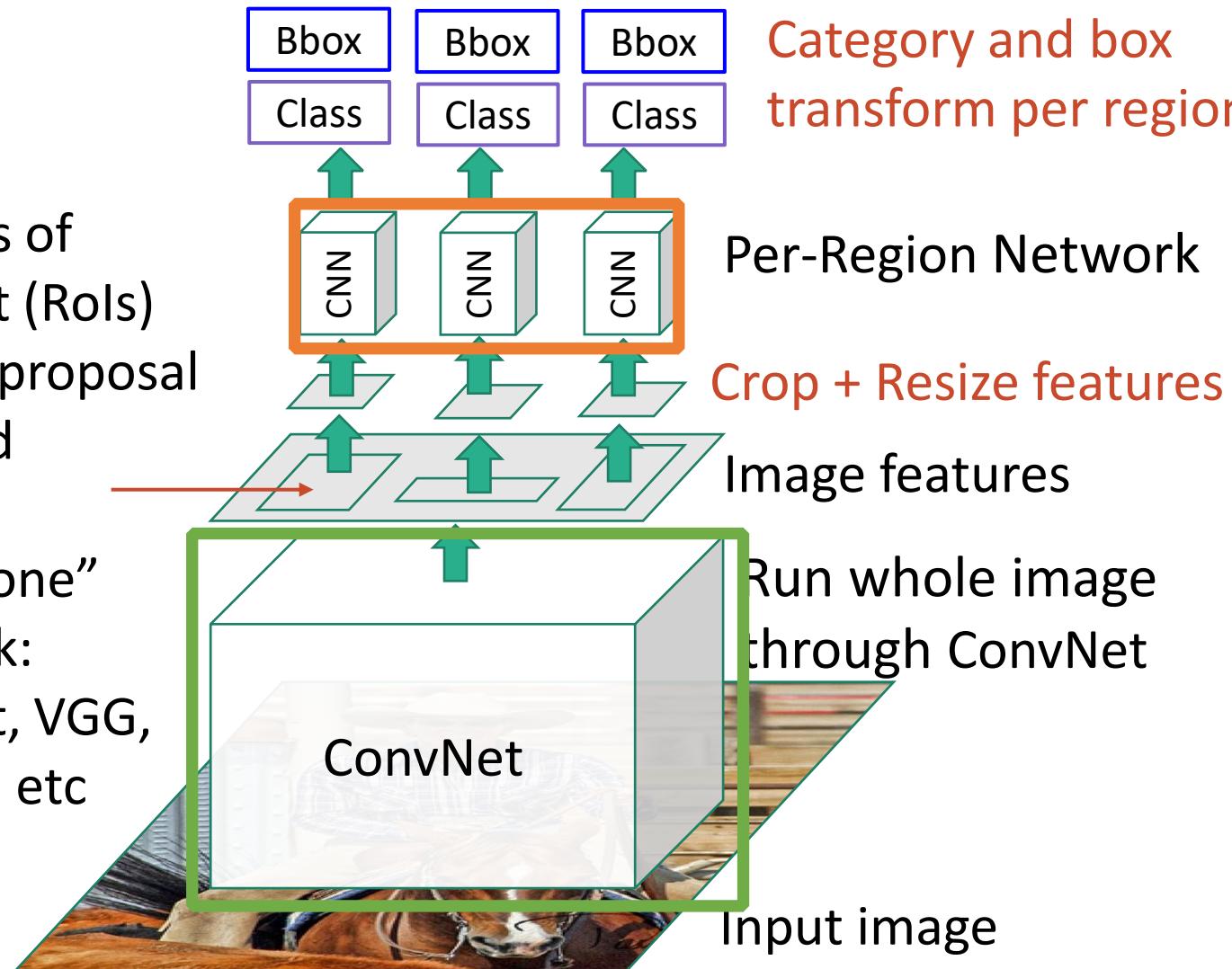
Per-Region network is relatively lightweight

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

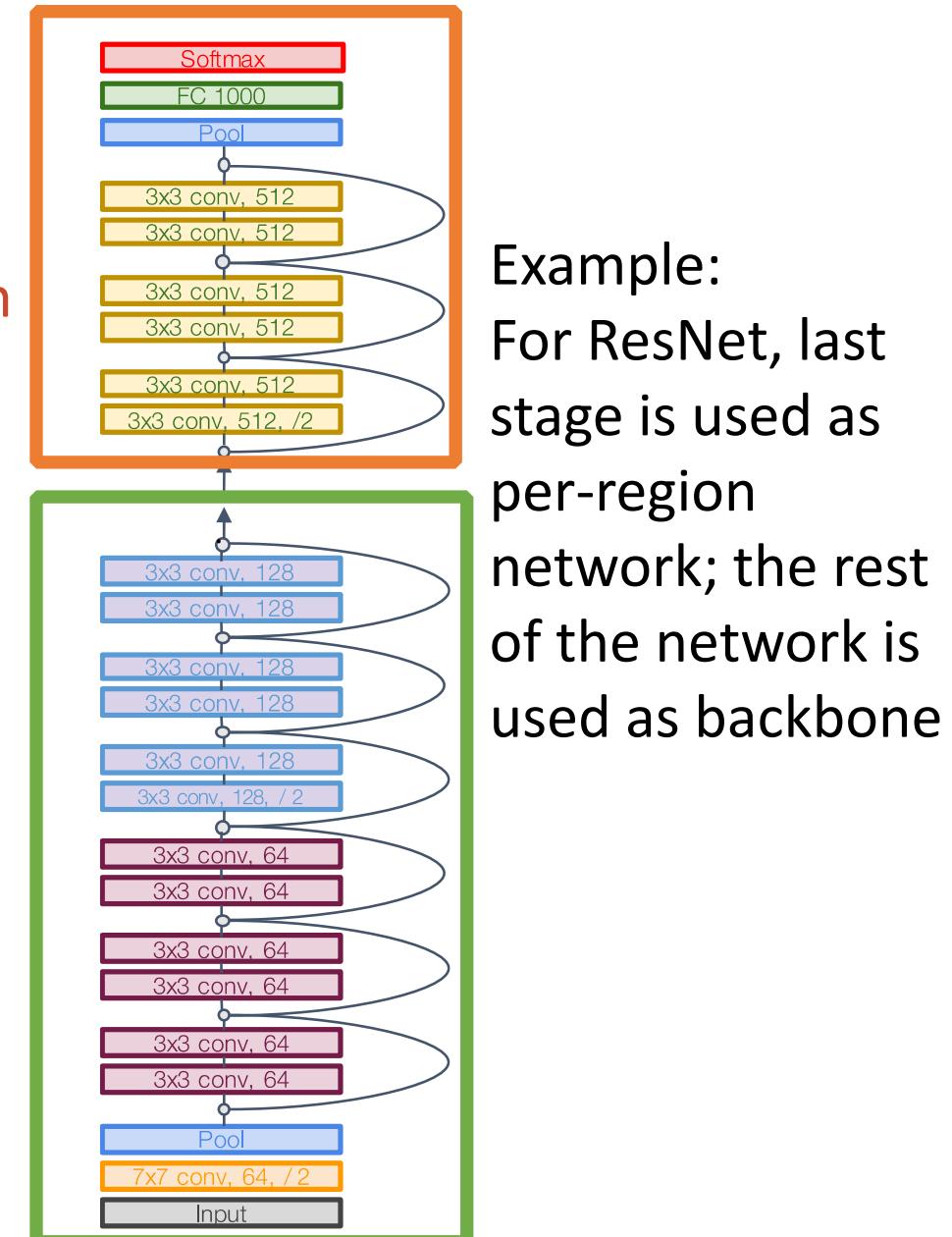
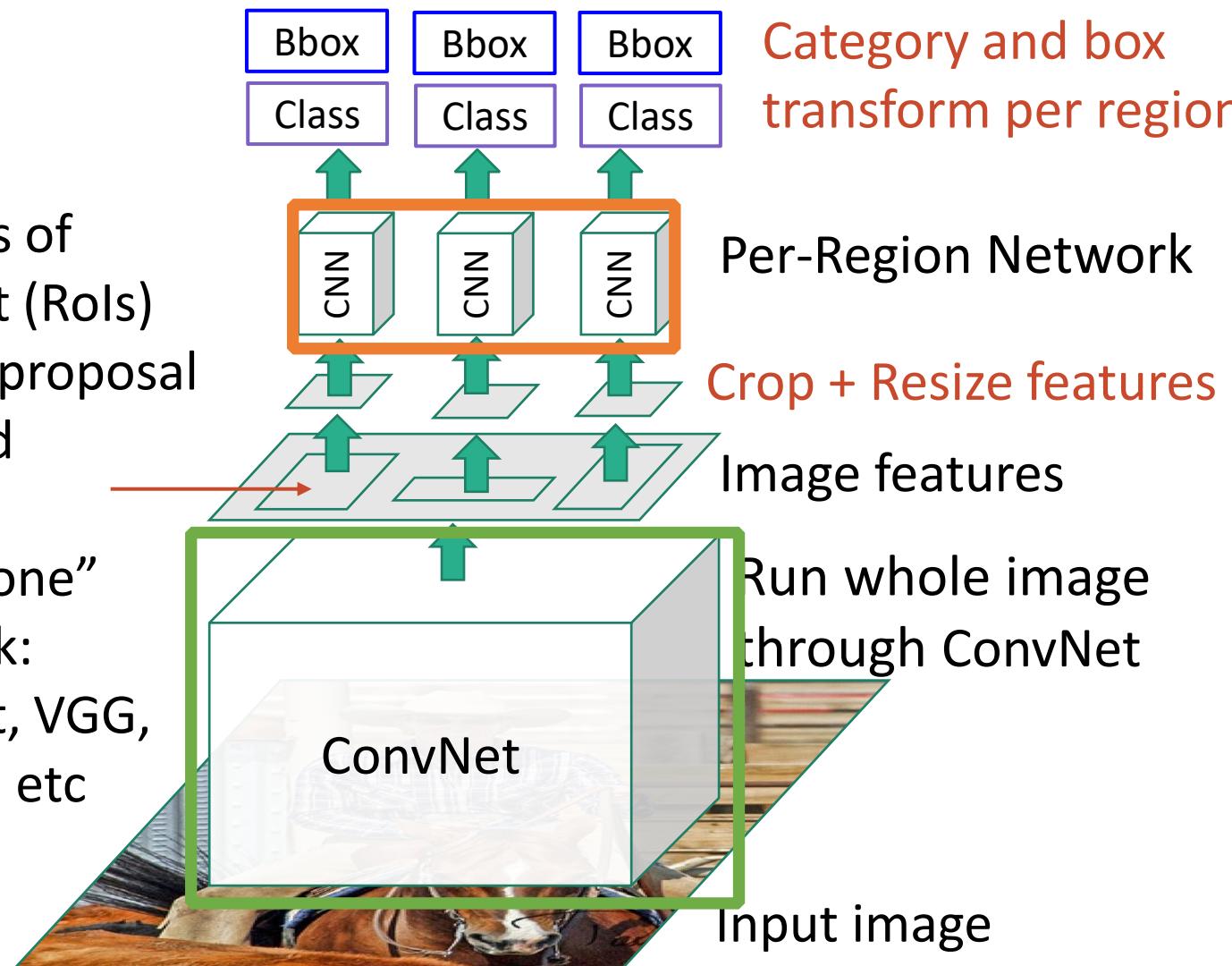


Example:
When using
AlexNet for
detection, five
conv layers are
used for
backbone and
two FC layers are
used for per-
region network

Fast R-CNN

Regions of Interest (Rois)
from a proposal
method

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

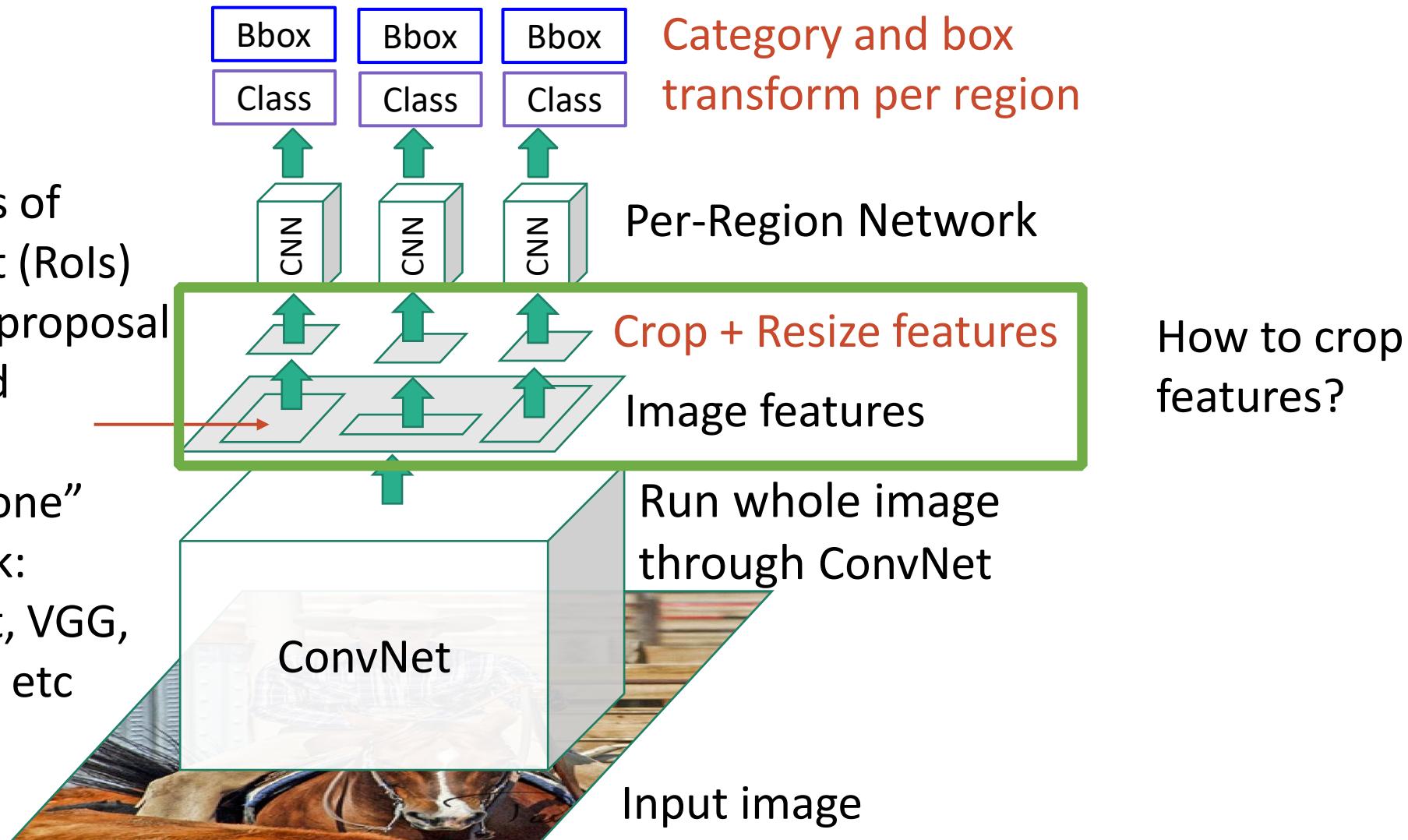


Example:
For ResNet, last
stage is used as
per-region
network; the rest
of the network is
used as backbone

Fast R-CNN

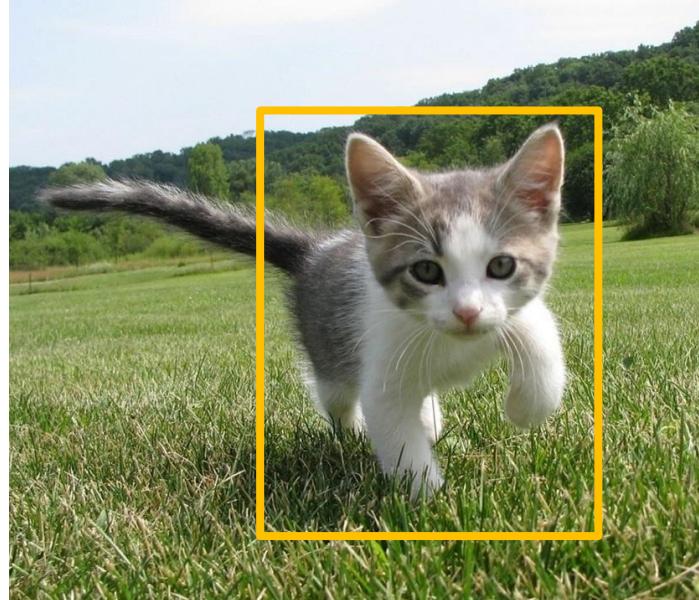
Regions of Interest (Rois)
from a proposal
method

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



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

Cropping Features: RoI Pool



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

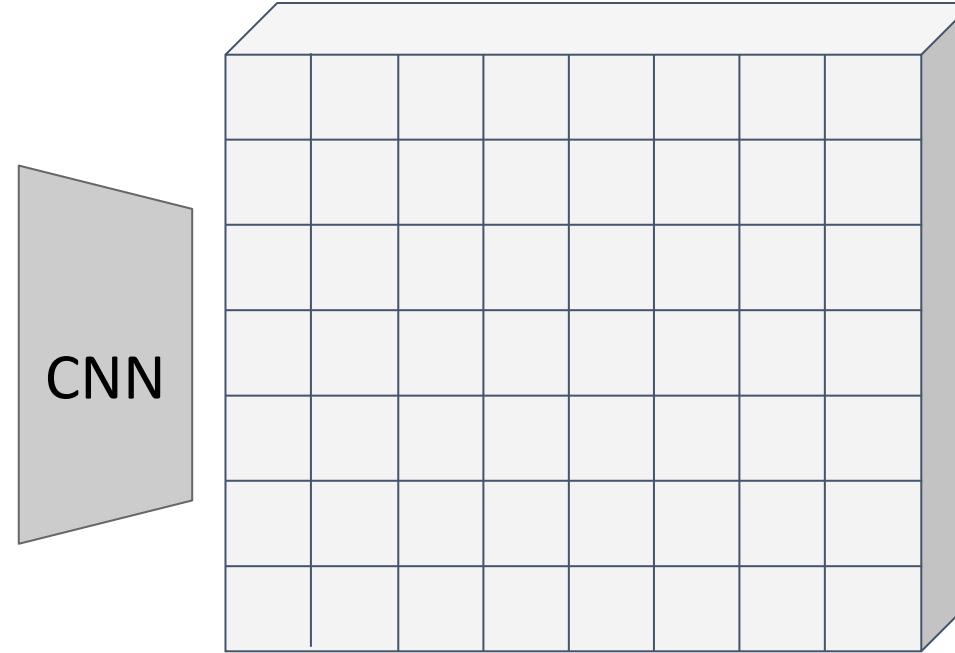
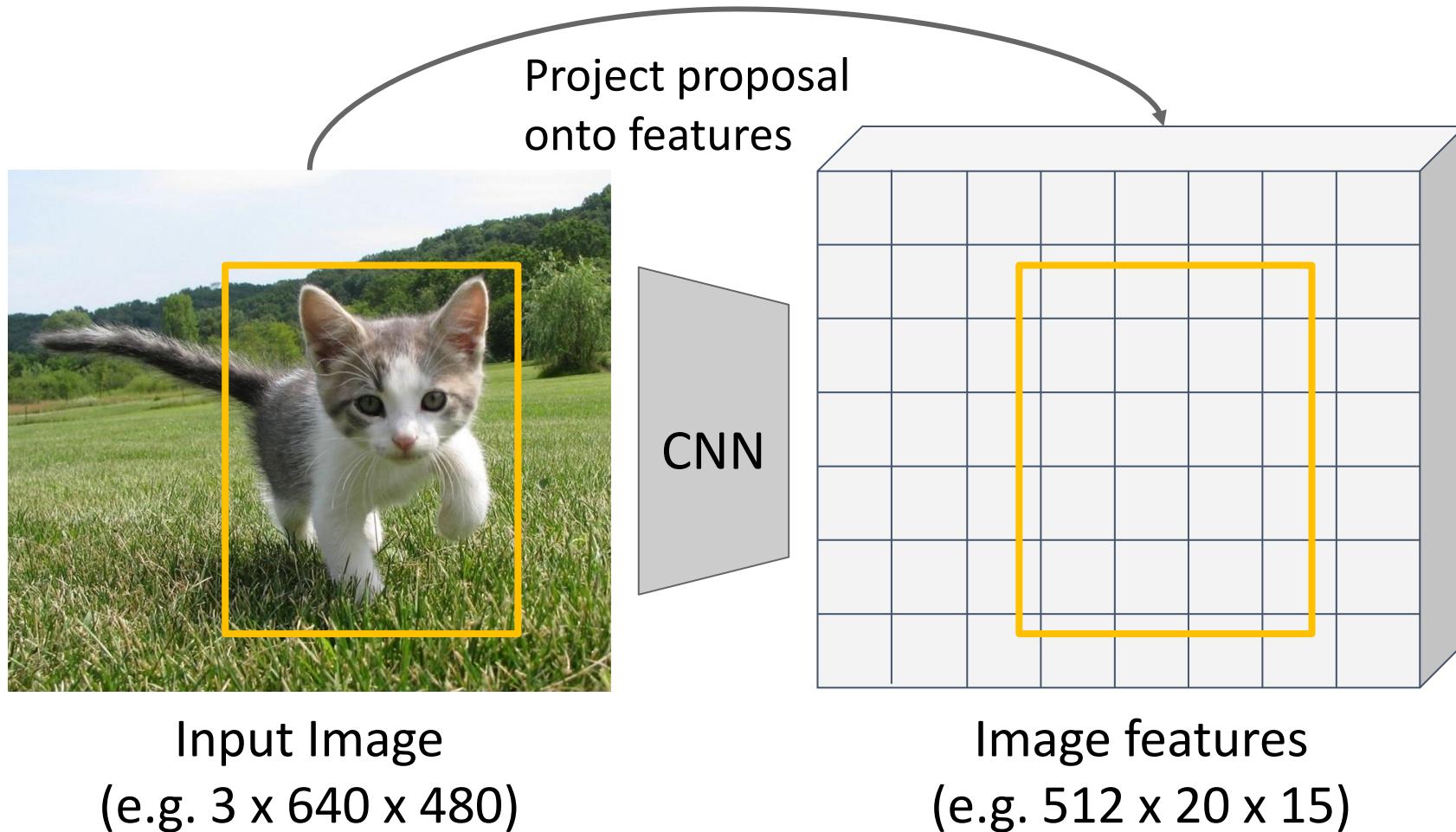


Image features
(e.g. $512 \times 20 \times 15$)

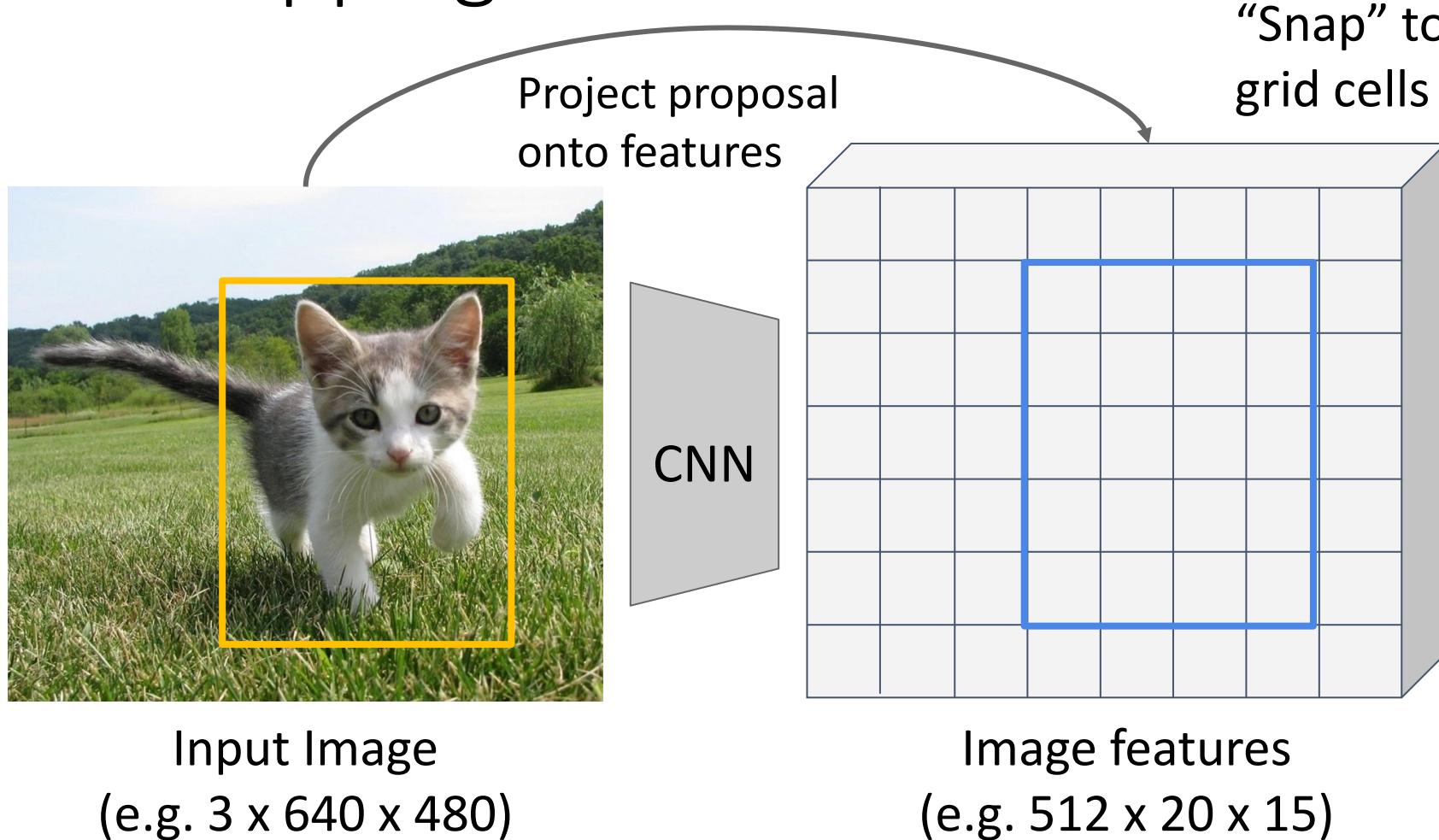
Want features for the
box of a fixed size
(2×2 in this example,
 7×7 or 14×14 in practice)

Cropping Features: RoI Pool



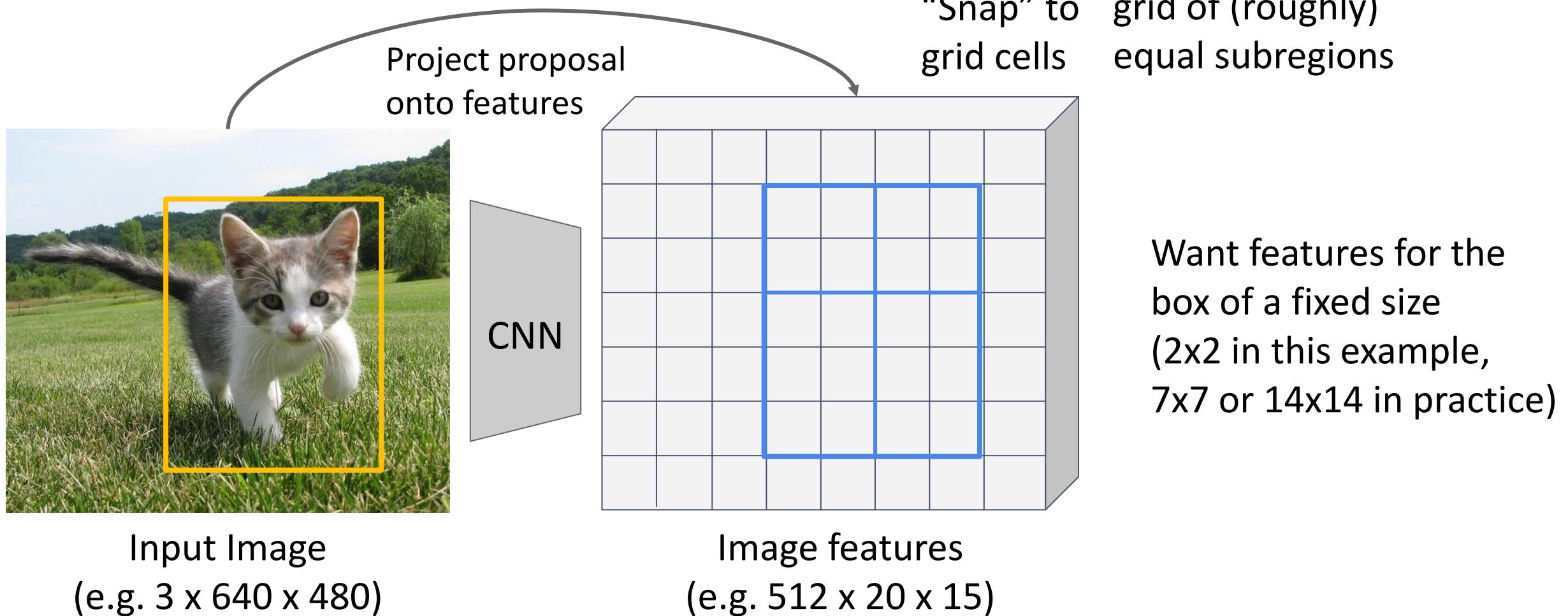
Want features for the
box of a fixed size
(2×2 in this example,
 7×7 or 14×14 in practice)

