

# Object Detection (Plus some bonuses)

EECS 442 – David Fouhey

Winter 2023, University of Michigan

[http://web.eecs.umich.edu/~fouhey/teaching/EECS442\\_W23/](http://web.eecs.umich.edu/~fouhey/teaching/EECS442_W23/)

# Last Time

“Semantic Segmentation”: Label each pixel with the object category it belongs to.

Input



Target



# Today – Object Detection

“Object Detection”: Draw a box around each instance of a list of categories

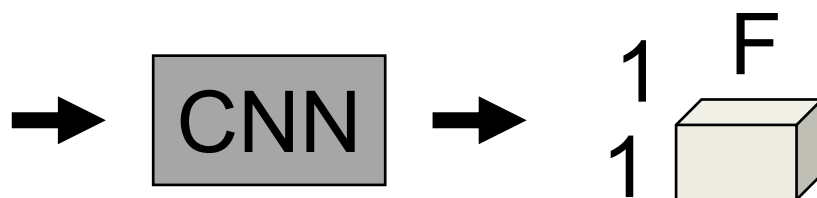
Input



Target



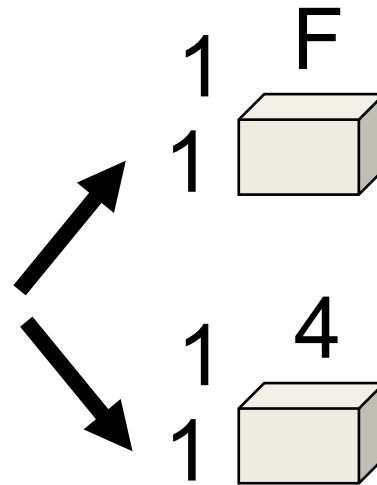
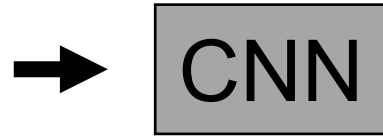
# The Wrong Way To Do It



## Starting point:

Can predict the probability of F classes  
 $P(\text{cat}), P(\text{goose}), \dots P(\text{tractor})$

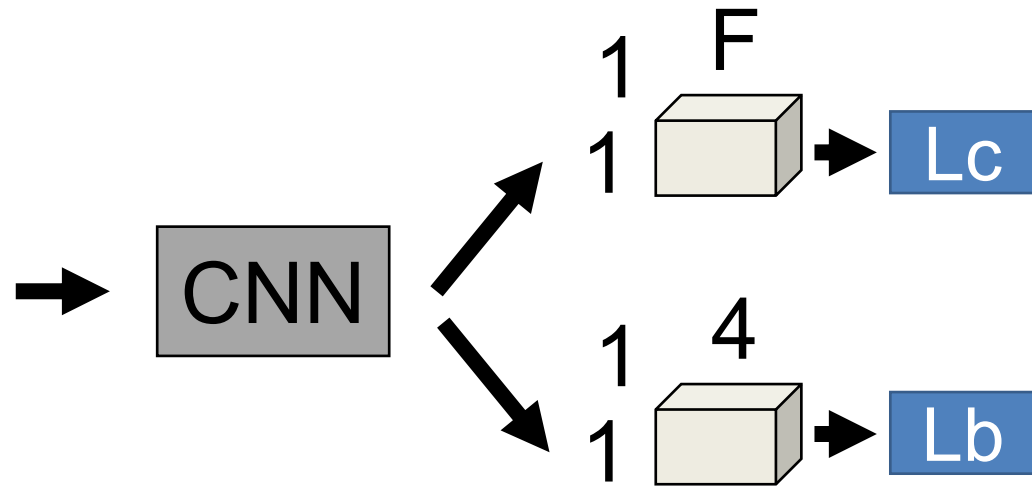
# The Wrong Way To Do It



**Add another output (why not):**

Predict the *bounding box* of the object  
[x,y,width,height] or [minX,minY,maxX,maxY]

# The Wrong Way To Do It



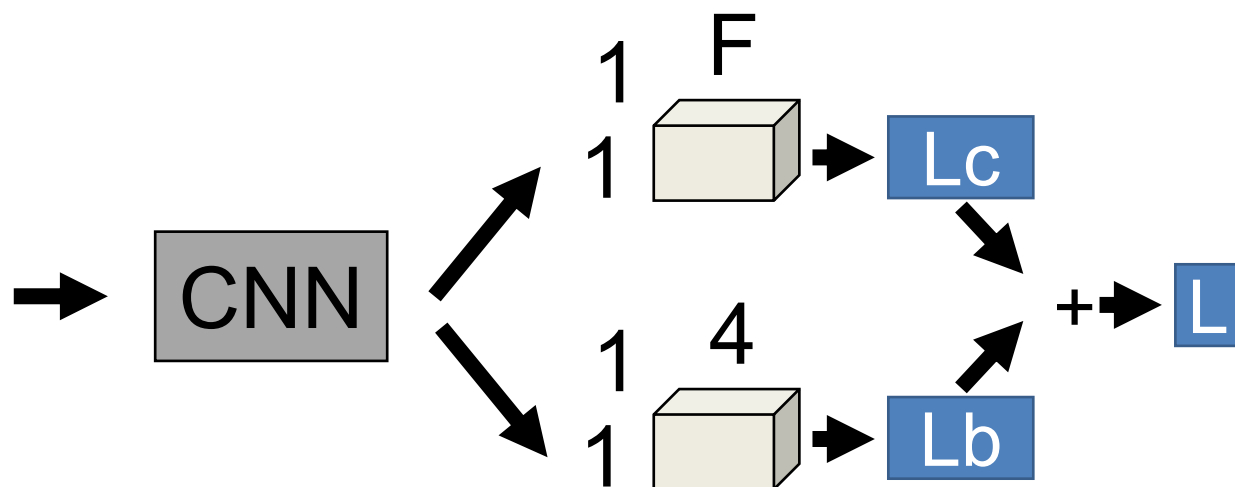
**Put a loss on it:**

Penalize mistakes on the classes with

$L_c = \text{negative log-likelihood}$

$L_b = \text{L2 loss}$

# The Wrong Way To Do It

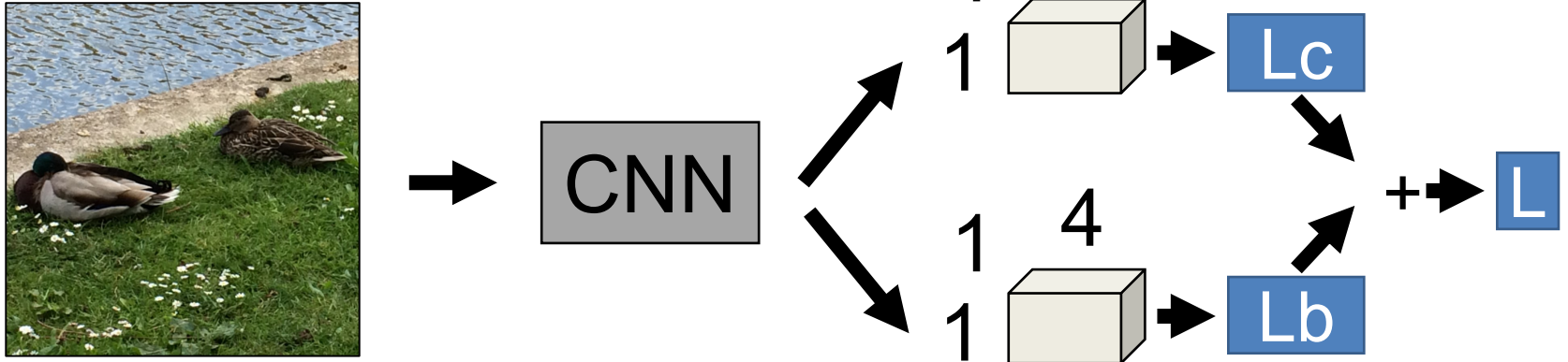


**Add losses, backpropagate**

$$\text{Final loss: } L = L_c + \lambda L_b$$

**Why do we need the  $\lambda$ ?**

# The Wrong Way To Do It



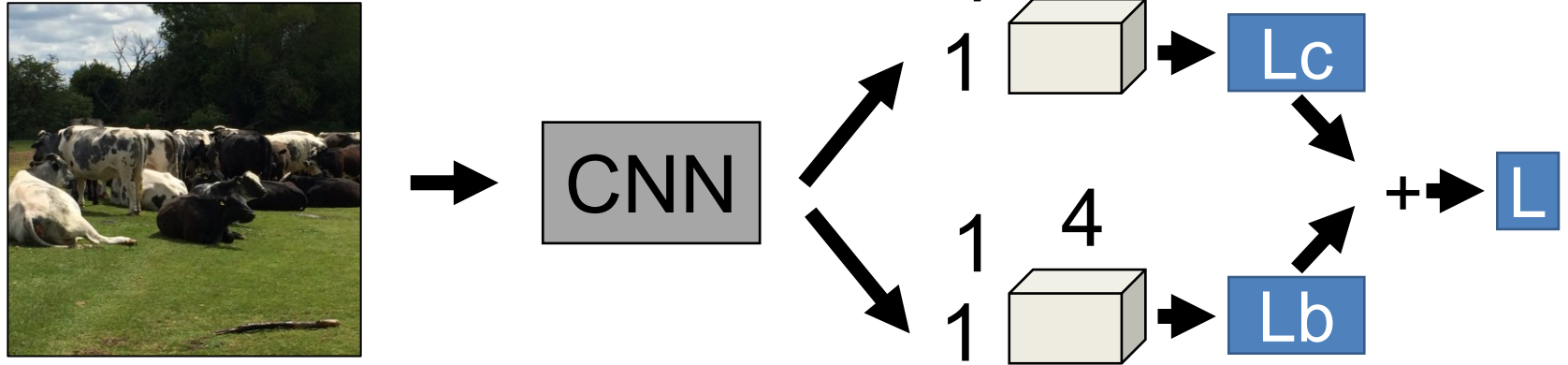
Now there are two ducks.

**How many outputs do we need?**

$$F, 4, F, 4 = 2*(F+4)$$



# The Wrong Way To Do It



Now it's a herd of cows.

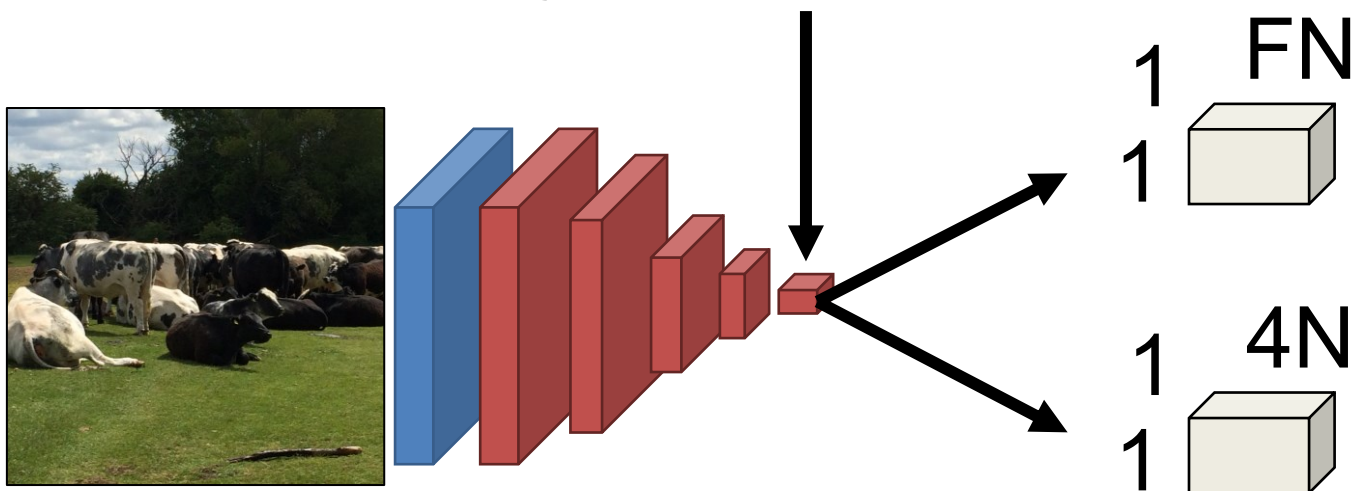
We need *lots* of outputs

(in fact the precise number of objects that are in the image, which is circular reasoning).

# In General

- Usually can't do varying-size outputs.
- Even if we could, think about how *you* would solve it if you were a network.

Bottleneck has to *encode* where the objects are for all objects and all N

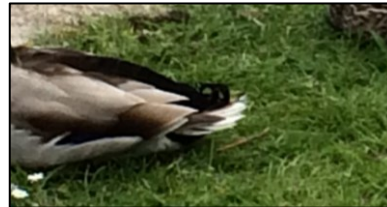


# An Alternate Approach

Examine every sub-window and determine if it is a tight box around an object

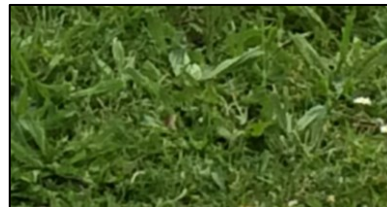


Yes



No?

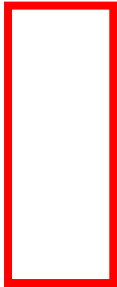
Hold this thought



No

# Sliding Window Classification

Let's assume we're looking for pedestrians in a box with a fixed aspect ratio.



# Sliding Window

Key idea – just try all the subwindows in the image at all positions.



# Generating hypotheses

Key idea – just try all the subwindows in the image at all positions **and scales**.



Note – Template did not change size

# Each window classified separately



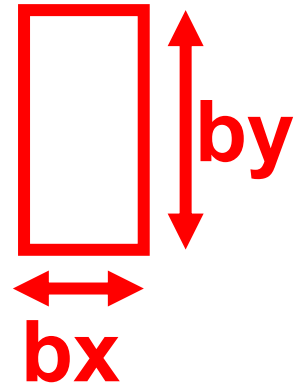
Slide credit: J. Hays

# How Many Boxes Are There?

Given a HxW image and a “template” of size by, bx.

**Q. How many sub-boxes are there of size (by,bx)?**

**A.  $(H-by)*(W-bx)$**



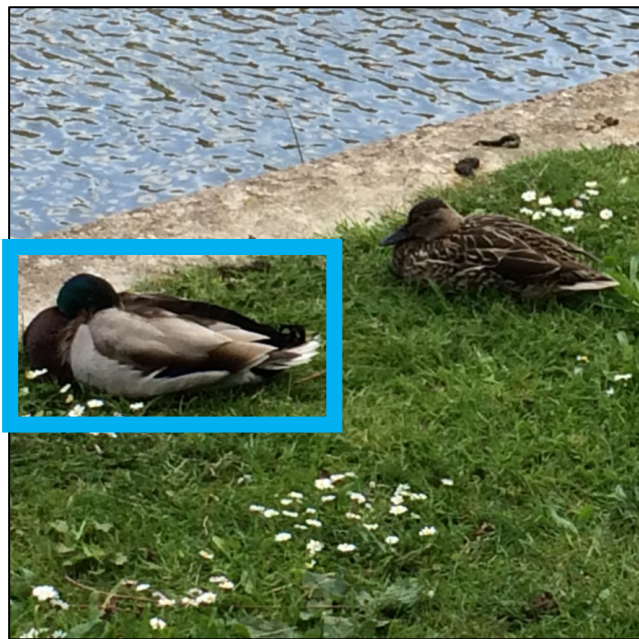
This is before considering adding:

- *scales* ( $by*s, bx*s$ )
- *aspect ratios* ( $by*sy, bx*sx$ )



# Challenges of Object Detection

- Have to evaluate *tons* of boxes
- Positive instances of objects are *extremely* rare



How many ways can we get the box wrong?

1. Wrong left x
2. Wrong right x
3. Wrong top y
4. Wrong bottom y

# Prime-time TV



Are You Smarter Than  
A 5<sup>th</sup> Grader?

Adults compete with  
5<sup>th</sup> graders on  
elementary school  
facts.

Adults often not  
smarter.

# Computer Vision TV



Are You Smarter Than  
A Random Number  
Generator?

Models trained on data  
compete with making  
random guesses.

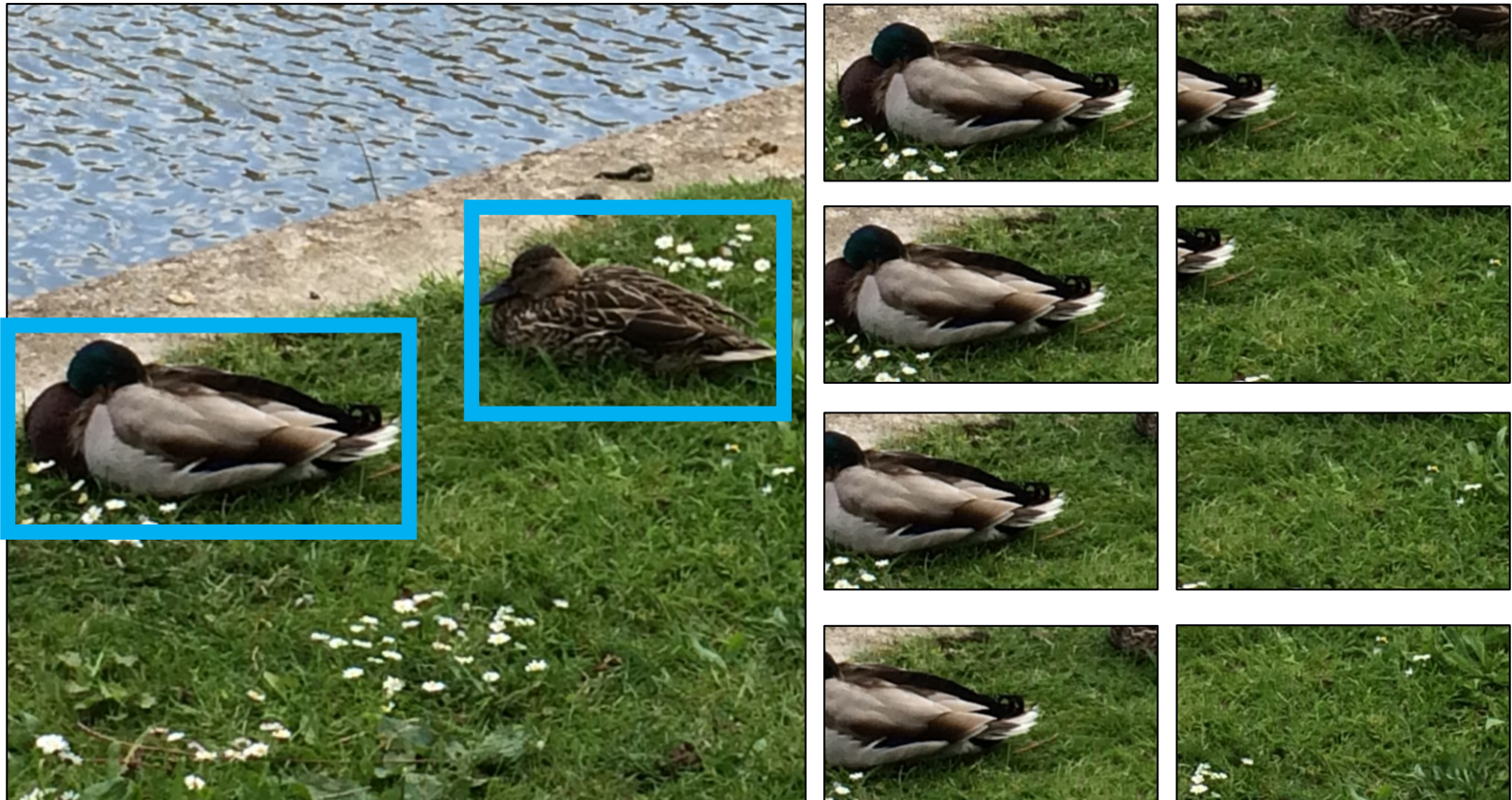
Models often not  
better.

