

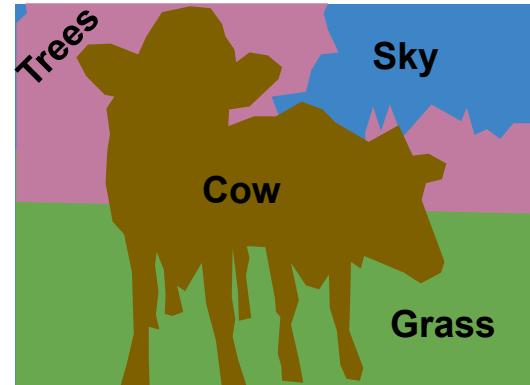
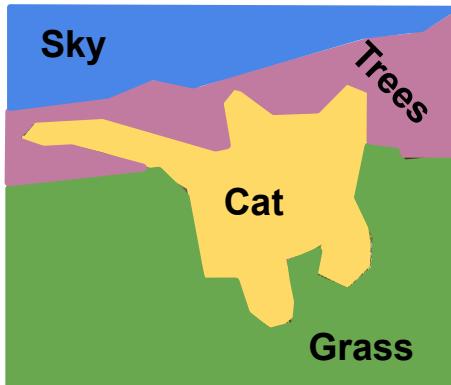
# Semantic Segmentation

Label each pixel in the image with a category label

Don't differentiate instances, only care about pixels



[This image is CC0 public domain](#)



# Object Detection

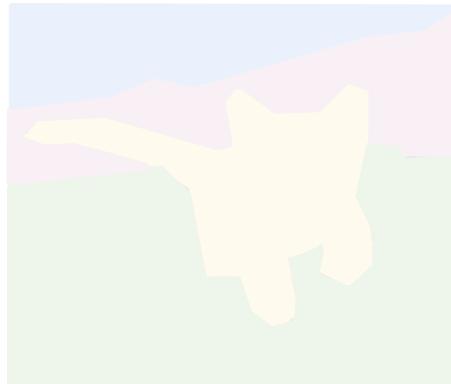
Classification



CAT

No spatial extent

Semantic Segmentation



GRASS, CAT,  
TREE, SKY

No objects, just pixels

Object  
Detection



DOG, DOG, CAT

Multiple Object

Instance  
Segmentation



DOG, DOG, CAT

# Object Detection

$x, y, h, w$  object

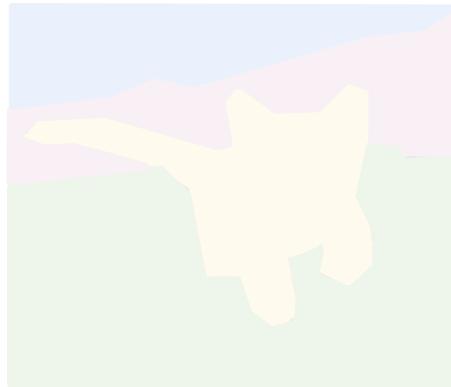
Classification



CAT

No spatial extent

Semantic Segmentation



GRASS, CAT,  
TREE, SKY

No objects, just pixels

Object  
Detection



DOG, DOG, CAT

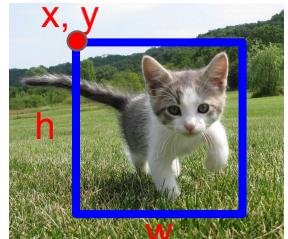
Instance  
Segmentation



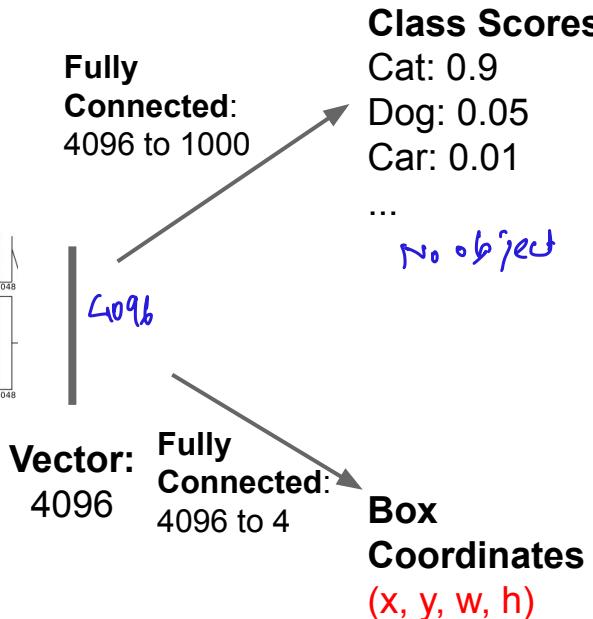
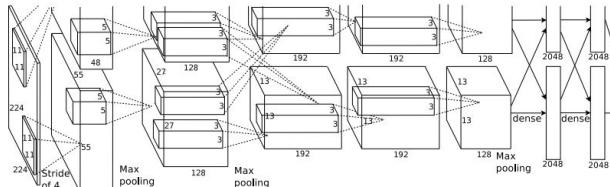
DOG, DOG, CAT

Multiple Object

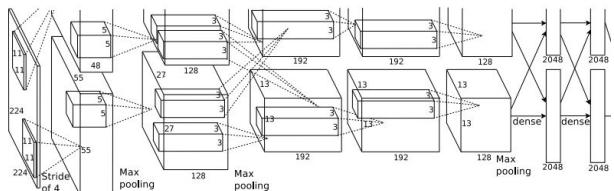
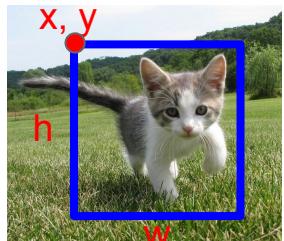
# Object Detection: Single Object (Classification + Localization)



This image is CC0 public domain



# Object Detection: Single Object (Classification + Localization)



Fully  
Connected:  
4096 to 1000

Vector:  
Fully  
Connected:  
4096 to 4

Treat localization as a  
regression problem!

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

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

Correct label:  
Cat

Softmax  
Loss

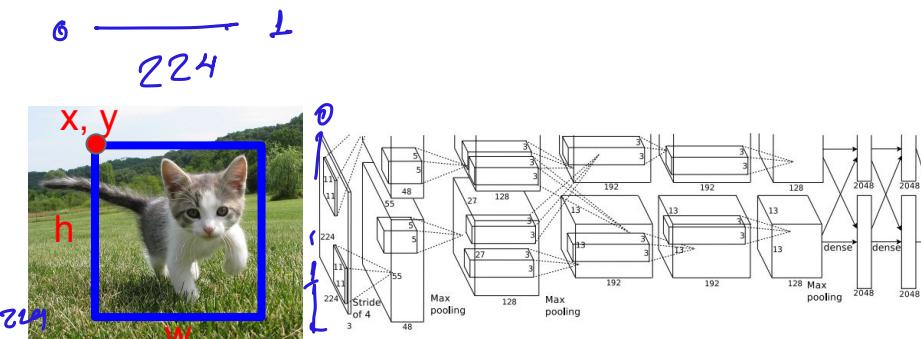
?