Cropping Features: RoI Pool



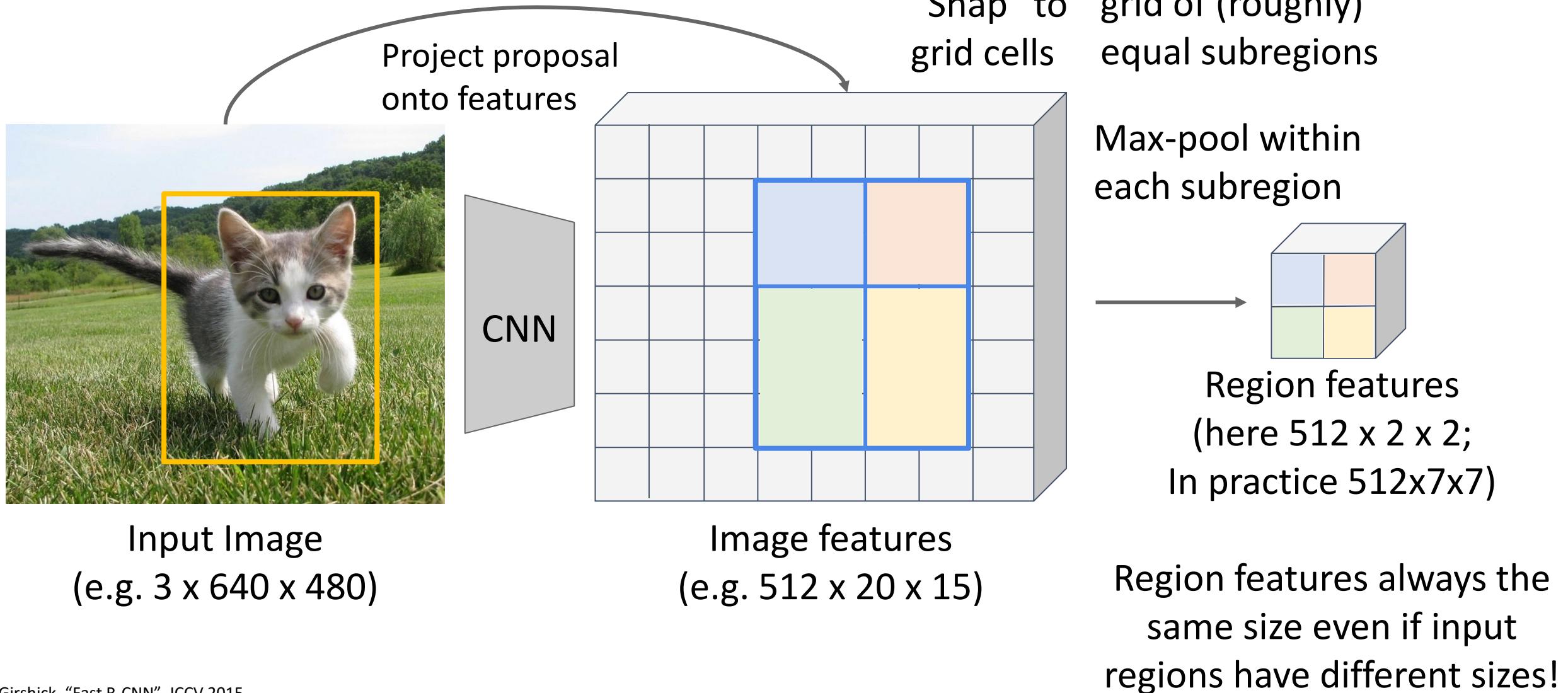
Want features for the box of a fixed size (2x2 in this example, 7x7 or 14x14 in practice)

Cropping Features: RoI Pool



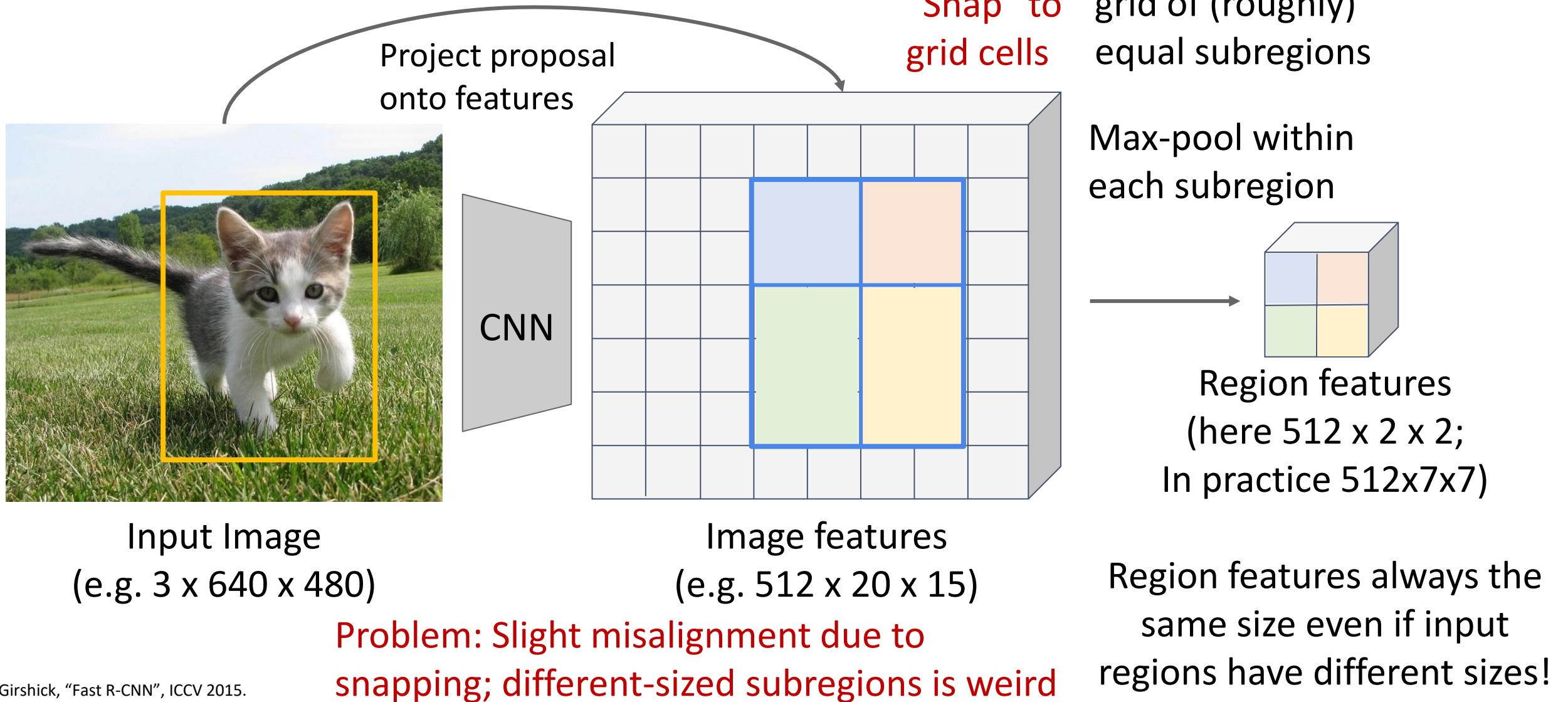
Girshick, "Fast R-CNN", ICCV 2015.

Cropping Features: RoI Pool



Girshick, “Fast R-CNN”, ICCV 2015.

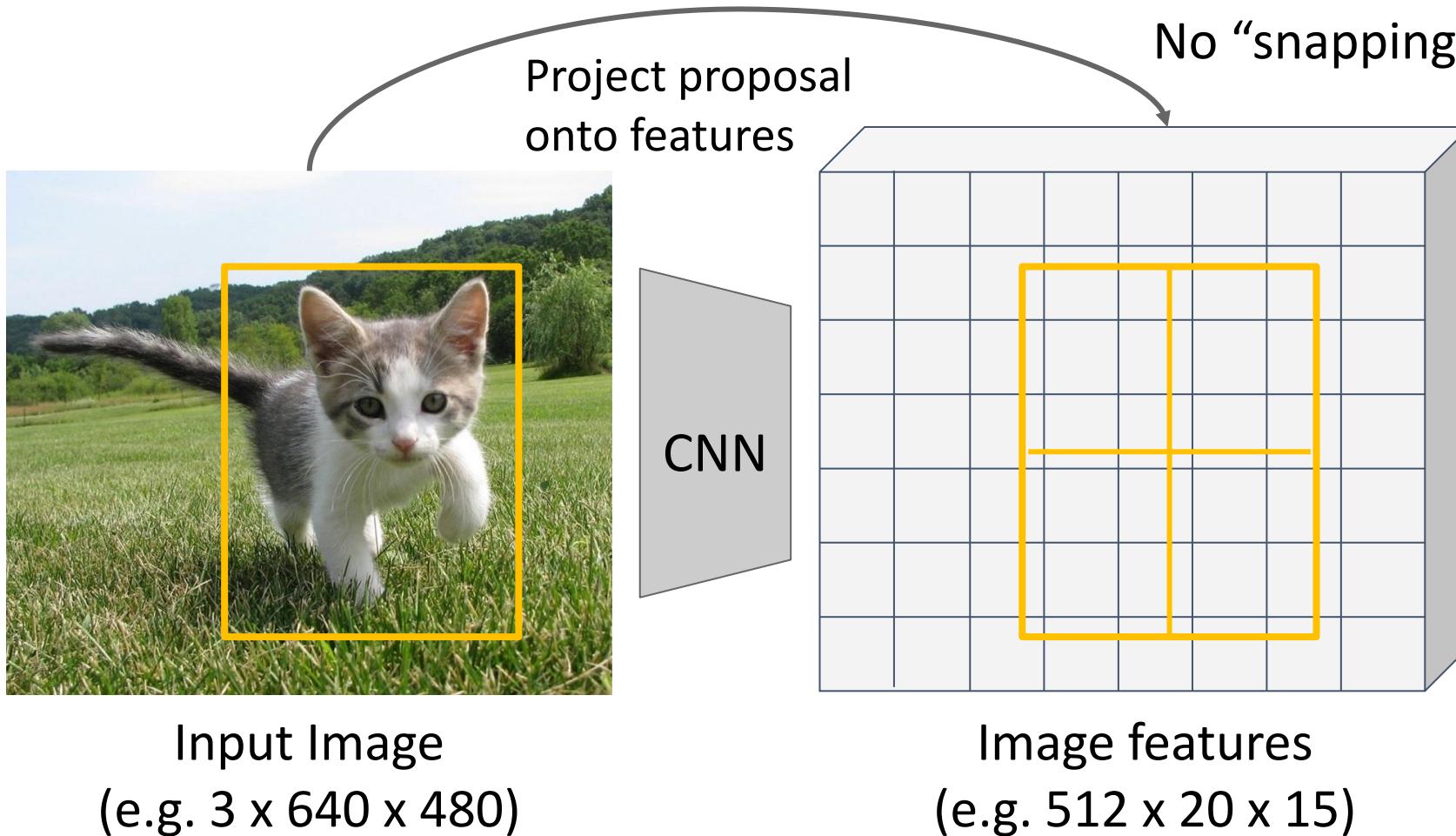
Cropping Features: RoI Pool



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

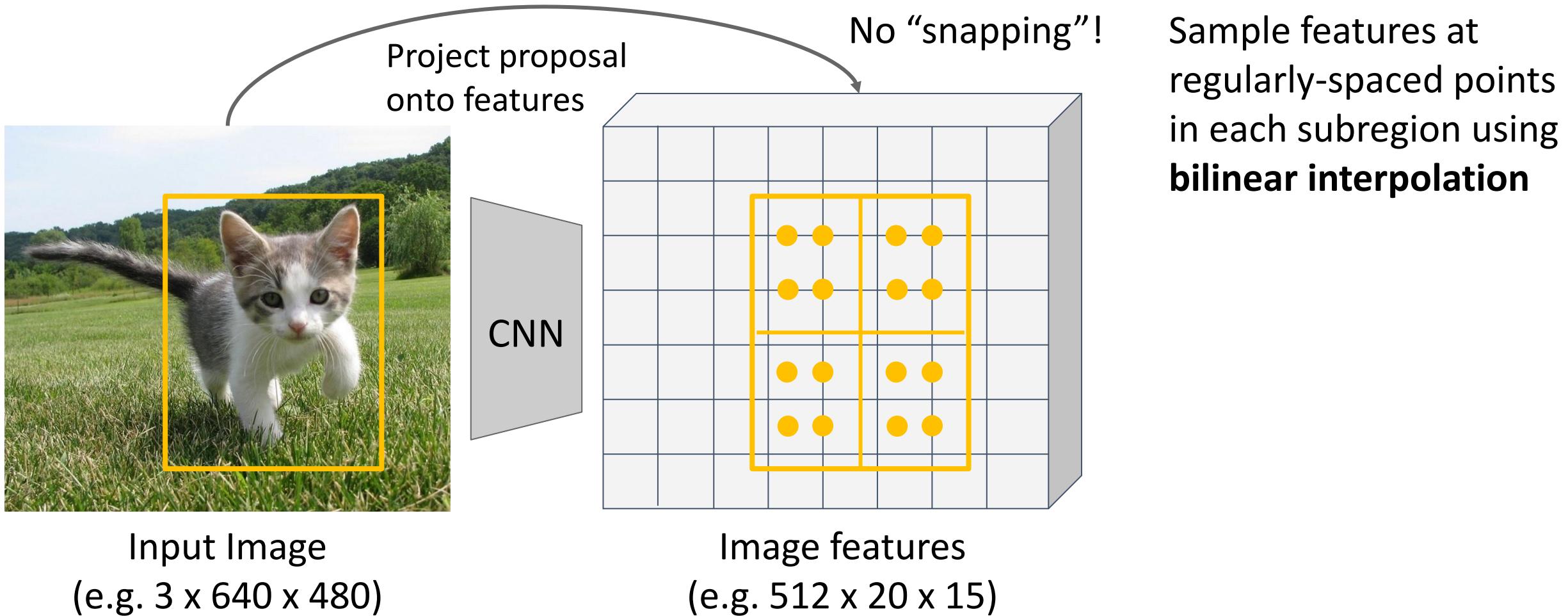
Cropping Features: RoI Align

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



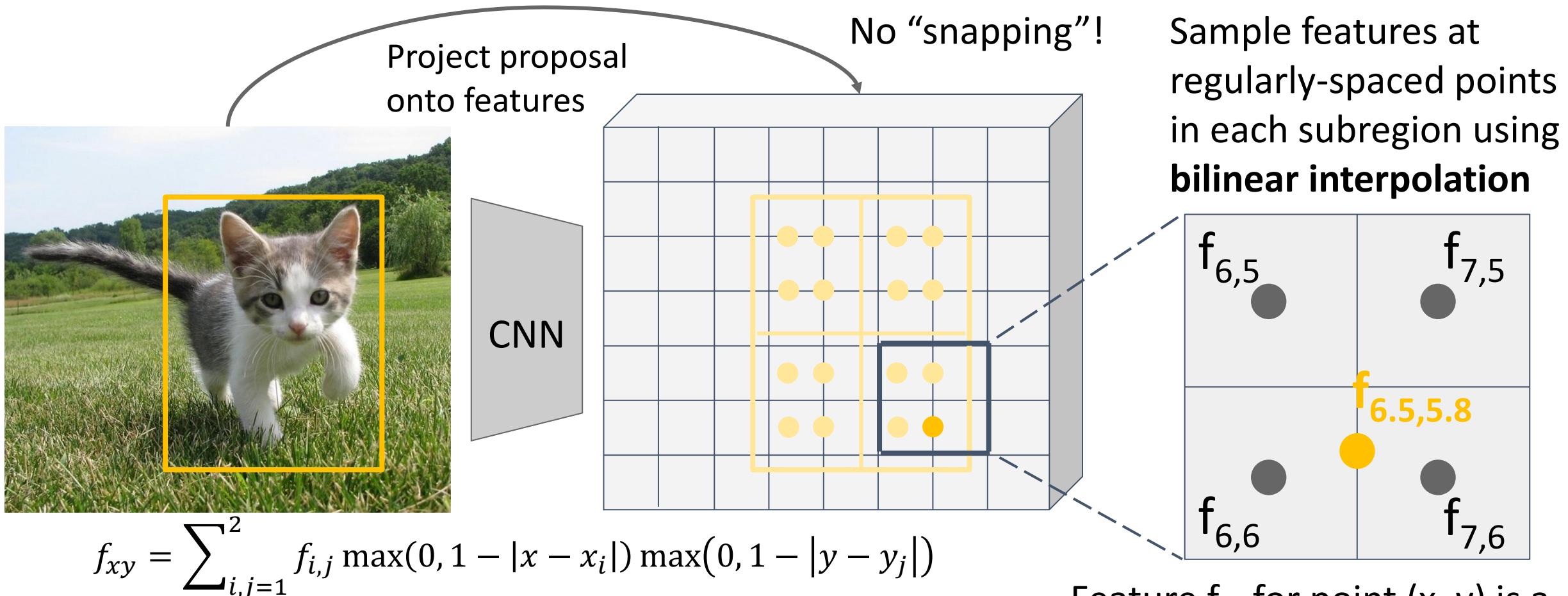
Cropping Features: RoI Align

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



Cropping Features: RoI Align

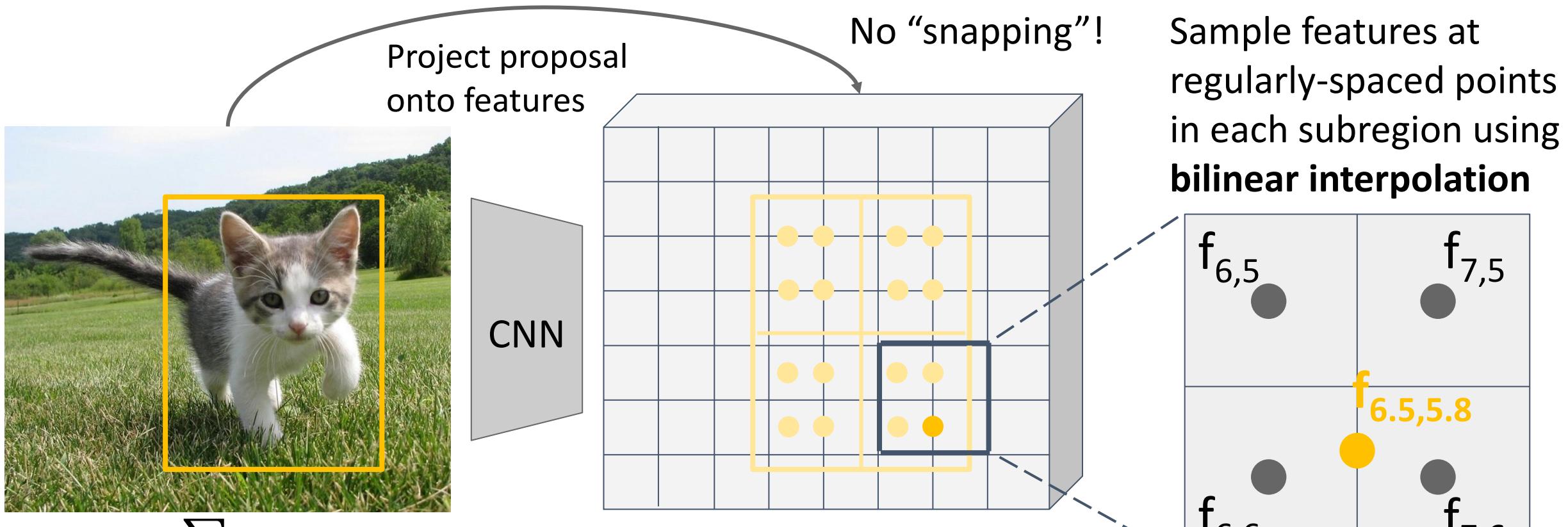
Divide into equal-sized subregions
(may not be aligned to grid!)



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

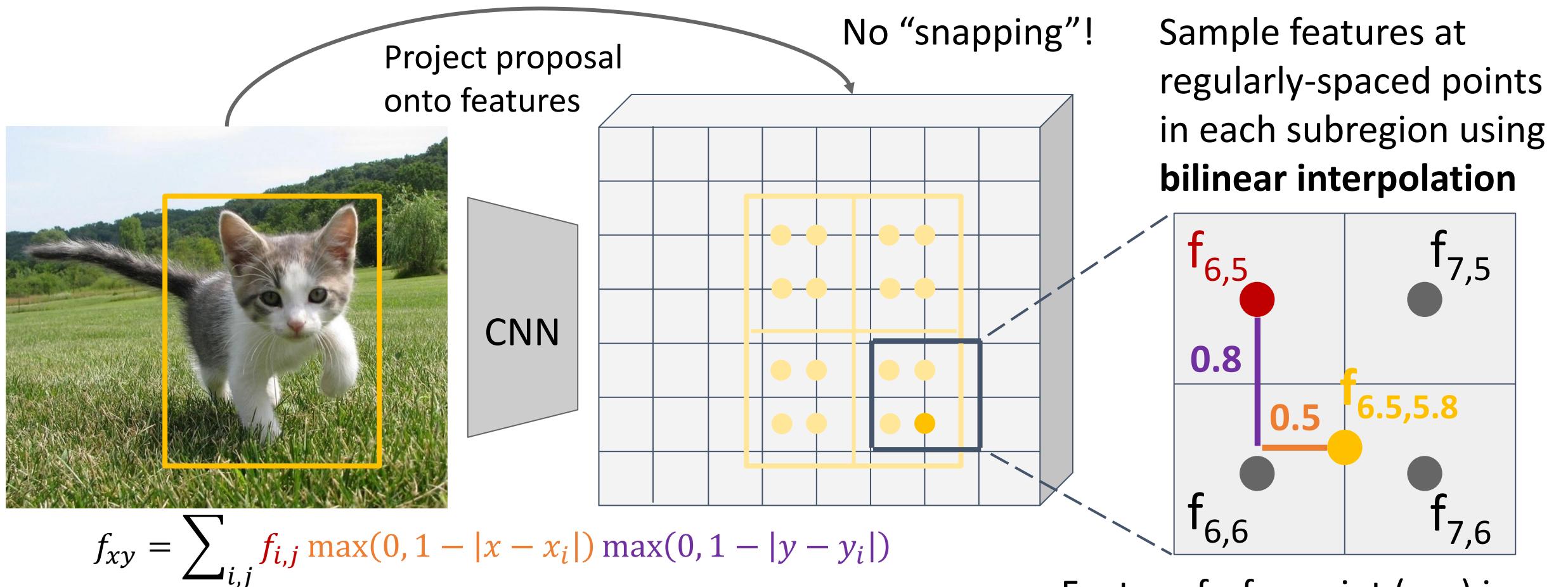
Cropping Features: RoI Align

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



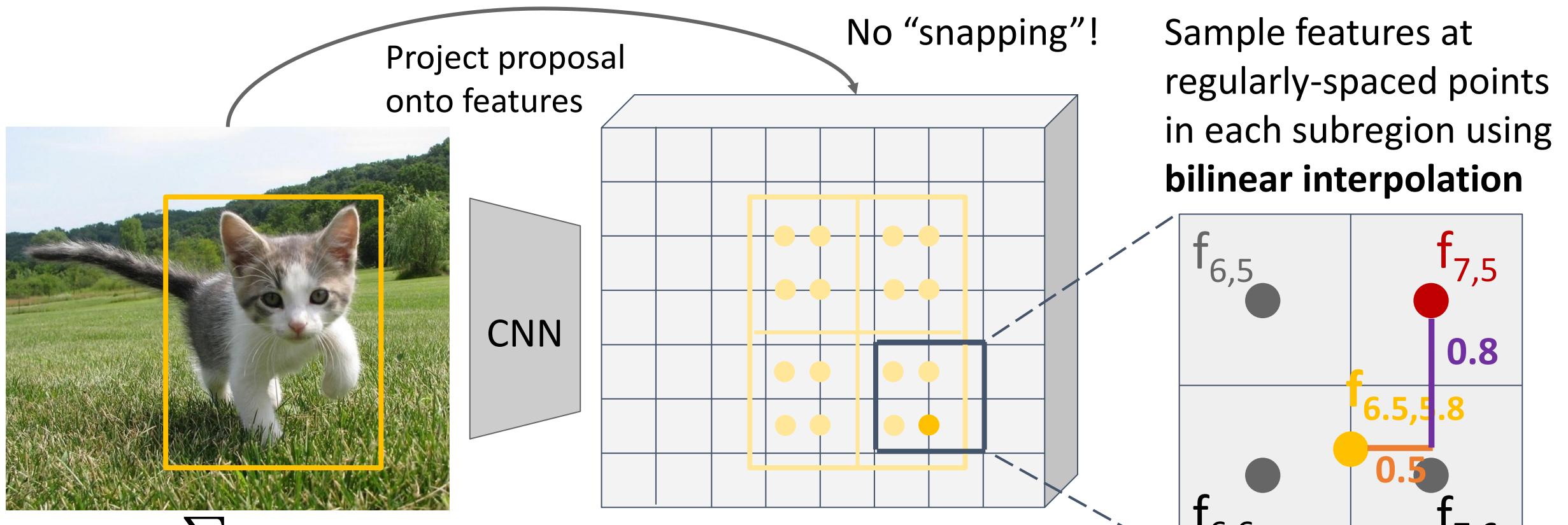
Feature f_{xy} for point (x, y) is a linear combination of features at its four neighboring grid cells:

Cropping Features: RoI Align



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

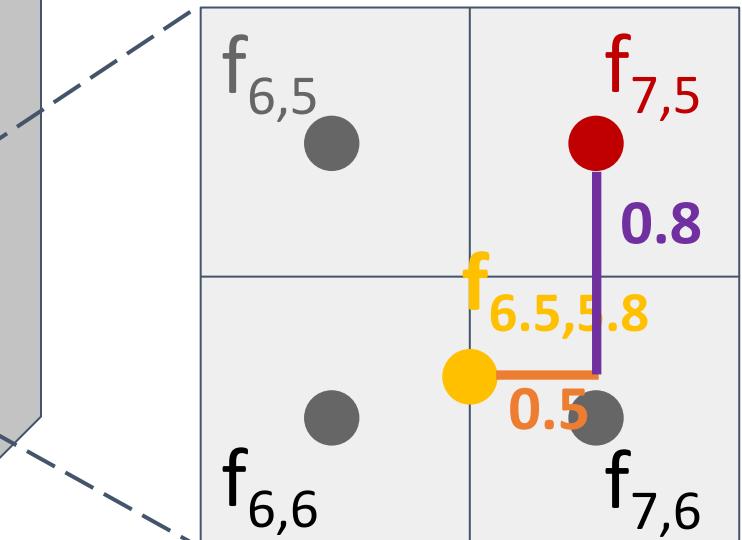
Cropping Features: RoI Align



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

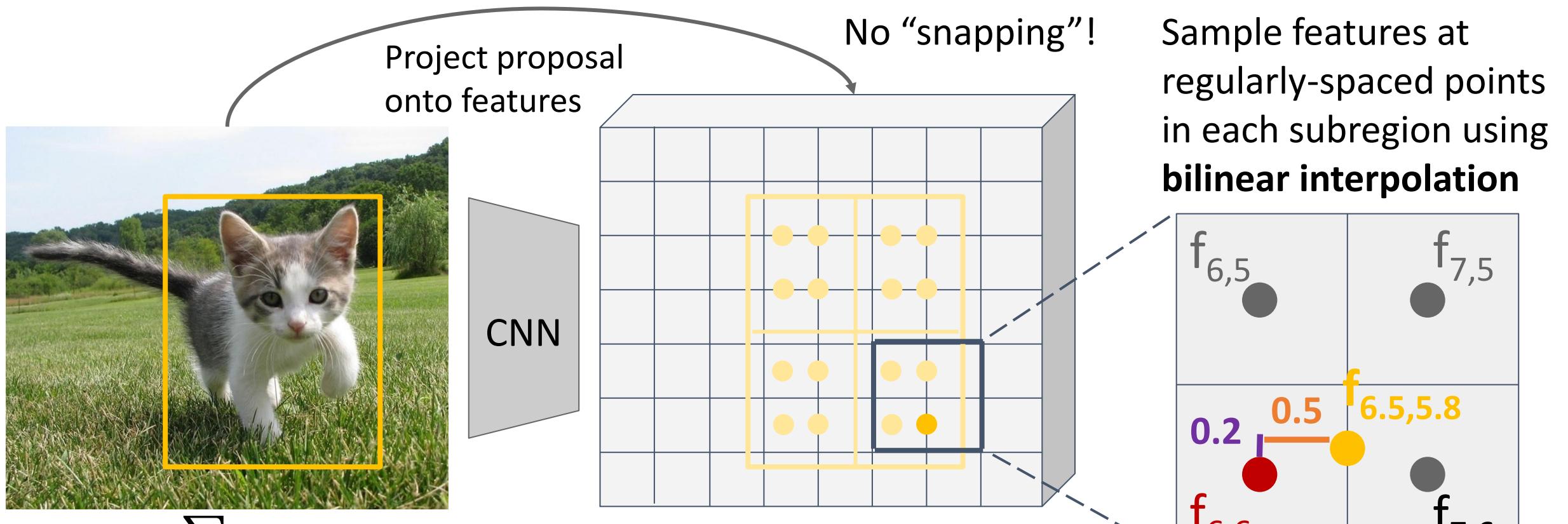
$$\begin{aligned} f_{6.5,5.8} &= (f_{6,5} * 0.5 * 0.2) + (f_{7,5} * 0.5 * 0.2) \\ &\quad + (f_{6,6} * 0.5 * 0.8) + (f_{7,6} * 0.5 * 0.8) \end{aligned}$$

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



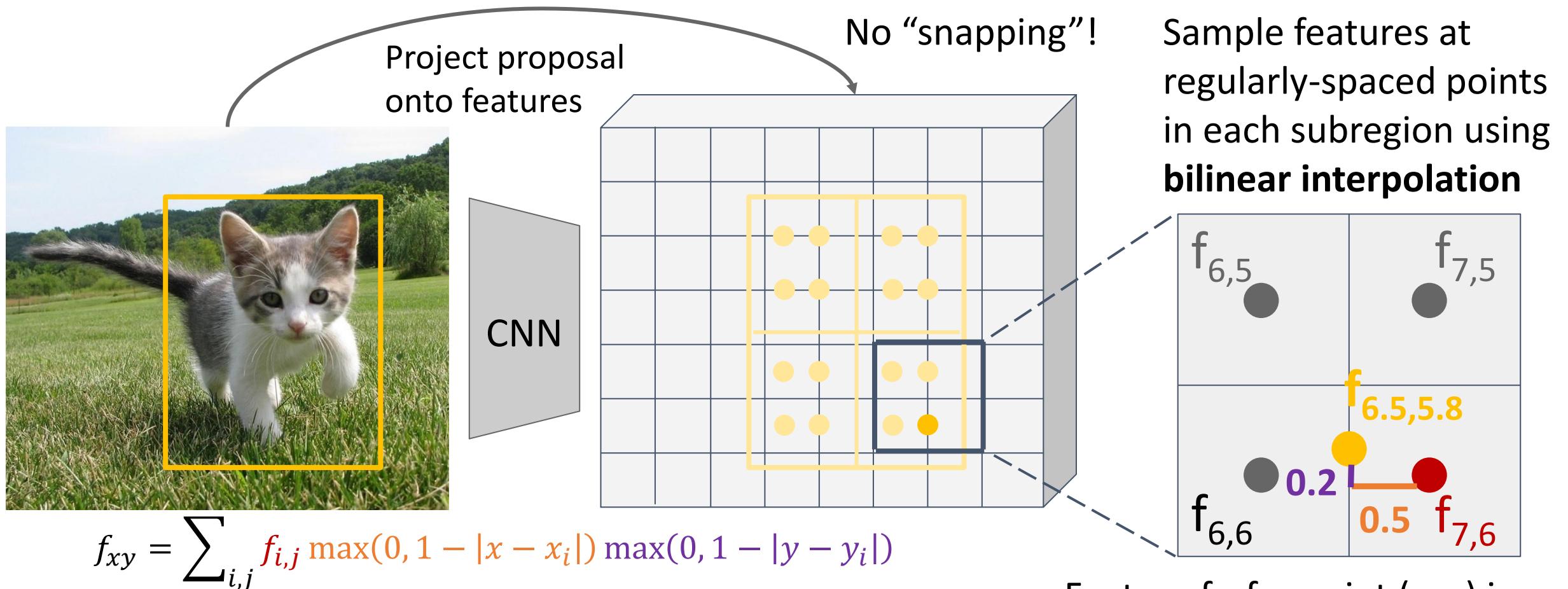
$$f_{xy} = \sum_{i,j} f_{i,j} \max(0, 1 - |x - x_i|) \max(0, 1 - |y - y_i|)$$