# Are You Smarter than a Random Number Generator?

- **Prob. of guessing 1k-way classification?**
  - 1/1,000
- **Prob. of guessing all 4 bounding box corners within 10% of image size?**
  - $(1/10)*(1/10)*(1/10)*(1/10)=1/10,000$
- Probability of guessing both: 1/10,000,000
- Detection is **hard** (via guessing and in general)
- Should *always compare* against guessing or picking most likely output label

# Evaluating – Bounding Boxes

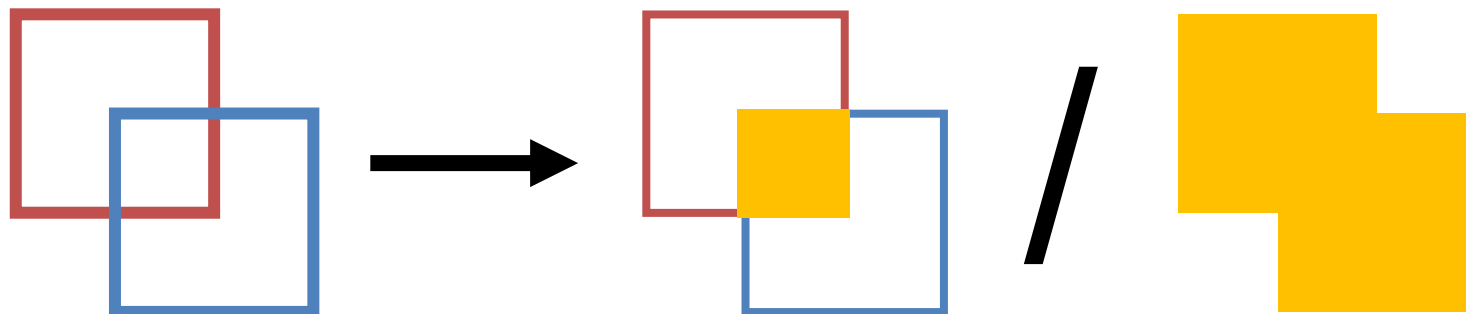
**Raise your hand when you think the detection stops being correct.**



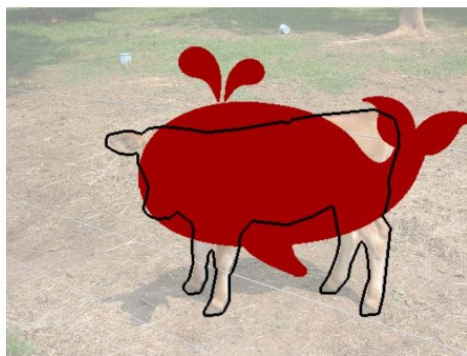
# Evaluating – Bounding Boxes

**Standard metric for two boxes:**

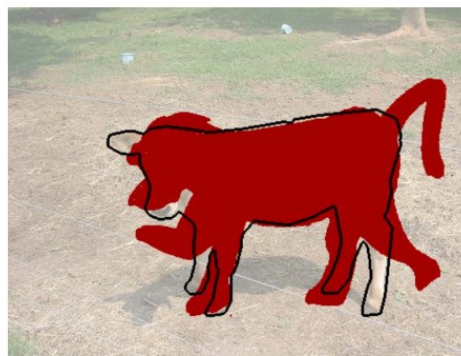
Intersection over union/IoU/Jaccard coefficient



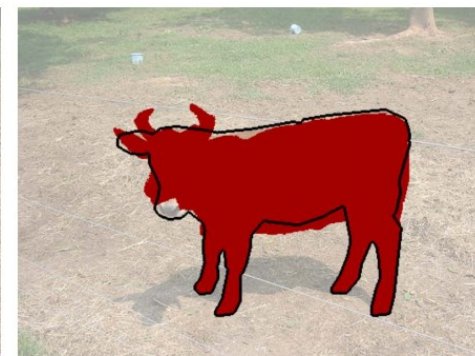
(a) Ground truth



(b)  $\mathcal{J} = 0.554$



(c)  $\mathcal{J} = 0.703$



(d)  $\mathcal{J} = 0.910$

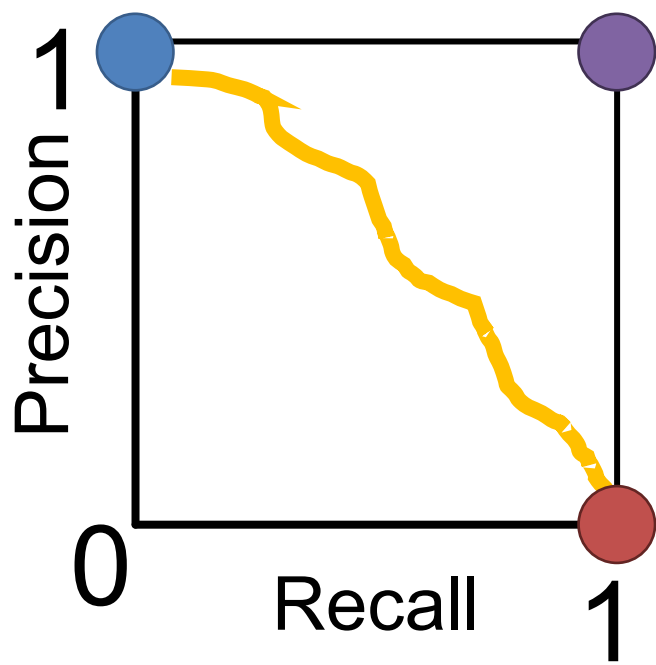
# Evaluating Performance

- Remember: accuracy = average of whether prediction is correct
- Suppose I have a system that gets 99% accuracy in person detection.
- **What's wrong?**
- I can get that by just saying no object everywhere!

# Evaluating Performance

- True detection aka true positive: high IoU
- Precision:  $\# \text{true detections} / \# \text{detections by detector}$
- Recall:  $\# \text{true detections} / \# \text{ground truth positives}$

**Reject everything: no mistakes**



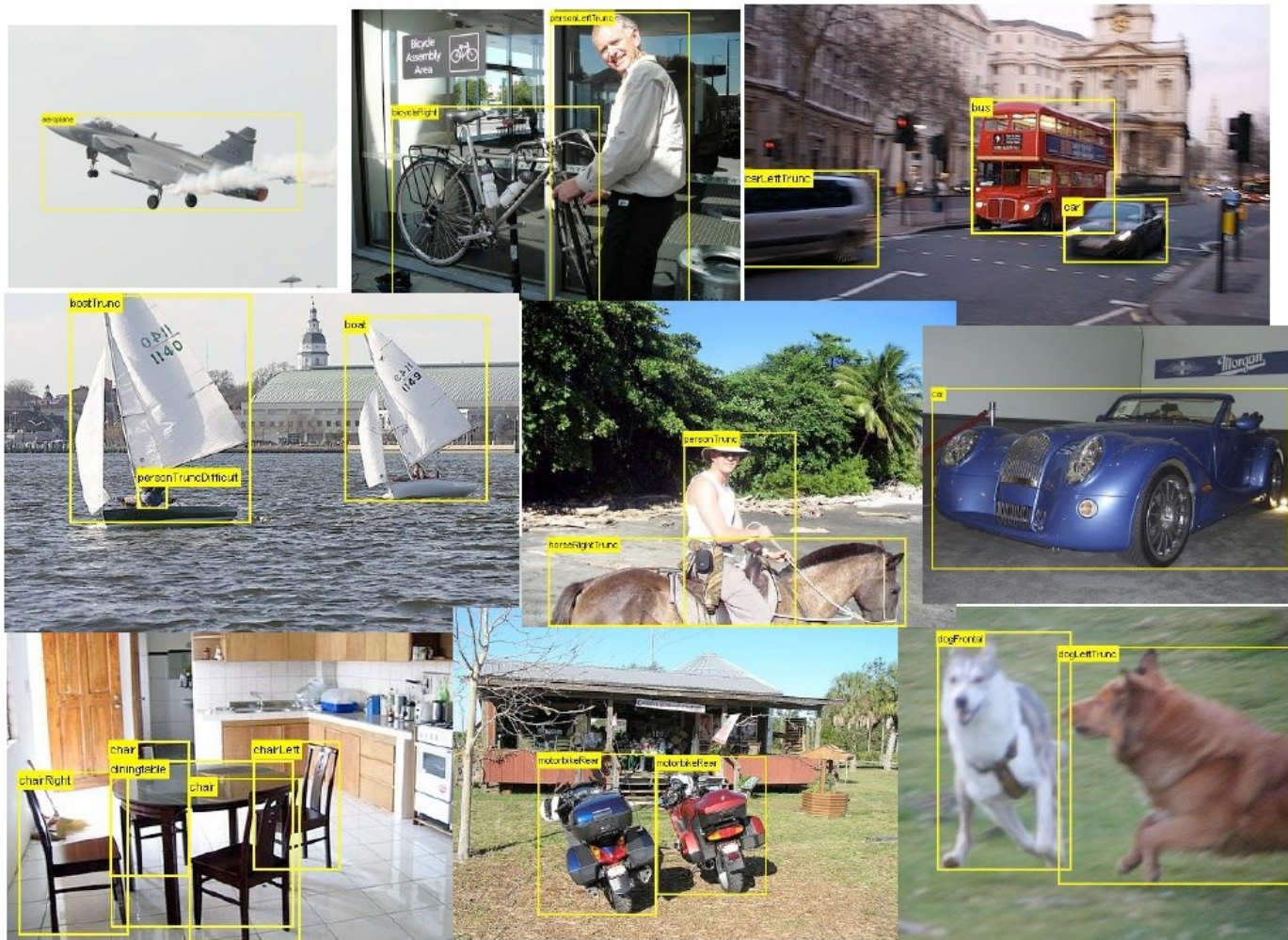
**Ideal!**

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

**Accept everything:  
Miss nothing**



# Generic object detection



# Histograms of oriented gradients (HOG)

Partition image into blocks and compute histogram of gradient orientations in each block

$H \times W \times 3$  Image



$H' \times W' \times C'$  Image

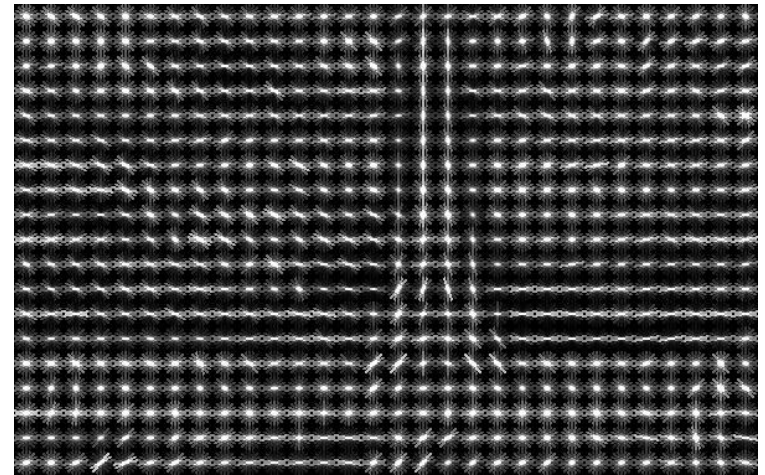


Image credit: N. Snavely

# Pedestrian detection with HOG

- Train a pedestrian template using a linear support vector machine

positive training examples



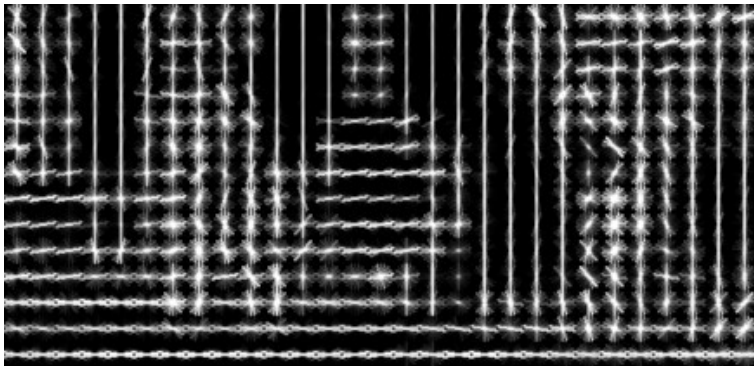
negative training examples



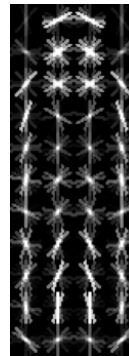
# Pedestrian detection with HOG

- Train pedestrian “template” using a linear svm
- At test time, convolve feature map with template
- Find local maxima of response
- For multi-scale detection, repeat over multiple levels of a HOG *pyramid*

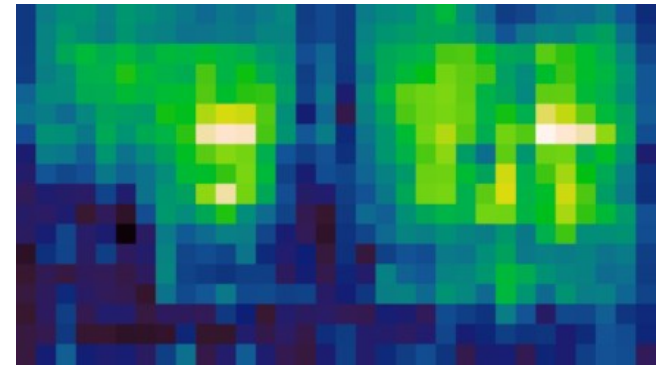
HOG feature map



Template



Detector response map



# Example detections



[Dalal and Triggs, CVPR 2005]

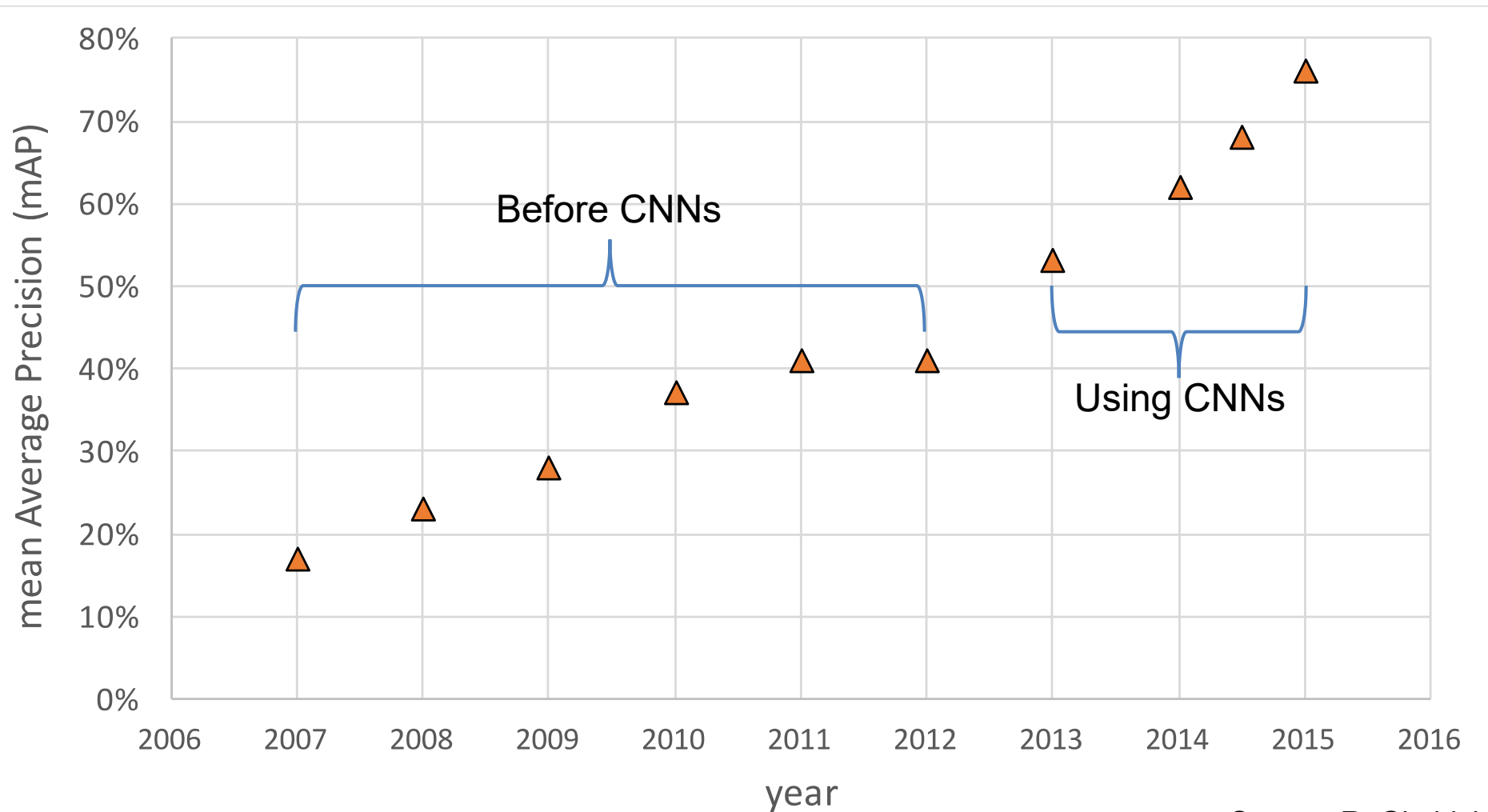
# PASCAL VOC Challenge (2005-2012)