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

L2 Loss

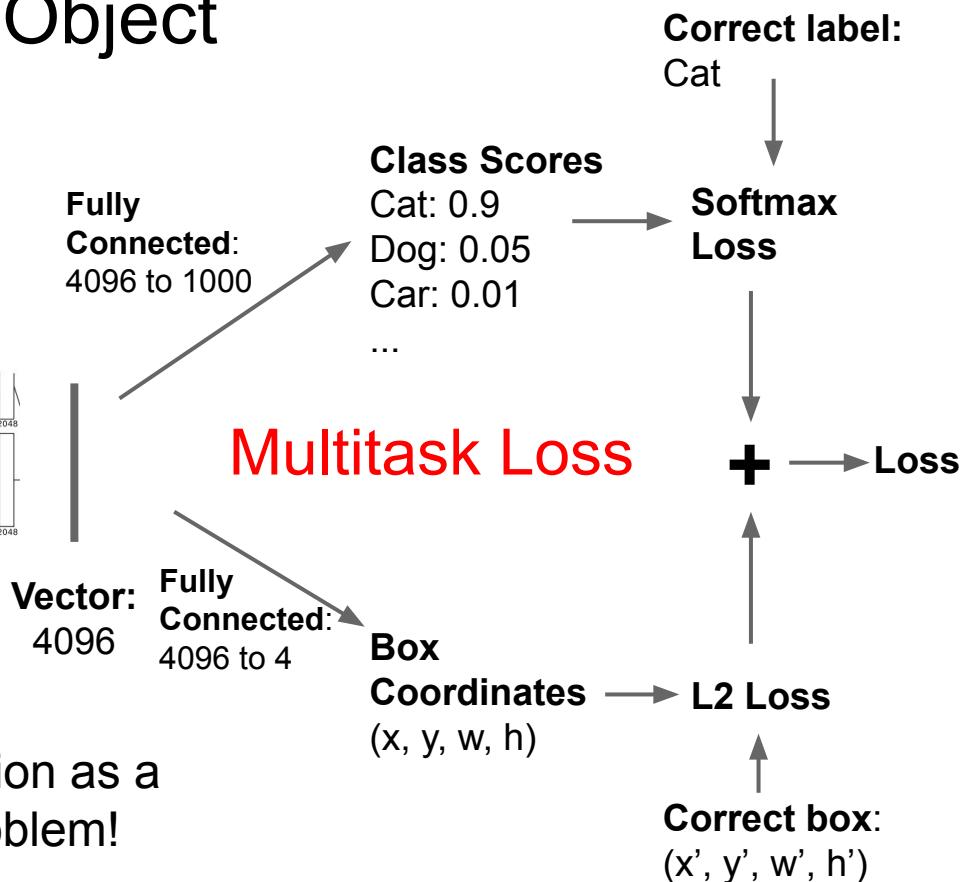
?

# Object Detection: Single Object (Classification + Localization)

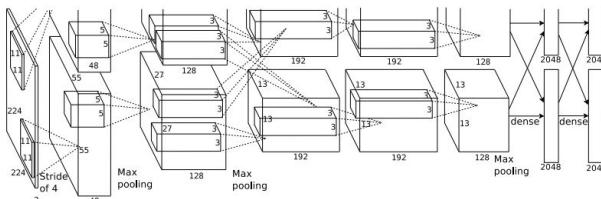


This image is CC0 public domain

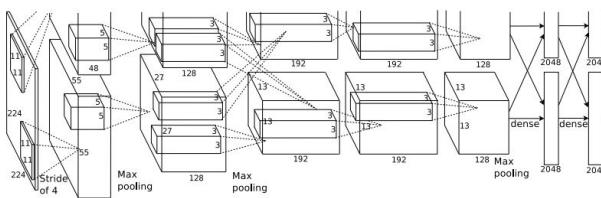
Treat localization as a regression problem!



# Object Detection: Multiple Objects



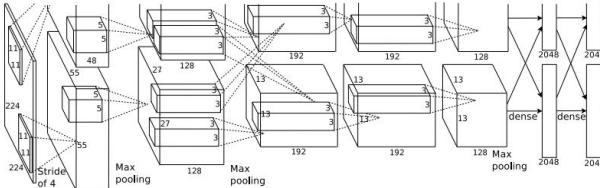
CAT: (x, y, w, h)



DOG: (x, y, w, h)

DOG: (x, y, w, h)

CAT: (x, y, w, h)



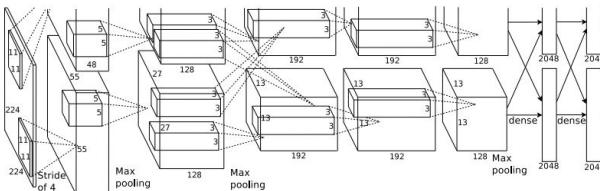
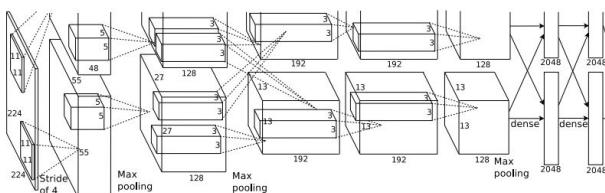
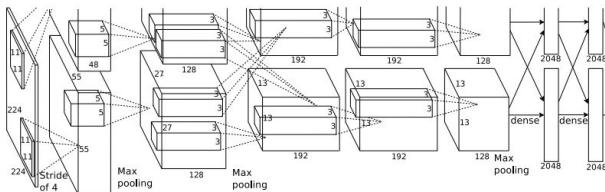
DUCK: (x, y, w, h)

DUCK: (x, y, w, h)

...

# Object Detection: Multiple Objects

Each image needs a different number of outputs!



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)

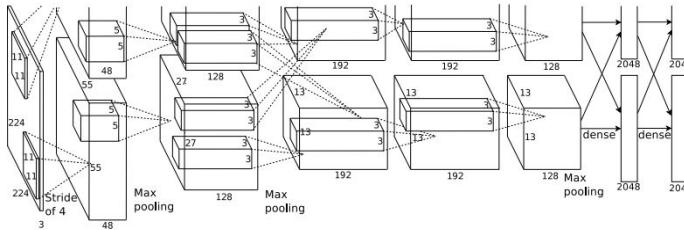
Many  
numbers!

DUCK: (x, y, w, h)

...

# Object Detection: Multiple Objects

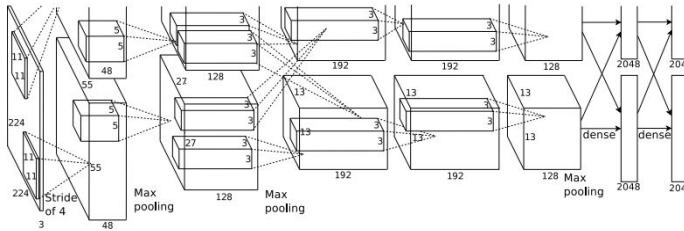
Apply a CNN to many different crops of the image, CNN classifies each crop as object or background



Dog? NO  
Cat? NO  
Background? YES

# Object Detection: Multiple Objects

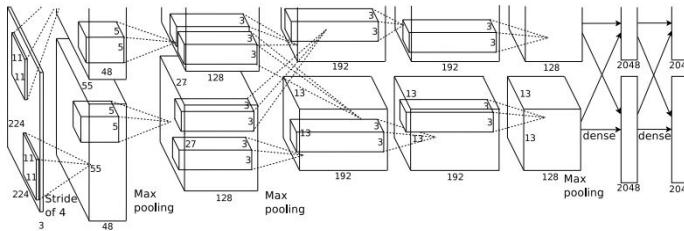
Apply a CNN to many different crops of the image, CNN classifies each crop as object or background