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

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

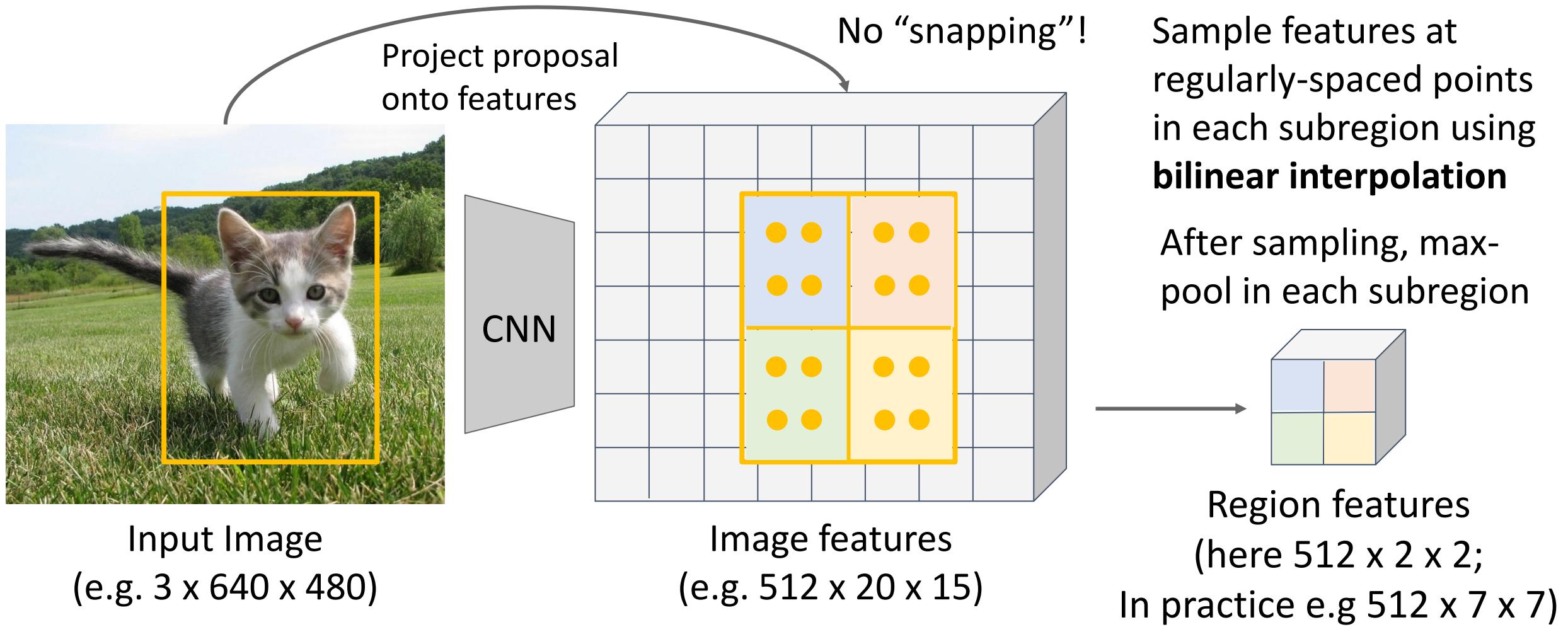
Feature f_{xy} for point (x, y) is a linear combination of features at its four neighboring grid cells:

Cropping Features: RoI Align



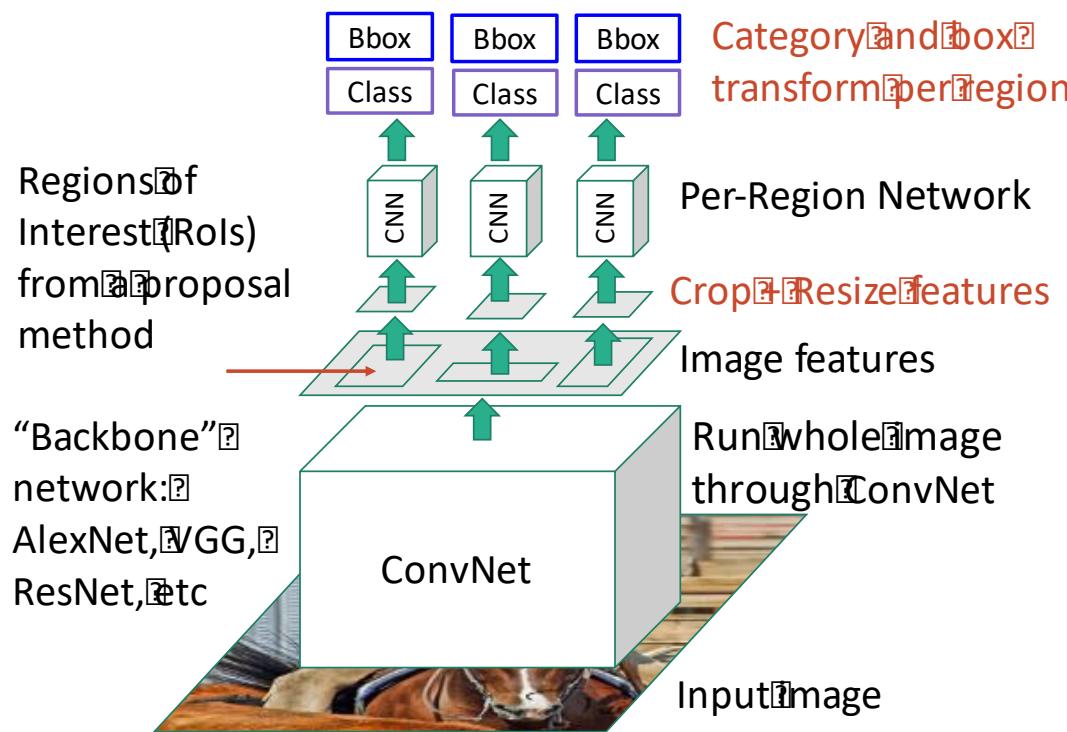
Feature f_{xy} for point (x, y) is a linear combination of features at its four neighboring grid cells:

Cropping Features: RoI Align

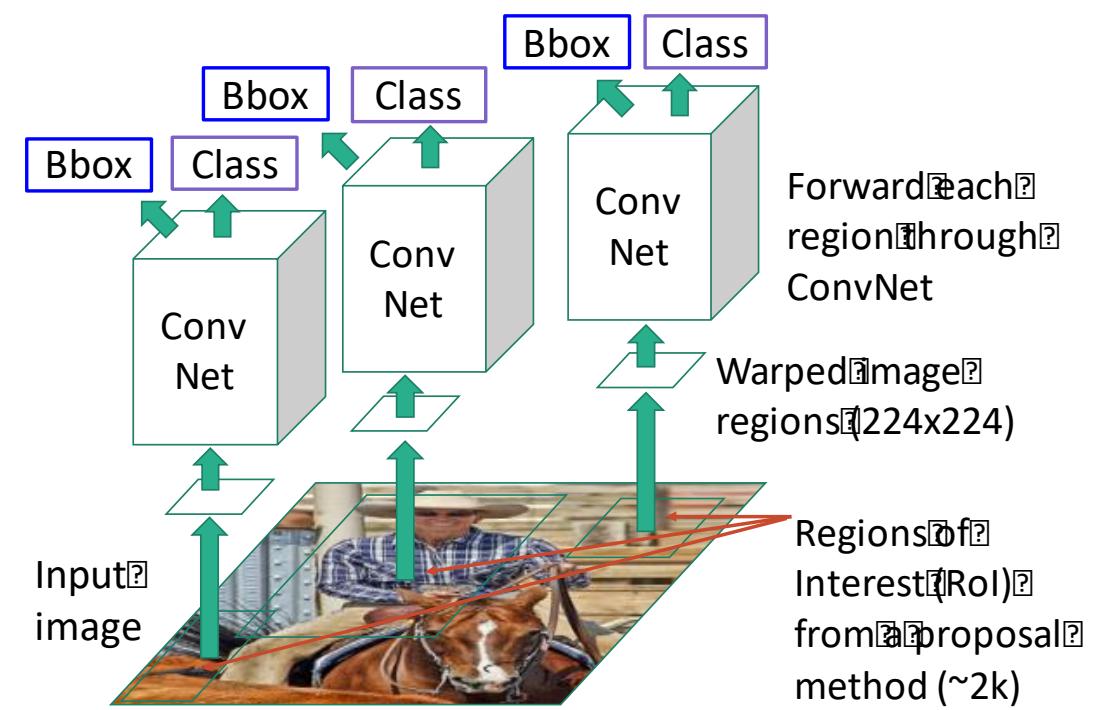


Fast R-CNN vs. “Slow” R-CNN

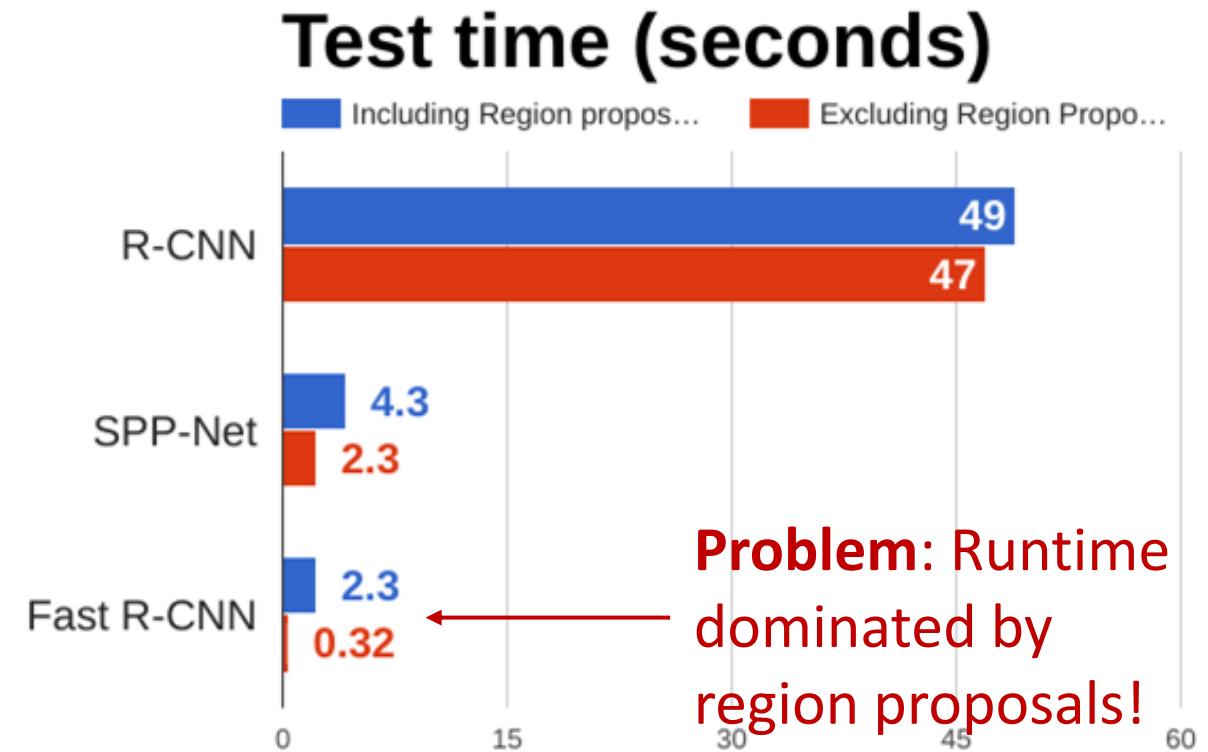
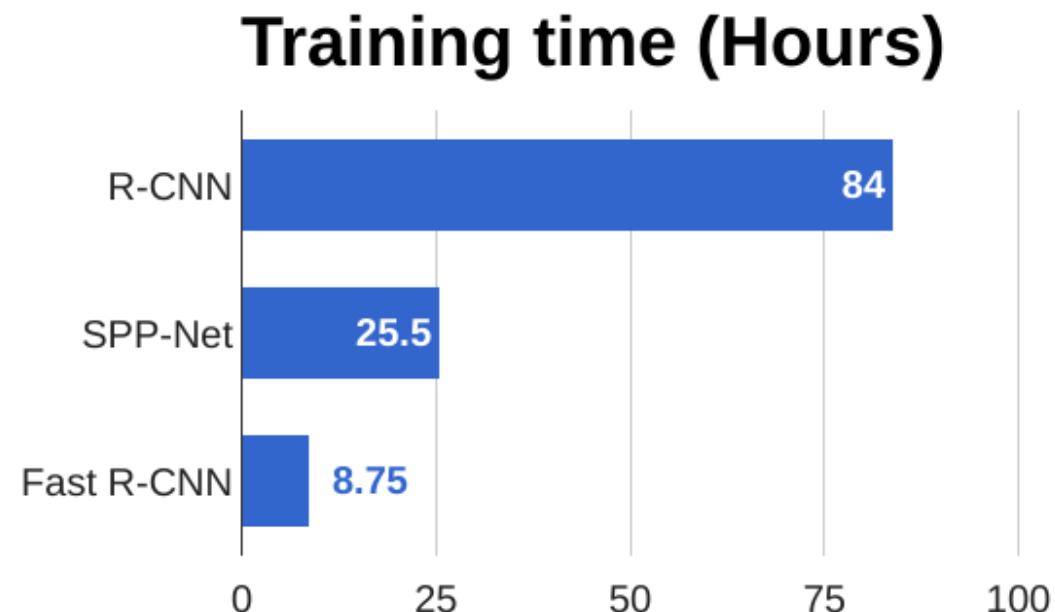
Fast R-CNN: Apply differentiable cropping to shared image features



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



Fast R-CNN vs. “Slow” R-CNN



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.

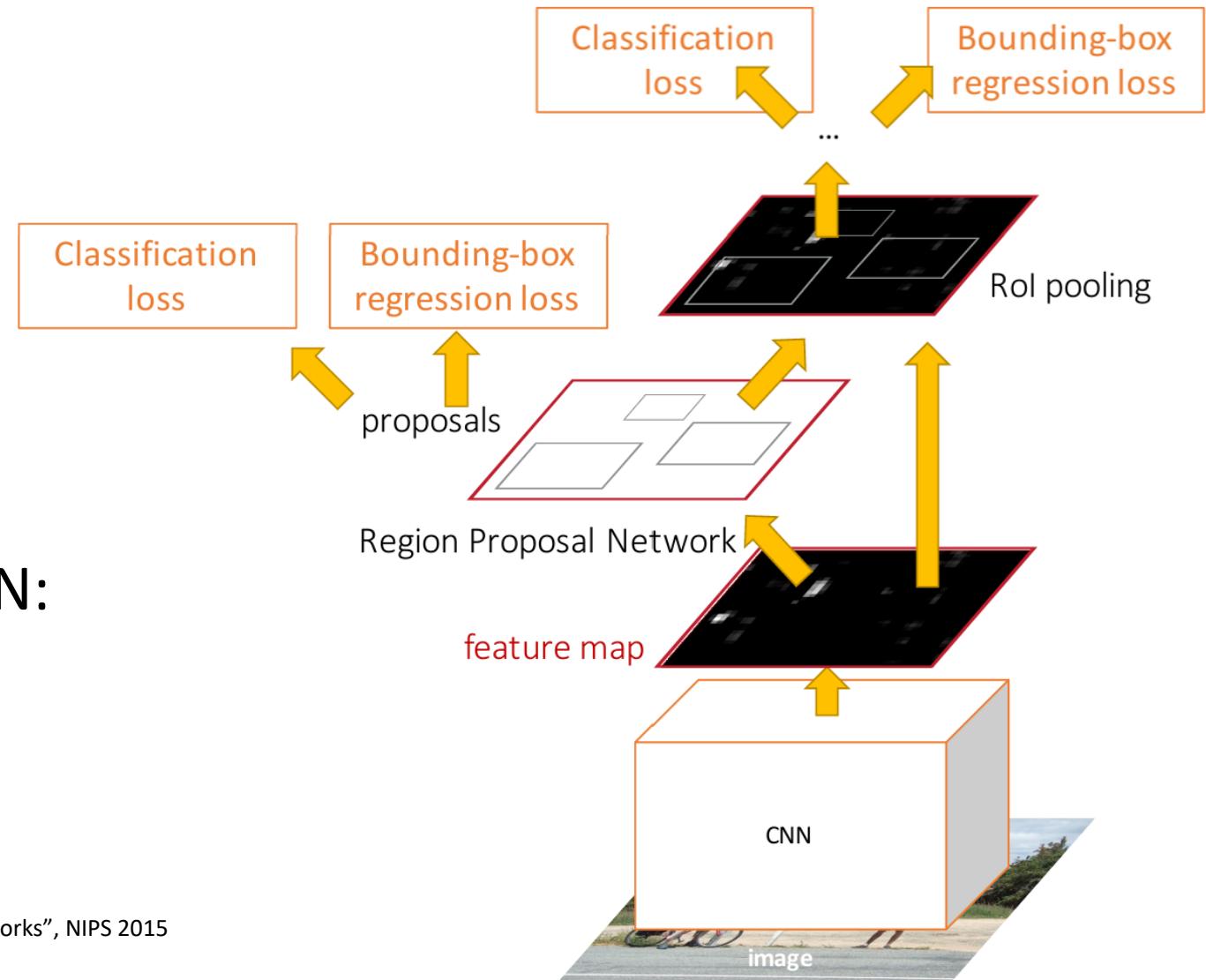
He et al, “Spatial pyramid pooling in deep convolutional networks for visual recognition”, ECCV 2014

Girshick, “Fast R-CNN”, ICCV 2015

Fasterer R-CNN: Learnable Region Proposals

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

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



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

Region Proposal Network (RPN)

Run backbone CNN to get
features aligned to input image



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

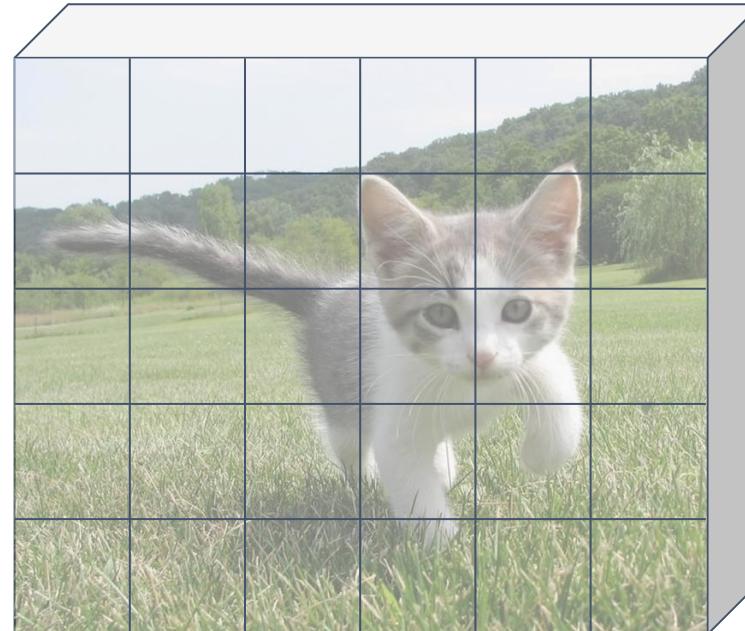


Image features
(e.g. $512 \times 20 \times 15$)

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

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

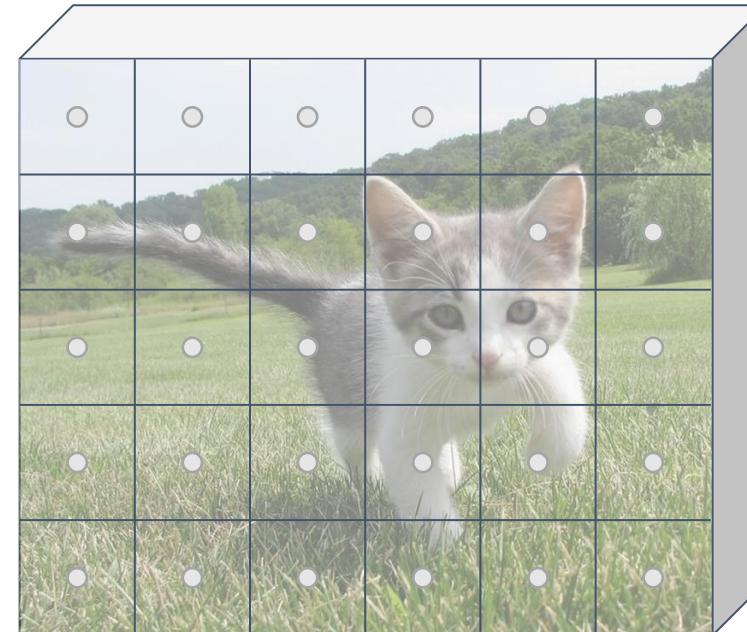
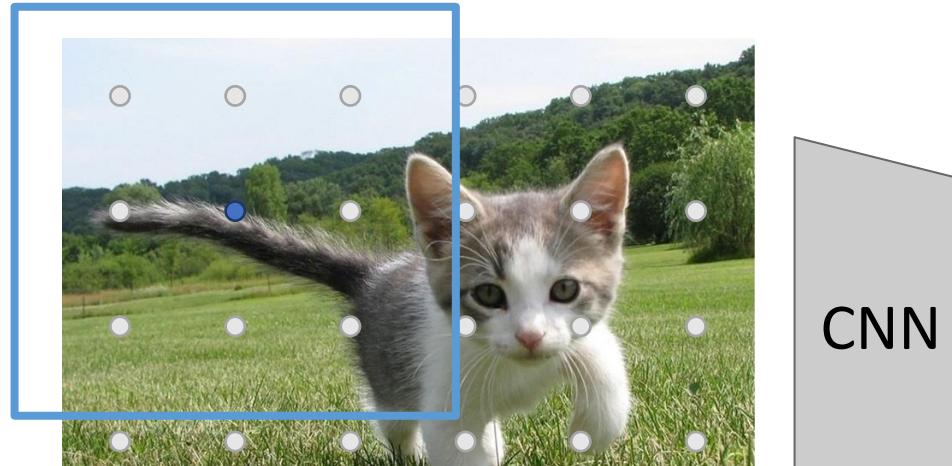


Image features
(e.g. $512 \times 20 \times 15$)

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

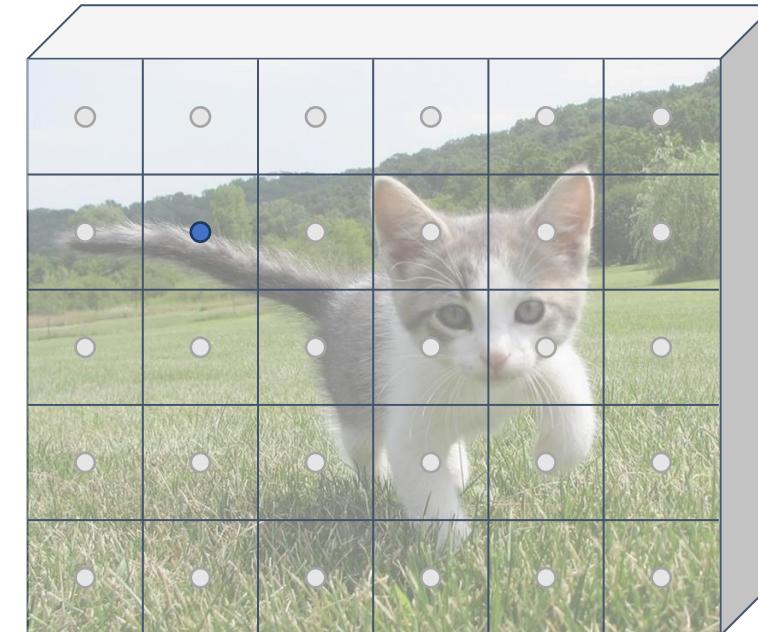
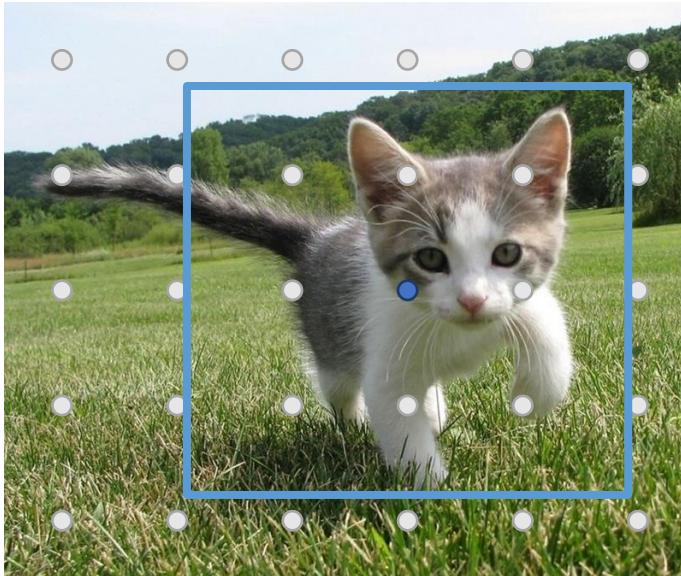


Image features
(e.g. $512 \times 20 \times 15$)

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

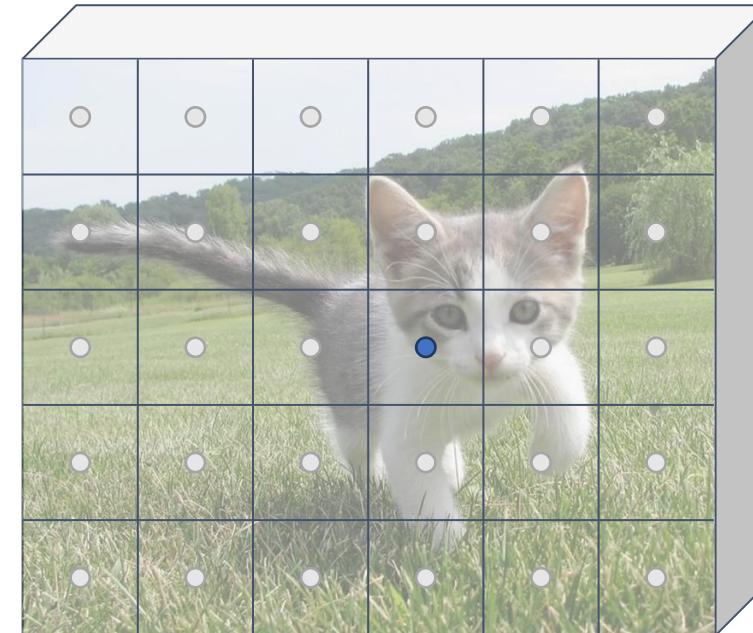
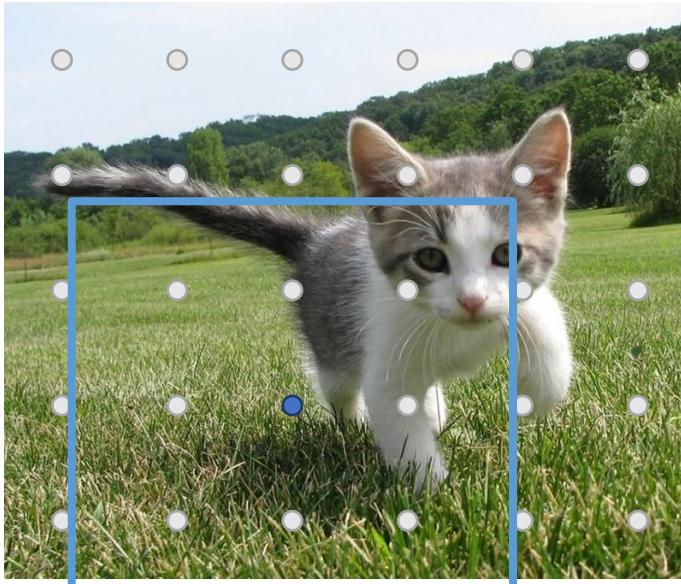


Image features
(e.g. $512 \times 20 \times 15$)

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

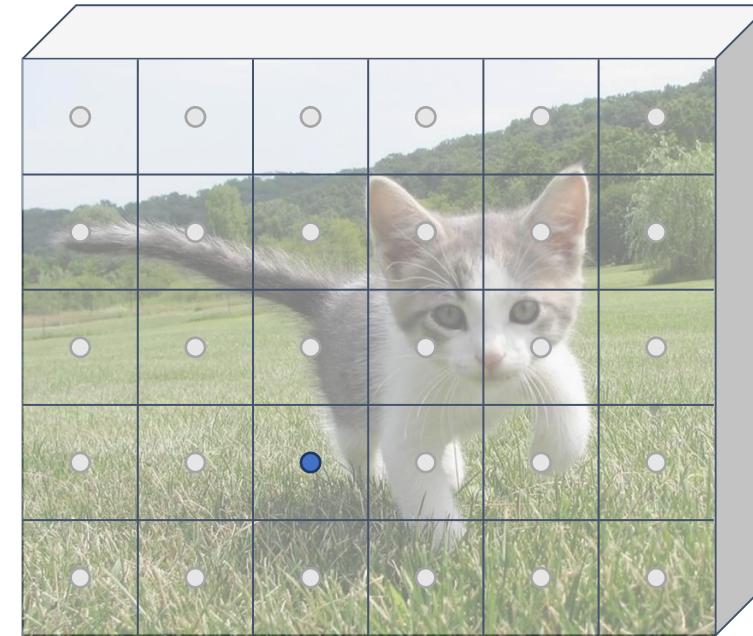
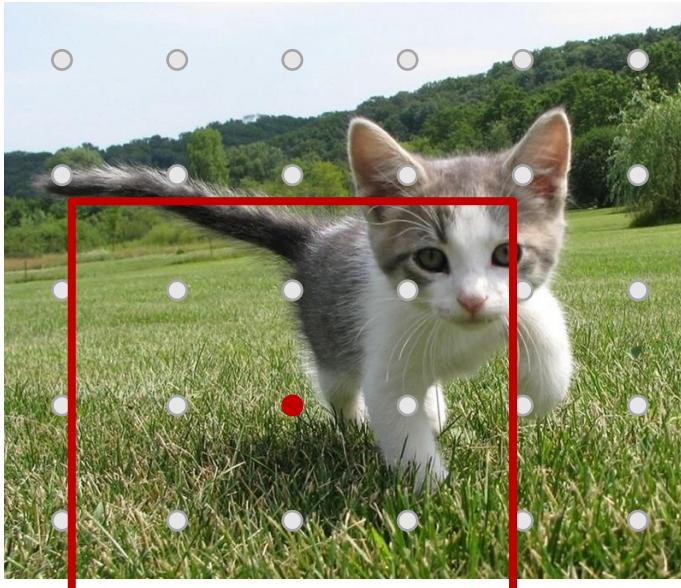