- 20 challenge classes:
- *Person*
- *Animals*: bird, cat, cow, dog, horse, sheep
- *Vehicles*: aeroplane, bicycle, boat, bus, car, motorbike, train
- *Indoor*: bottle, chair, dining table, potted plant, sofa, tv/monitor
- Dataset size (by 2012): 11.5K training/validation images, 27K bounding boxes, 7K segmentations

# Object detection progress

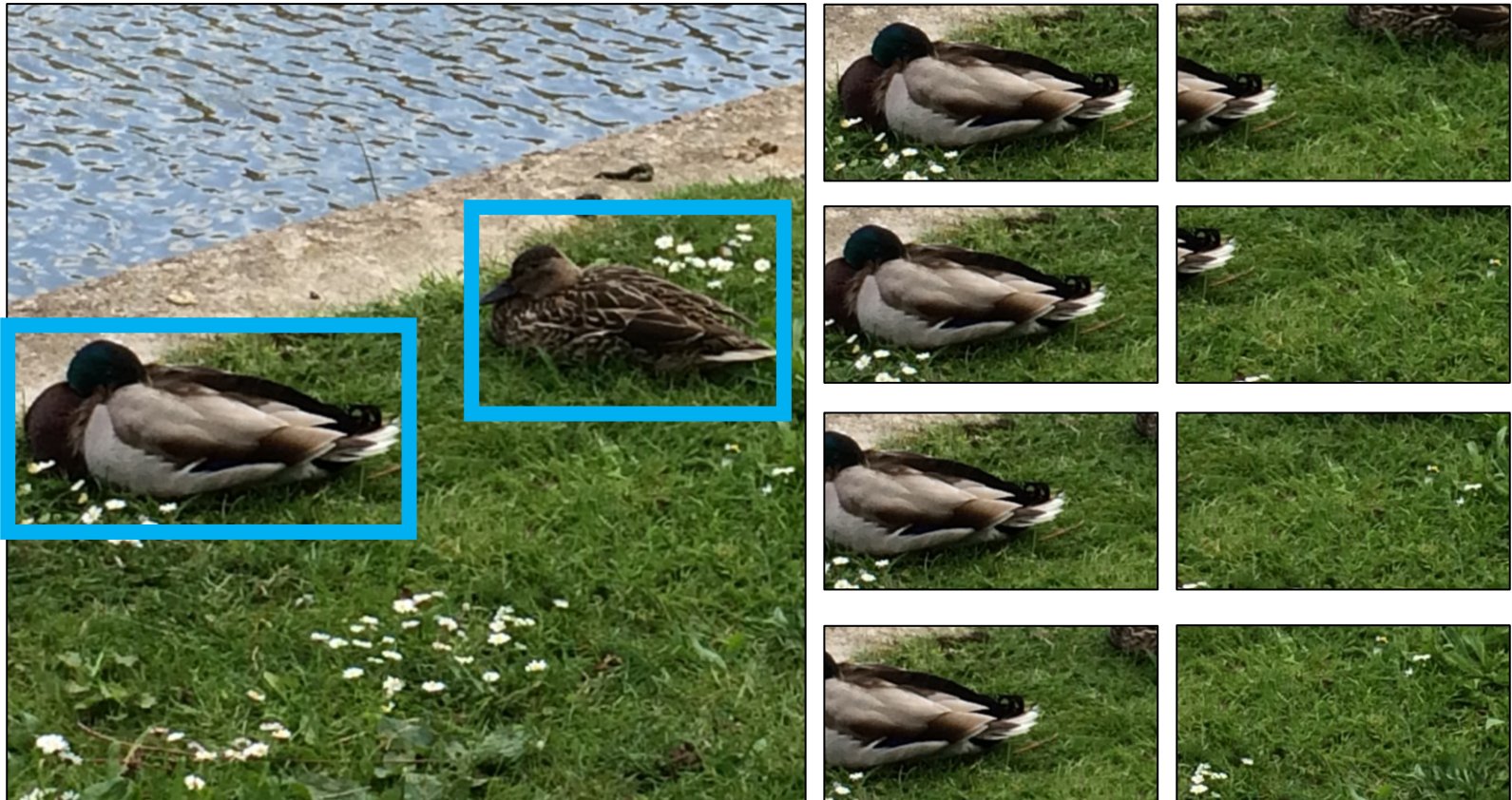
PASCAL VOC



Source: R. Girshick

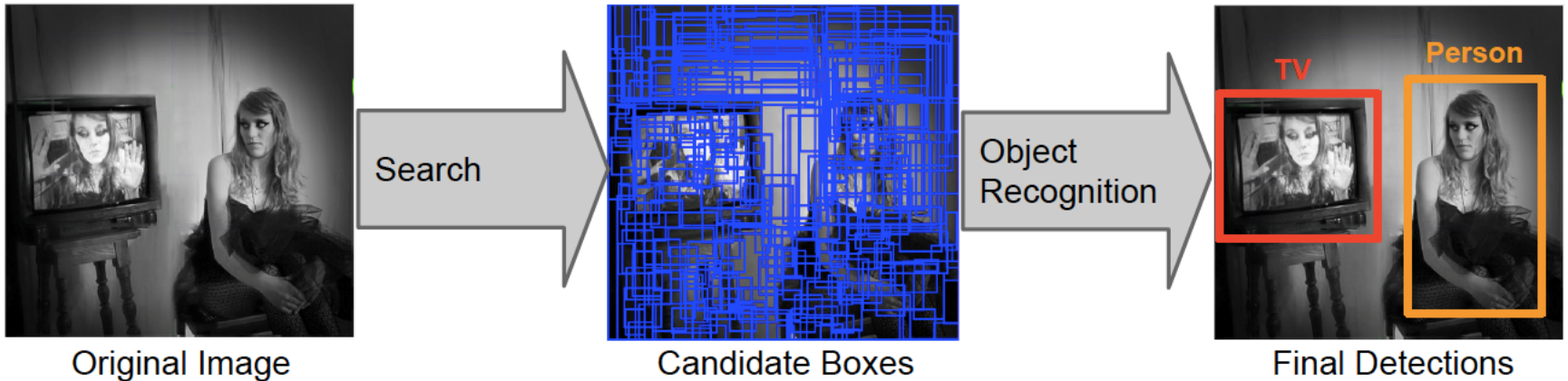
# Region Proposals

Do I need to spend a lot of time filtering all the boxes covering grass?





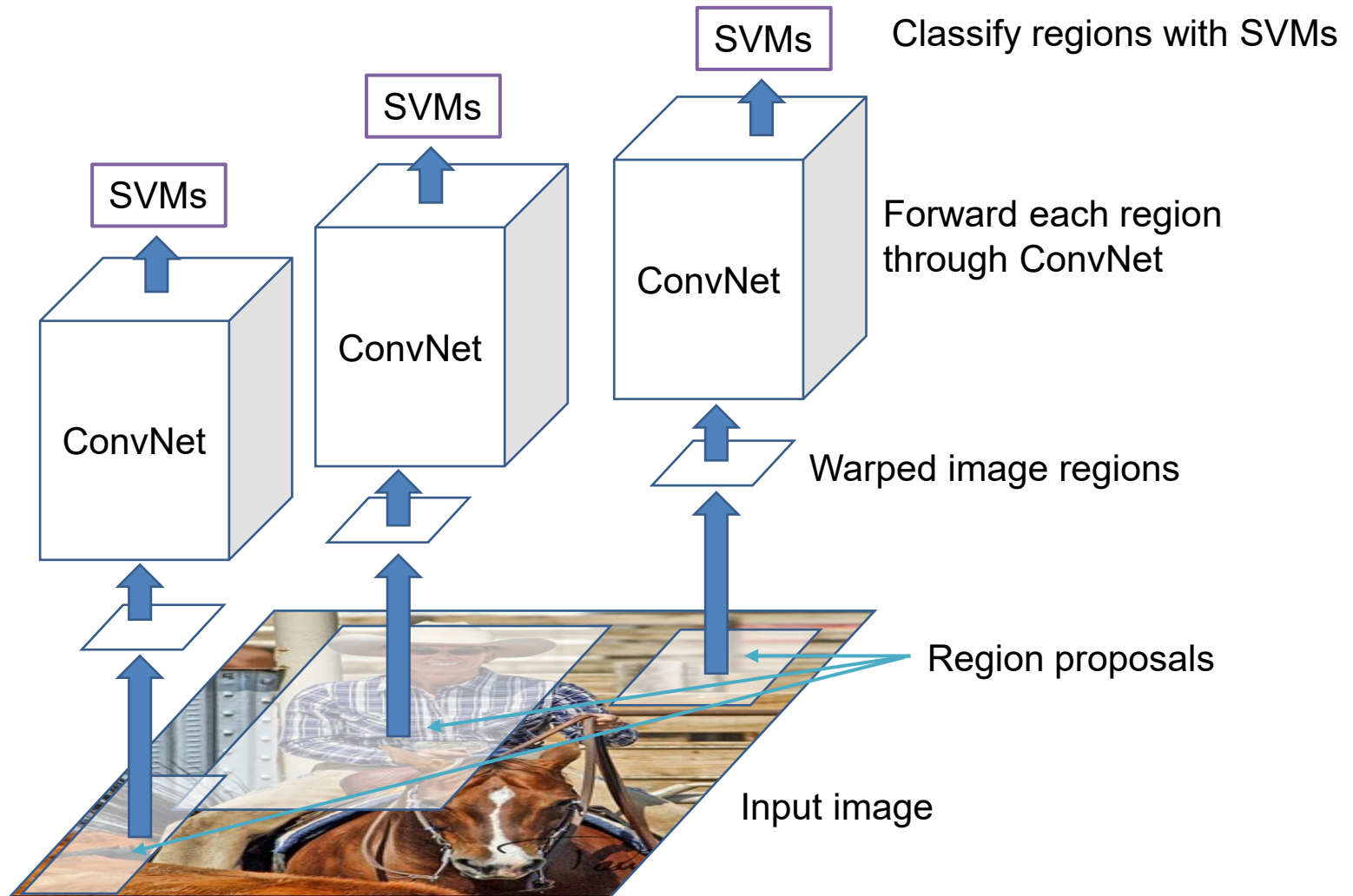
# Region Proposals



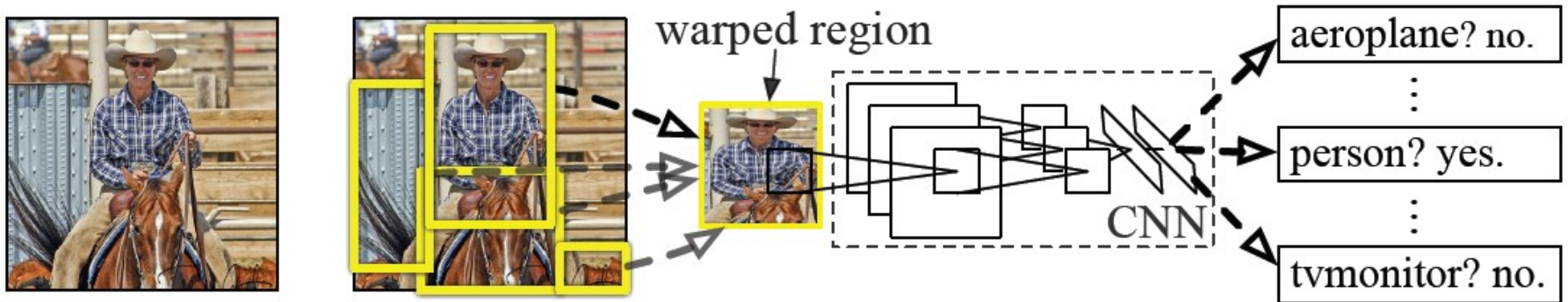
- As an alternative to sliding window search, evaluate a few hundred *region proposals*
  - Can use slower but more powerful features and classifiers
  - Proposal mechanism can be category-independent
  - Proposal mechanism can be trained

# R-CNN: Region proposals + CNN features

Source: R. Girshick



# R-CNN details

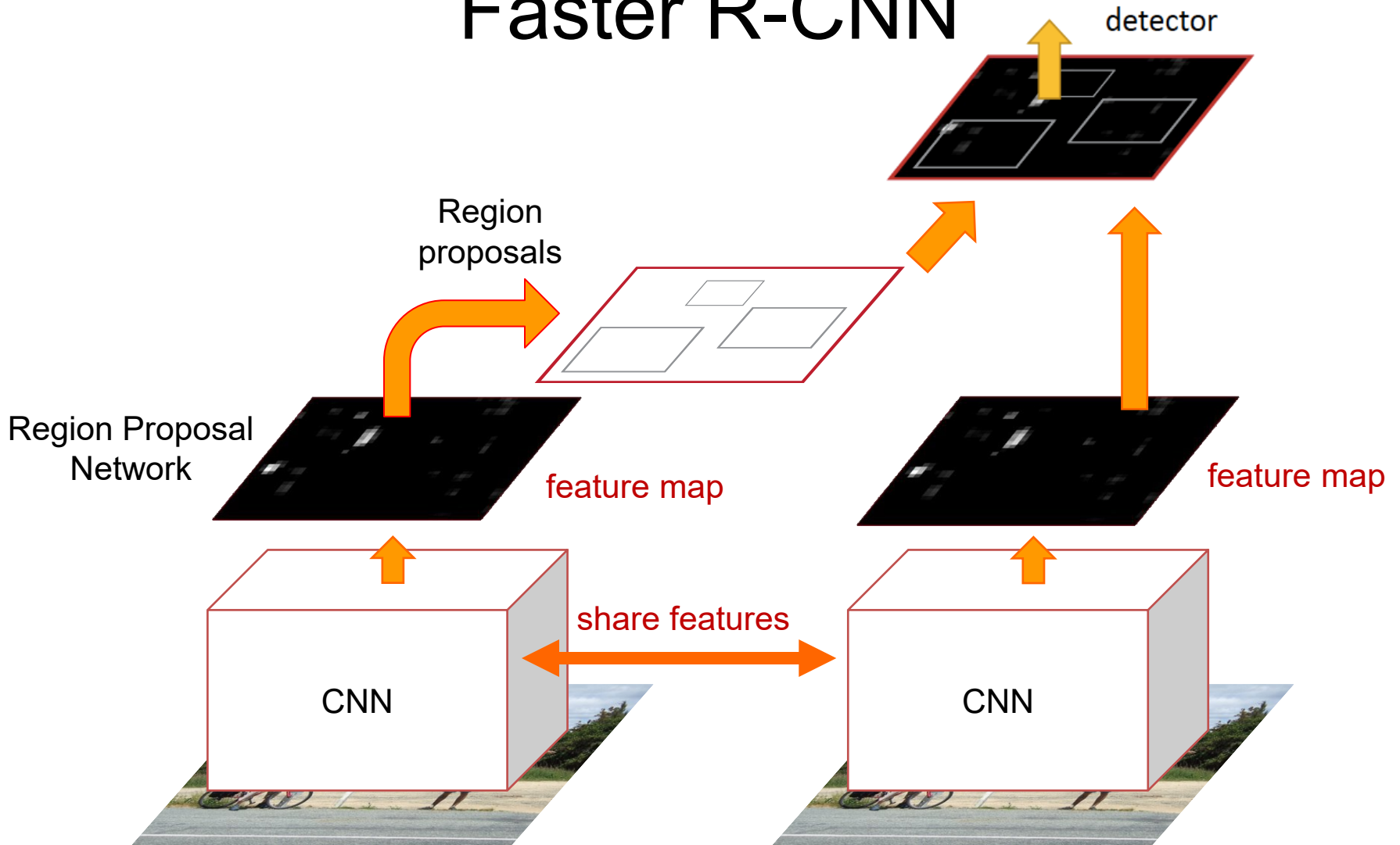


- **Regions:** ~2000 Selective Search proposals
- **Network:** AlexNet *pre-trained* on ImageNet (1000 classes), *fine-tuned* on PASCAL (21 classes)
- **Final detector:** warp proposal regions, extract fc7 network activations (4096 dimensions), classify with linear SVM
- **Bounding box regression** to refine box locations
- **Performance:** mAP of 53.7% on PASCAL 2010 (vs. 35.1% for Selective Search and 33.4% for DPM).