Dog? YES  
Cat? NO  
Background? NO

# Object Detection: Multiple Objects

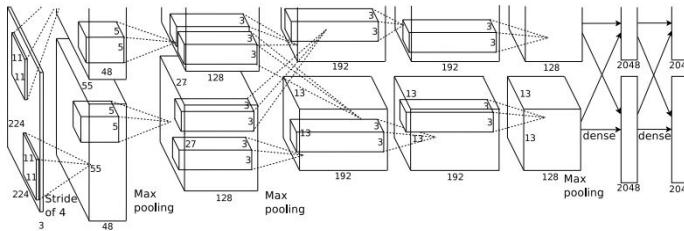
Apply a CNN to many different crops of the image, CNN classifies each crop as object or background



Dog? YES  
Cat? NO  
Background? NO

# Object Detection: Multiple Objects

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

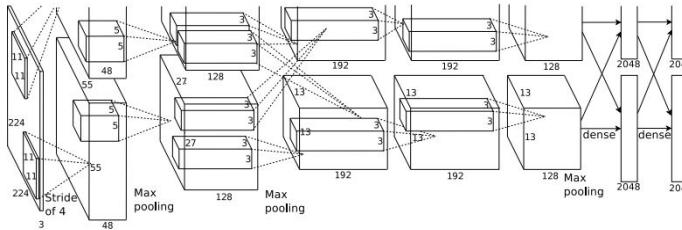
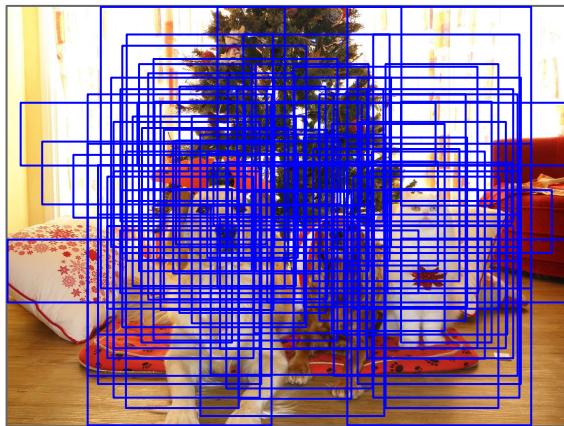


Dog? NO  
Cat? YES  
Background? NO

Q: What's the problem with this approach?

# Object Detection: Multiple Objects

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



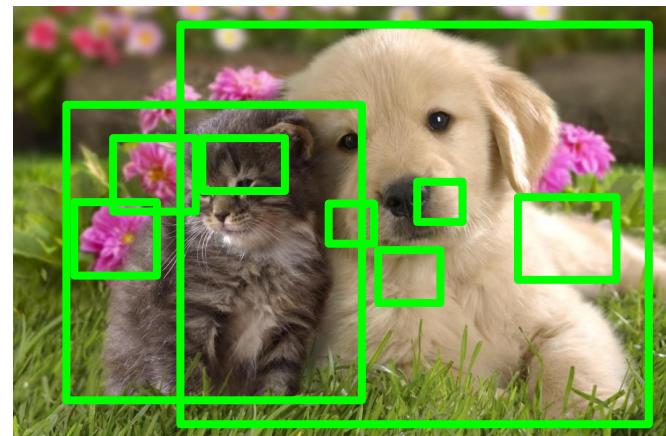
Dog? NO  
Cat? YES  
Background? NO

Problem: Need to apply CNN to huge number of locations, scales, and aspect ratios, very computationally expensive!

# Region Proposals: Selective Search

RPN

- Find “blobby” image regions that are likely to contain objects
- 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

# Faster R-CNN:

Make CNN do proposals!

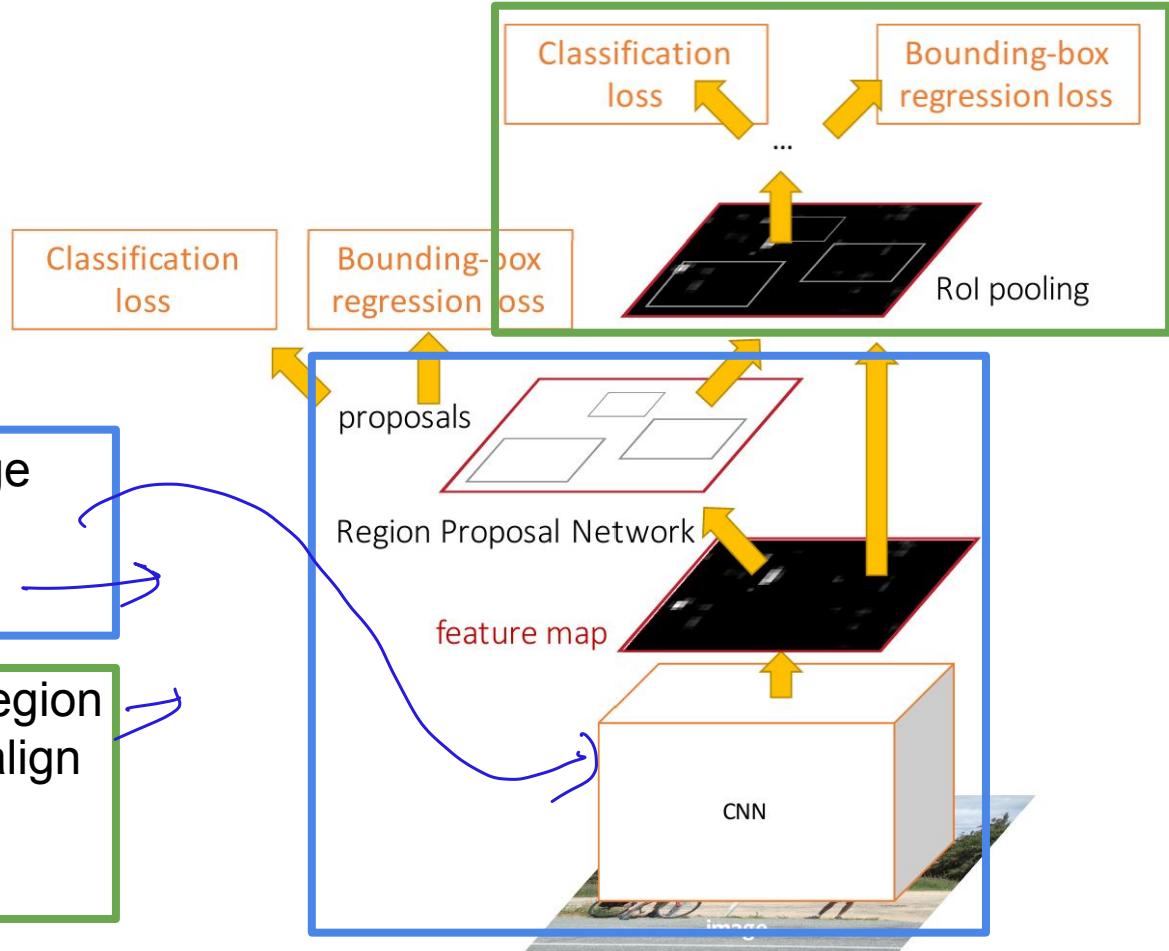
Faster R-CNN is a  
**Two-stage object detector**