Image features
(e.g. $512 \times 20 \times 15$)

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

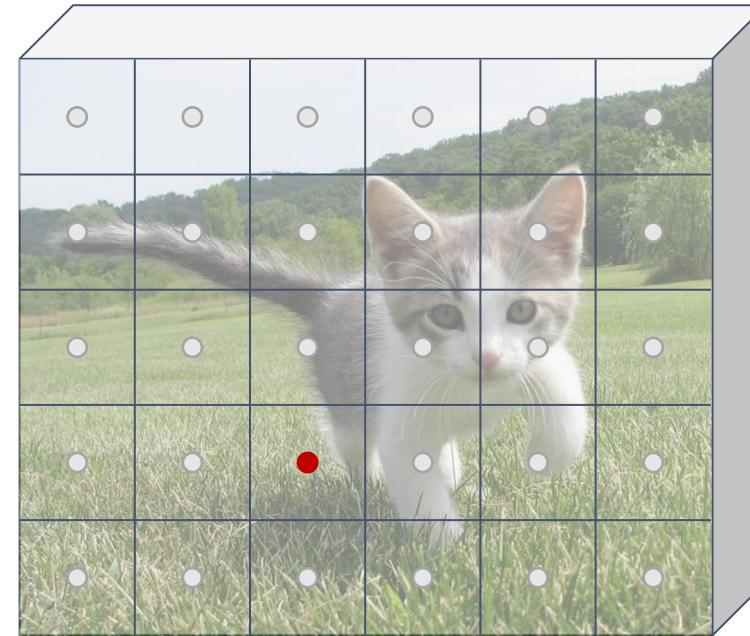


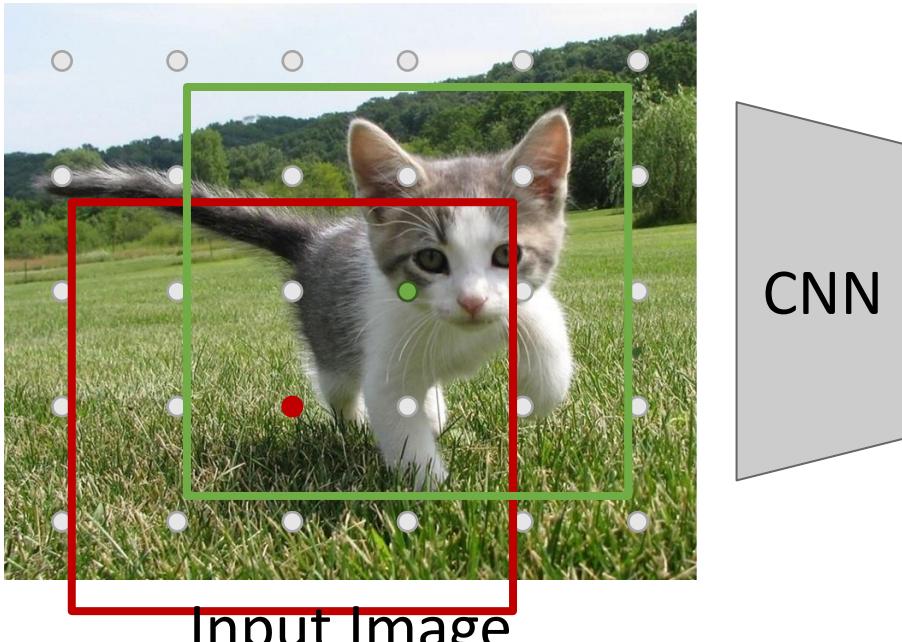
Image features
(e.g. $512 \times 20 \times 15$)

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

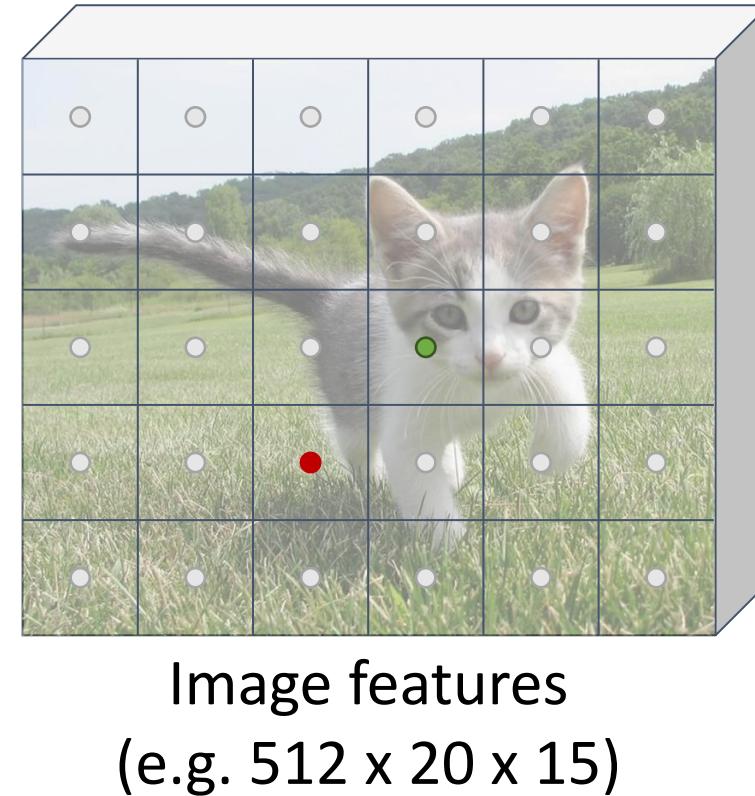
Classify each anchor as **positive (object)** or **negative (no object)**

Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



Each feature corresponds to a point in the input

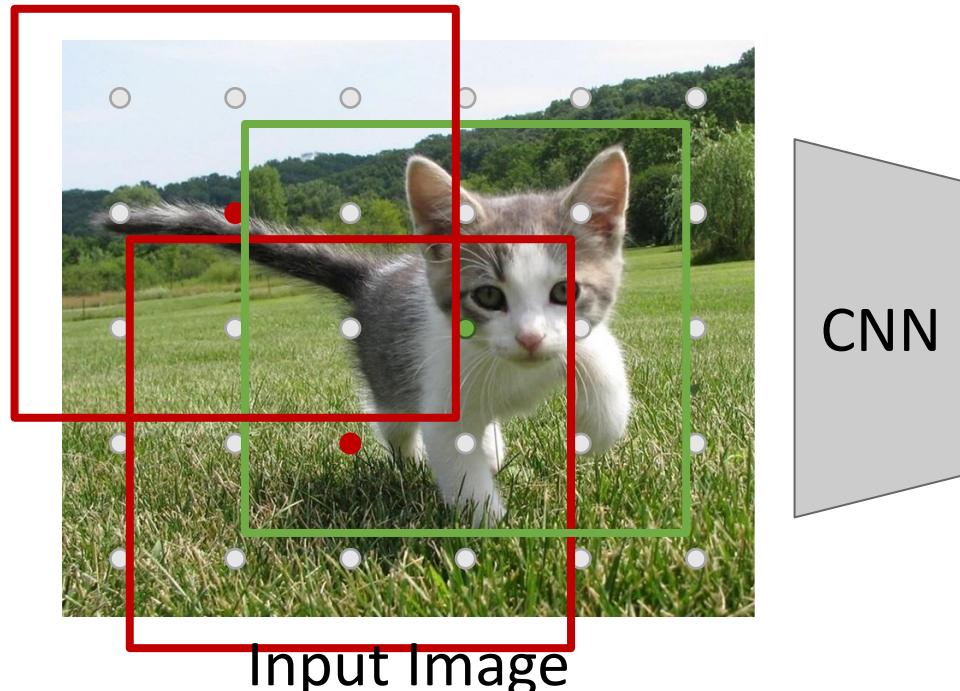


Imagine an **anchor box** of fixed size at each point in the feature map

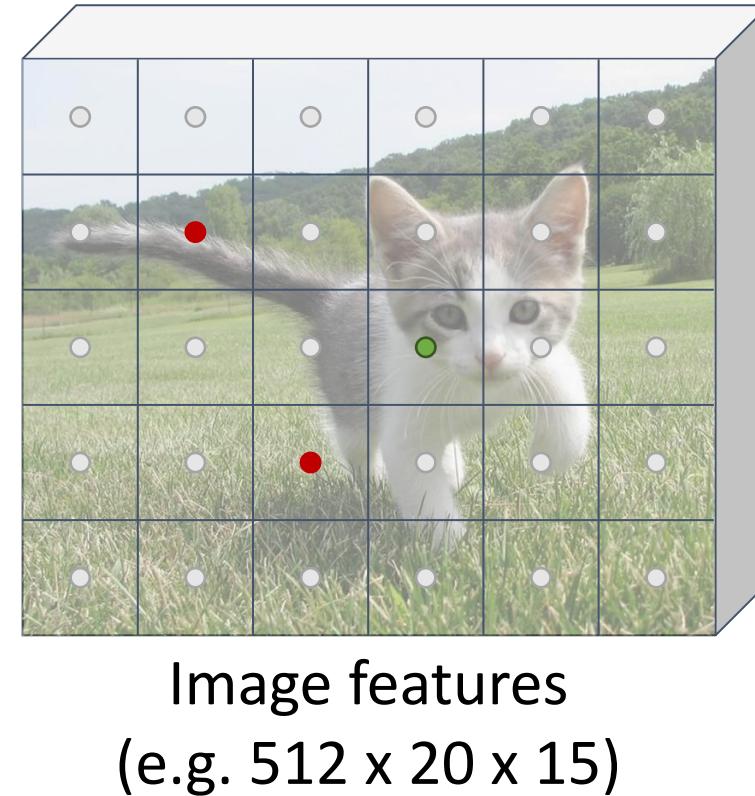
Classify each anchor as **positive (object)** or **negative (no object)**

Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



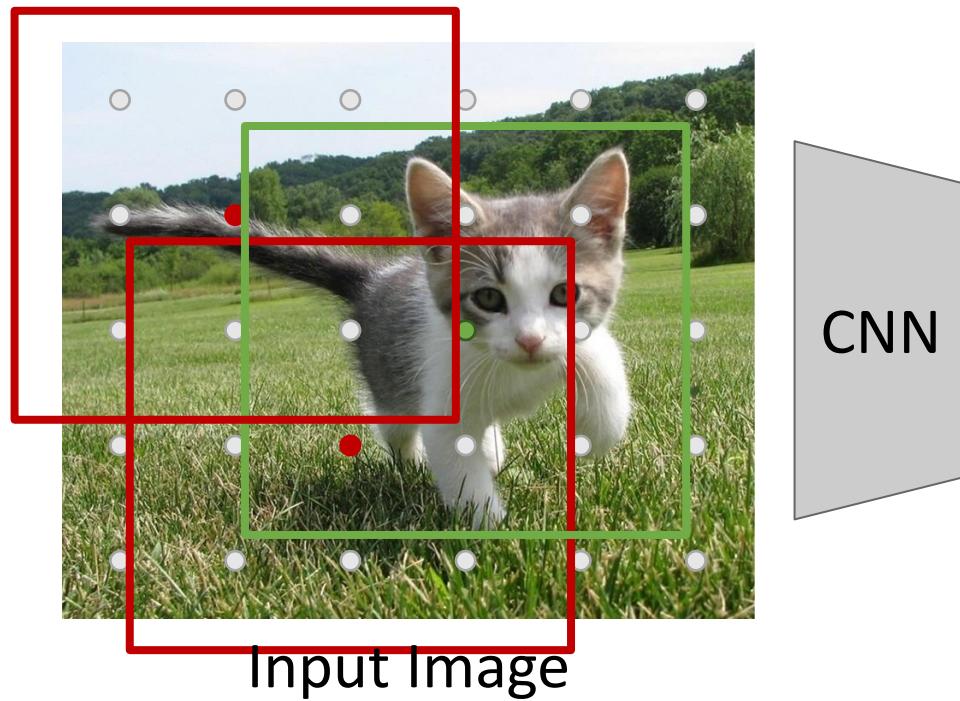
Each feature corresponds to a point in the input



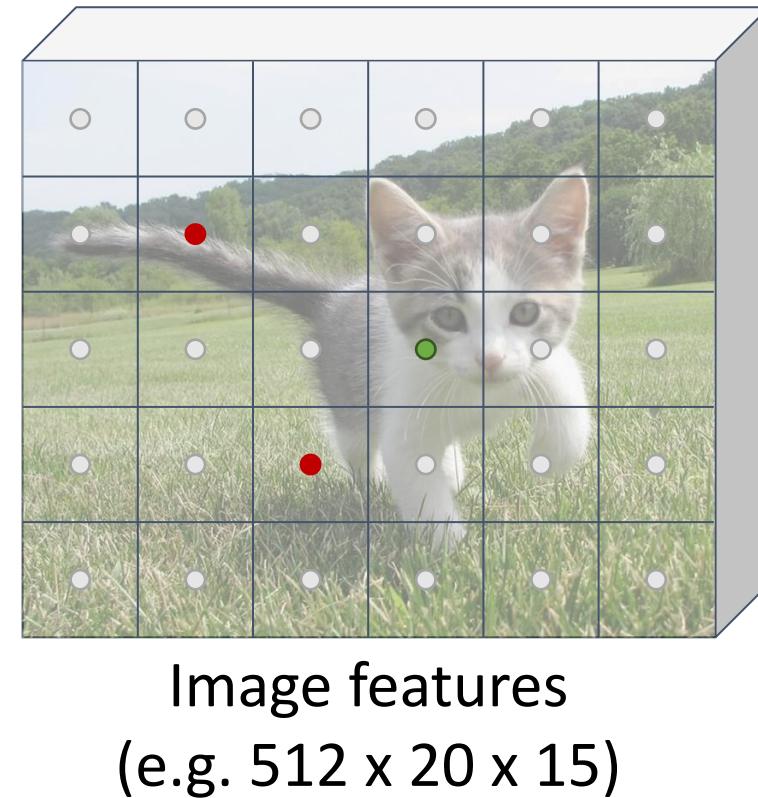
Imagine an **anchor box** of fixed size at each point in the feature map

Region Proposal Network (RPN)

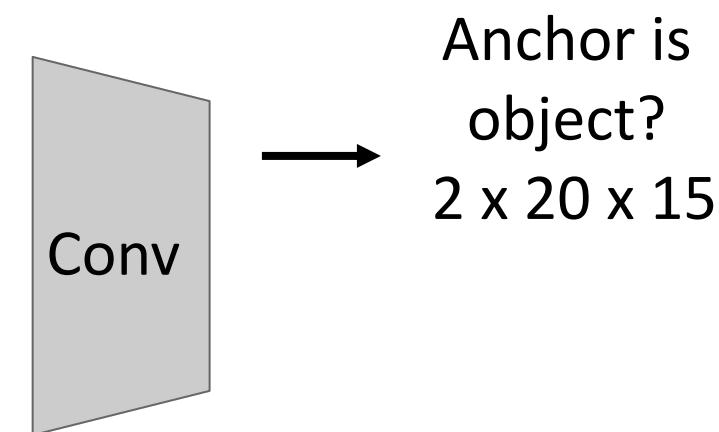
Run backbone CNN to get features aligned to input image



Each feature corresponds to a point in the input



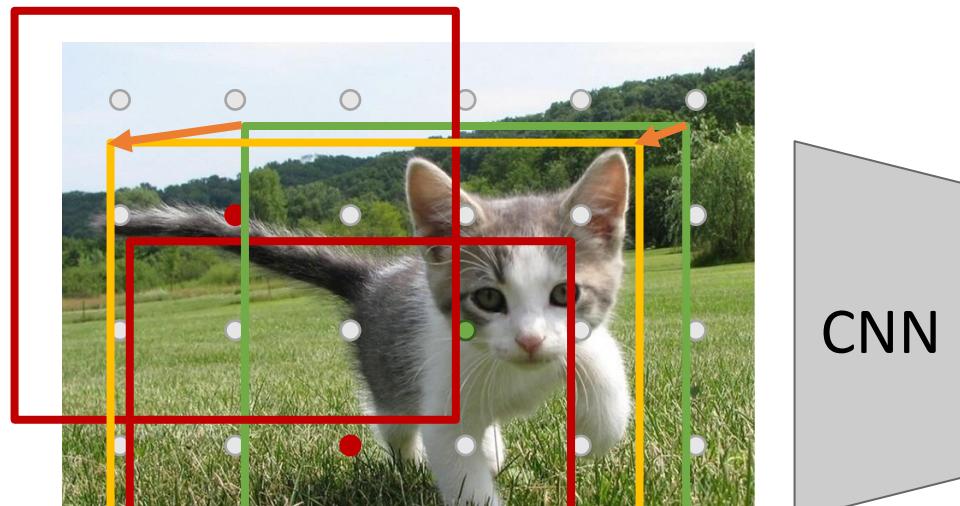
Predict object vs. not object scores for all anchors with a conv layer (512 input filters, 2 output filters)



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

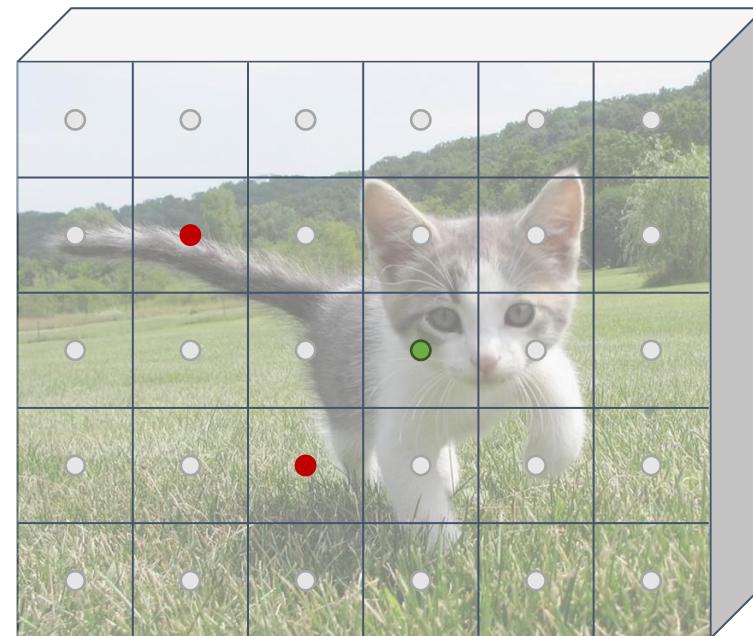
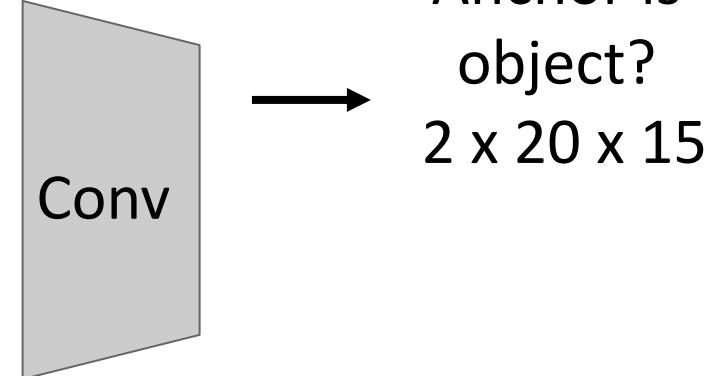


Image features
(e.g. $512 \times 20 \times 15$)

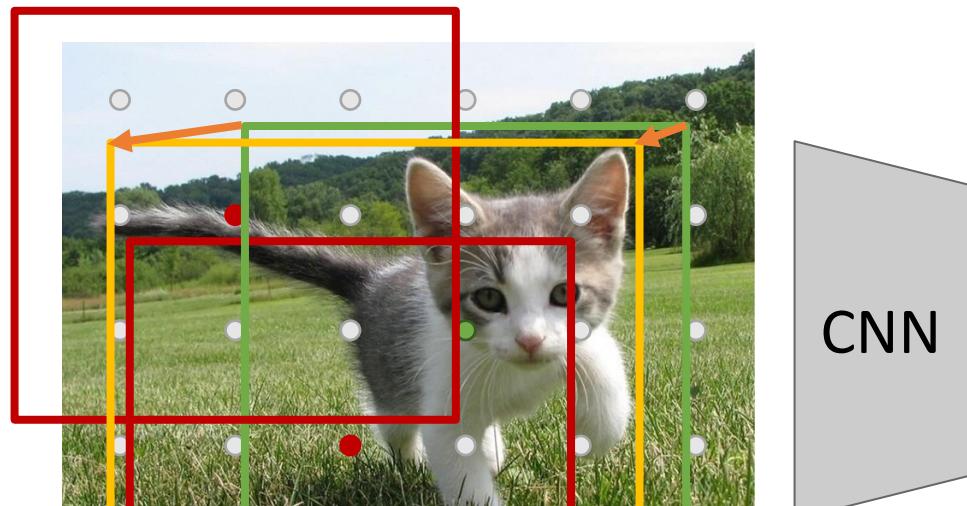
For **positive anchors**, also predict a **transform** that converting the anchor to the **GT box** (like R-CNN)



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

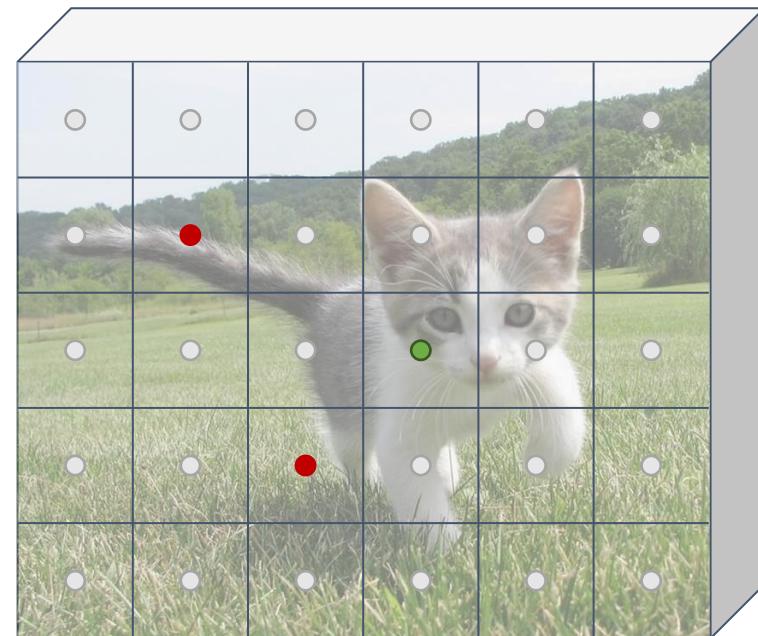
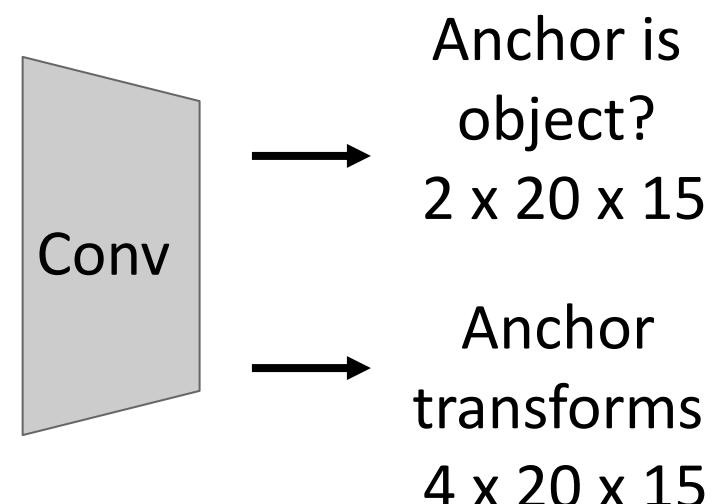


Image features
(e.g. $512 \times 20 \times 15$)

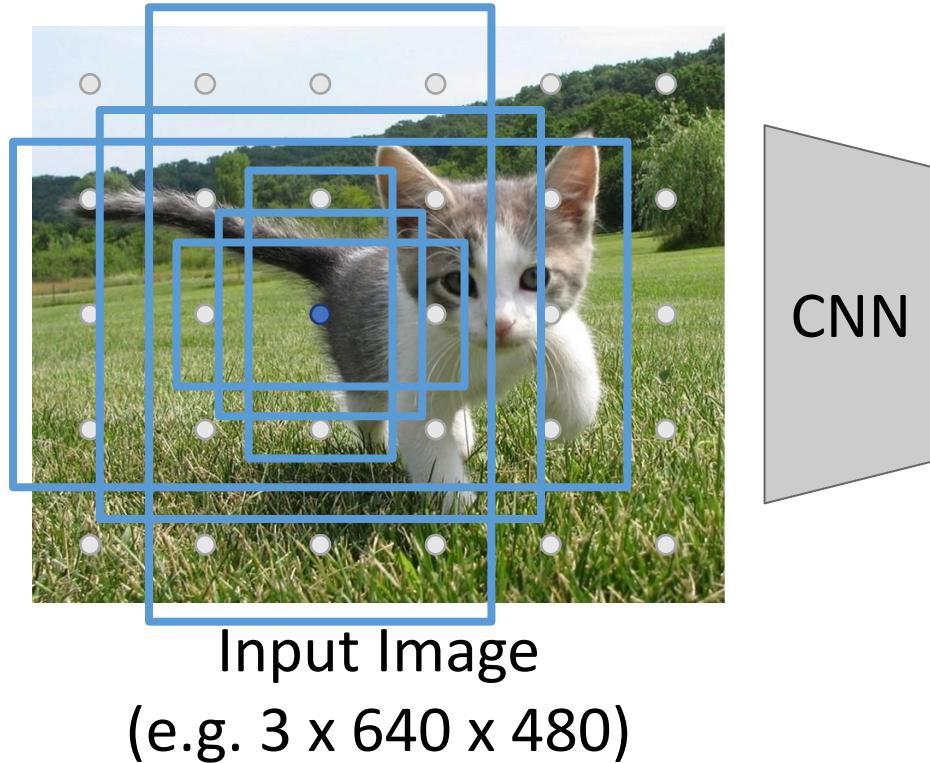
For **positive anchors**, also predict a **transform** that converting the anchor to the **GT box** (like R-CNN)
Predict transforms with conv



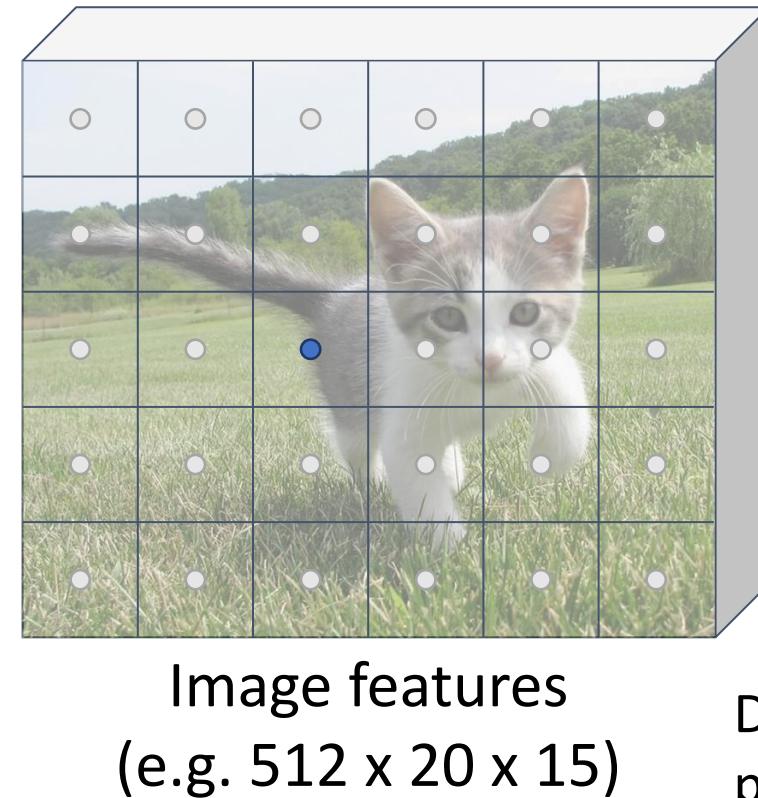
Classify each anchor as **positive (object)** or **negative (no object)**

Region Proposal Network (RPN)

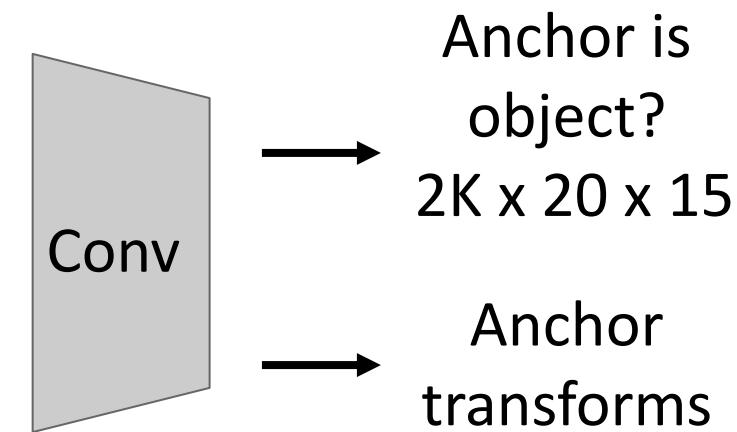
Run backbone CNN to get features aligned to input image



Each feature corresponds to a point in the input



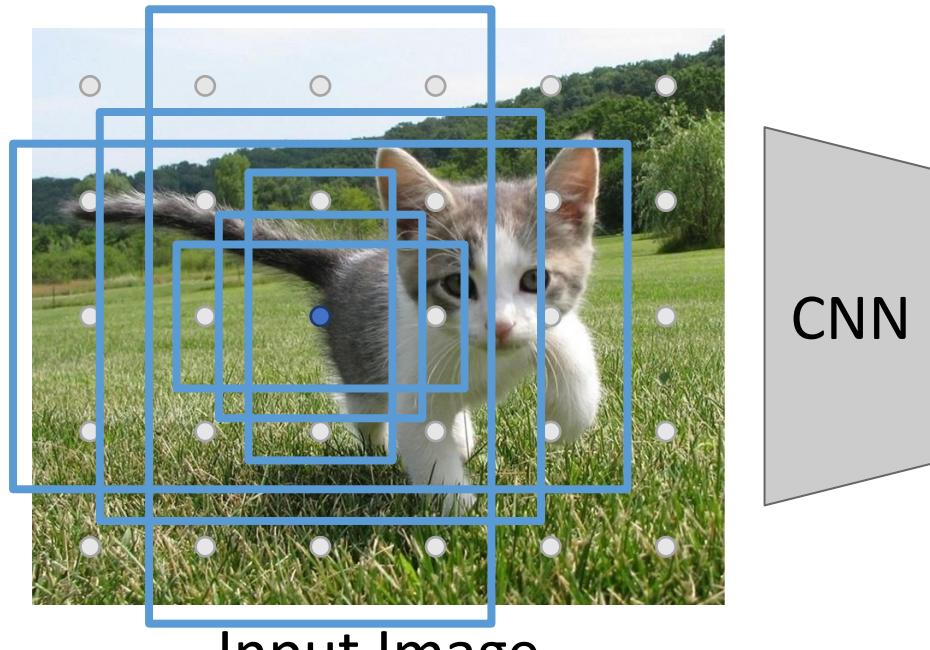
In practice: Rather than using one anchor per point, instead consider K different anchors with different size and scale (here $K = 6$)