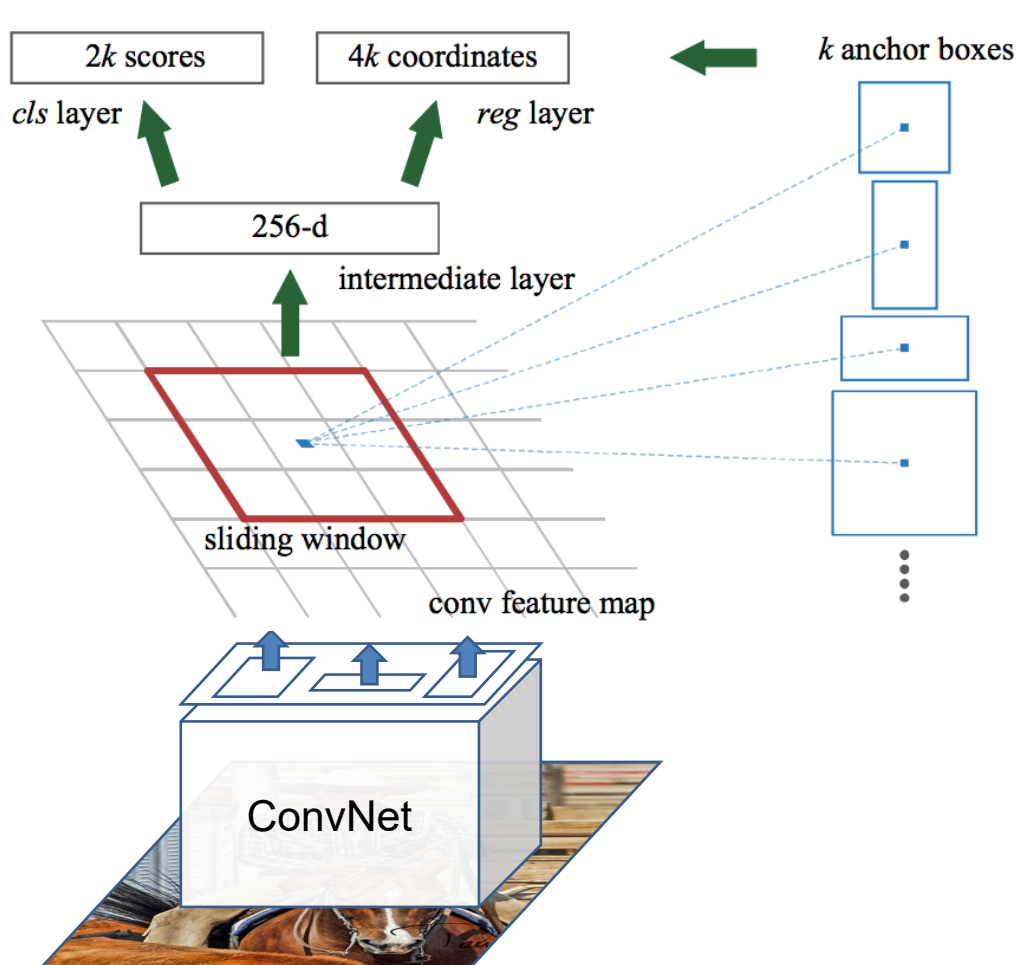
# R-CNN pros and cons

- **Pros**
  - Accurate!
  - Any deep architecture can immediately be “plugged in”
- **Cons**
  - Ad hoc training objectives
    - Fine-tune network with softmax classifier (log loss)
    - Train post-hoc linear SVMs (hinge loss)
    - Train post-hoc bounding-box regressions (least squares)
  - Training is slow (84h), takes a lot of disk space
    - 2000 CNN passes per image
  - Inference (detection) is slow (47s / image with VGG16)

# Faster R-CNN



# Region Proposal Network (RPN)



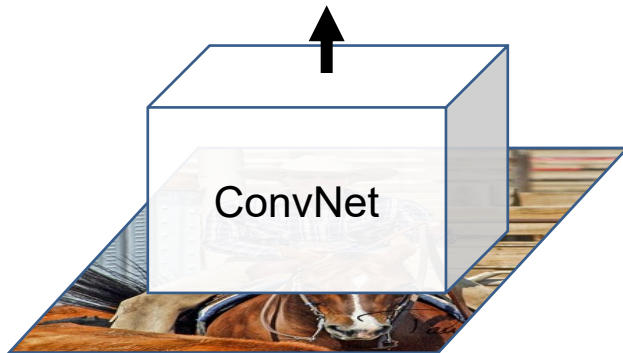
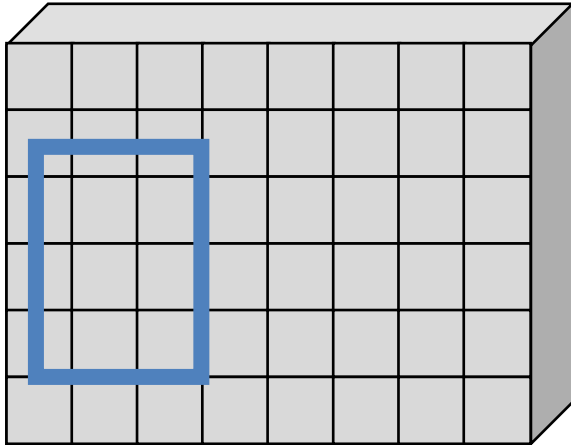
Small network applied to conv5 feature map.

Predicts:

- good box or not (classification),
  - how to modify box (regression)
- for k “anchors” or boxes relative to the position in feature map.

# ROI Pooling/Align

Feature Map  
(e.g.,  $6 \times 8 \times 256$ )

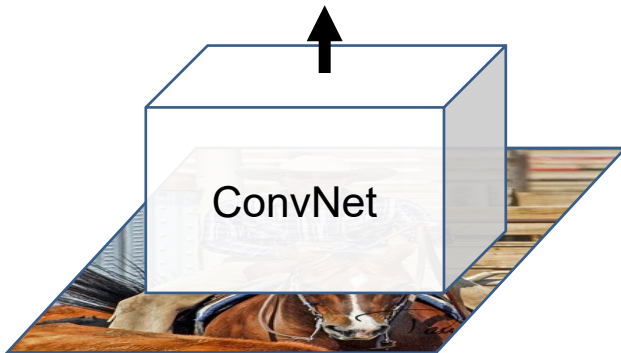
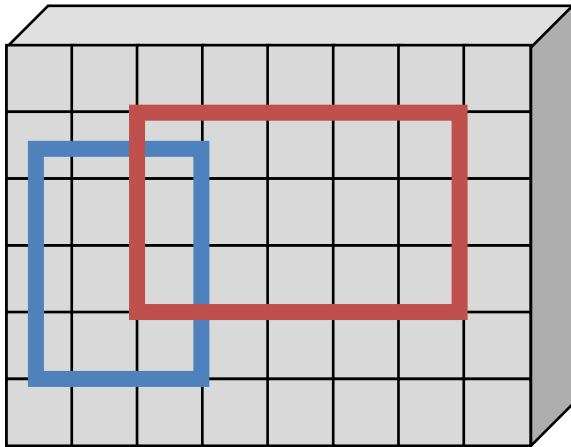


Given box in original image, calculate where the box goes.

- Example:  $H=600, W=800$
- Feature map is  $H'=6, W'=8$
- Box: left  $x=50$ , top  $y=150$ , width=250, height=350
- Feature map box: left  $x=0.5$ ,  $y=1.5$ , width=2.5, height=3.5

# ROI Pooling/Align

Feature Map  
(e.g.,  $6 \times 8 \times 256$ )



Given box in original image, calculate where the box goes.

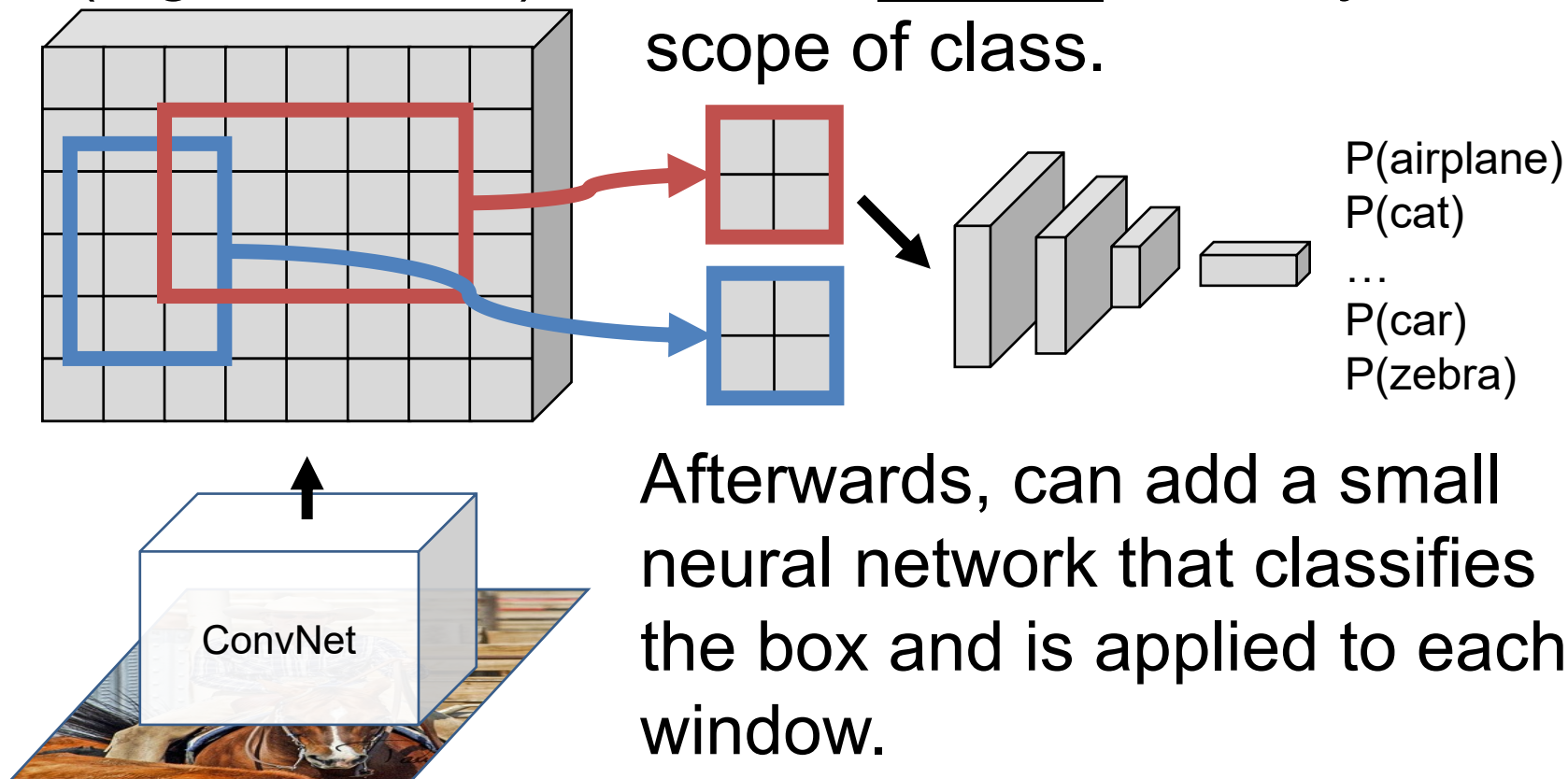
- Example:  $H=600, W=800$
- Feature map is  $H'=6, W'=8$
- Box: left  $x=50$ , top  $y=150$ , width=250, height=350
- Feature map box: left  $x=0.5$ ,  $y=1.5$ , width=2.5, height=3.5
- Other feature map box: left  $x=2$ , top  $y=1$ , width=5, height=3



# ROI Pooling/Align

Feature Map  
(e.g., 6x8x256)

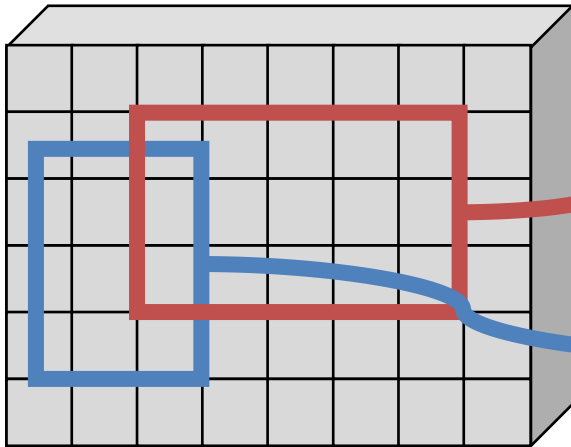
Resize to fixed size (e.g., 7x7)  
Details **critical**, but beyond  
scope of class.



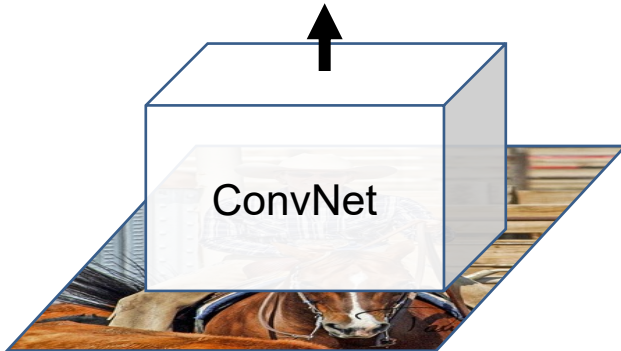
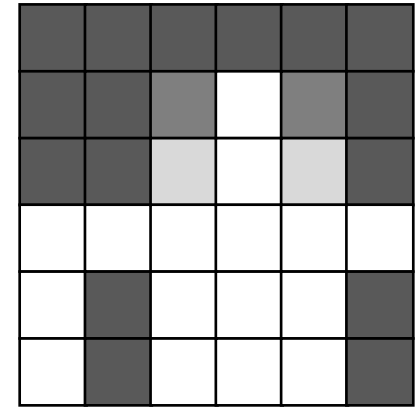
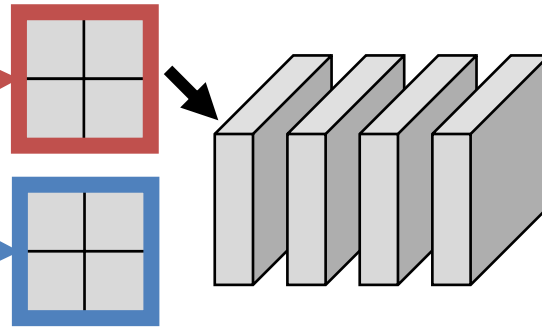
Afterwards, can add a small  
neural network that classifies  
the box and is applied to each  
window.

# Mask RCNN

Feature Map  
(e.g., 6x8x256)



Resize to fixed size (e.g., 7x7)  
Details **critical**, but beyond  
scope of class.



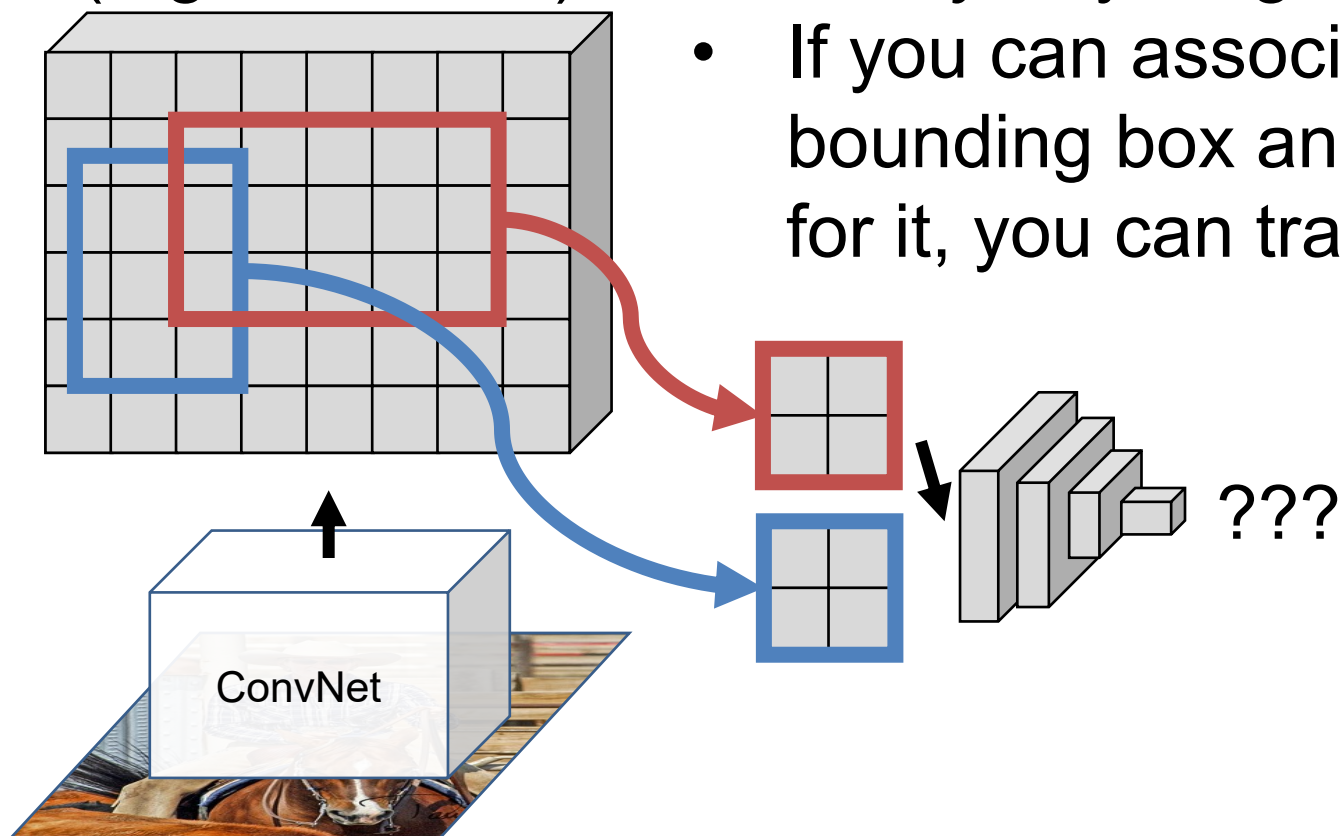
Can also predict a mask!  
Everything is learned together  
Simple and effective; details critical

# MaskRCNN – Results



# Extending Object Detection

Feature Map  
(e.g., 6x8x256)



- Can ask the network to predict *nearly* anything.
- If you can associate it with a bounding box and can get data for it, you can train the model

# Extending Object Detection

Example: RGB image input, detect planar surfaces



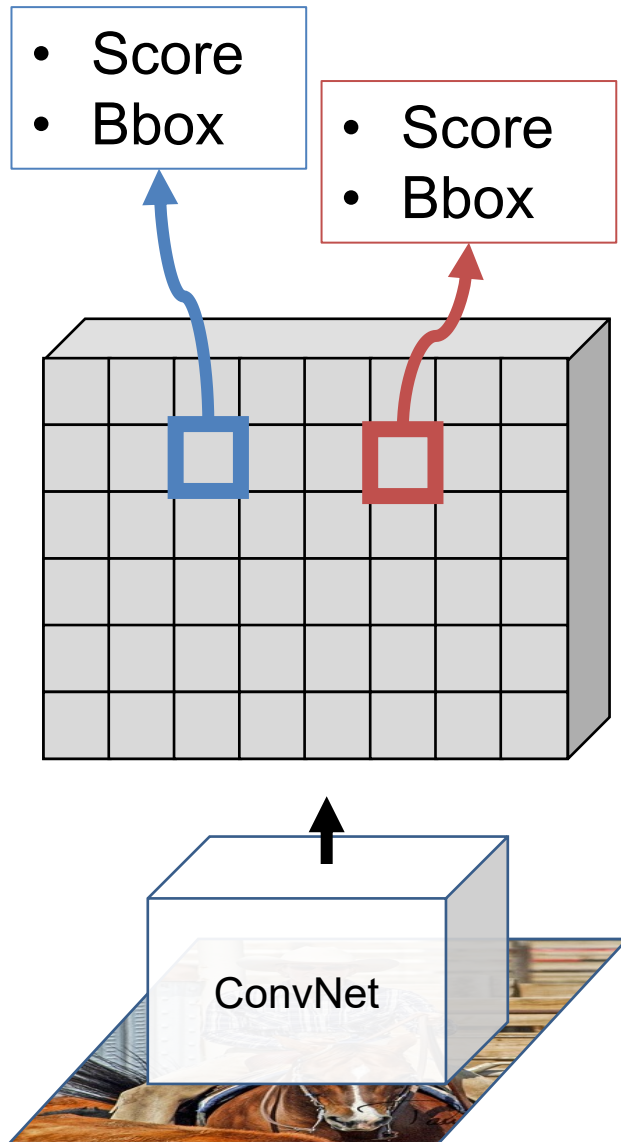
# Extending Object Detection



Core building block is detecting plane in image.