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 ✓



# Region Proposal Network



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

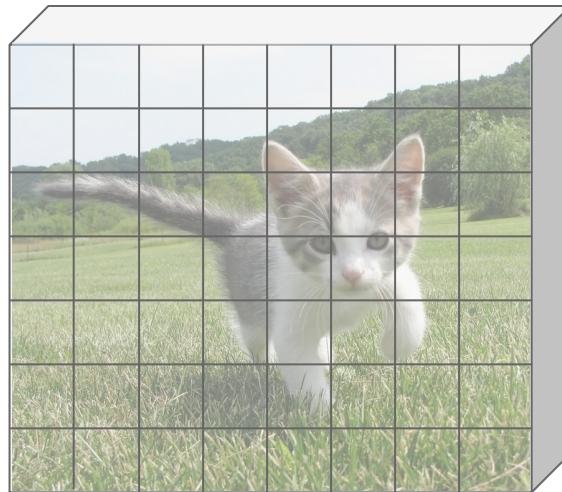
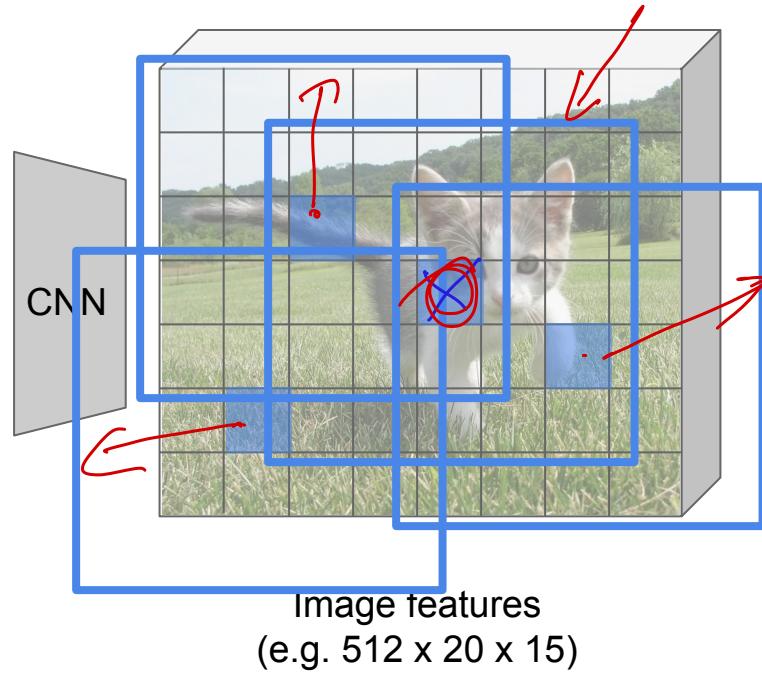


Image features  
(e.g.  $512 \times 20 \times 15$ )  
300

# Region Proposal Network



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



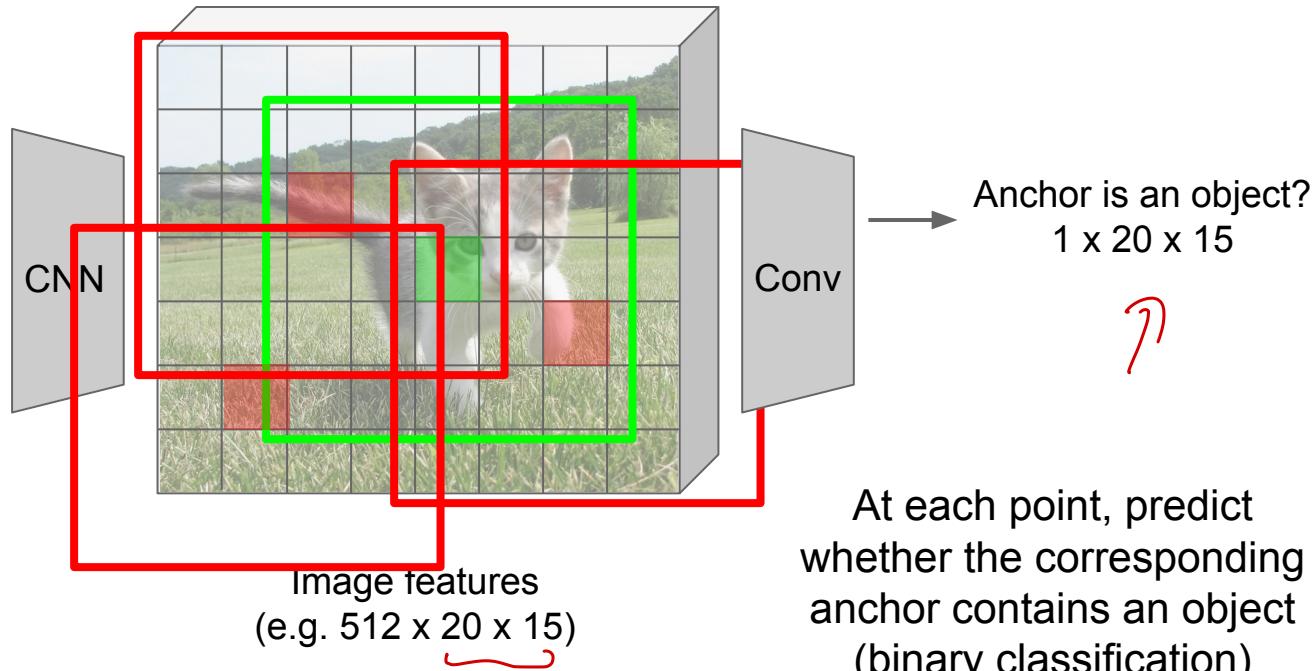
Imagine an **anchor box**  
of fixed size at each  
point in the feature map

IOU Threshold

# Region Proposal Network



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



# Region Proposal Network



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

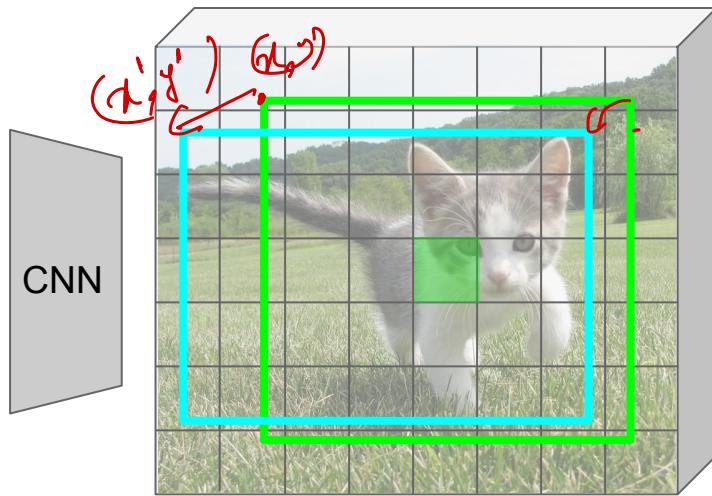
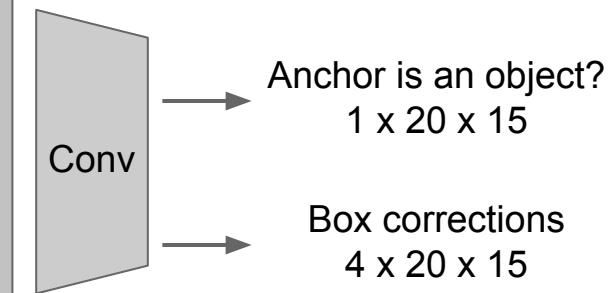


Image features  
(e.g.  $512 \times 20 \times 15$ )