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. $3 \times 640 \times 480$)

Each feature corresponds to a point in the input

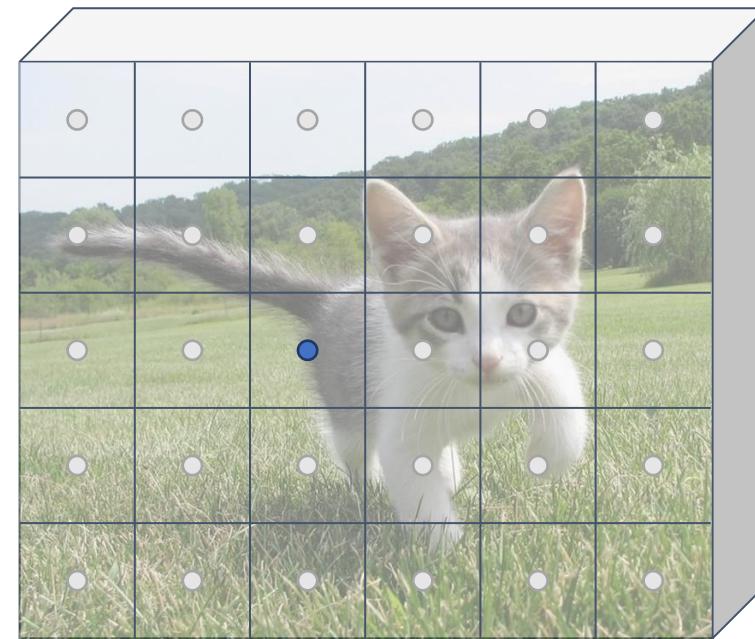
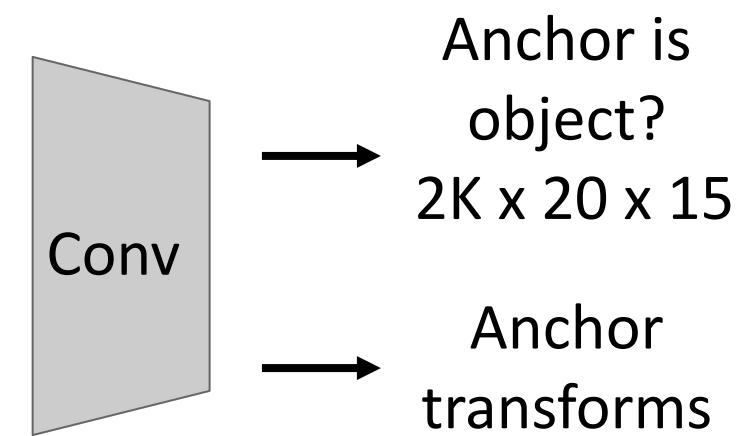


Image features
(e.g. $512 \times 20 \times 15$)

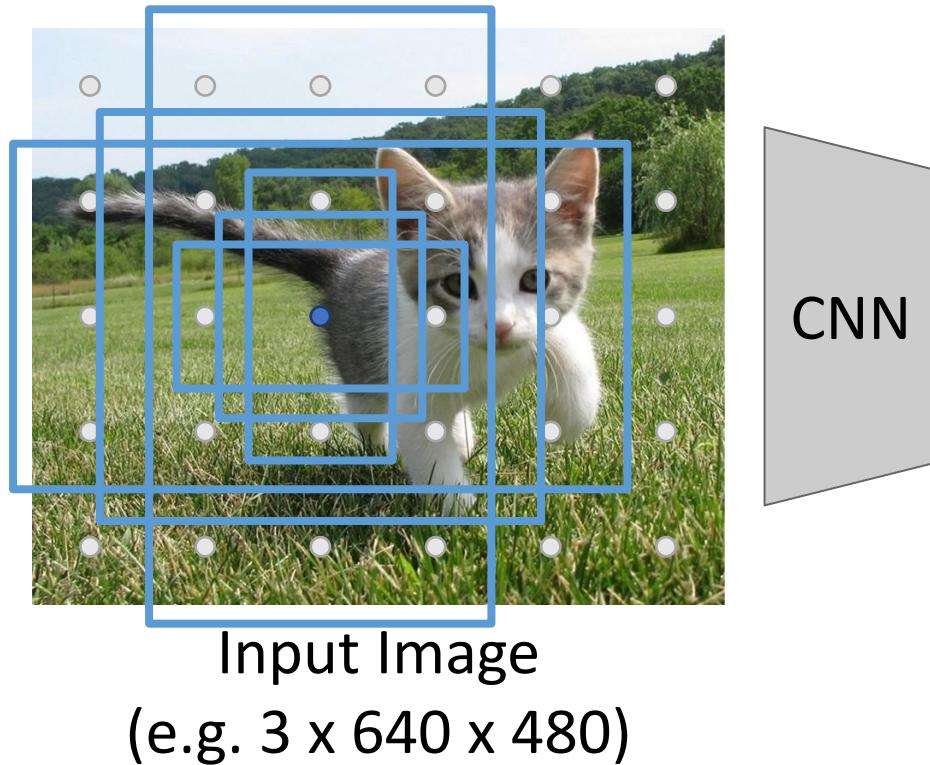
In practice: Rather than using one anchor per point, instead consider K different anchors with different size and scale (here $K = 6$)



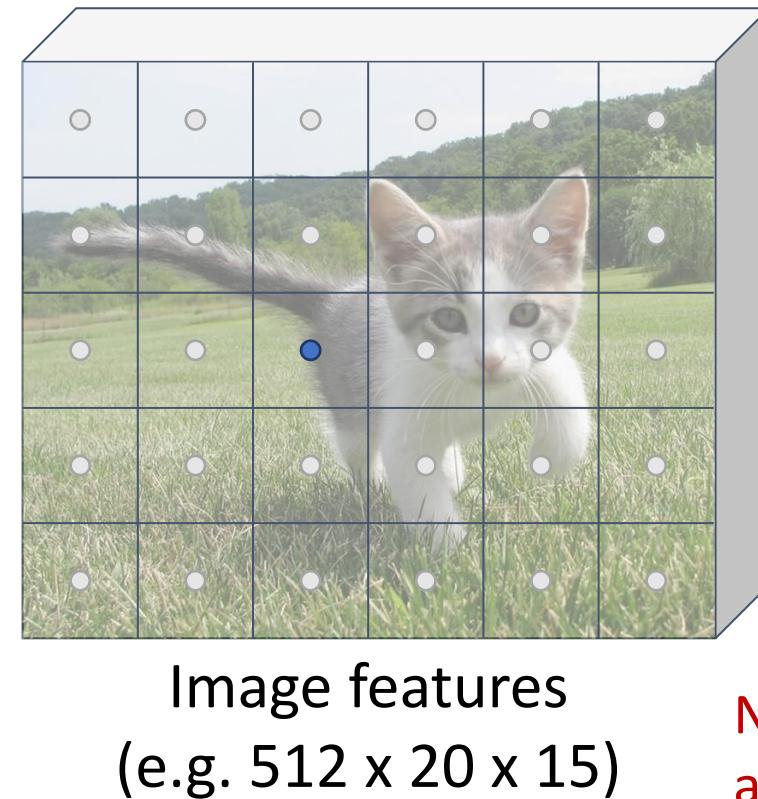
Positive anchors: ≥ 0.7 IoU with some GT box (plus highest IoU to each GT)

Region Proposal Network (RPN)

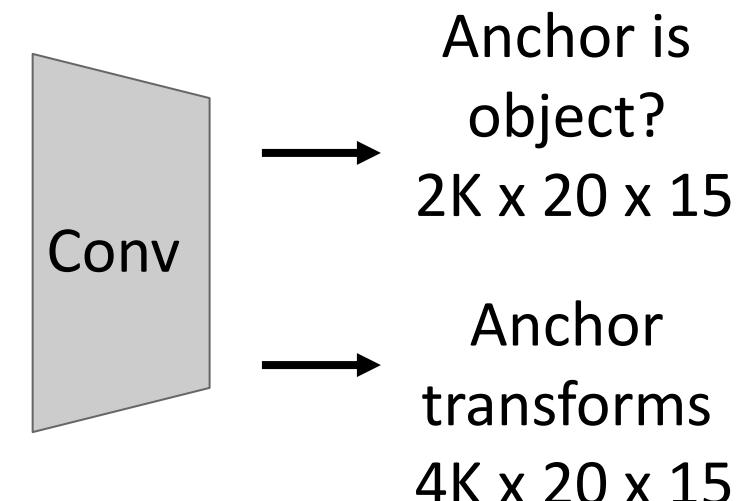
Run backbone CNN to get features aligned to input image



Each feature corresponds to a point in the input



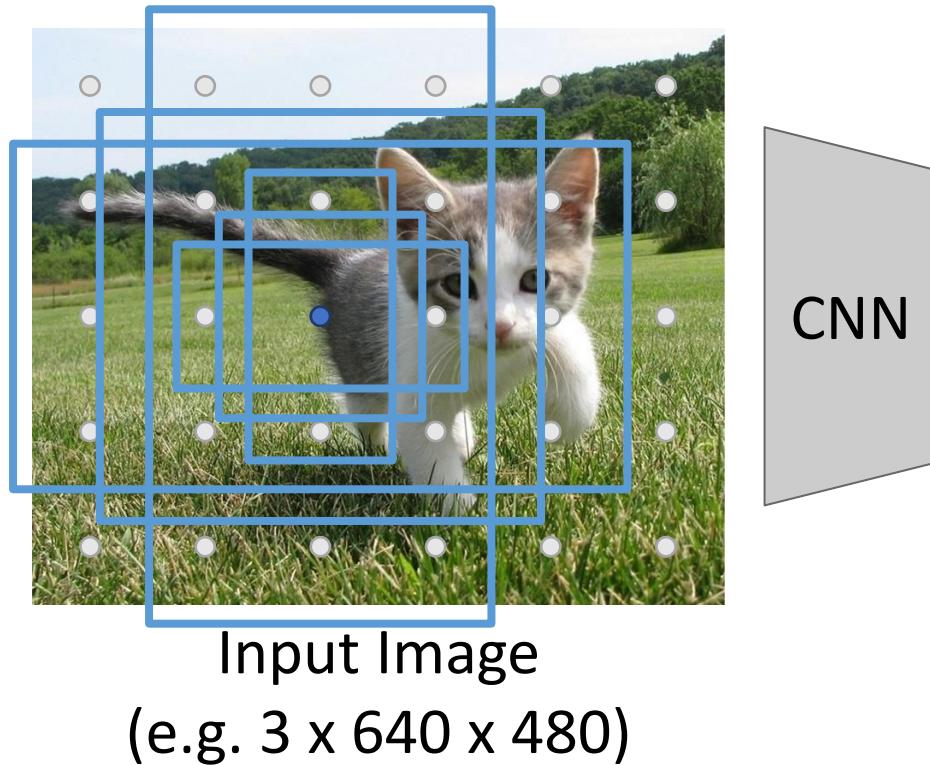
In practice: Rather than using one anchor per point, instead consider K different anchors with different size and scale (here $K = 6$)



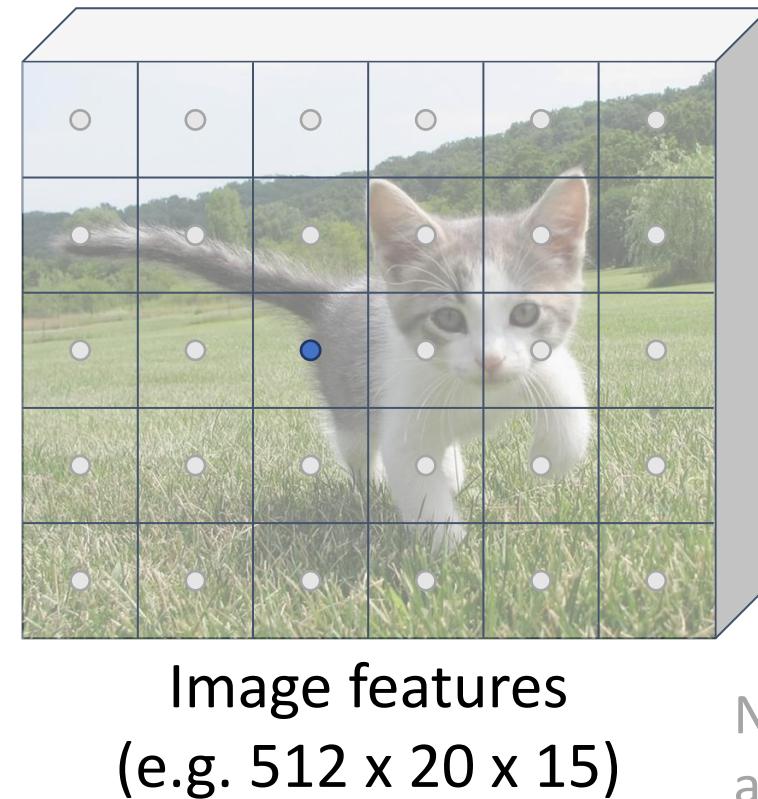
Negative anchors: < 0.3 IoU with all GT boxes. Don't supervise transforms for negative boxes.

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 $K = 6$)

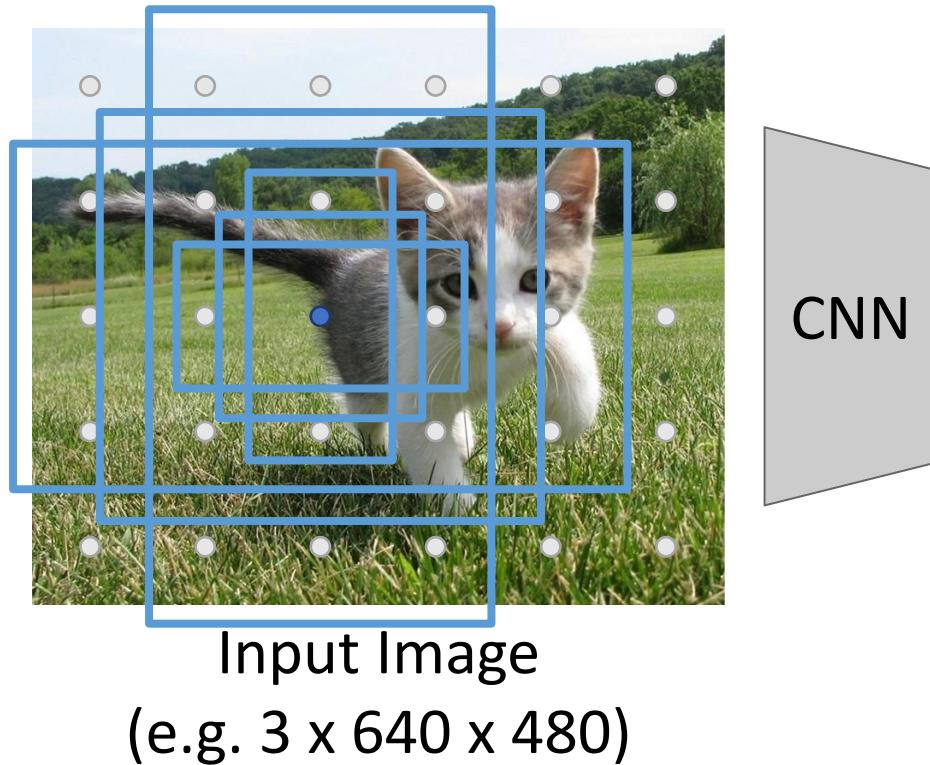
Anchor is object?
 $2K \times 20 \times 15$

Anchor transforms
 $4K \times 20 \times 15$

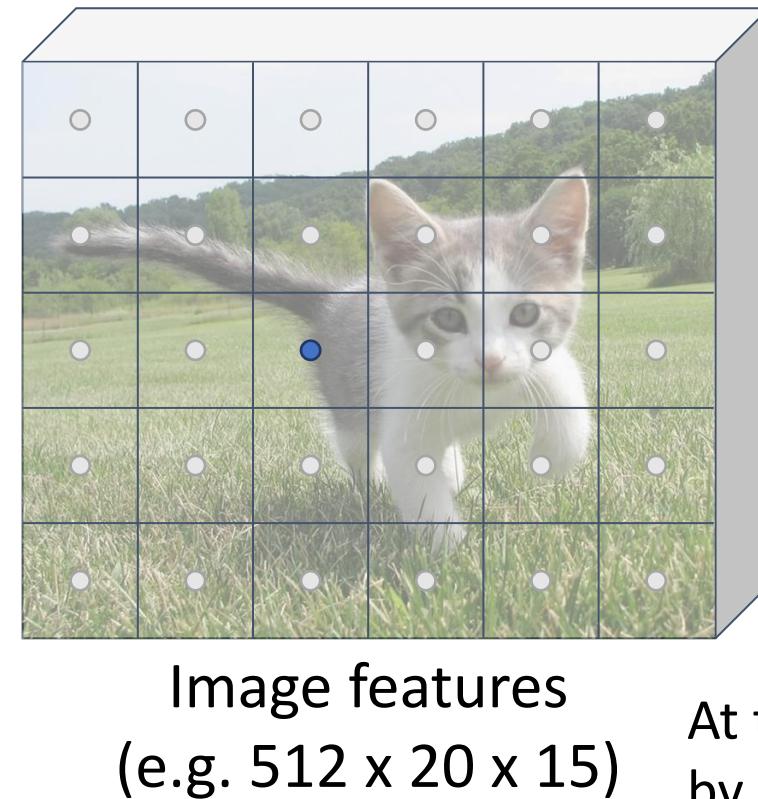
Neutral anchors: between 0.3 and 0.7 IoU with all GT boxes; ignored during training

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 $K = 6$)

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

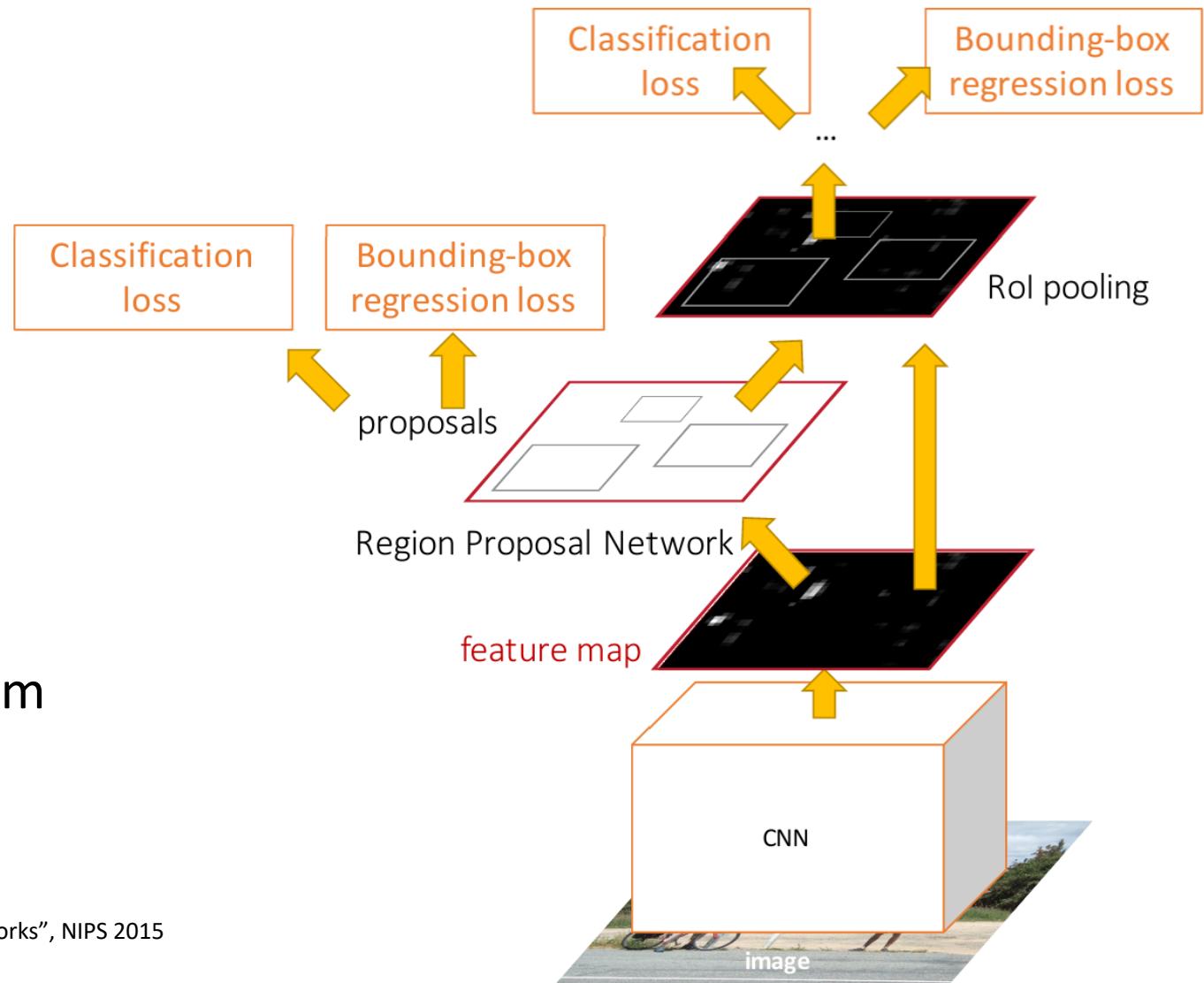
Anchor transforms
 $4K \times 20 \times 15$

At test-time, sort all $K \times 20 \times 15$ boxes by their positive score, take top 300 as our region proposals

Fasterer R-CNN: Learnable Region Proposals

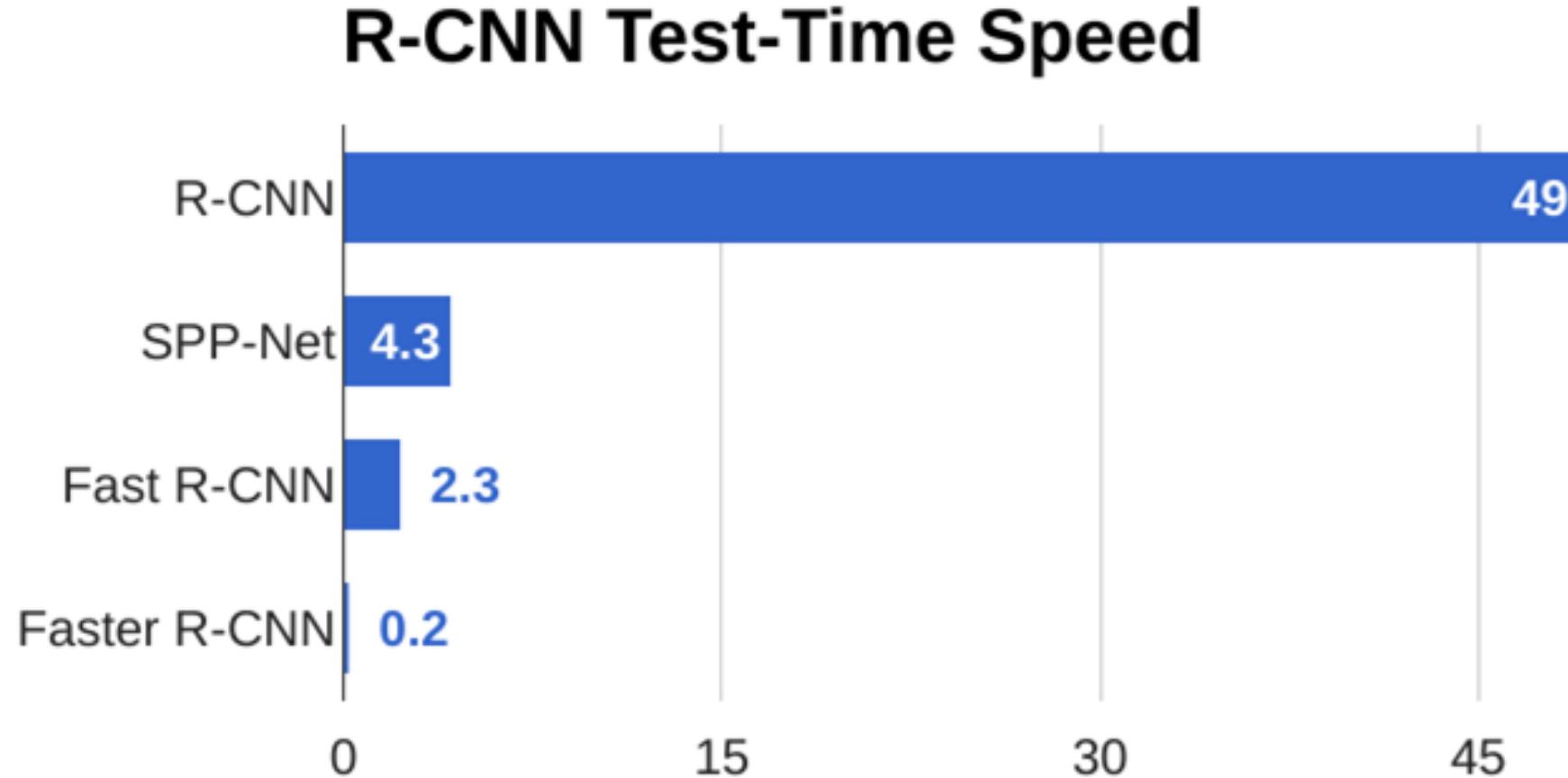
Jointly train with 4 losses:

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



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

Fasterer R-CNN: Learnable Region Proposals



Object Detection Progress

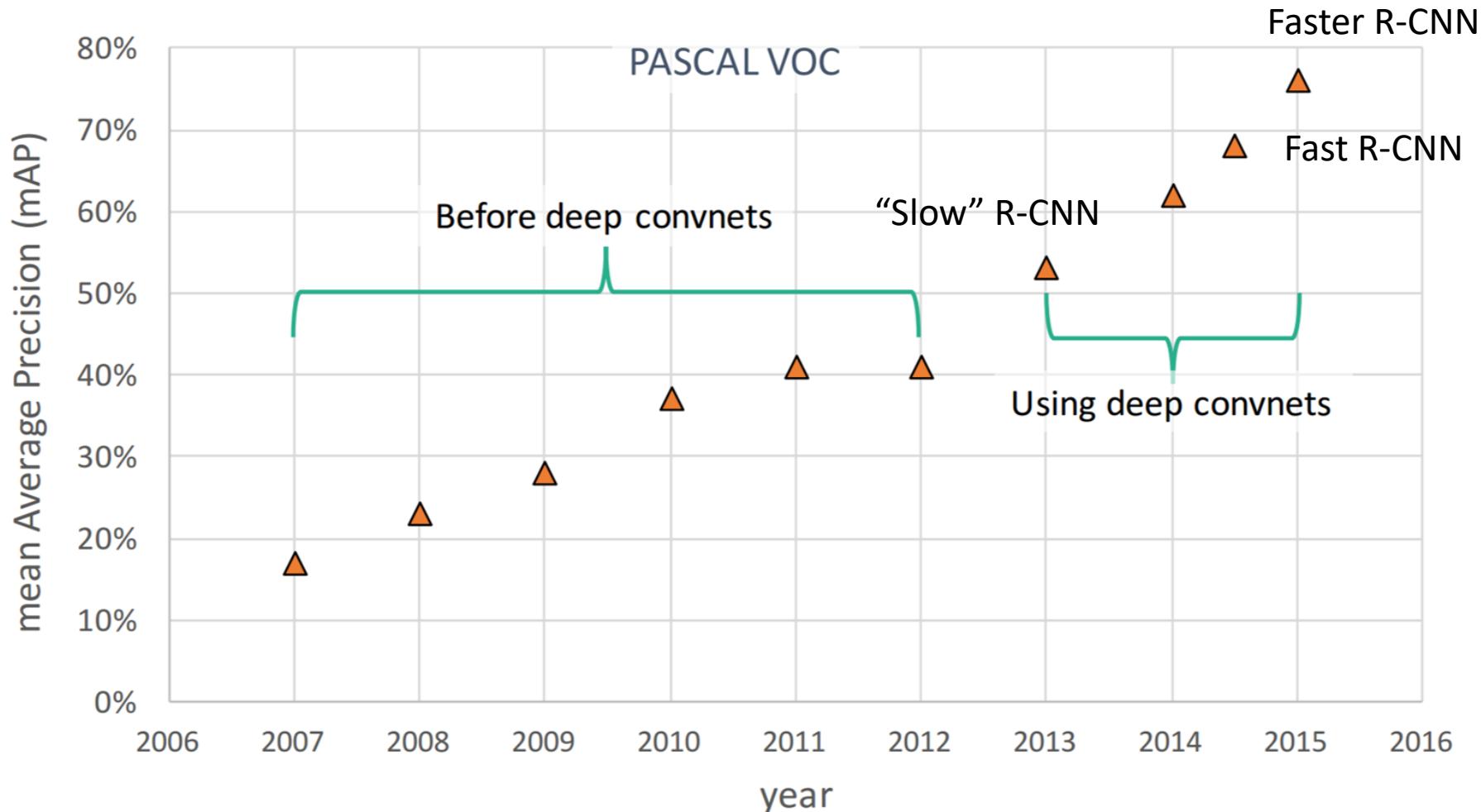


Figure copyright Ross Girshick, 2015.
Reproduced with permission.