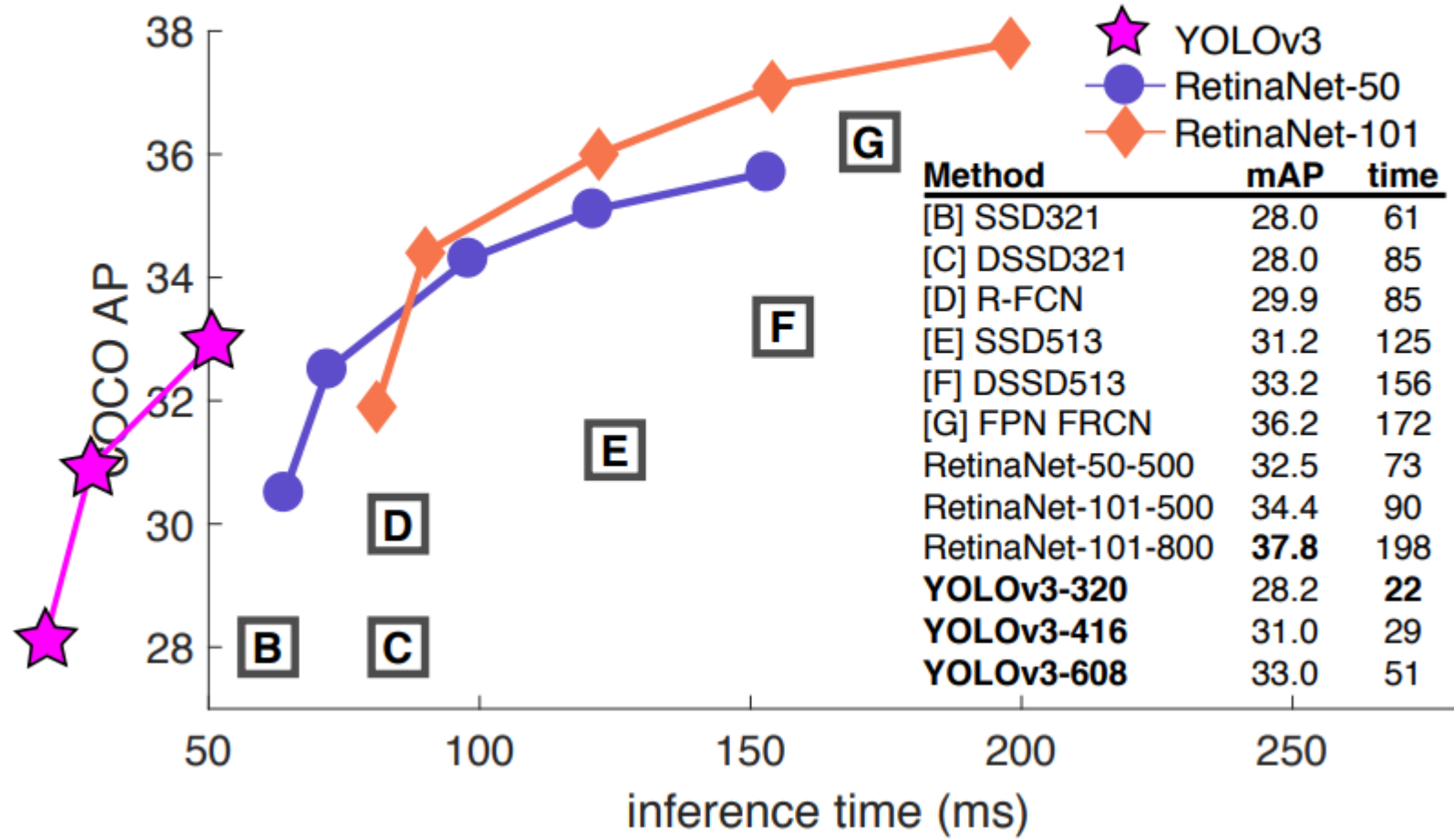
# YOLO



- No proposals
- Predict at each location in 7x7 feature map, score for each class + 2 bboxes
- 7x faster than Faster-RCNN, but worse accuracy, precision
- Immensely popular in robotics
- Loads of similar methods (YOLOv2, YOLOv3)

# YOLO

Score



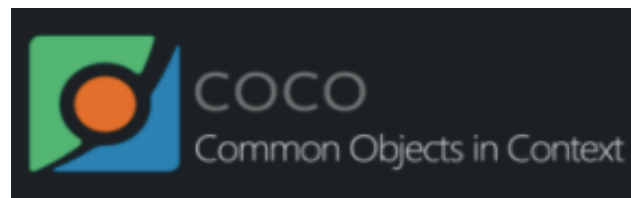
7

1  
S



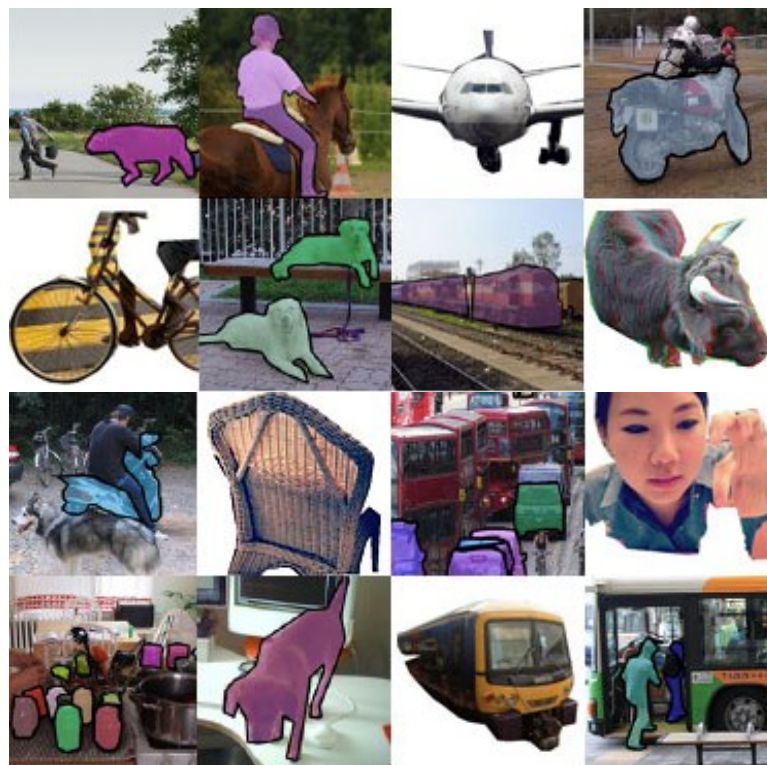
# New detection benchmark: COCO (2014)

- 80 categories instead of PASCAL's 20
- Current best mAP: 52%



COCO is a large-scale object detection, segmentation, and captioning dataset. COCO has several features:

- ✓ Object segmentation
- ✓ Recognition in context
- ✓ Superpixel stuff segmentation
- ✓ 330K images (>200K labeled)
- ✓ 1.5 million object instances
- ✓ 80 object categories
- ✓ 91 stuff categories
- ✓ 5 captions per image
- ✓ 250,000 people with keypoints



<http://cocodataset.org/#home>

# A Few Caveats

- Flickr images come from a really weird process
- Step 1: user takes a picture
- Step 2: user decides to upload it
- Step 3: user decides to write something like “refrigerator” somewhere in the description
- Step 4: a vision person stumbles on it while searching Flickr for refrigerators for a dataset



**Who takes  
photos of open  
refrigerators  
?????**



# Guess the category!

These were detected with  $>90\%$  confidence, corresponding to 99% precision on original dataset



(1) Person

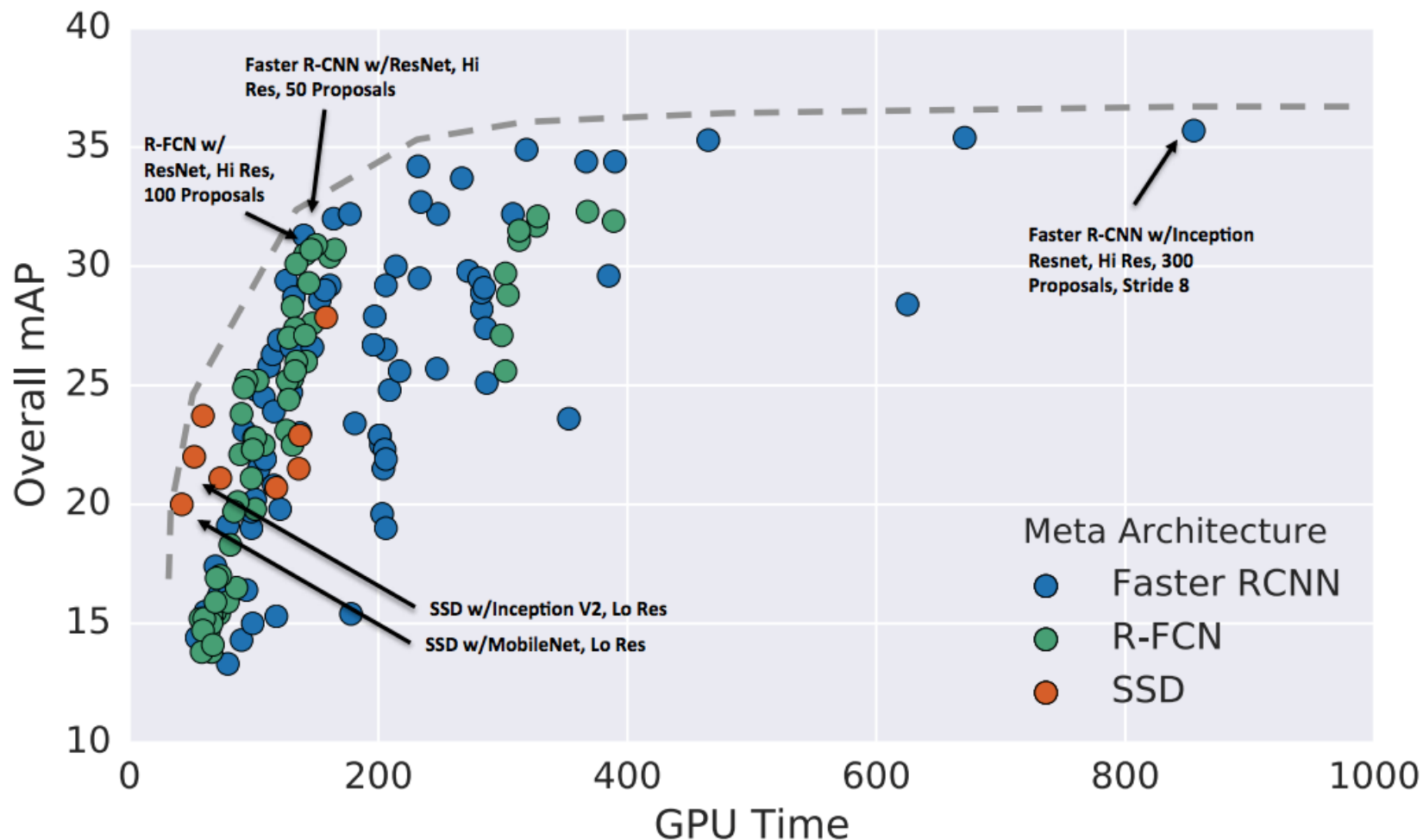
(2) Giraffe

(3) Bicycle

# Kitchens from Googling



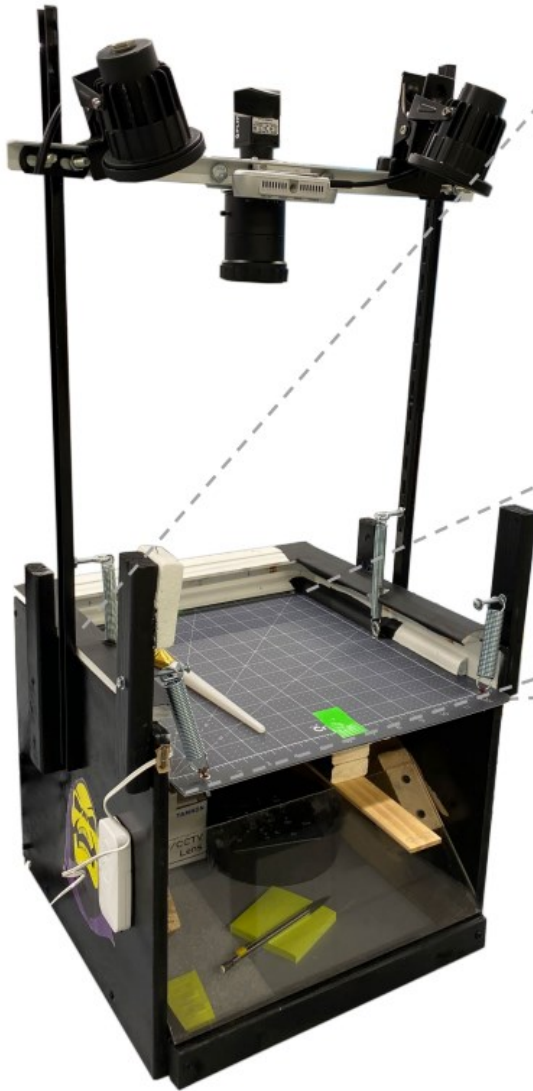
# New detection benchmark: COCO (2014)



J. Huang et al., [Speed/accuracy trade-offs for modern convolutional object detectors](#), CVPR 2017

# What's It Good For?

## Imaging Setup



## Instance Segmentation



**Skull**



**36.08 mm**

**Femur**

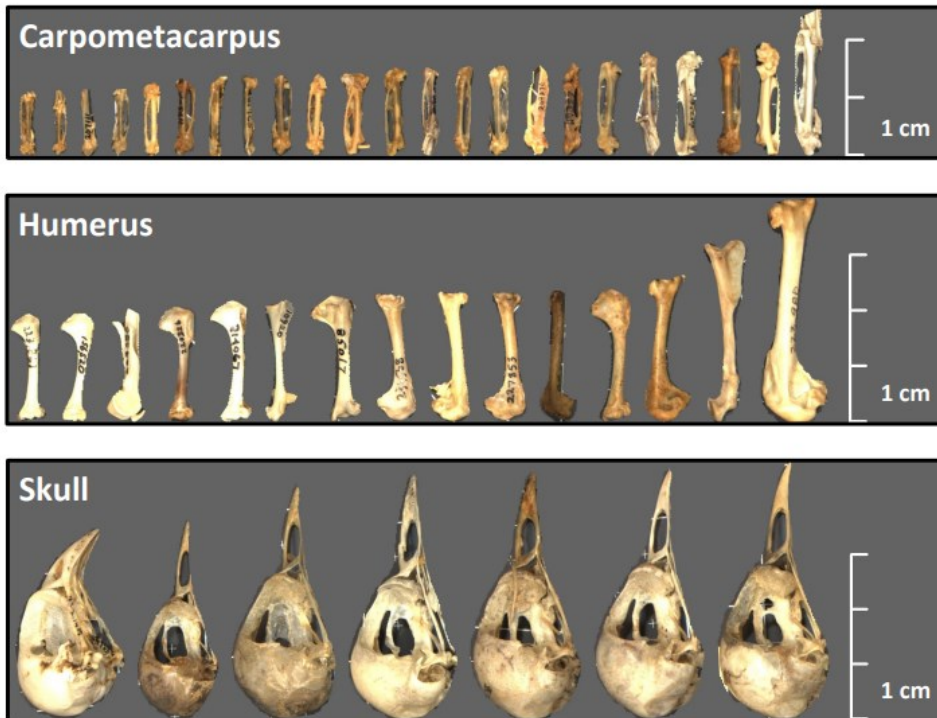


**18.97 mm**

Z. Zhou, G. Hassena, B.C. Weeks, D.F. Fouhey. *Quantifying Bird Skeletons*. CVPR Workshops, 2021.

B.C. Weeks, Z. Zhou, B.K. O'Brien, R. Darling, M. Dean, T. Dias, G. Hassena, M. Zhang, D.F. Fouhey. *A deep neural network for high throughput measurement of functional traits on museum skeletal specimens*. *Methods in Ecology and Evolution*, 2022.

# What's It Good For?



- UM has 25K bird skeletons but measuring bird bones by hand is hard and tedious.
- New solution: dump the bones, take a picture, use a deep network.
- Now enabling testing hypotheses about birds at huge scale

Z. Zhou, G. Hassena, B.C. Weeks, D.F. Fouhey. *Quantifying Bird Skeletons*. CVPR Workshops, 2021.

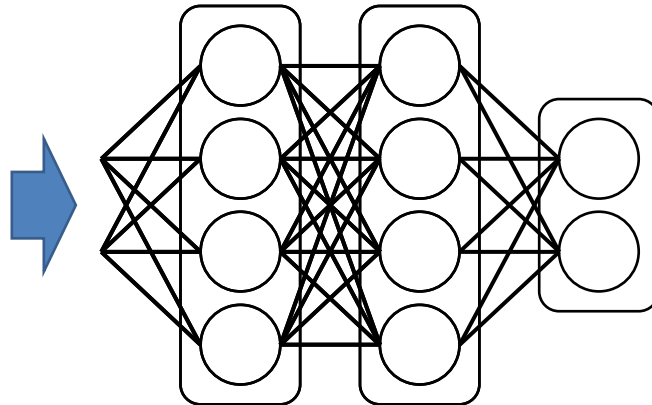
B.C. Weeks, Z. Zhou, B.K. O'Brien, R. Darling, M. Dean, T. Dias, G. Hassena, M. Zhang, D.F. Fouhey. *A deep neural network for high throughput measurement of functional traits on museum skeletal specimens*. *Methods in Ecology and Evolution*, 2022.



# And Now For Something Completely Different

- I planned to cover this, but I'll cover it in the next class.