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

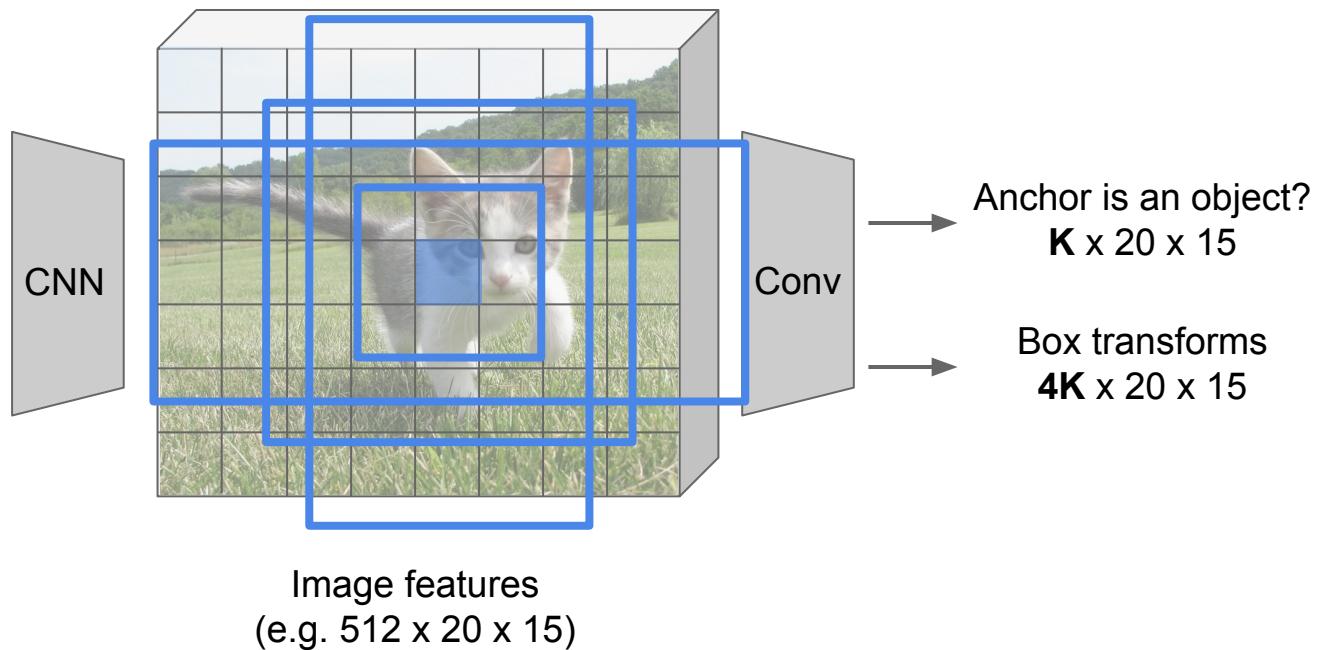


For positive boxes, also predict a corrections from the anchor to the ground-truth box (regress 4 numbers per pixel)

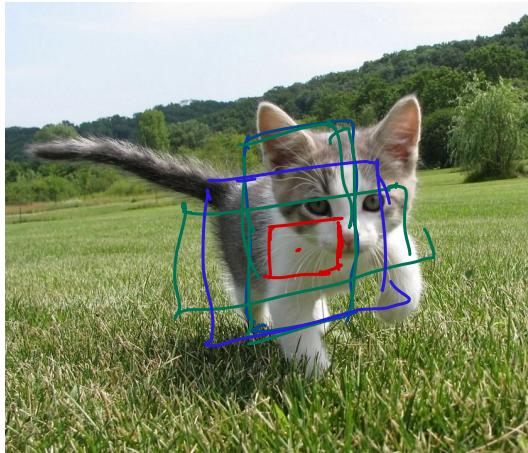
# Region Proposal Network



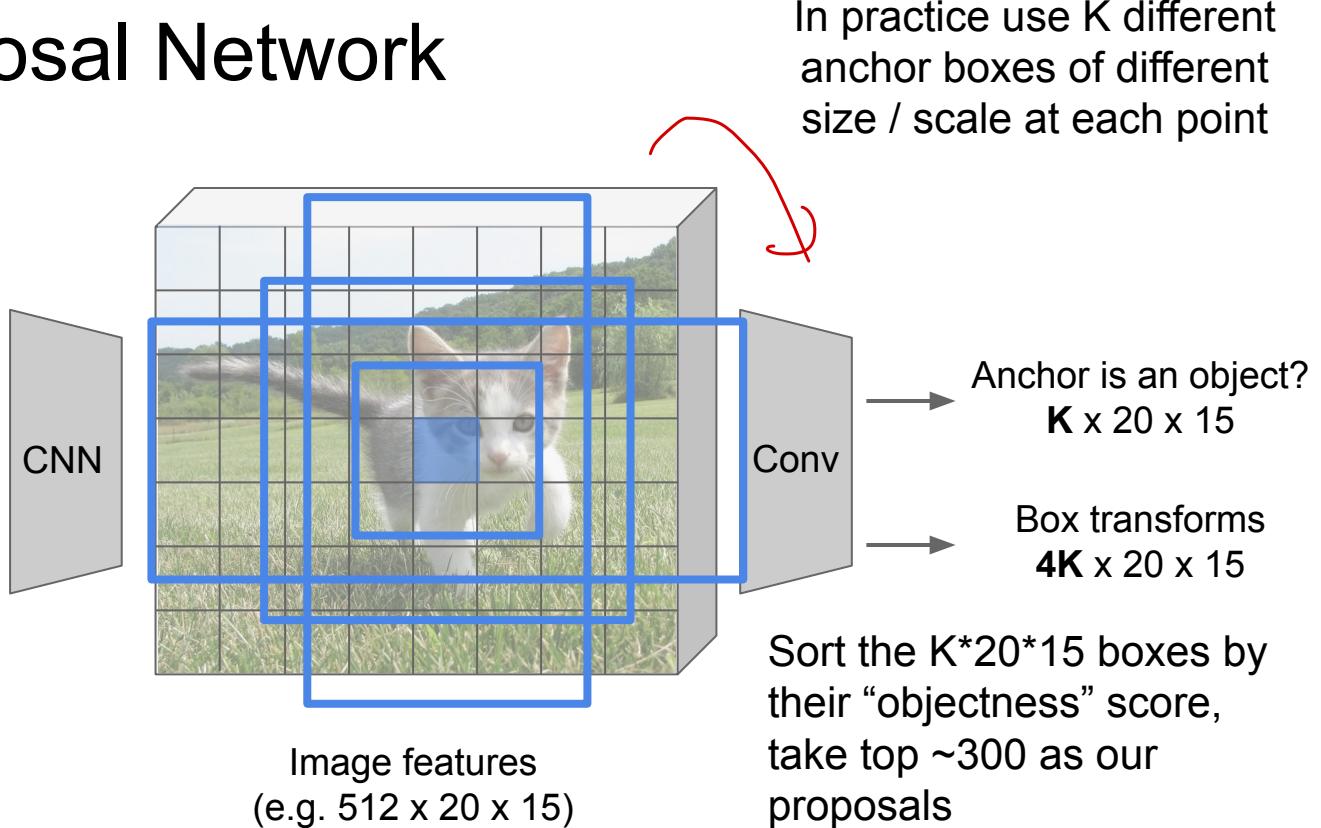
Input Image  
(e.g.  $3 \times 640 \times 480$ )