Fasterer R-CNN: Two-Stage Object Detector

Faster R-CNN is a
Two-stage object detector

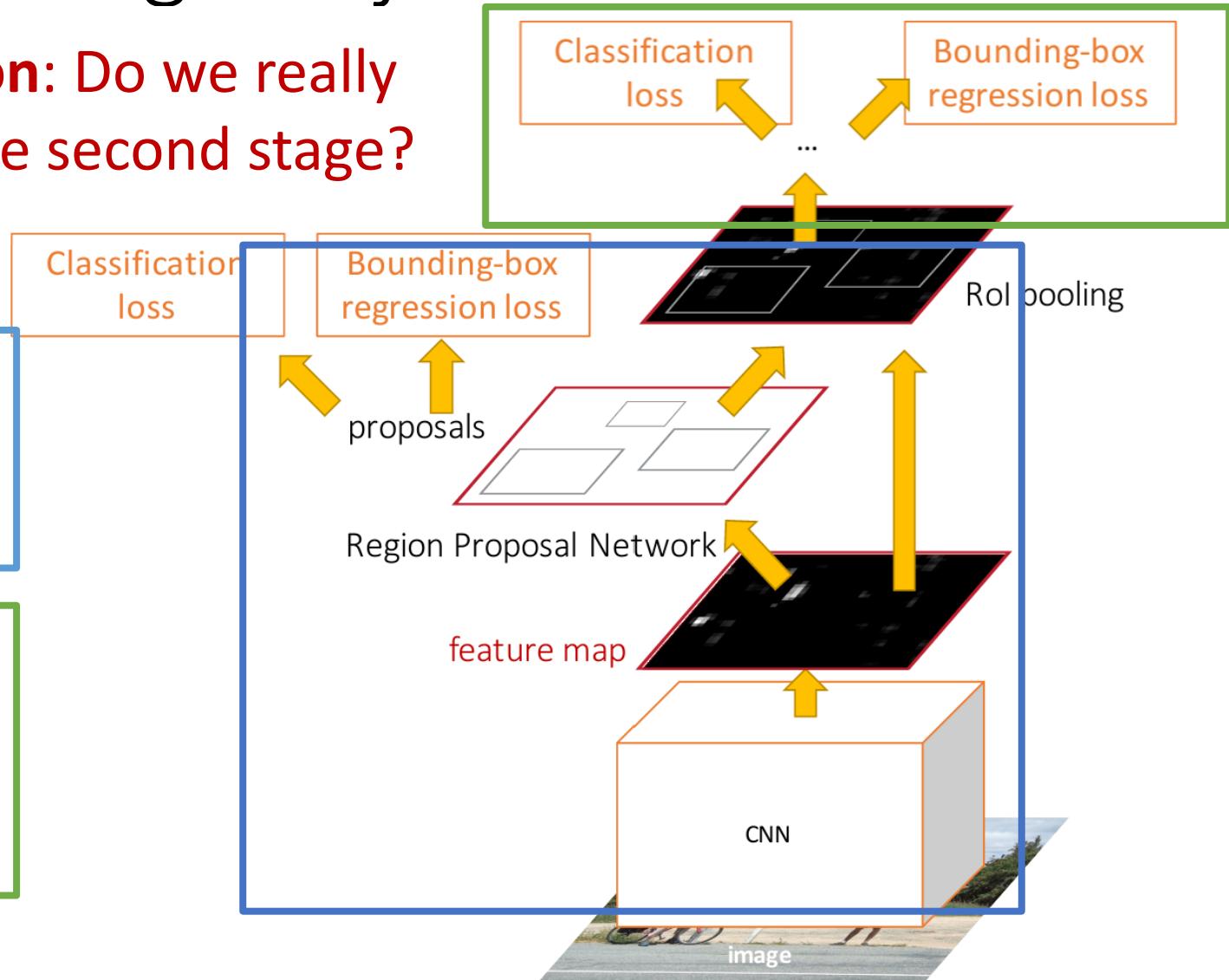
First stage: Run once per image

- Backbone network
- Region proposal network

Second stage: Run once per region

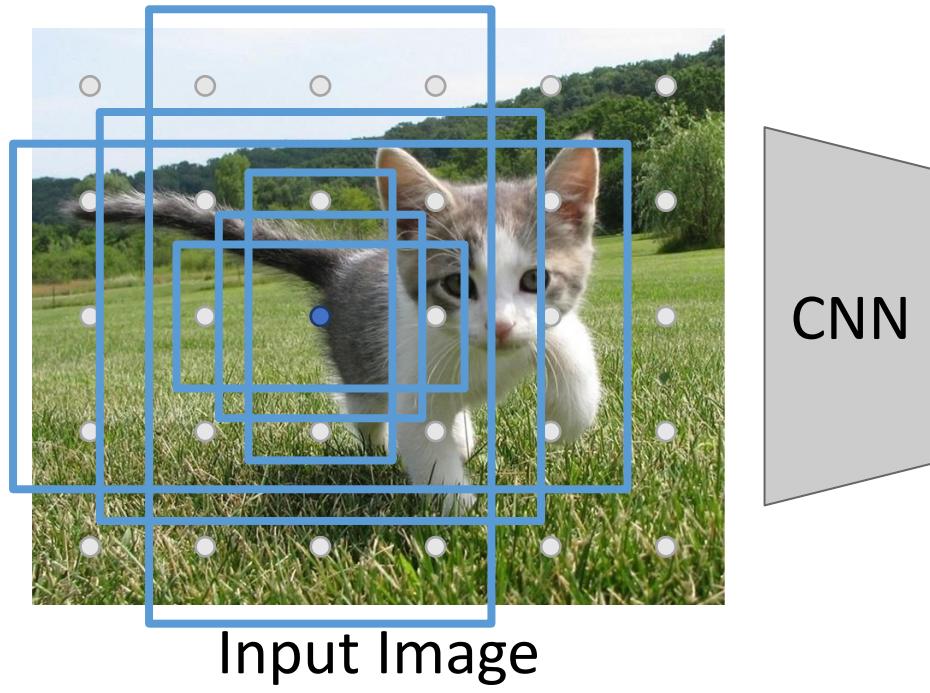
- Crop features: RoI pool / align
- Predict object class
- Prediction bbox offset

Question: Do we really
need the second stage?

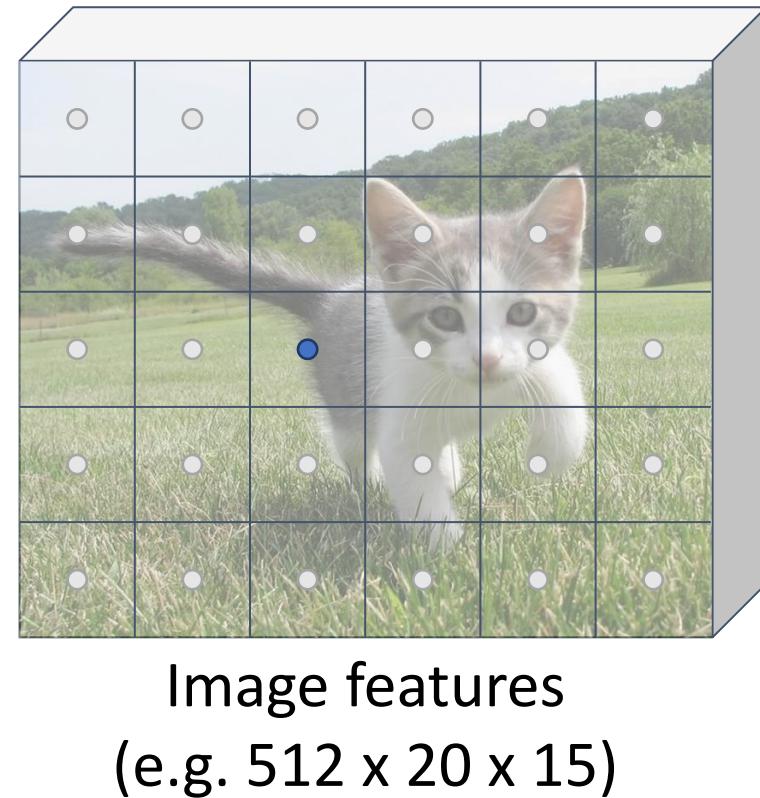


Single-Stage Detectors: RetinaNet

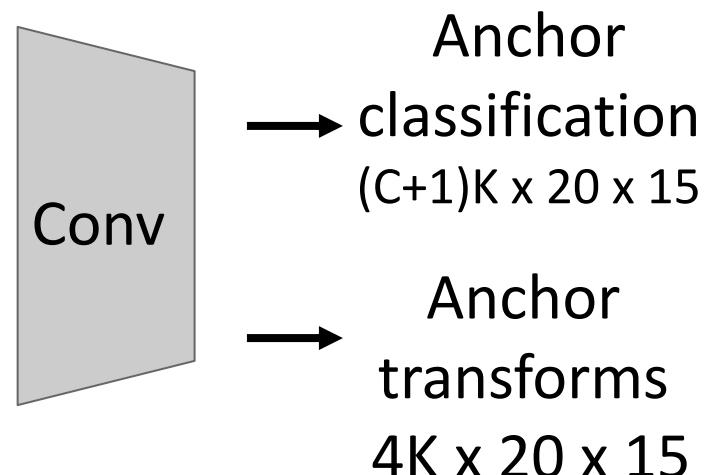
Run backbone CNN to get features aligned to input image



Each feature corresponds to a point in the input

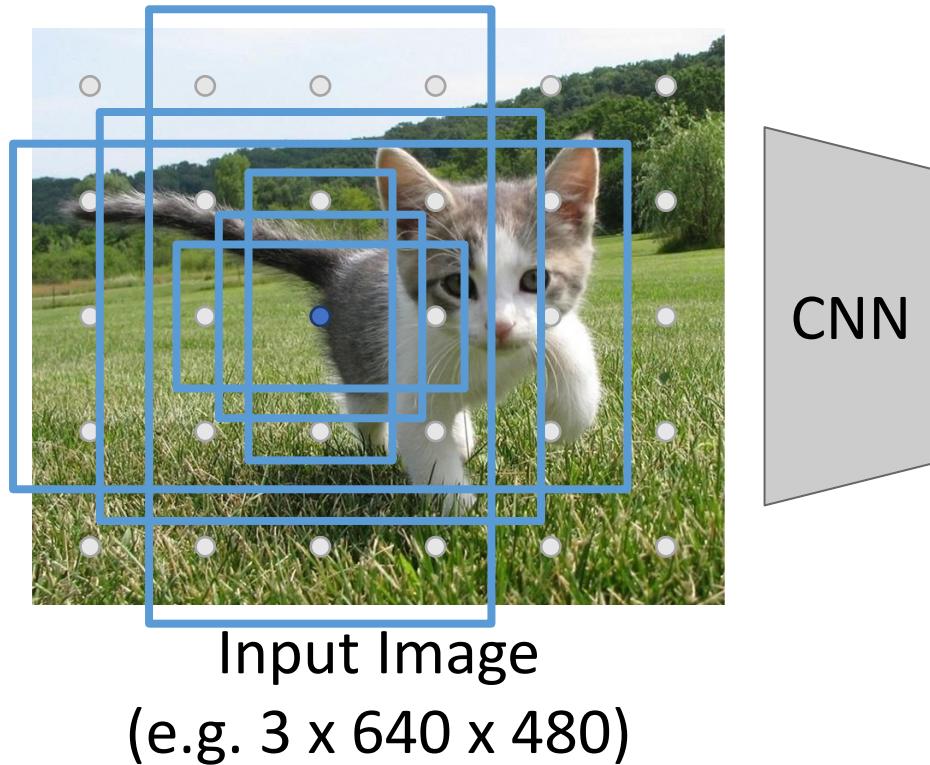


Similar to RPN – but rather than classify anchors as object/no object, directly predict object category (among C categories) or background

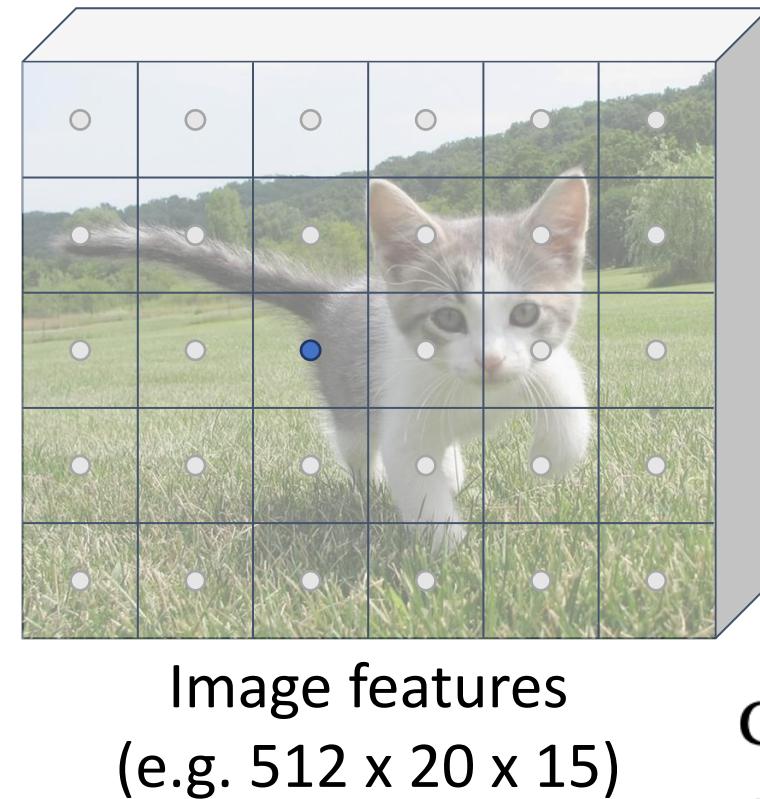


Single-Stage Detectors: RetinaNet

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: “Focal Loss”; see the paper for details

Anchor classification
 $(C+1)K \times 20 \times 15$

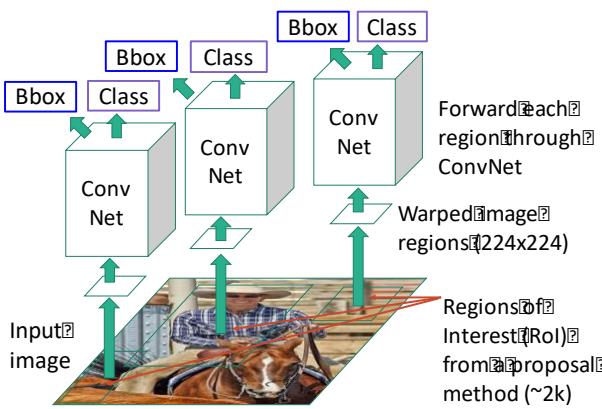
Anchor transforms
 $4K \times 20 \times 15$

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

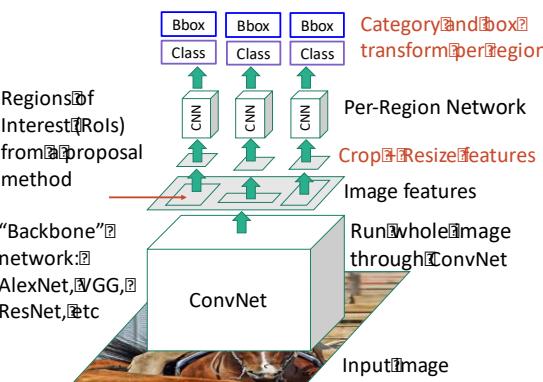
$$\text{FL}(p_t) = -(1 - p_t)^\gamma \log(p_t)$$

Summary: Object Detection

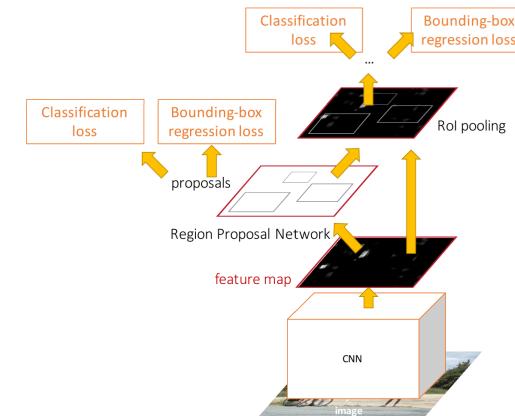
“Slow” R-CNN: Run CNN independently for each region



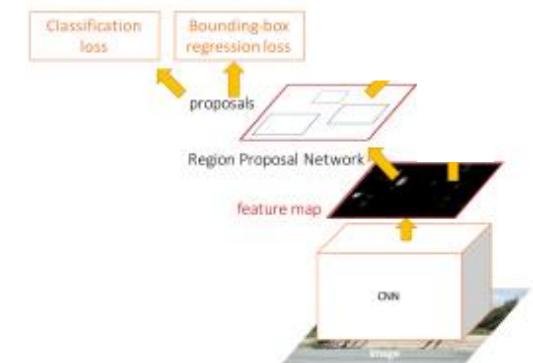
Fast R-CNN: Apply differentiable cropping to shared image features



Faster R-CNN: Compute proposals with CNN



Single-Stage: Fully convolutional detector (e.g., RetinaNet)



Computer Vision Task: Object Detection

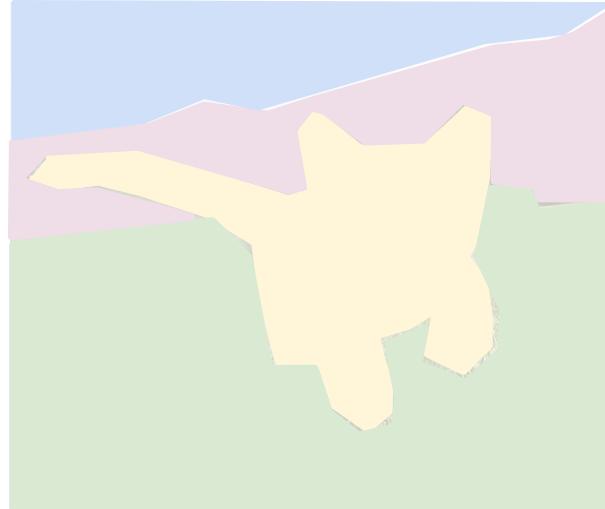
Classification



CAT

No spatial extent

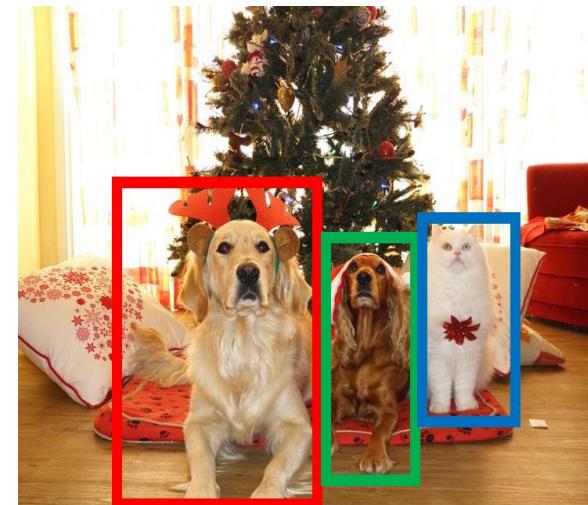
Semantic Segmentation



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](#)

Computer Vision Task: Instance Segmentation

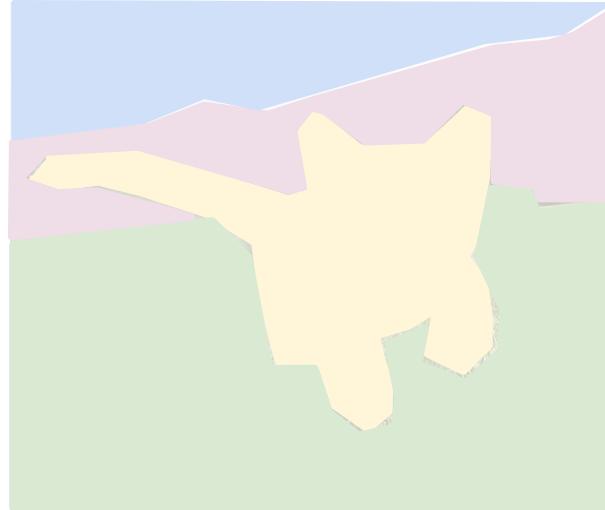
Classification



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

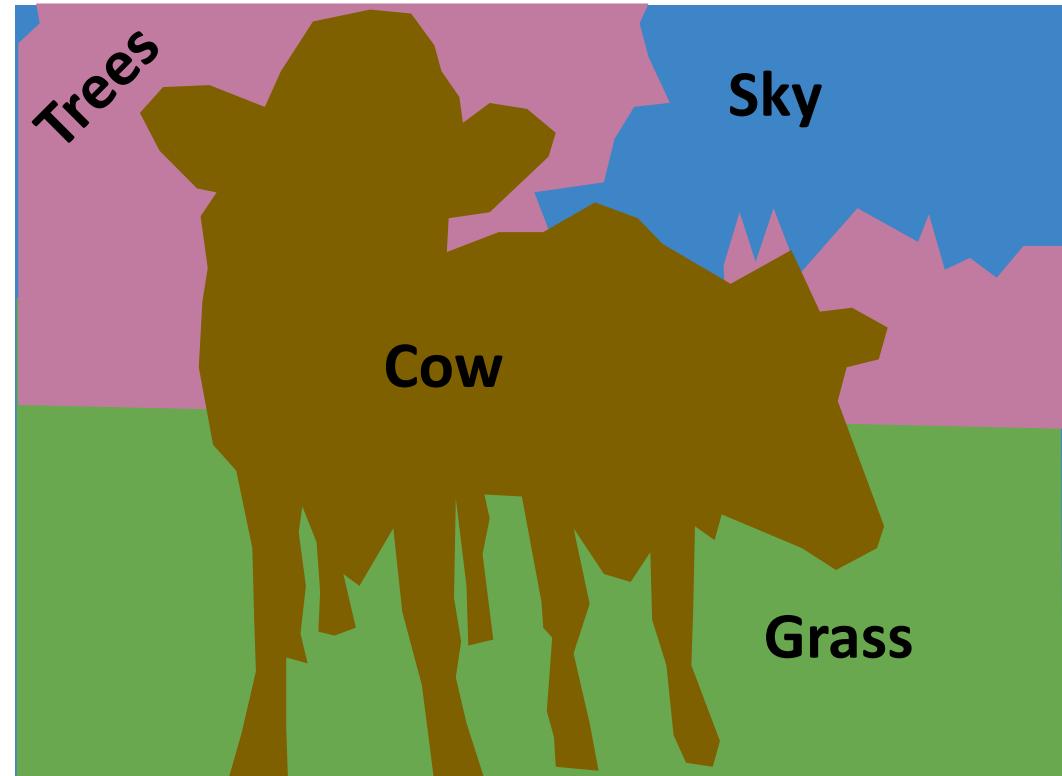
[This image is CC0 public domain](#)

Computer Vision Tasks

Object Detection: Detects individual object instances, but only gives box



Semantic Segmentation: Gives per-pixel labels, but merges instances

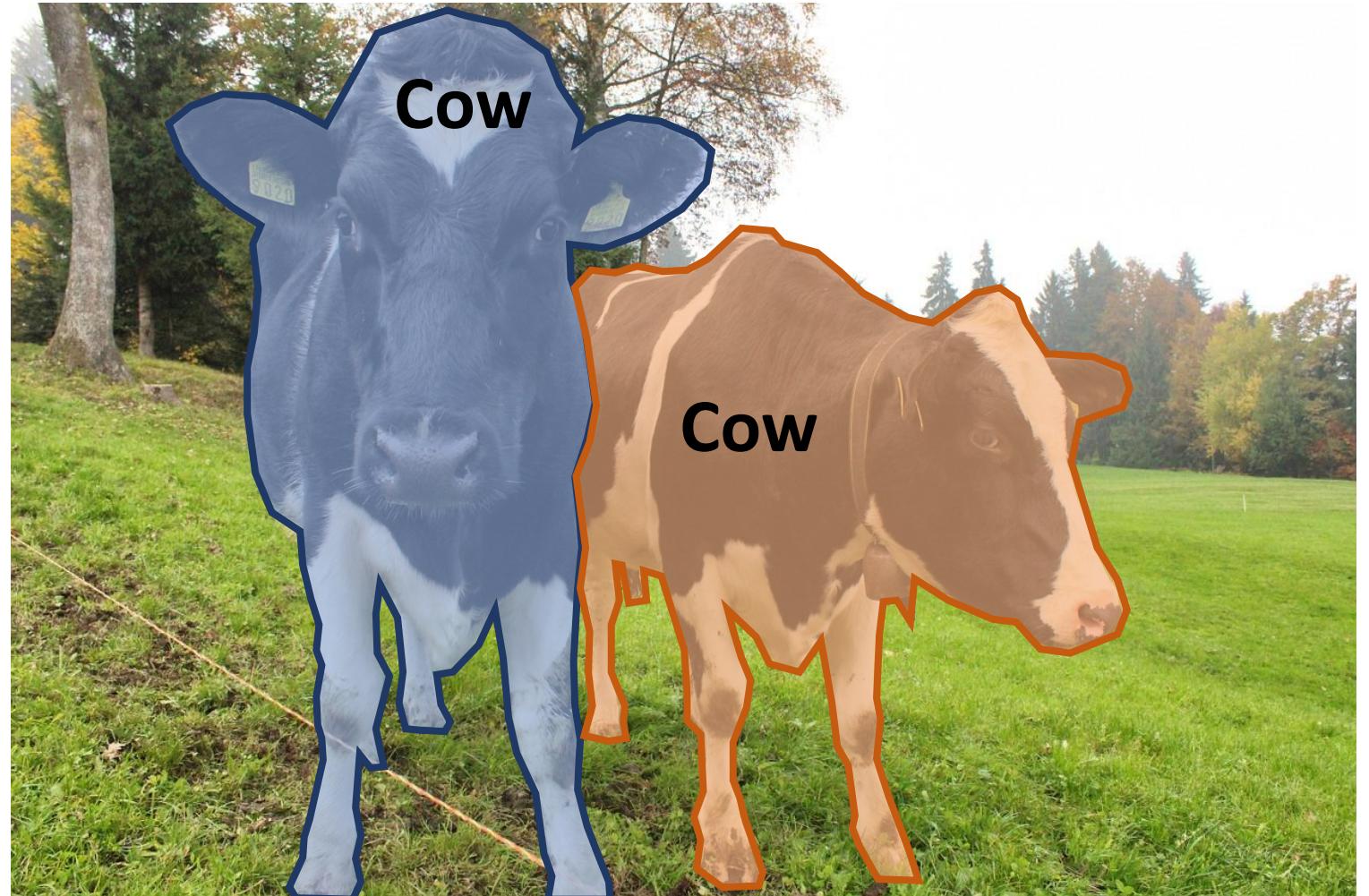


This image is CC0 public domain

Computer Vision Task: Instance Segmentation

Instance Segmentation:
Detect all objects in the image, and identify the pixels that belong to each object.

Approach: Perform object detection, then predict a segmentation mask for each object!