# ImageNet + Deep Learning

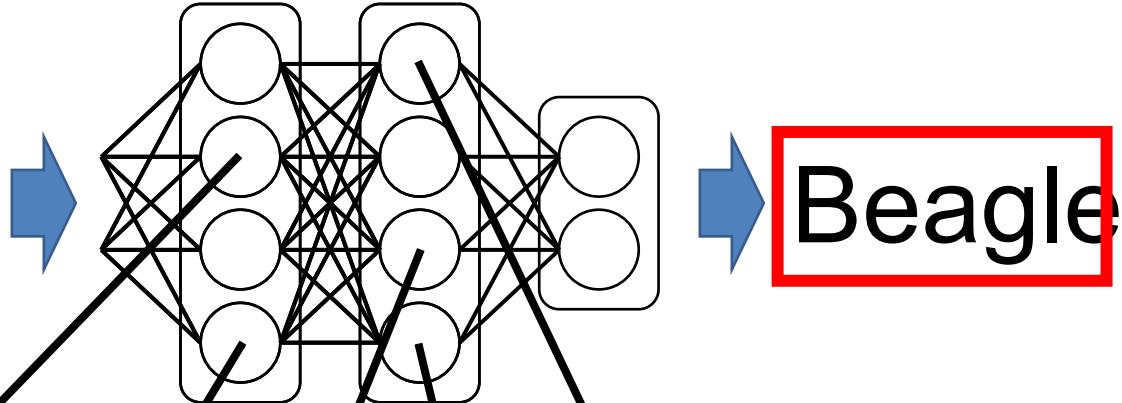


**Beagle**



- Image Retrieval
- Detection
- Segmentation
- Depth Estimation
- ...

# ImageNet + Deep Learning



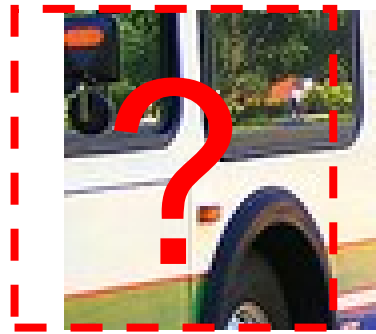
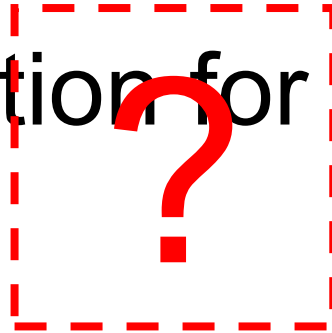
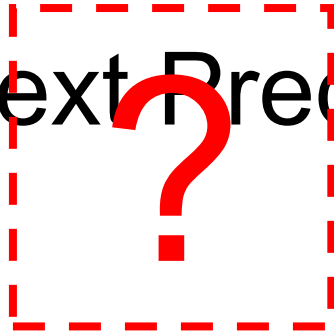
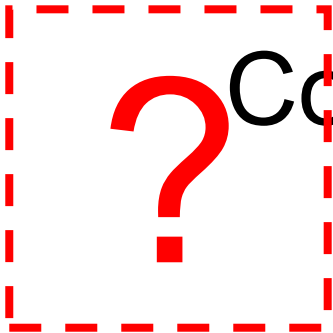
Materials?                      Pose?  
*Do we even need semantic*  
*Do we need this task?*  
Parts?  
*labels?*  
Geometry?                      Boundaries?

# To Make It Super Clear

- $w = \text{weights\_from\_somewhere\_else}$
- for batch in batches:
  - inputs, labels = batch
  - calculate gradient of loss function with respect to  $w$  applied to samples in inputs
  - $w += \text{gradient}$

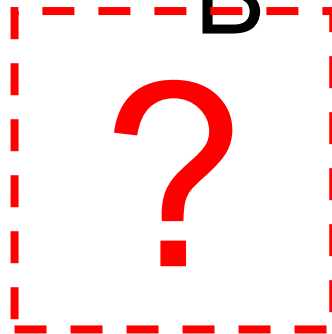
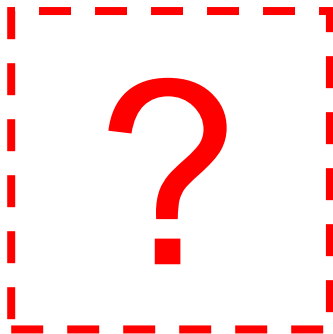


# Context Prediction for Images

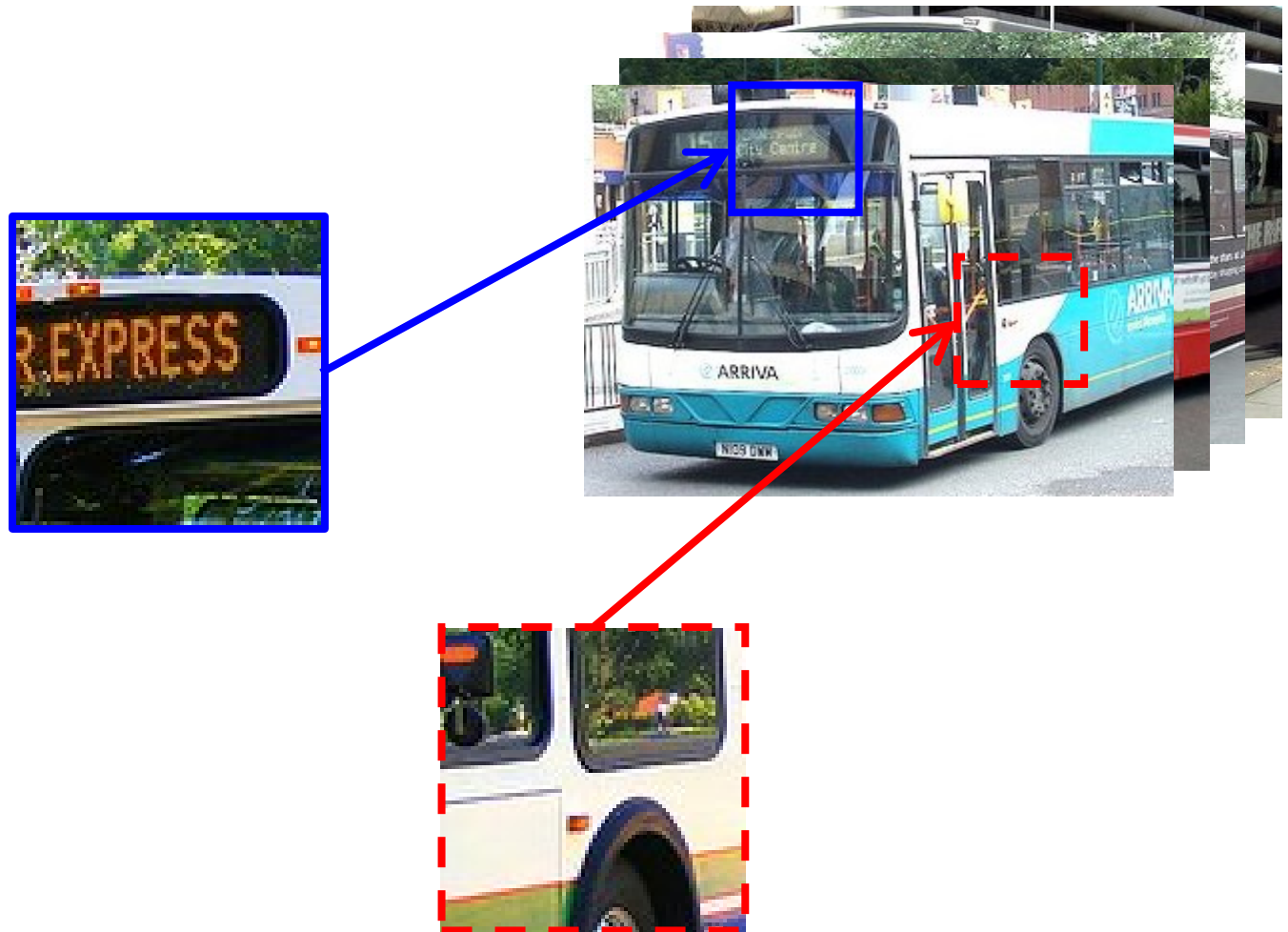


A

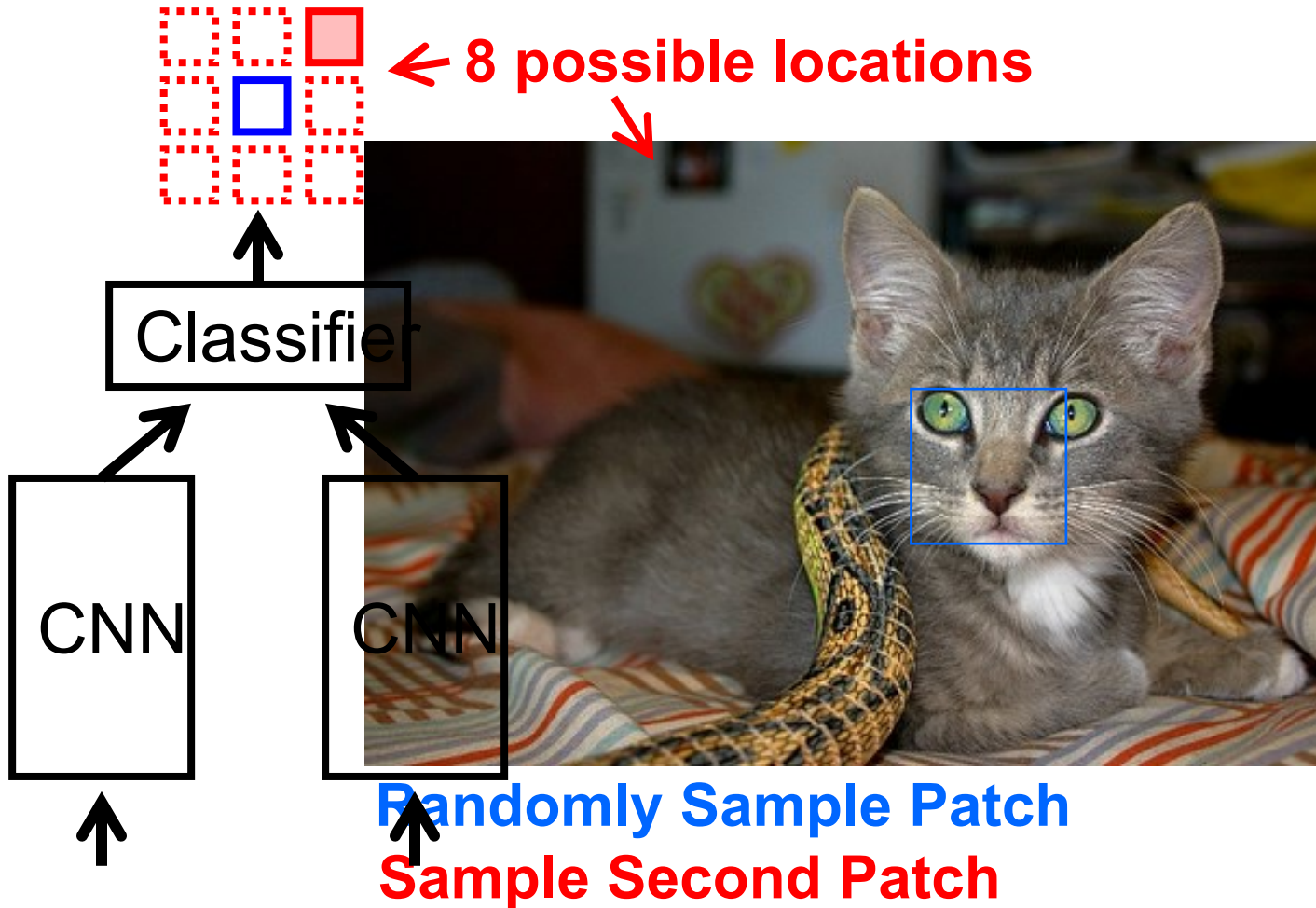
B



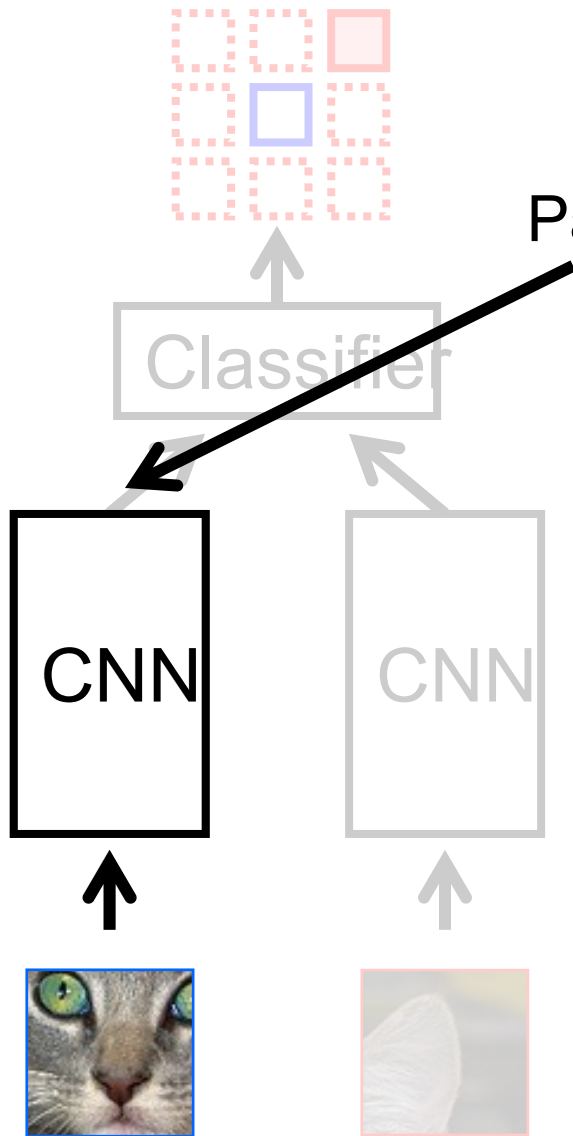
# Semantics from a non-semantic task



# Relative Position Task







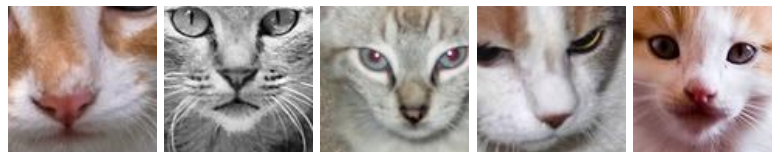
Patch Embedding

Input



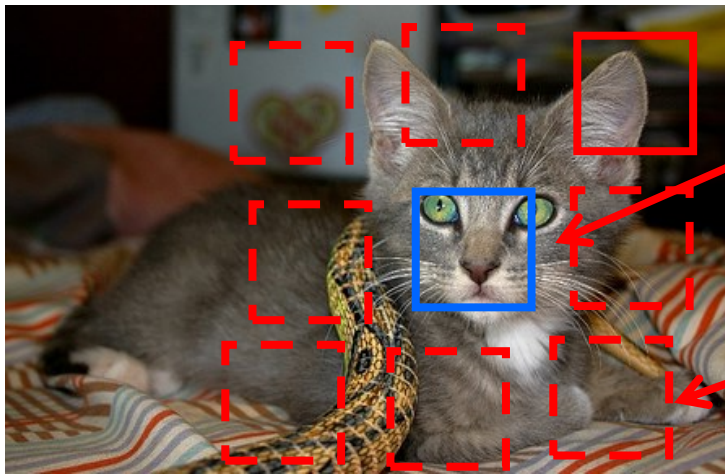
:

Nearest Neighbors



Note: connects ***across*** instances

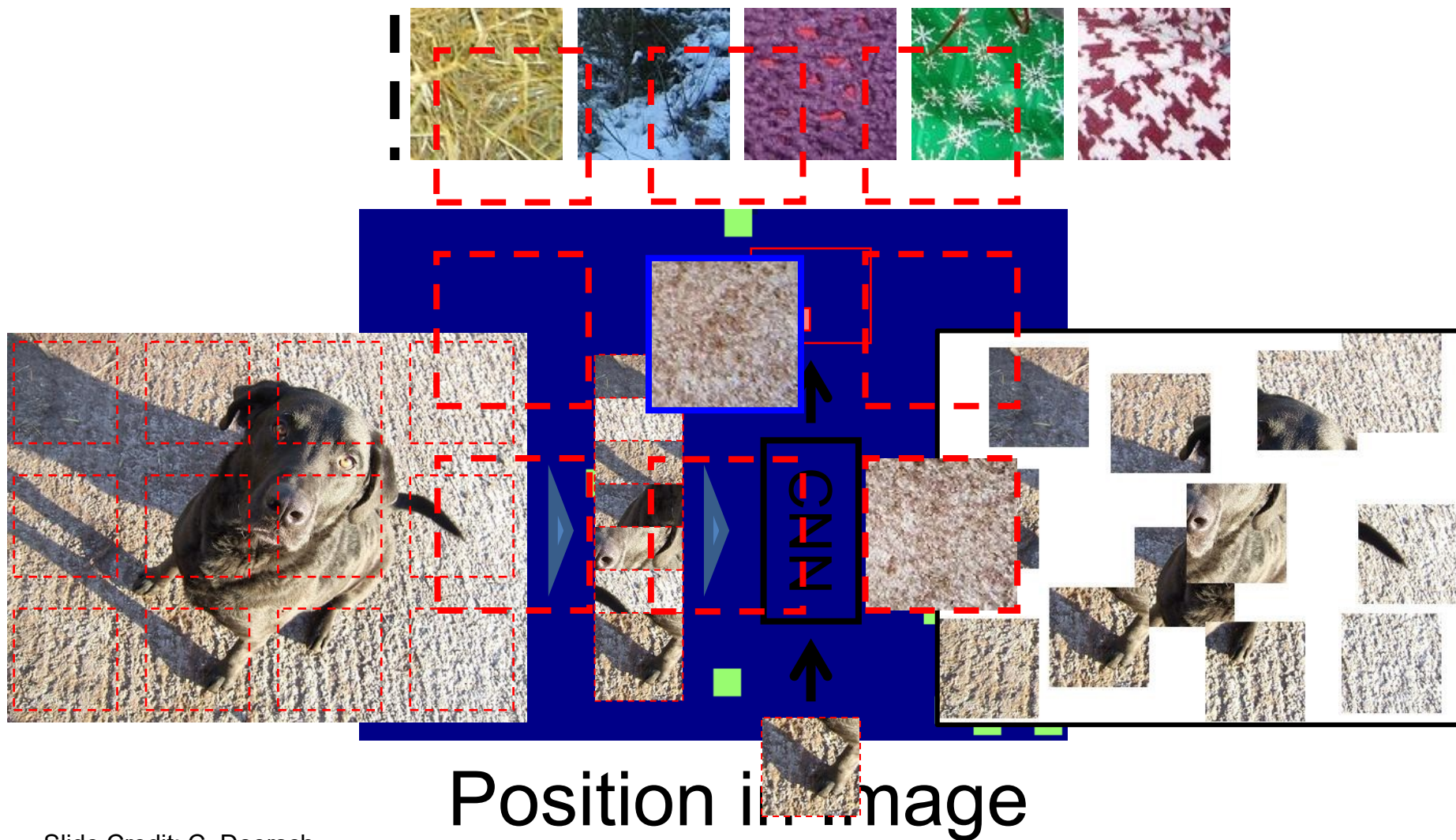
# Avoiding Trivial Shortcuts



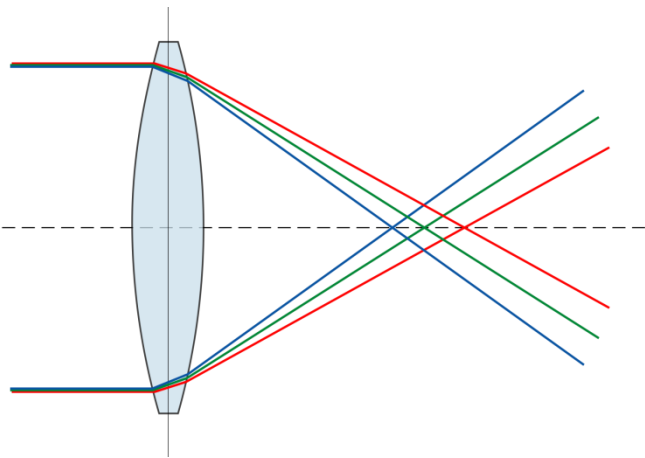
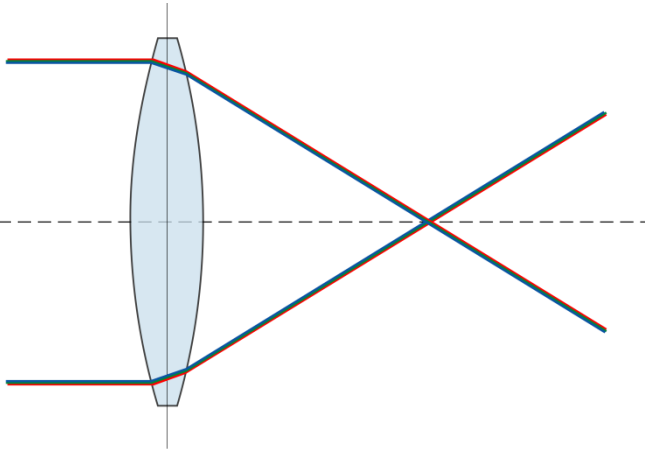
Include a gap