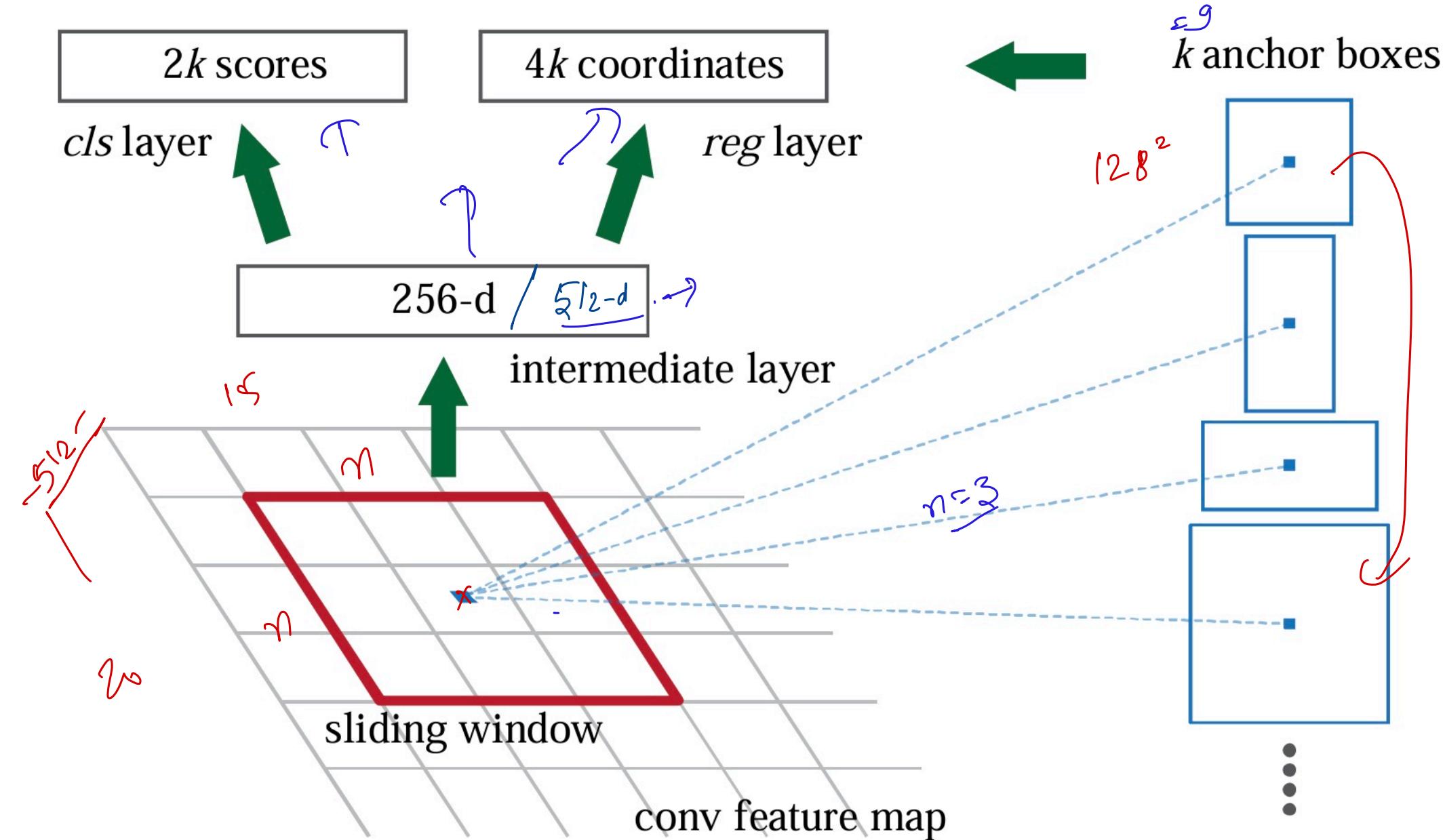


# Region Proposal Network



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



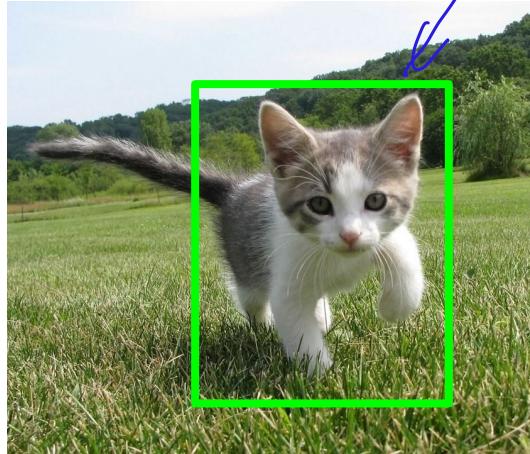


Suppose conv. feature map has 512 channels  
for a  $3 \times 3$  slidely window, and 512-d fix  
representation with 9 anchor boxes, what'll be  
the parameters for the proposal layer?



$$\overbrace{\left[ 512 \times 3 \times 3 \times 512 + 512 \times 6 \times 9 \right]}^{\text{parameters for the proposal layer}}$$

# Cropping Features: RoI Pool



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

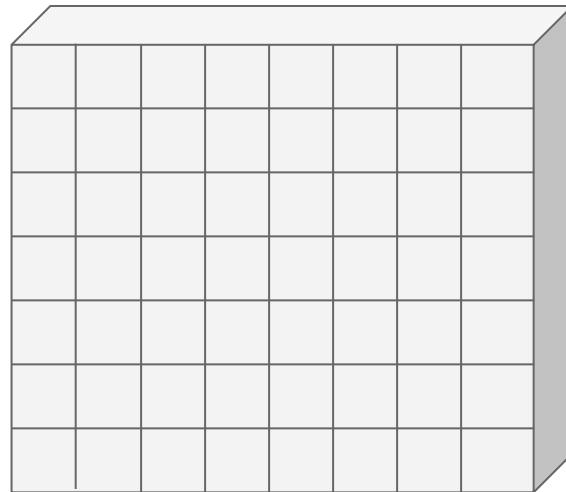
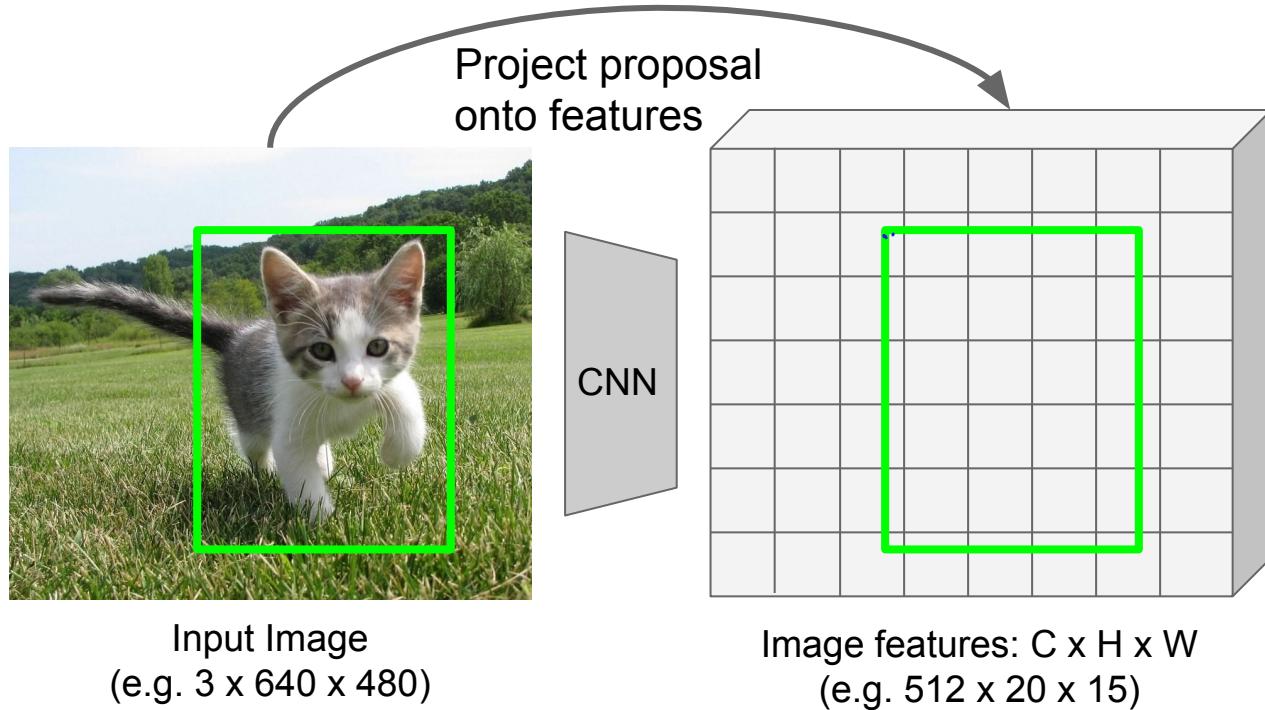