This image is CC0 public domain

Object Detection: Faster R-CNN

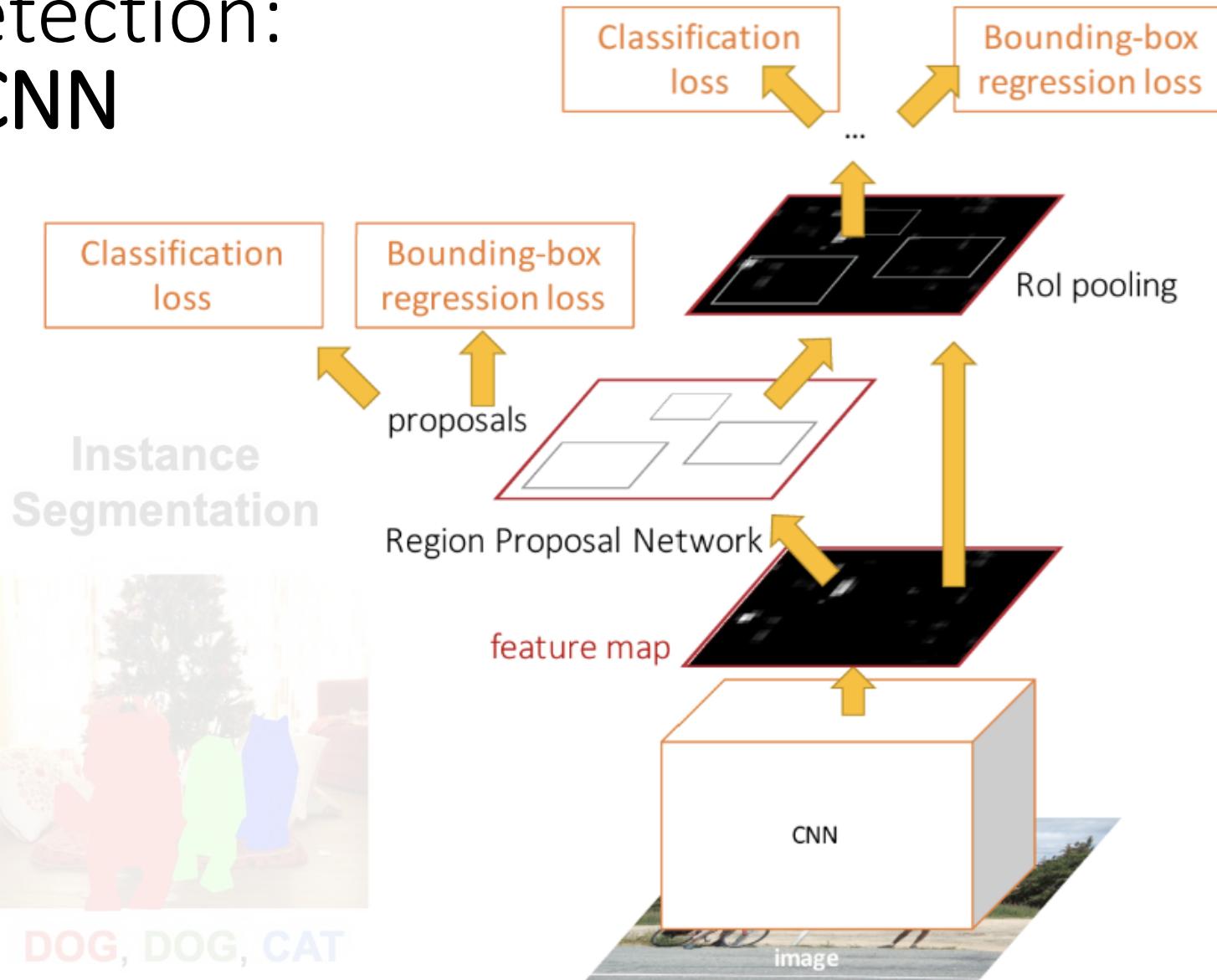
Object Detection



DOG, DOG, CAT



DOG, DOG, CAT



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

Instance Segmentation: Mask R-CNN

Object
Detection



DOG, DOG, CAT



DOG, DOG, CAT

Instance Segmentation

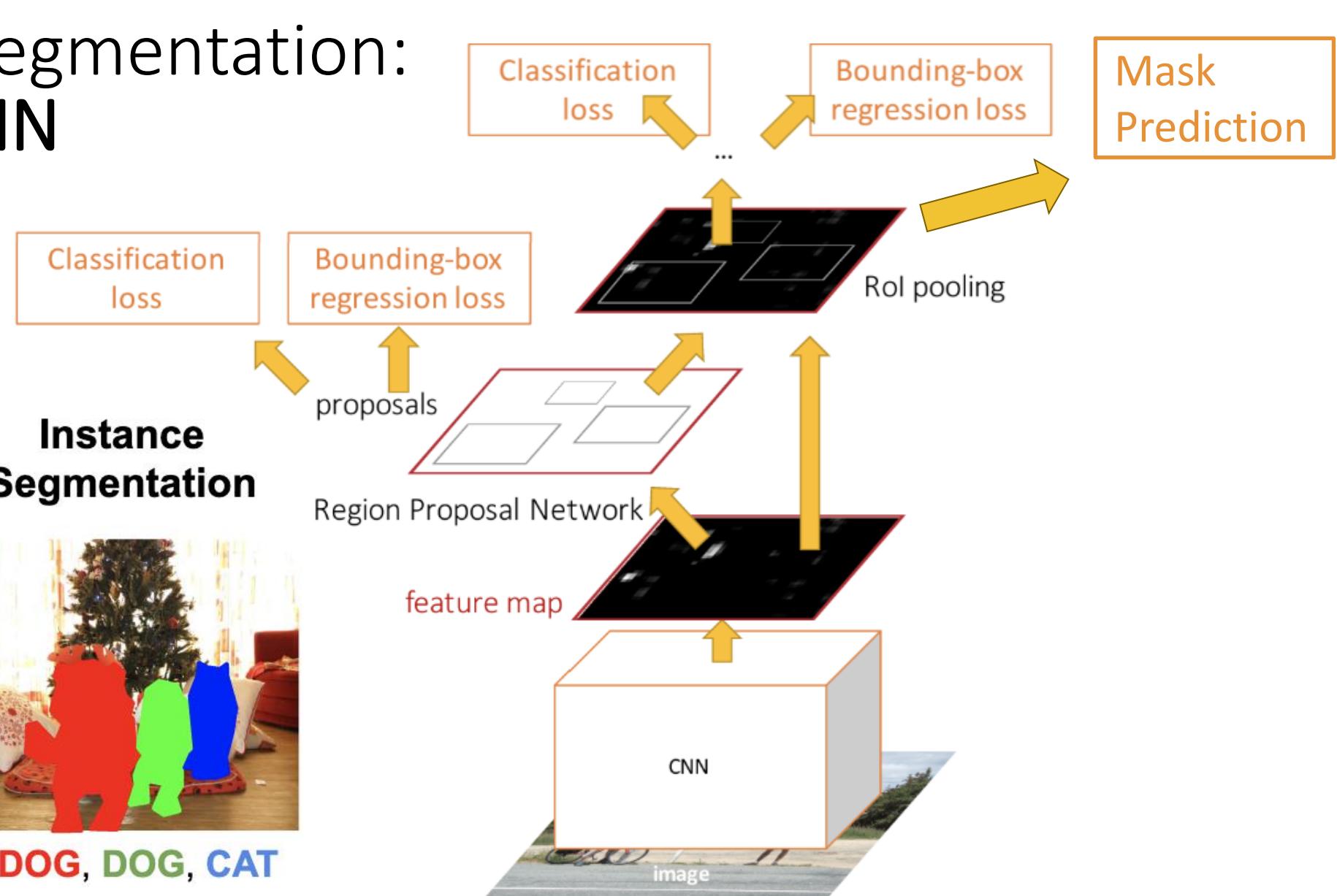
Classification
loss

Bounding-box
regression loss

Classification
loss

Bounding-box
regression loss

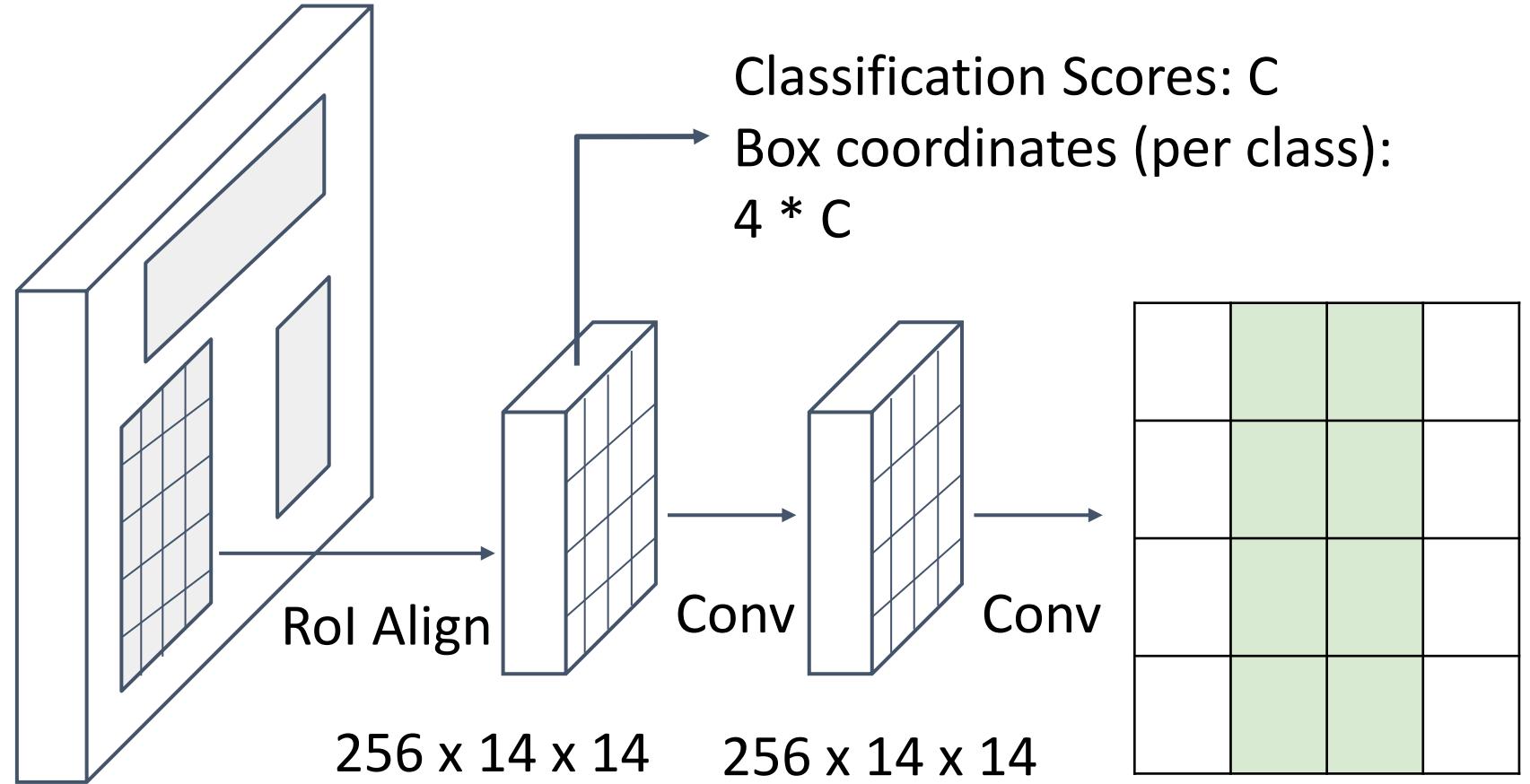
Mask
Prediction



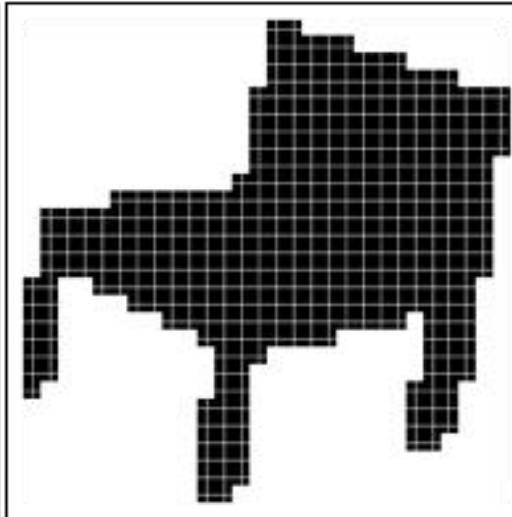
Mask R-CNN



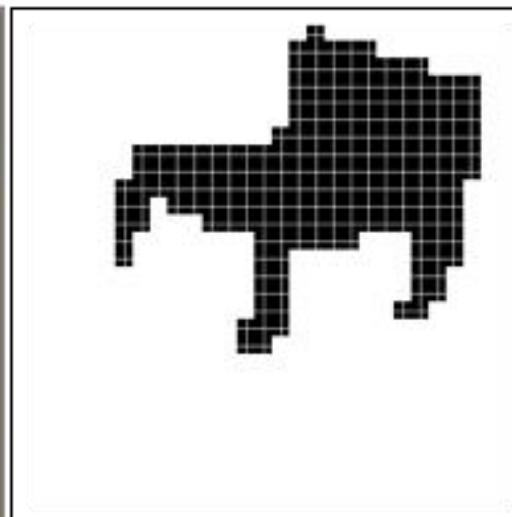
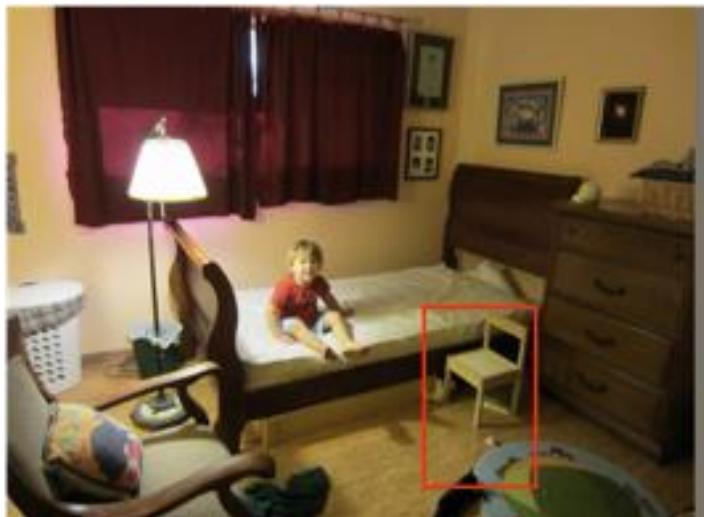
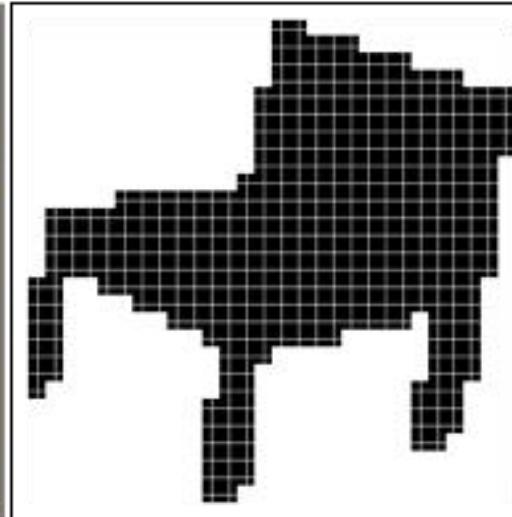
CNN
+RPN



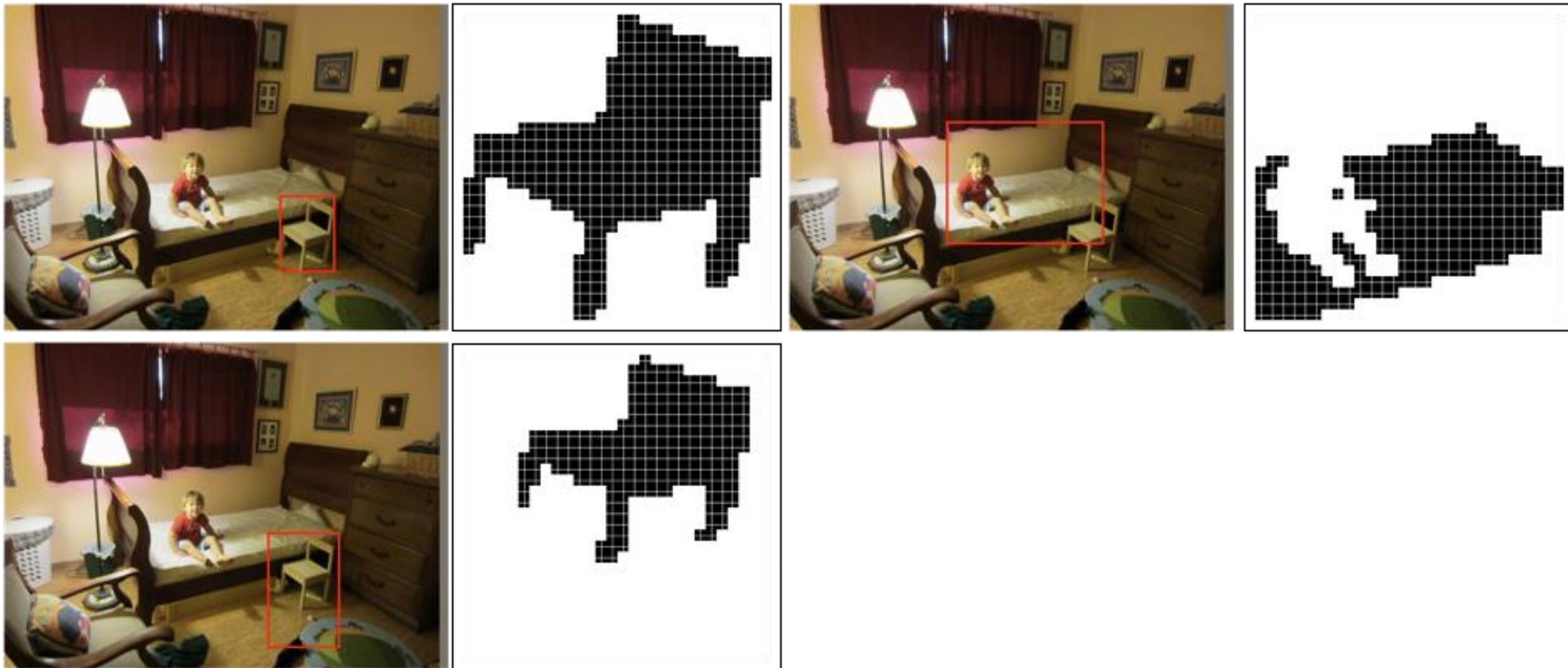
Mask R-CNN: Example Training Targets



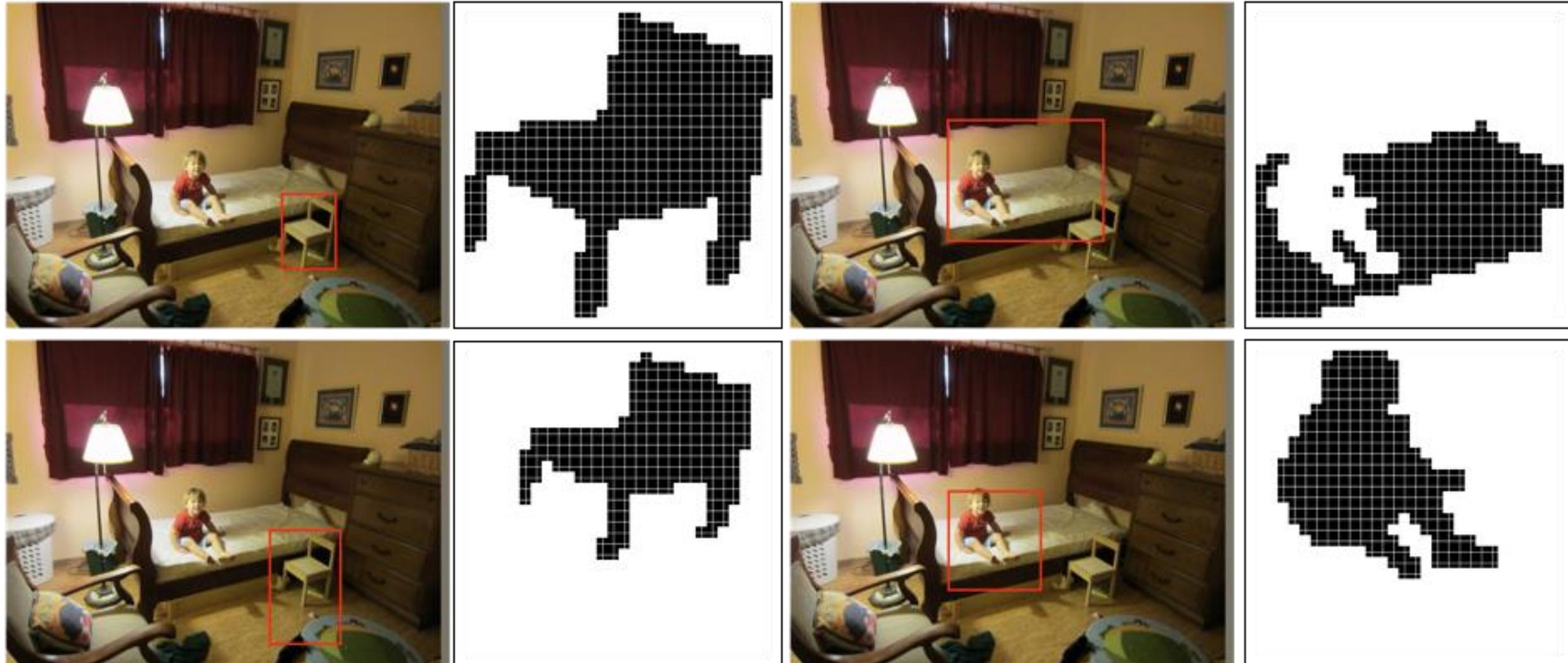
Mask R-CNN: Example Training Targets



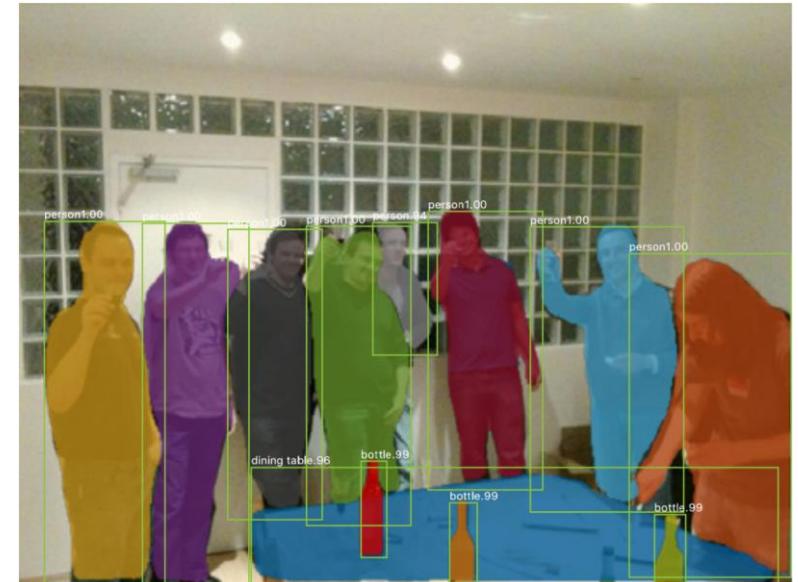
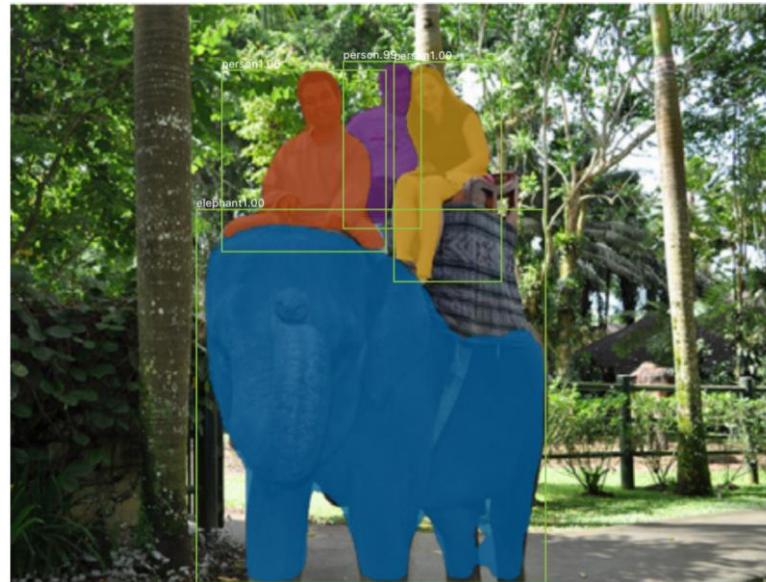
Mask R-CNN: Example Training Targets



Mask R-CNN: Example Training Targets

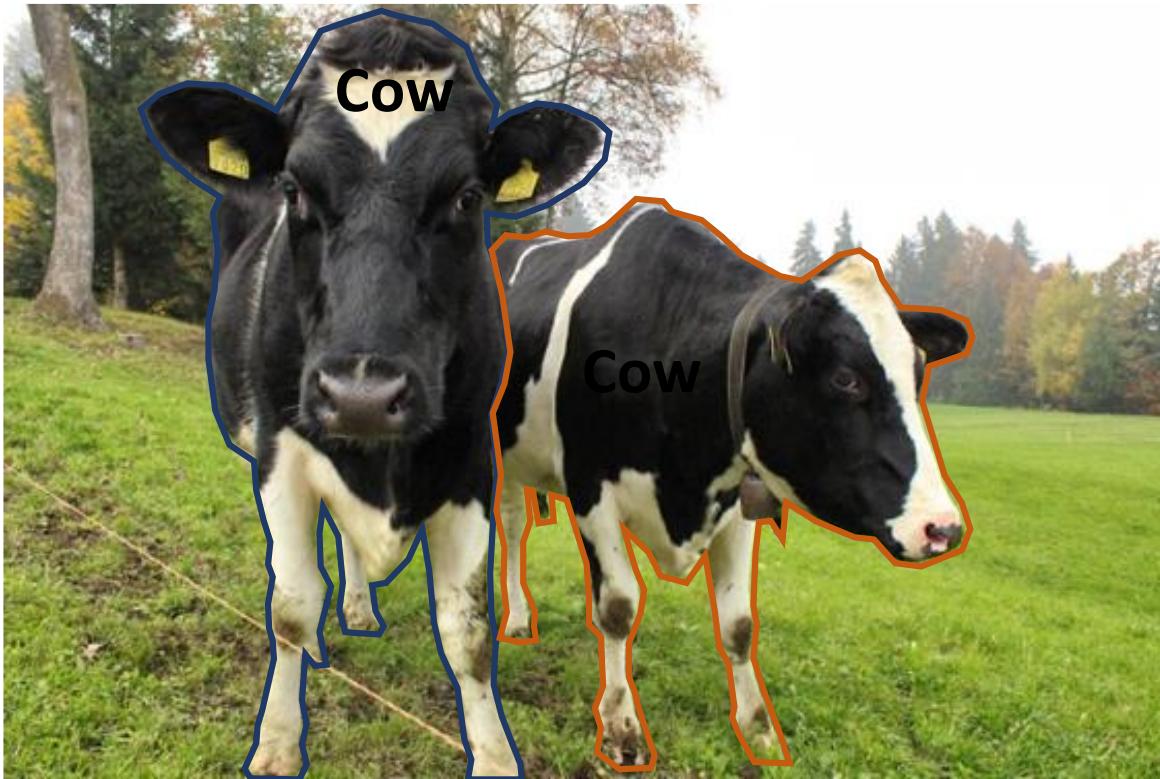


Mask R-CNN: Very Good Results!

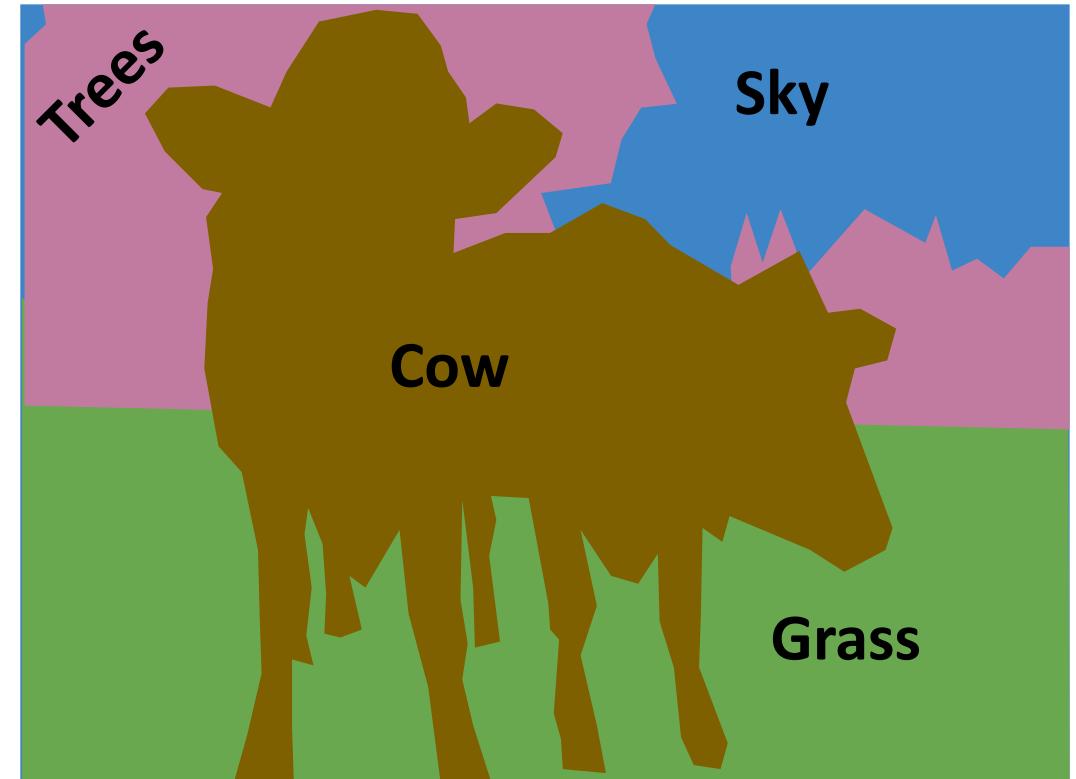


Beyond Instance Segmentation

Instance Segmentation: Separate object instances, but only things



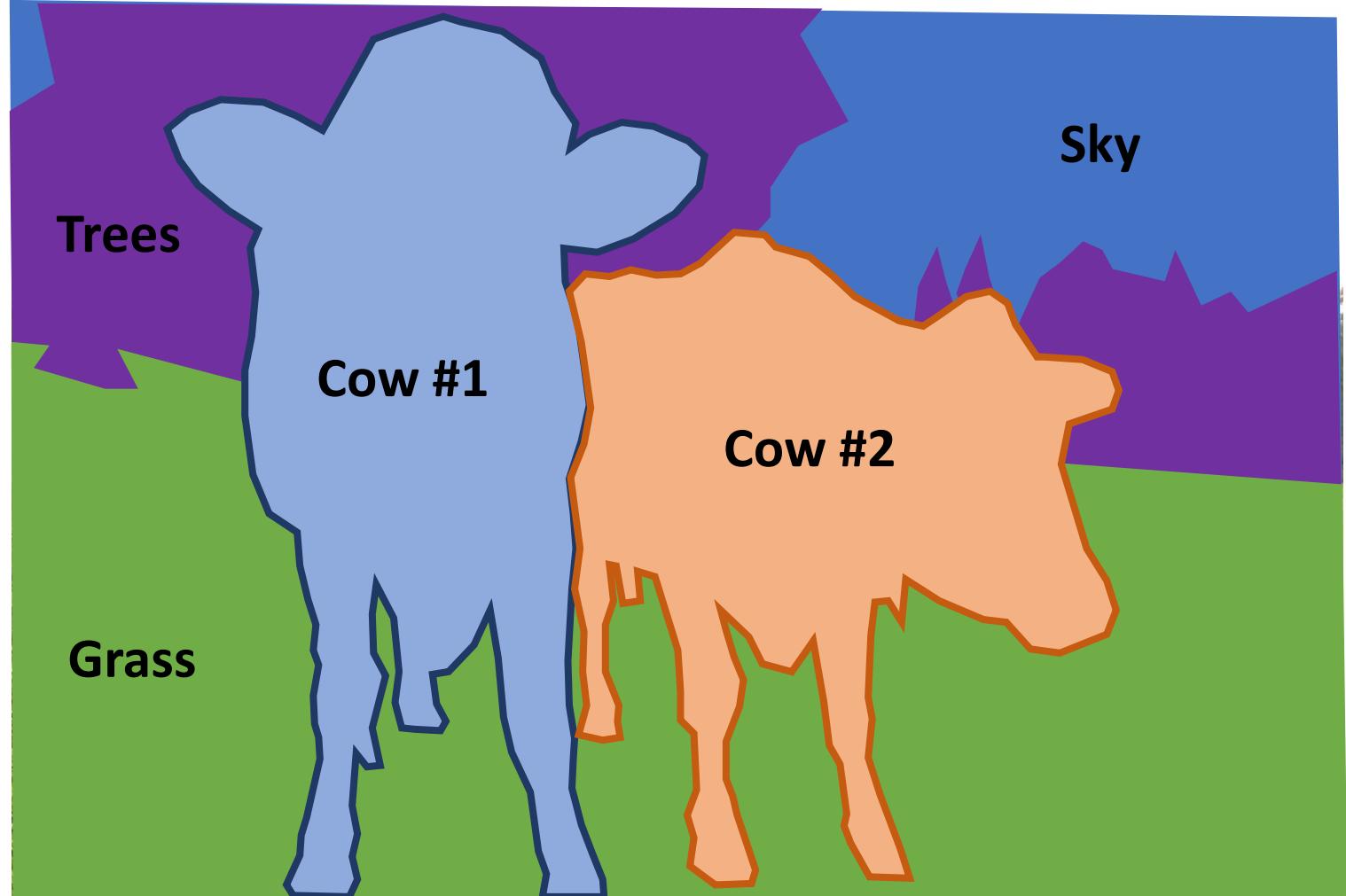
Semantic Segmentation: Identify both things and stuff, but doesn't separate instances



Beyond Instance Segmentation: Panoptic Segmentation

Label all pixels in the image (both things and stuff)