Jitter the patch locations

# A Not-So “Trivial” Shortcut



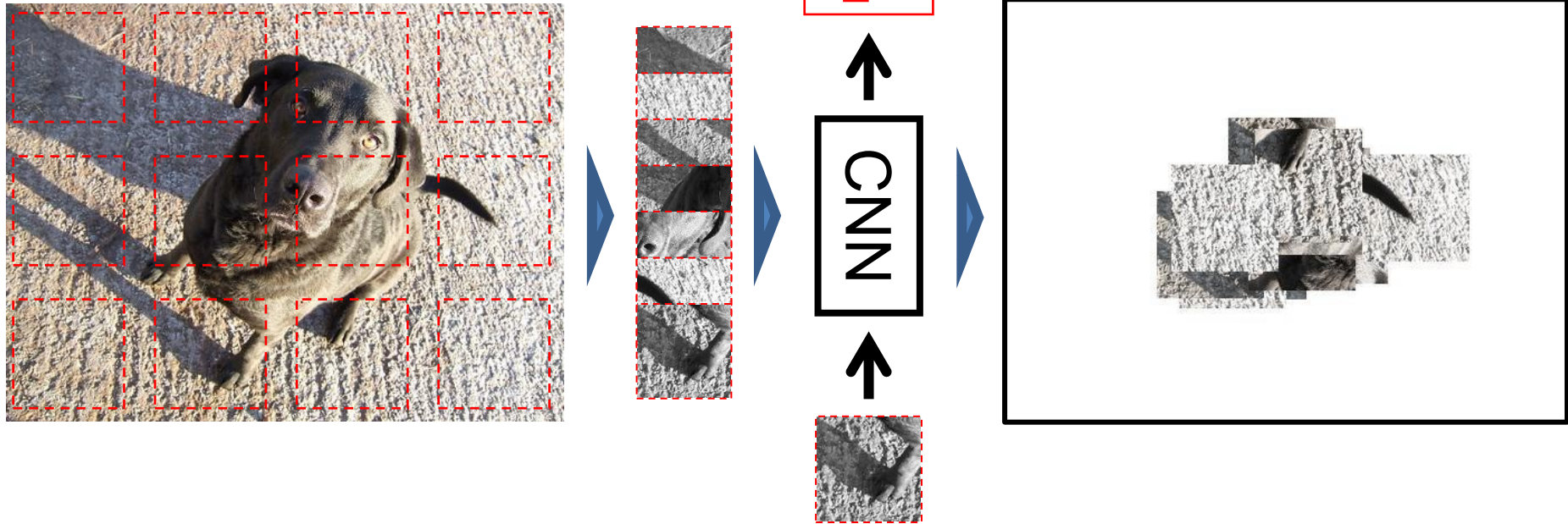
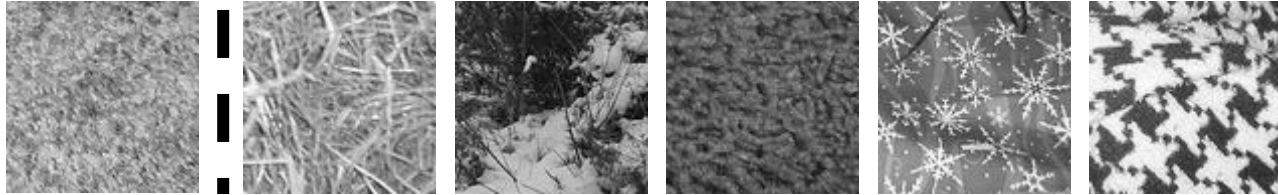
# Chromatic Aberration



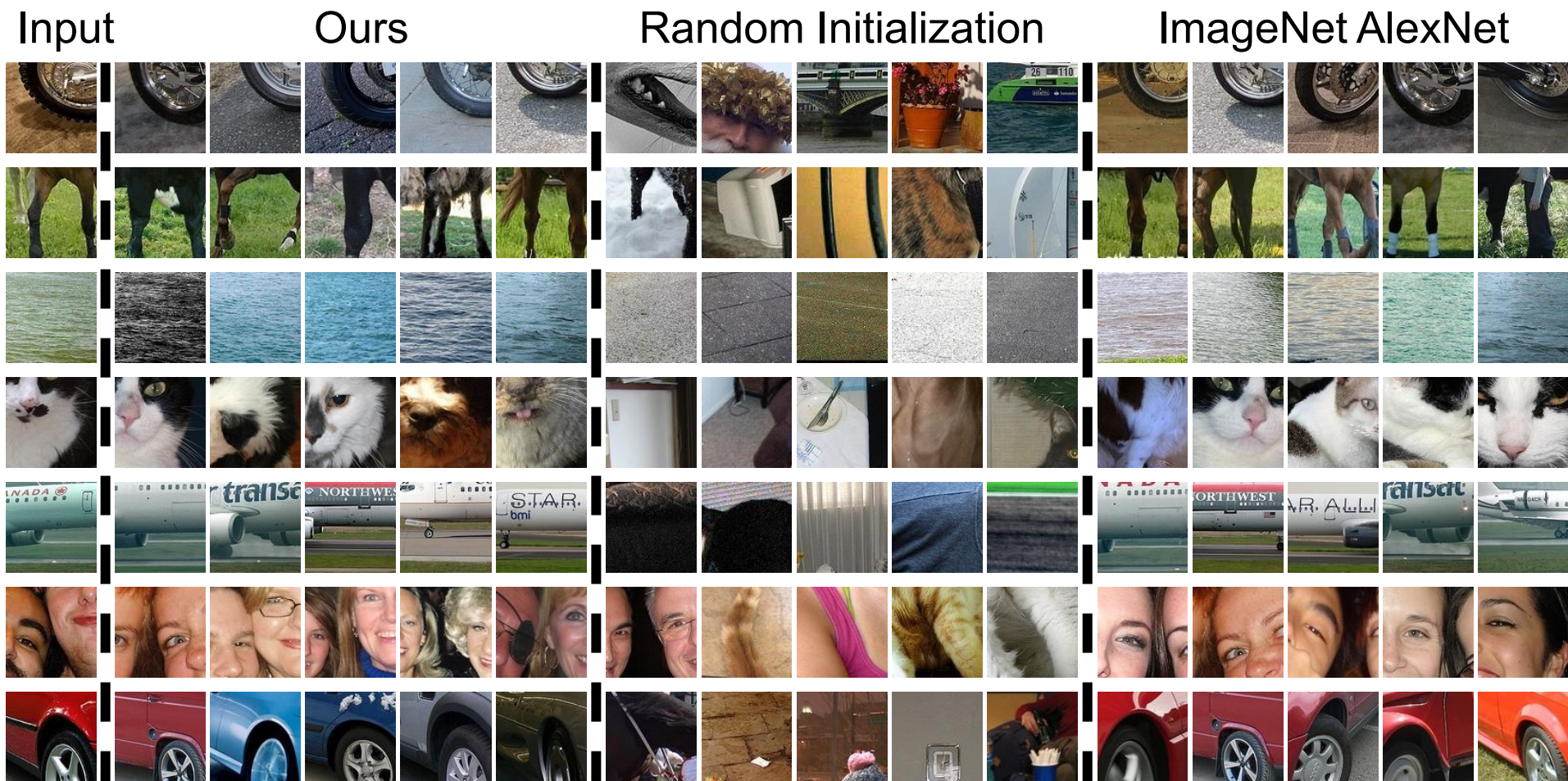
Slide Credit: C. Doersch



# Chromatic Aberration



# What is learned?



# Pre-Training for R-CNN

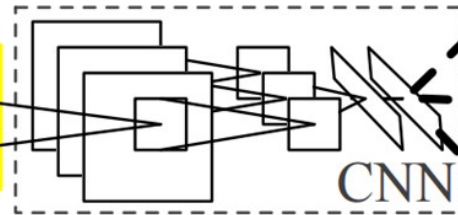


1. Input image



2. Extract region proposals (~2k)

warped region



CNN

3. Compute CNN features

aeroplane? no.

⋮

person? yes.

⋮

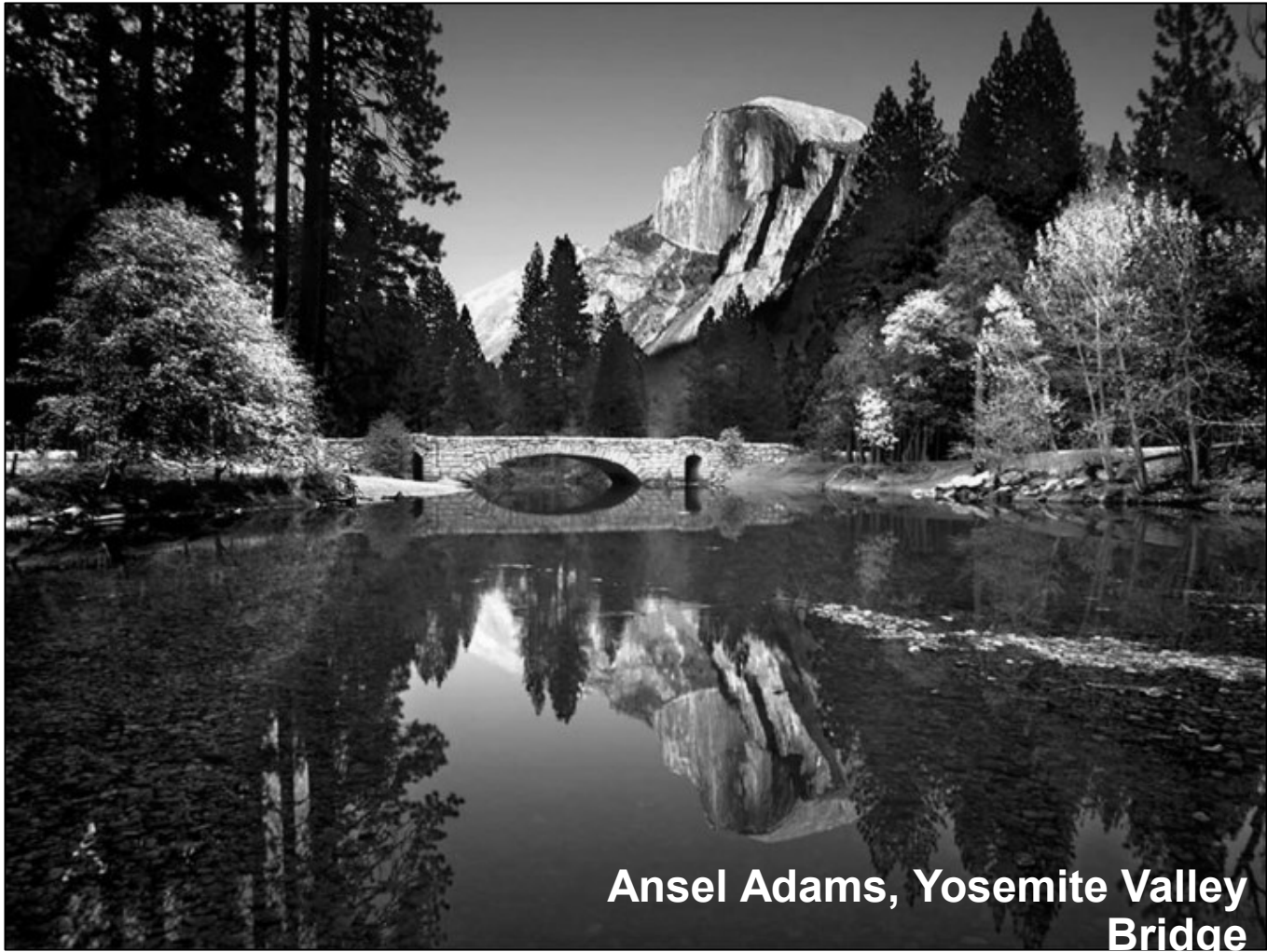
tvmonitor? no.

4. Classify regions

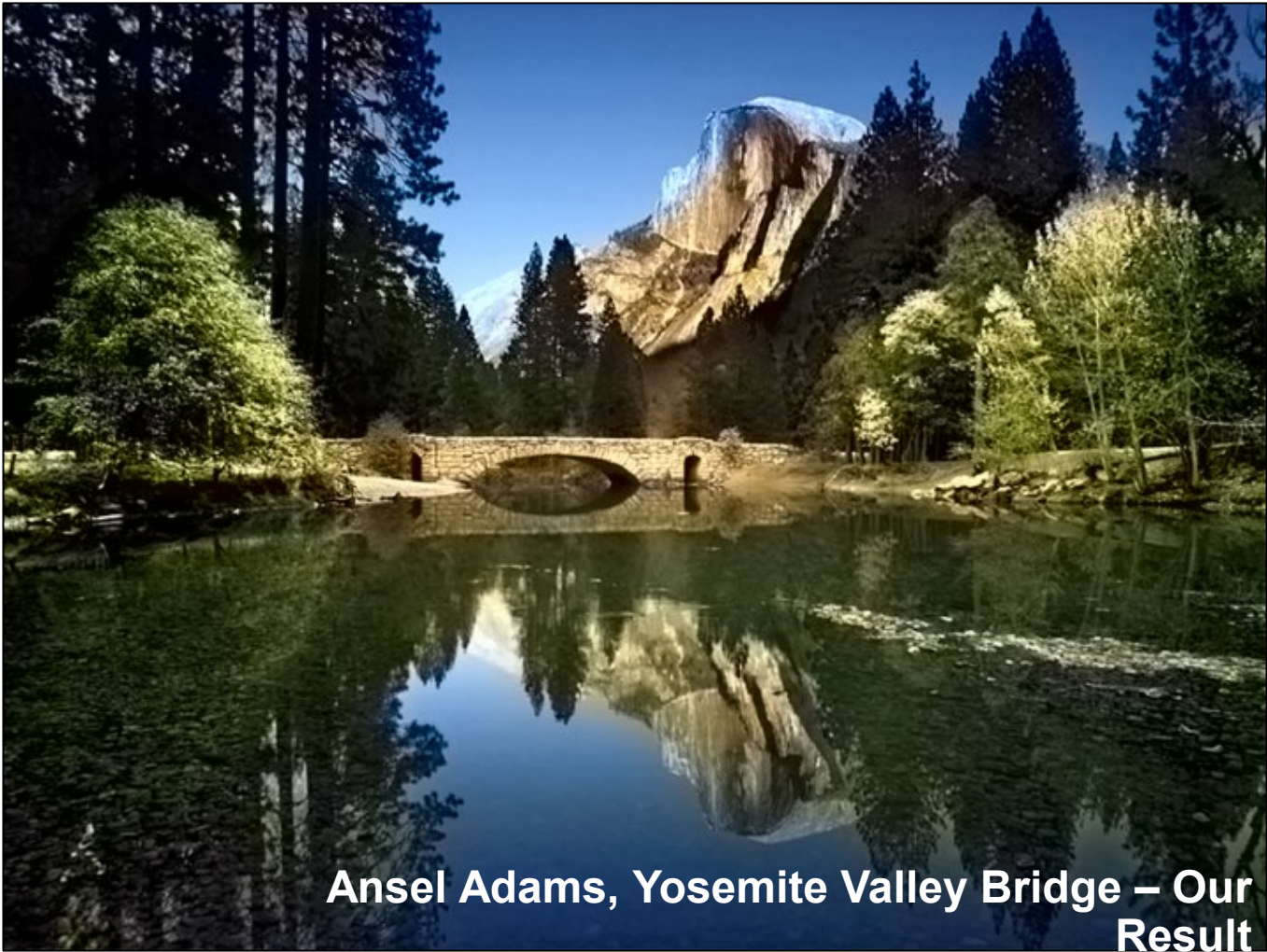
Pre-train on relative-position task, w/o labels