Image features:  $C \times H \times W$   
(e.g.  $512 \times 20 \times 15$ )

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

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

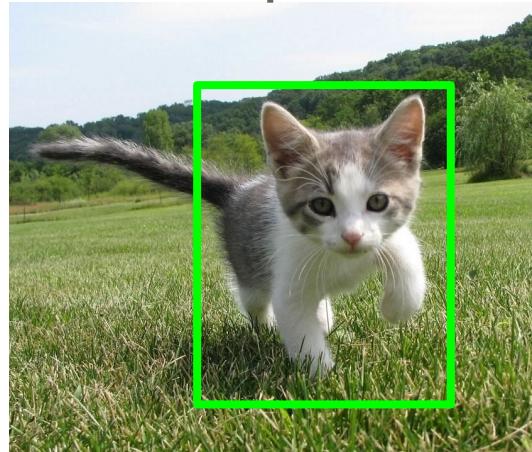
# Cropping Features: RoI Pool



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

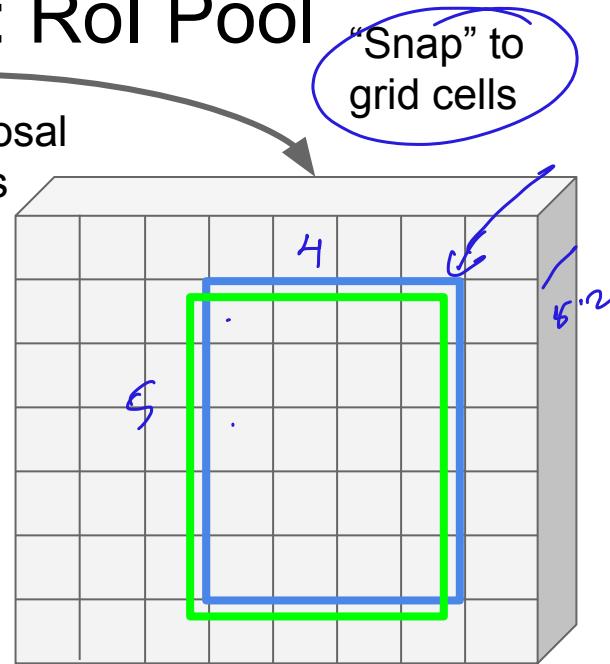
Girshick, "Fast R-CNN", ICCV 2015.

# Cropping Features: RoI Pool



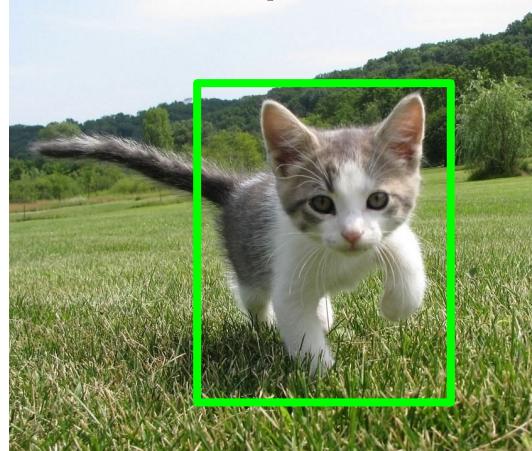
Input Image  
(e.g.  $3 \times 640 \times 480$ )

Project proposal  
onto features



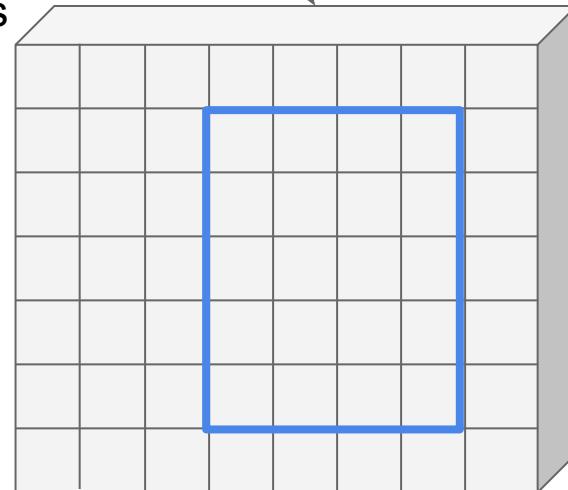
Girshick, “Fast R-CNN”, ICCV 2015.

# Cropping Features: ROI Pool



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

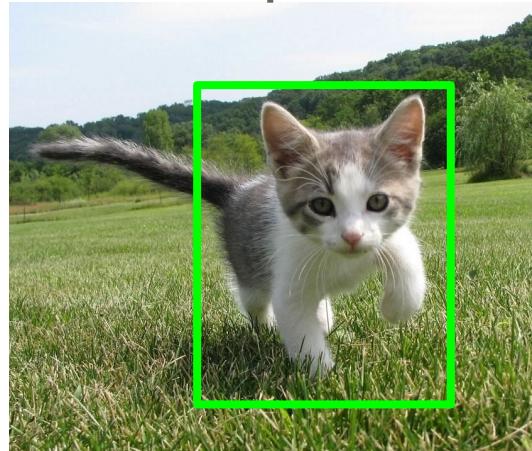
## Project proposal onto features



“Snap” to  
grid cells

Q: how do we resize the 512 x 5 x 4 region to, e.g., a 512 x 2 x 2 tensor?

# Cropping Features: RoI Pool



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

Project proposal  
onto features

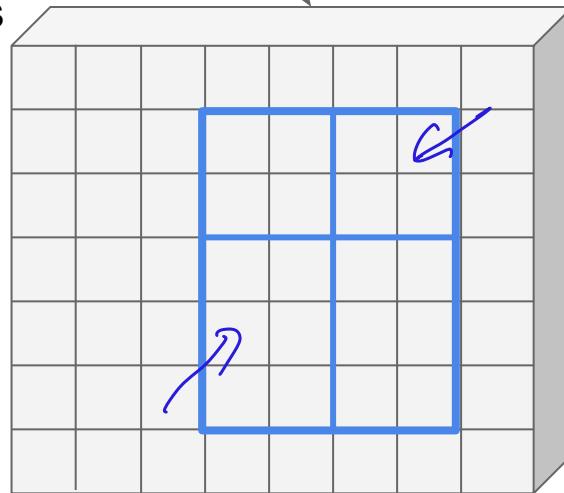


Image features:  $C \times H \times W$   
(e.g.  $512 \times 20 \times 15$ )

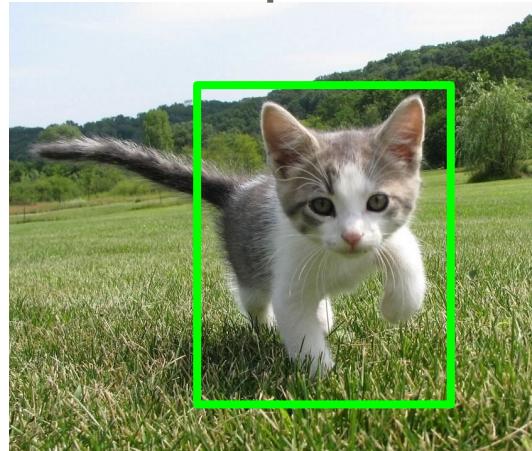
“Snap” to  
grid cells

Divide into  $2 \times 2$   
grid of (roughly)  
equal subregions

Q: how do we resize the  $512 \times 5 \times 4$  region to, e.g., a  $512 \times 2 \times 2$  tensor?