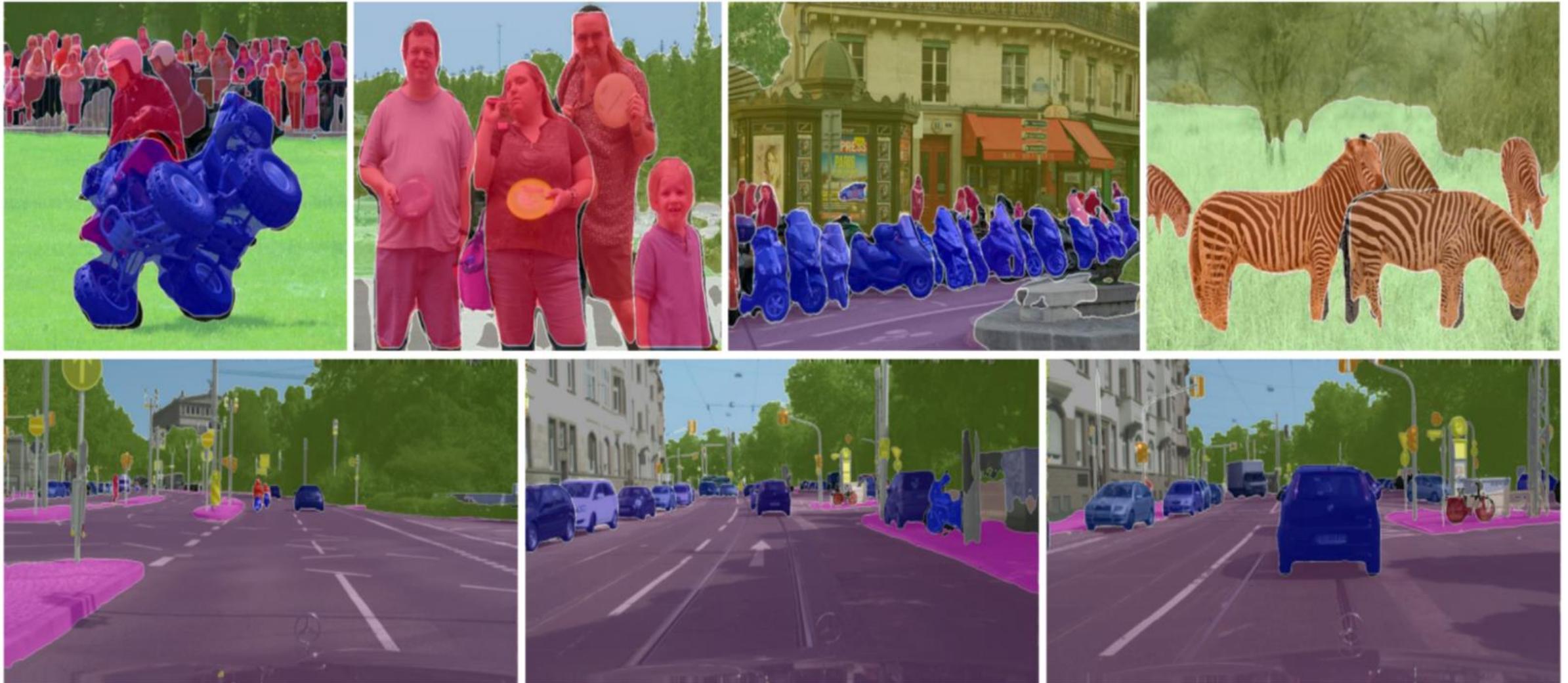
For “thing” categories also separate into instances



Kirillov et al, “Panoptic Segmentation”, CVPR 2019

Kirillov et al, “Panoptic Feature Pyramid Networks”, CVPR 2019

Beyond Instance Segmentation: Panoptic Segmentation



Kirillov et al, "Panoptic Feature Pyramid Networks", CVPR 2019

Instance Segmentation: Mask R-CNN

Object
Detection



DOG, DOG, CAT



DOG, DOG, CAT

Instance Segmentation

Classification
loss

Bounding-box
regression loss

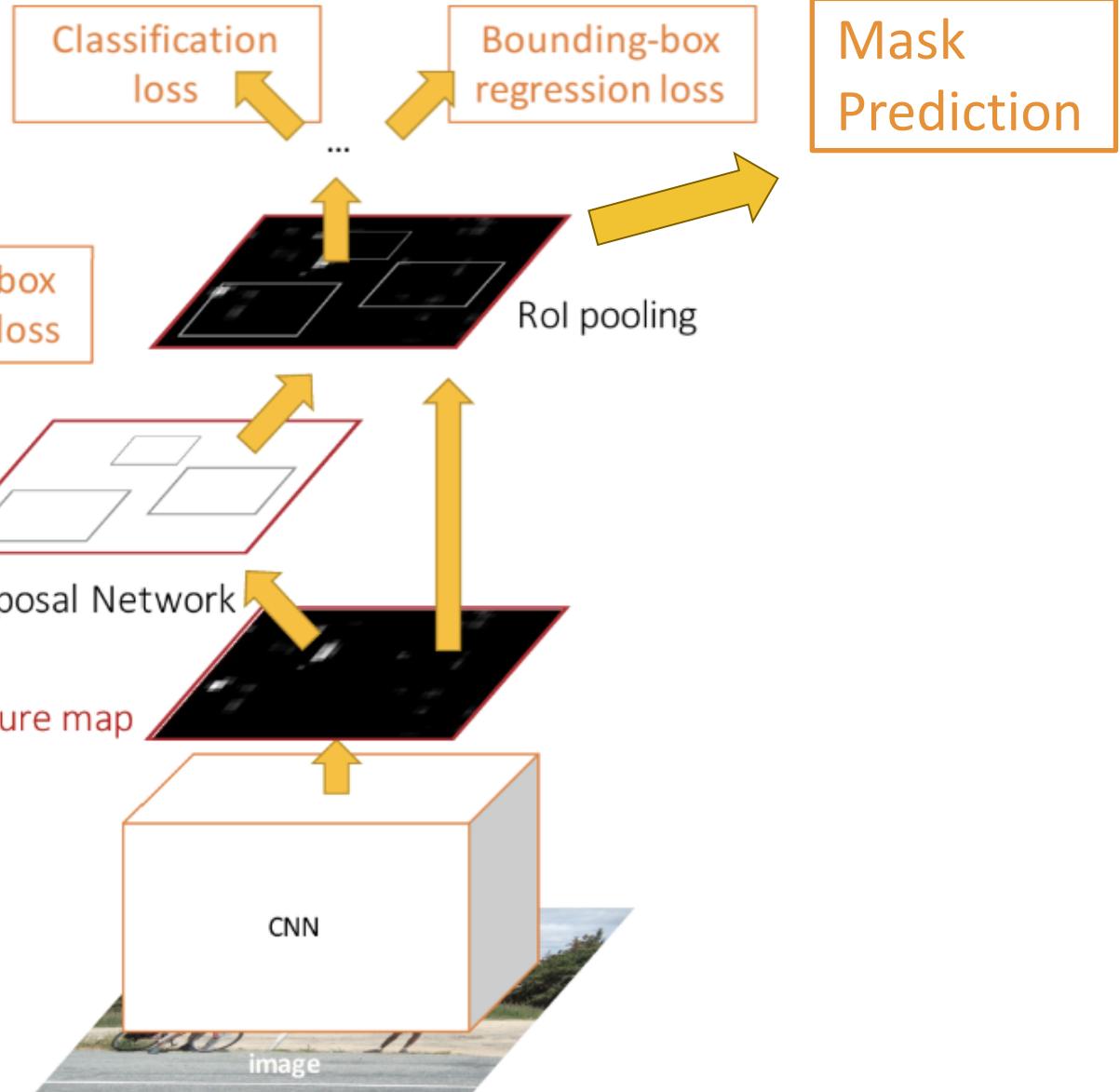
proposals

Region Proposal Network

feature map

CNN

image



General Idea: Add Per-Region “Heads”!

Object Detection



DOG, DOG, CAT

Instance Segmentation



DOG, DOG, CAT

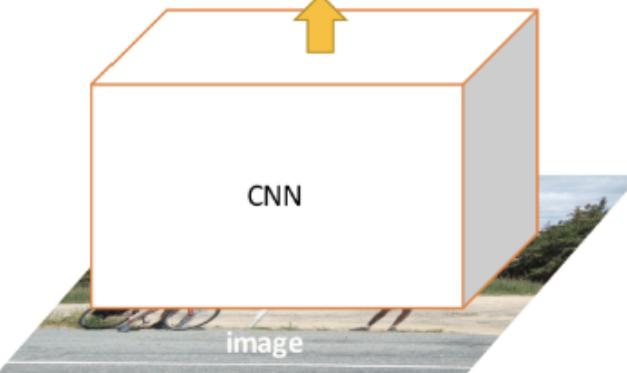
Classification loss

Bounding-box regression loss

proposals

Region Proposal Network

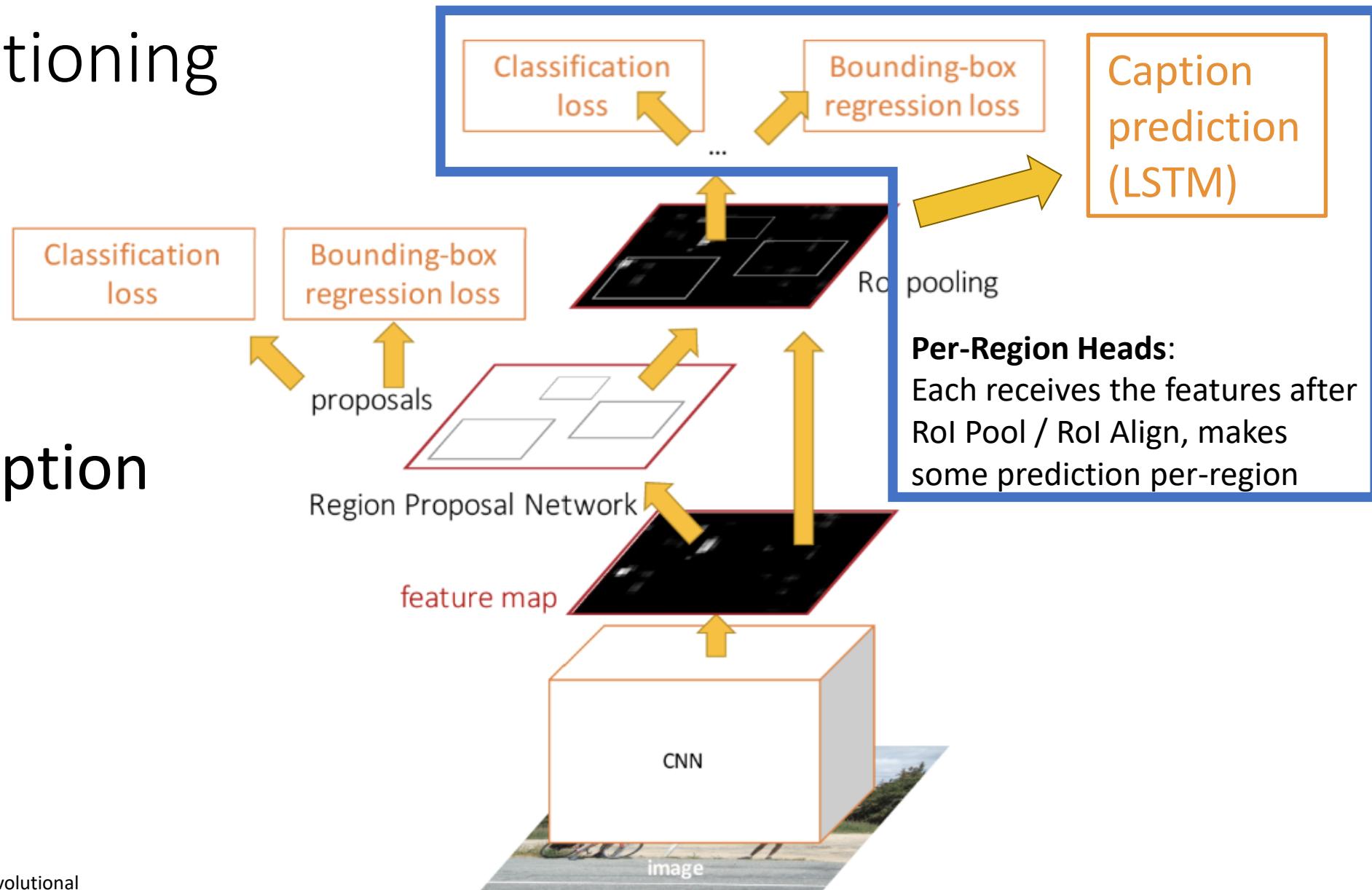
feature map



Per-Region Heads:
Each receives the features after
RoI Pool / RoI Align, makes
some prediction per-region

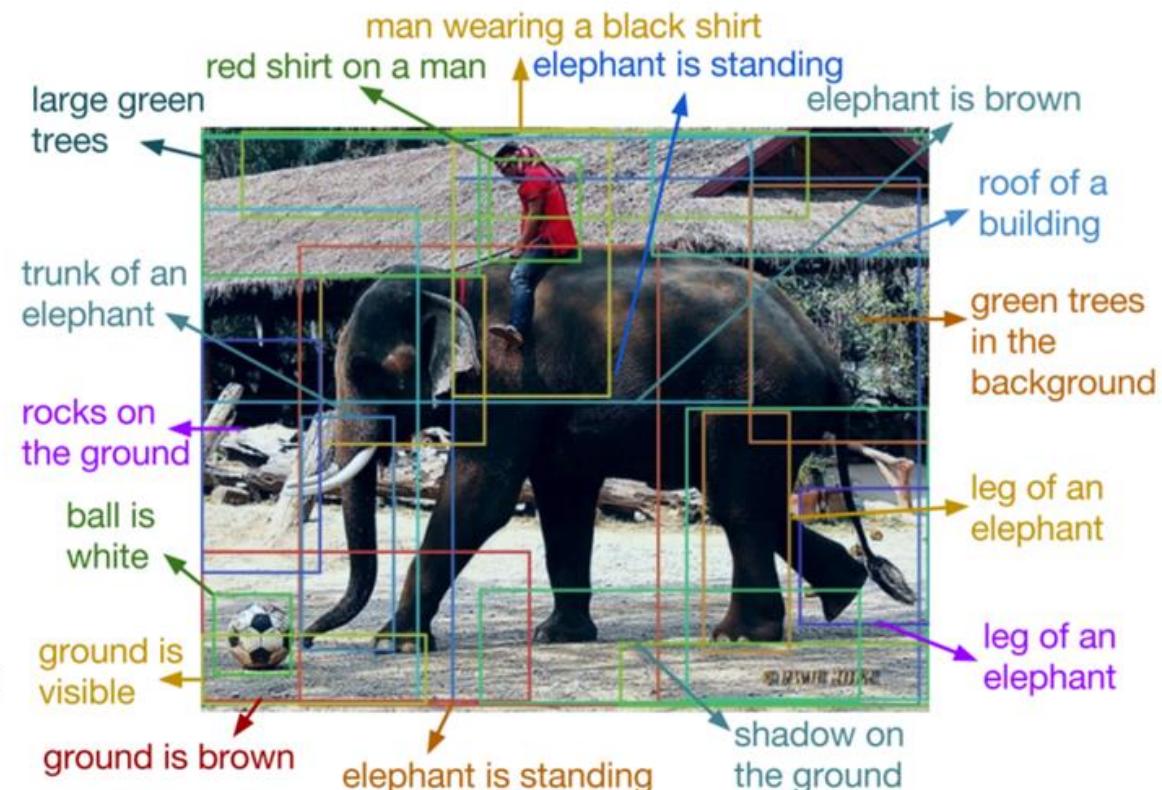
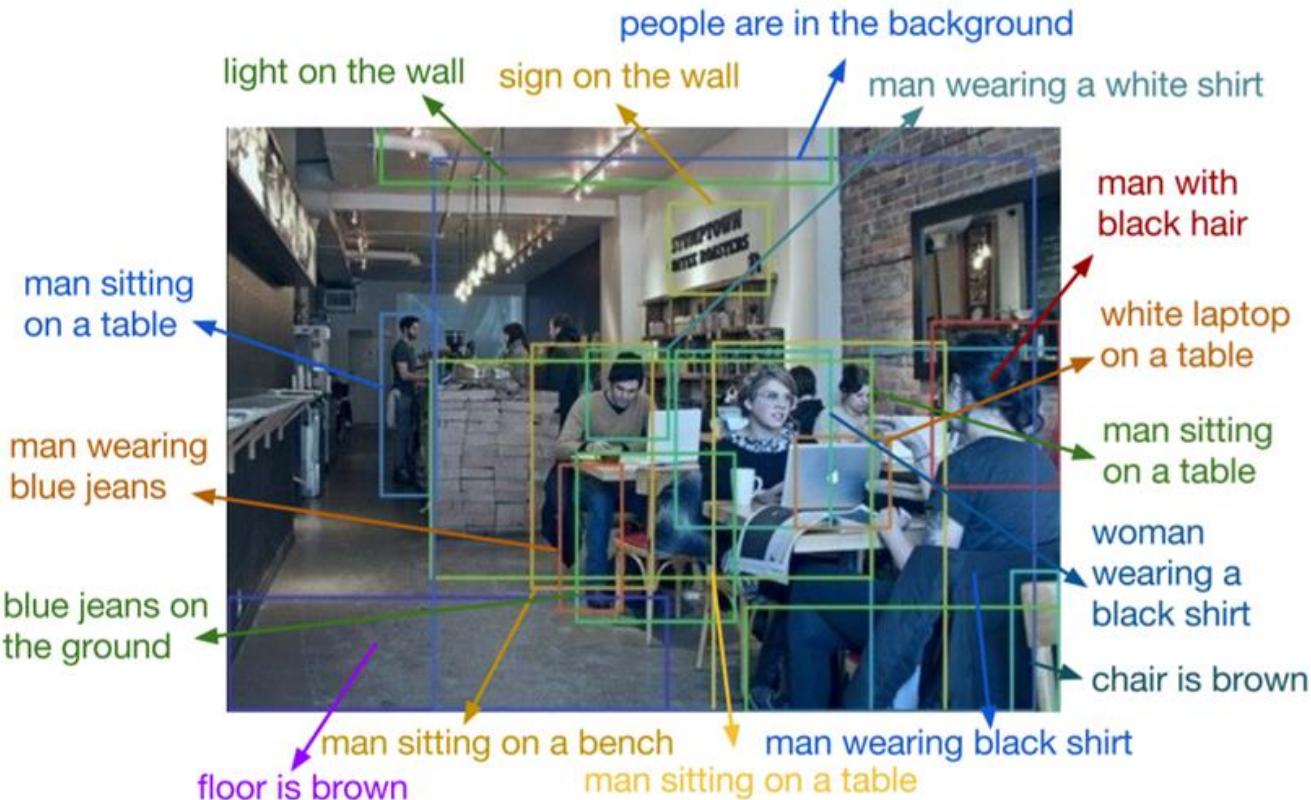
Dense Captioning

Predict a caption per region!

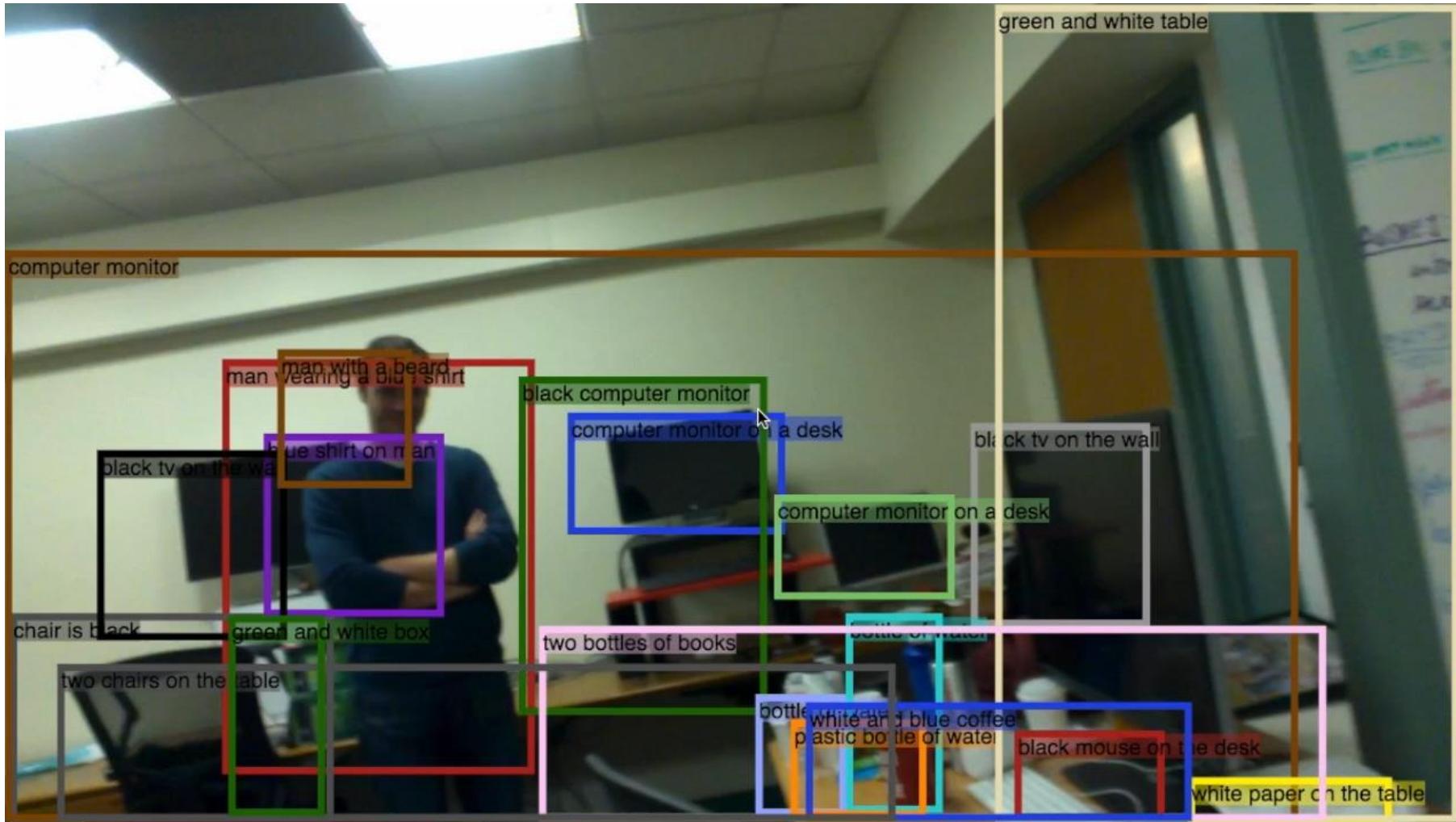


Johnson, Karpathy, and Fei-Fei, "DenseCap: Fully Convolutional Localization Networks for Dense Captioning", CVPR 2016

Dense Captioning



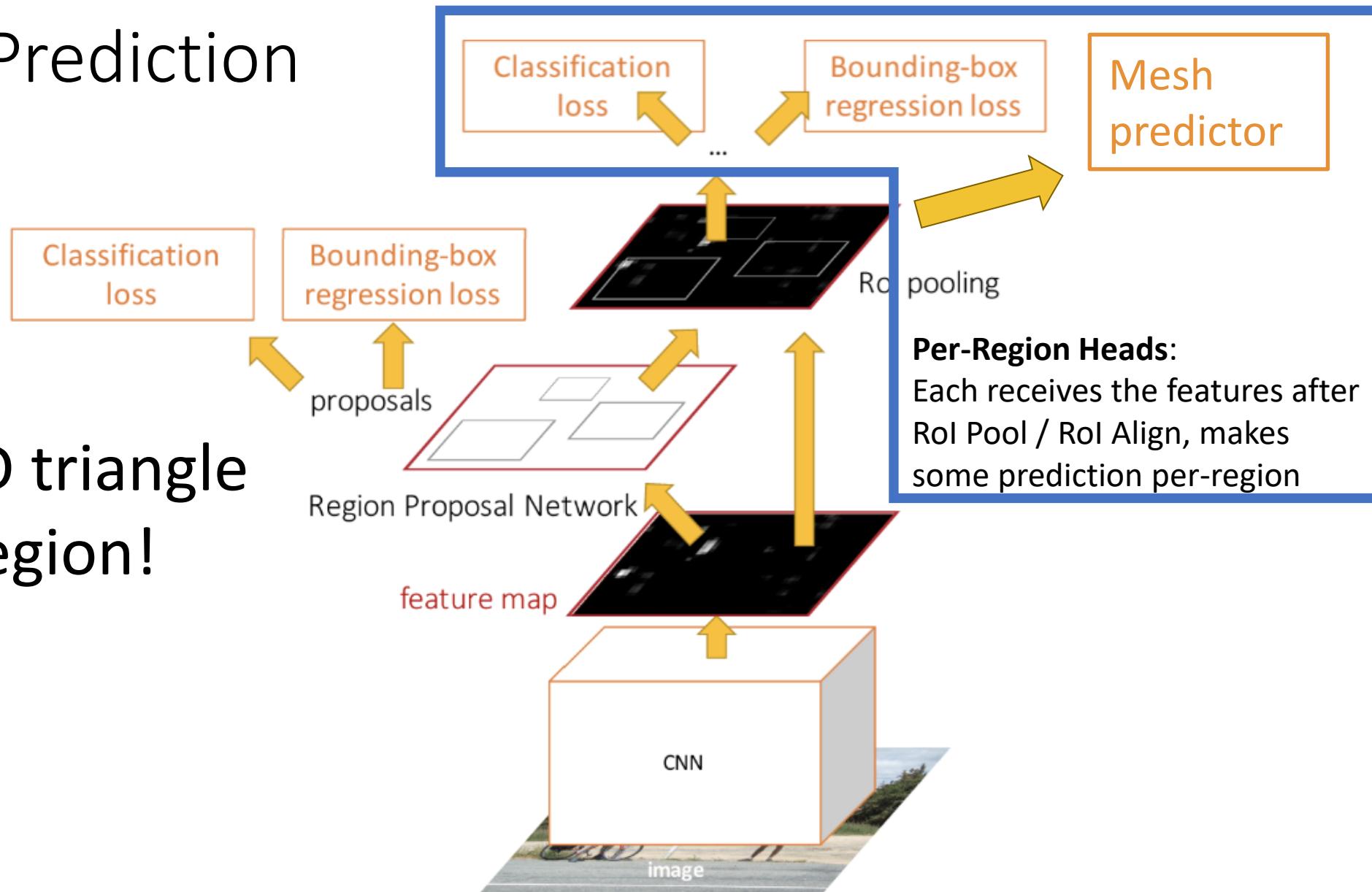
Dense Captioning



Johnson, Karpathy, and Fei-Fei, "DenseCap: Fully Convolutional Localization Networks for Dense Captioning", CVPR 2016

3D Shape Prediction

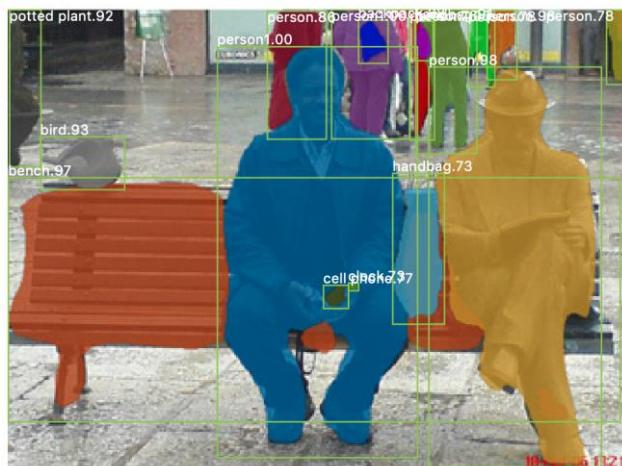
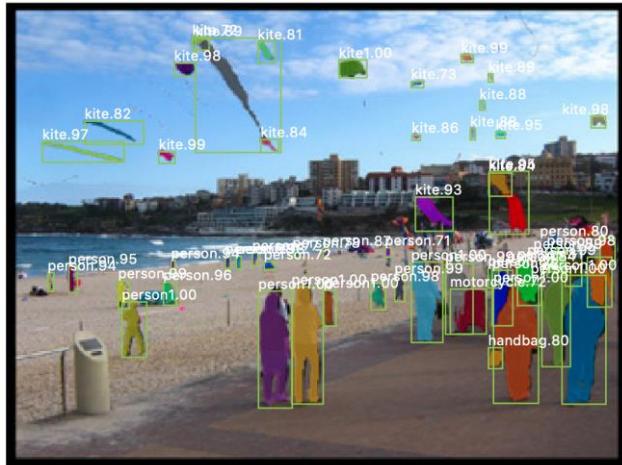
Predict a 3D triangle
mesh per region!



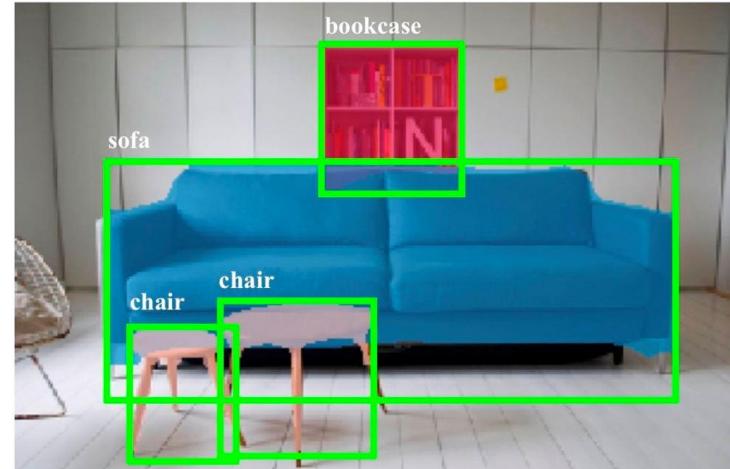
Gkioxari, Malik, and Johnson, "Mesh R-CNN", ICCV 2019

3D Shape Prediction: Mask R-CNN + Mesh Head

Mask R-CNN:
2D Image -> 2D shapes



Mesh R-CNN:
2D Image -> 3D shapes



He, Gkioxari, Dollár, and
Girshick, "Mask R-CNN",
ICCV 2017

Gkioxari, Malik, and Johnson,
"Mesh R-CNN", ICCV 2019

Summary: Many Computer Vision Tasks!

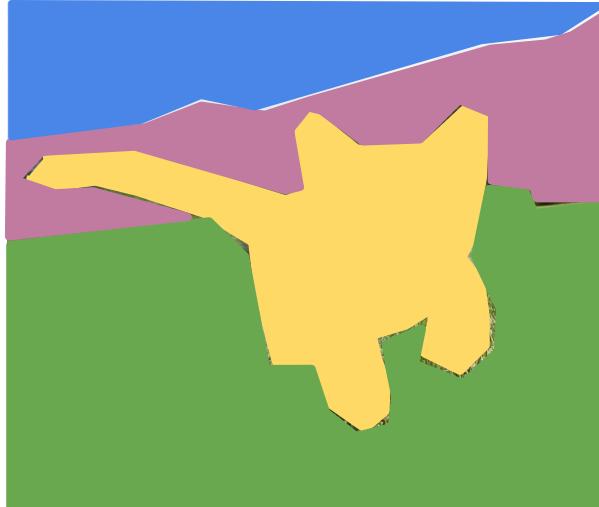
Classification



CAT

No spatial extent

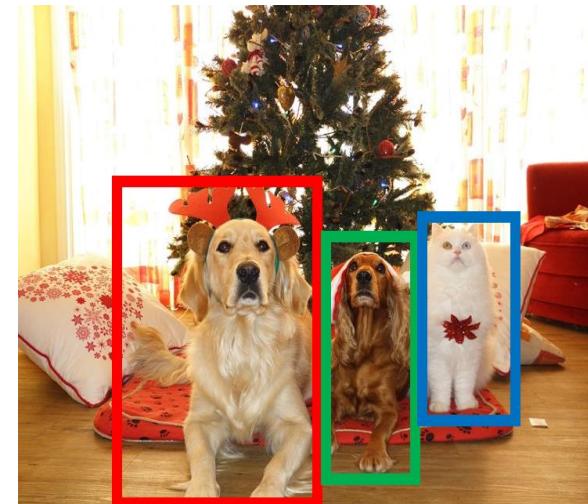
Semantic Segmentation



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](#)