# Cropping Features: RoI Pool



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

Project proposal  
onto features

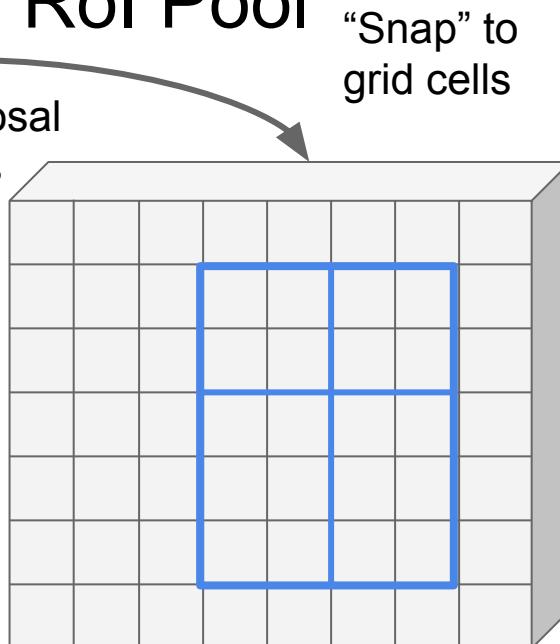
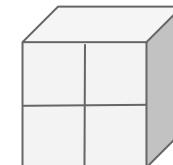


Image features:  $C \times H \times W$   
(e.g.  $512 \times 20 \times 15$ )

“Snap” to  
grid cells

Divide into  $2 \times 2$   
grid of (roughly)  
equal subregions

Max-pool within  
each subregion



Region features  
(here  $\underline{512} \times 2 \times 2$ ;  
In practice e.g.  $\underline{512} \times \underline{7} \times \underline{7}$ )

Region features always the  
same size even if input  
regions have different sizes!

Girshick, “Fast R-CNN”, ICCV 2015.

# Faster R-CNN:

Make CNN do proposals!

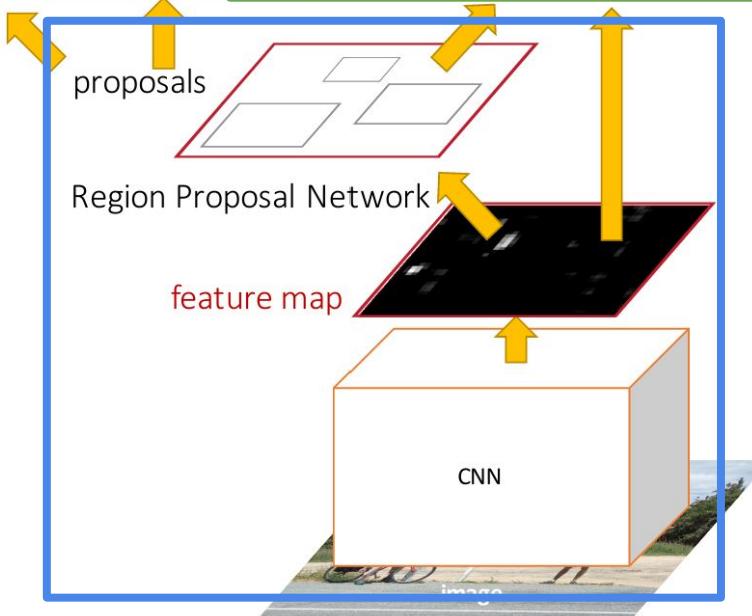
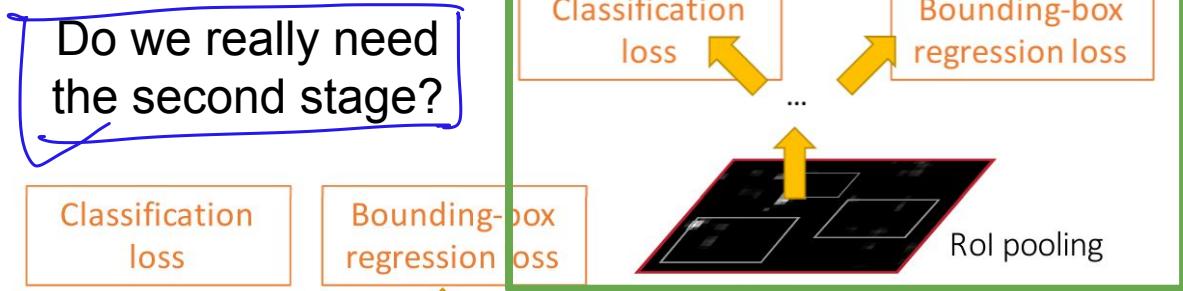
Faster R-CNN is a  
**Two-stage object detector**

First stage: Run once per image

- Backbone network
- Region proposal network ✓

Second stage: Run once per region

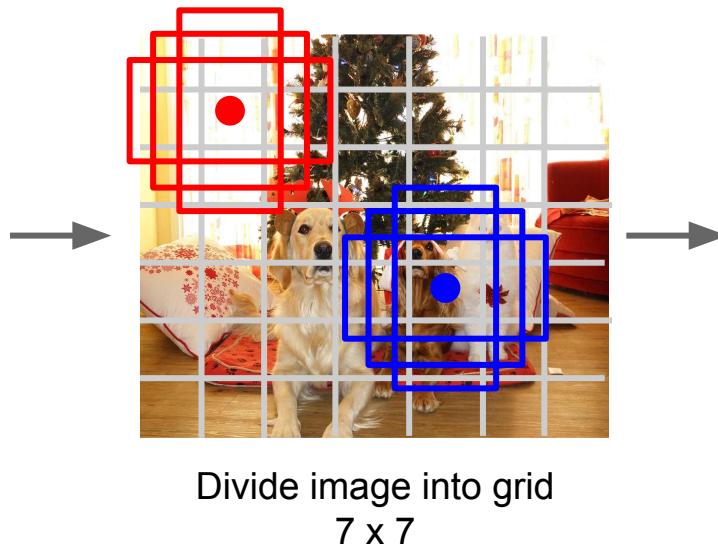
- Crop features: RoI pool / align
- Predict object class ↴
- Prediction bbox offset



# Single-Stage Object Detectors: YOLO / SSD / RetinaNet



Input image  
 $3 \times H \times W$



Divide image into grid  
 $7 \times 7$

Image a set of **base boxes**  
centered at each grid cell  
Here  $B = 3$

Within each grid cell:

- Regress from each of the  $B$  base boxes to a final box with 5 numbers:  
( $dx$ ,  $dy$ ,  $dh$ ,  $dw$ , confidence)
- Predict scores for each of  $C$  classes (including background as a class)
- Looks a lot like RPN, but category-specific!

Output:  $\nearrow^{10}$   
 $7 \times 7 \times (5 * B + C)$

Redmon et al, "You Only Look Once: Unified, Real-Time Object Detection", CVPR 2016  
Liu et al, "SSD: Single-Shot MultiBox Detector", ECCV 2016  
Lin et al, "Focal Loss for Dense Object Detection", ICCV 2017

# How does the output look like?