# Other Sources Of Signal

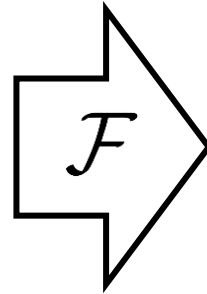




**Ansel Adams, Yosemite Valley  
Bridge**



**Ansel Adams, Yosemite Valley Bridge – Our Result**

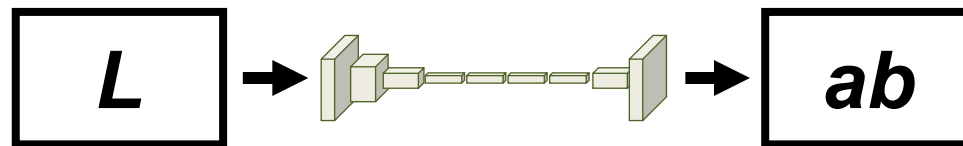


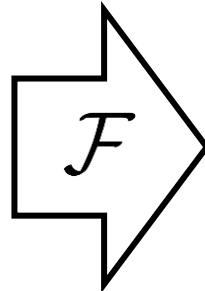
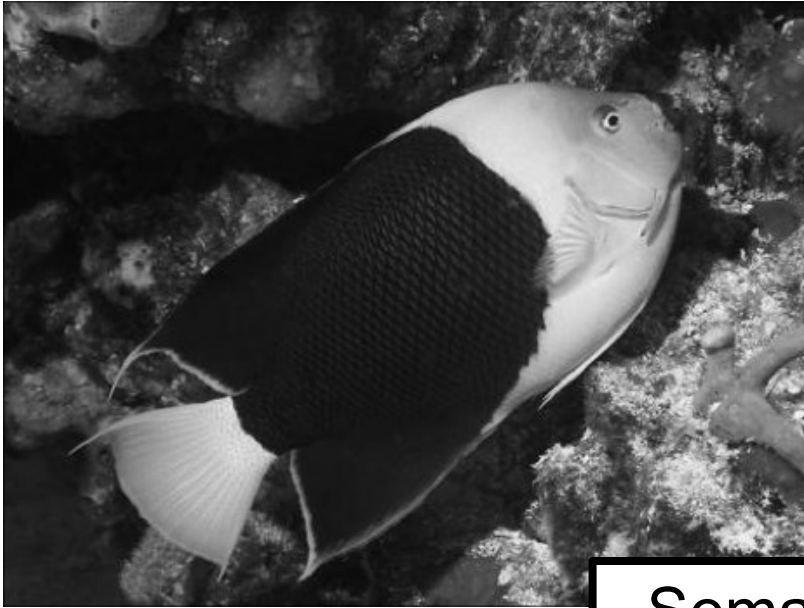
Grayscale image:  $L$  channel

$$\mathbf{X} \in \mathbb{R}^{H \times W \times 1}$$

Color information:  $ab$  channels

$$\hat{\mathbf{Y}} \in \mathbb{R}^{H \times W \times 2}$$

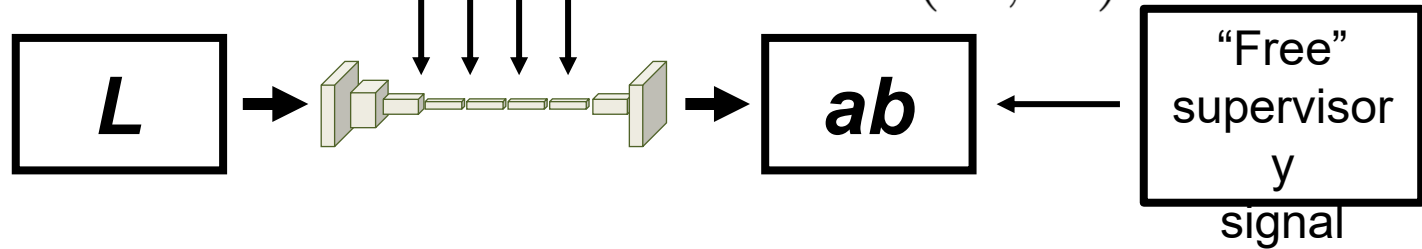




Semantics? Higher-level abstraction?

Grayscale image:  $L$   
 $\mathbf{X} \in \mathbb{R}^{H \times W \times 1}$

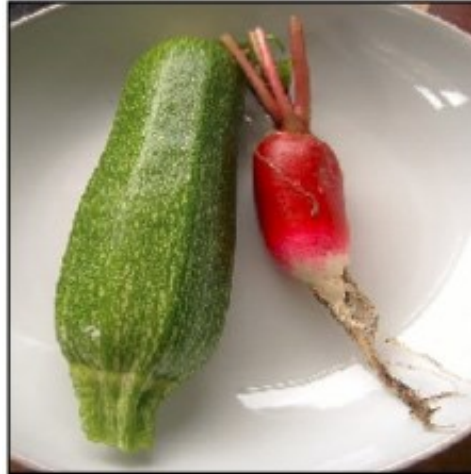
Concatenate  $(L, ab)$   
 $(\mathbf{X}, \hat{\mathbf{Y}})$



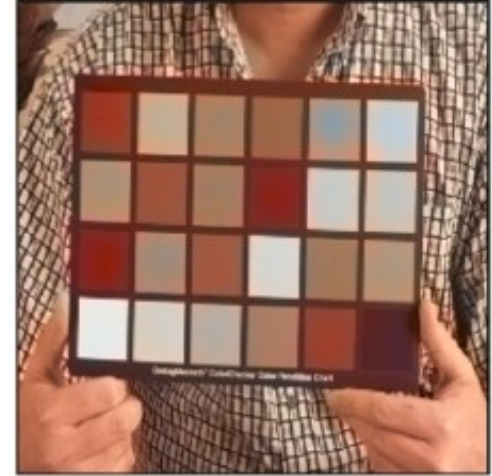
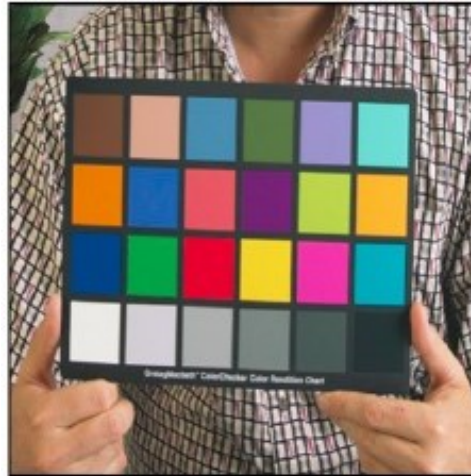
# Input



# Ground Truth



# Output





# Visually Indicated Sounds

Andrew Owens

Phillip Isola

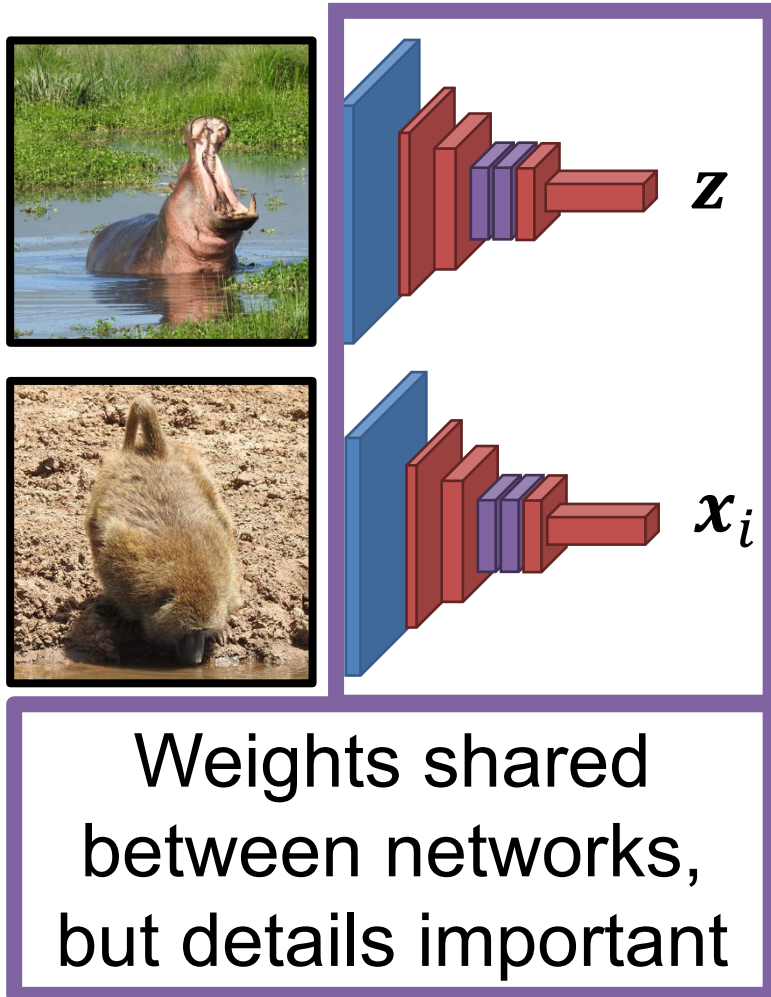
Josh McDermott

Antonio Torralba

Edward Adelson

William Freeman

# Contrastive Learning



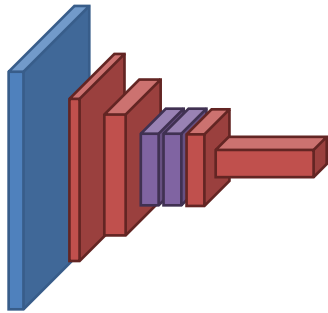
Given sample, construct feature  $\mathbf{z}$ ; take a bunch of other images, minimize:

$$\frac{\exp(\mathbf{z}^T \mathbf{z})}{\exp(\mathbf{z}^T \mathbf{z} + \sum_i \mathbf{z}^T \mathbf{x}_i)}$$

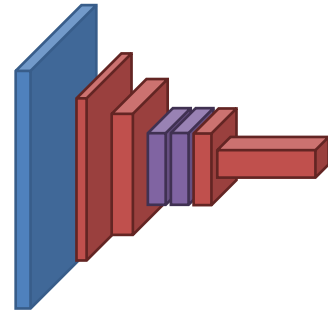
Basically, a scoring function with  $\mathbf{w} = \mathbf{z}$ :

$$\frac{\exp(\mathbf{w}^T \mathbf{z})}{\exp(\mathbf{w}^T \mathbf{z} + \sum_i \mathbf{w}^T \mathbf{x}_i)}$$

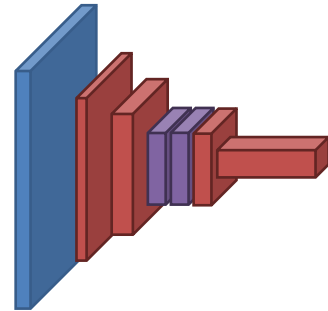
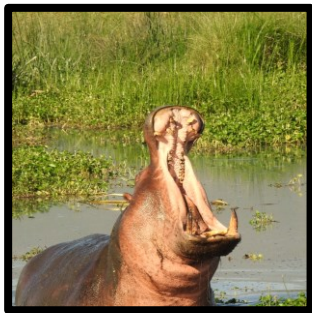
# Contrastive Learning



$z$



$x_i$



$z'$

Best performing methods measure distance to augmented sample  $z'$ :

$$\frac{\exp(\mathbf{z}^T \mathbf{z}')}{\exp(\mathbf{z}^T \mathbf{z}' + \sum_i \mathbf{z}^T \mathbf{x}_i)}$$

Goal: score augmented sample higher than everything else.

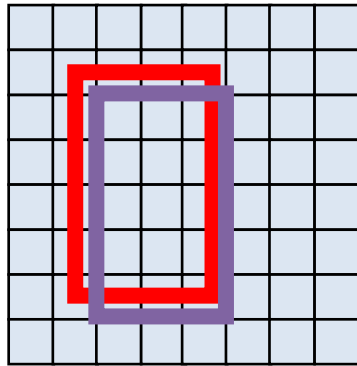


# Next Time

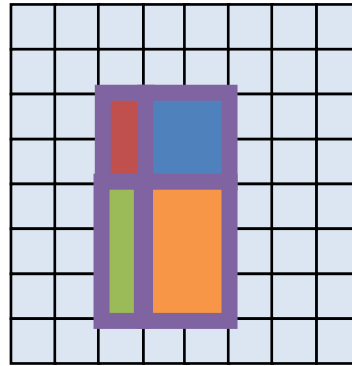
- Synthesizing Images

# Extra Stuff

# Fast R-CNN – ROI-Pool



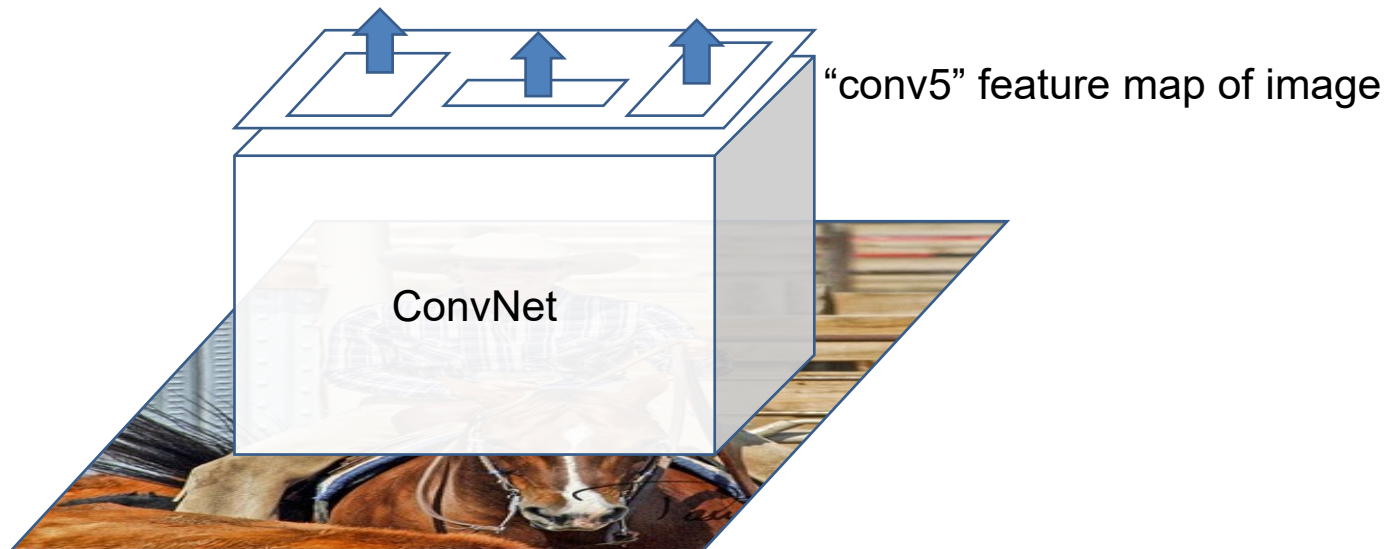
Line up



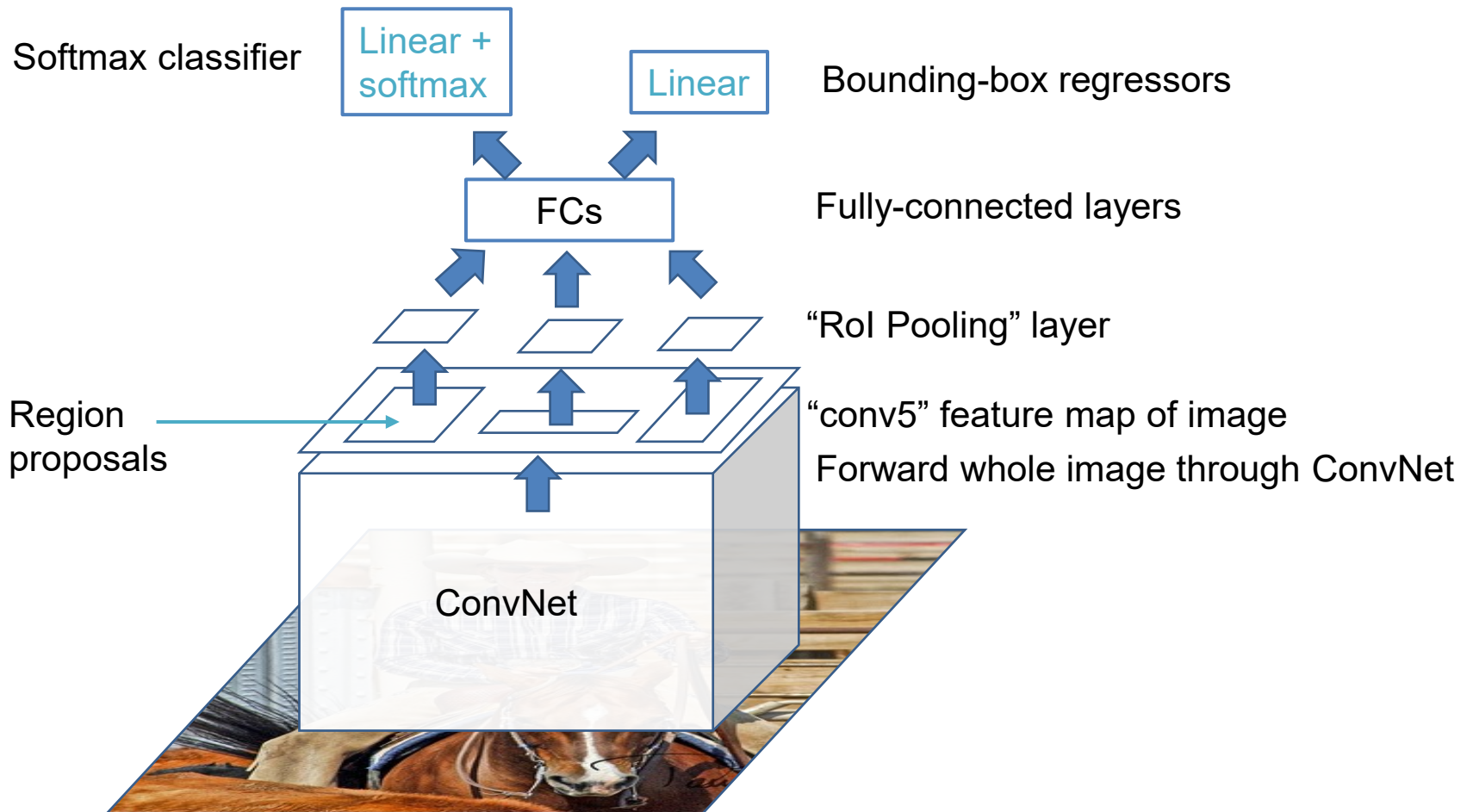
Divide



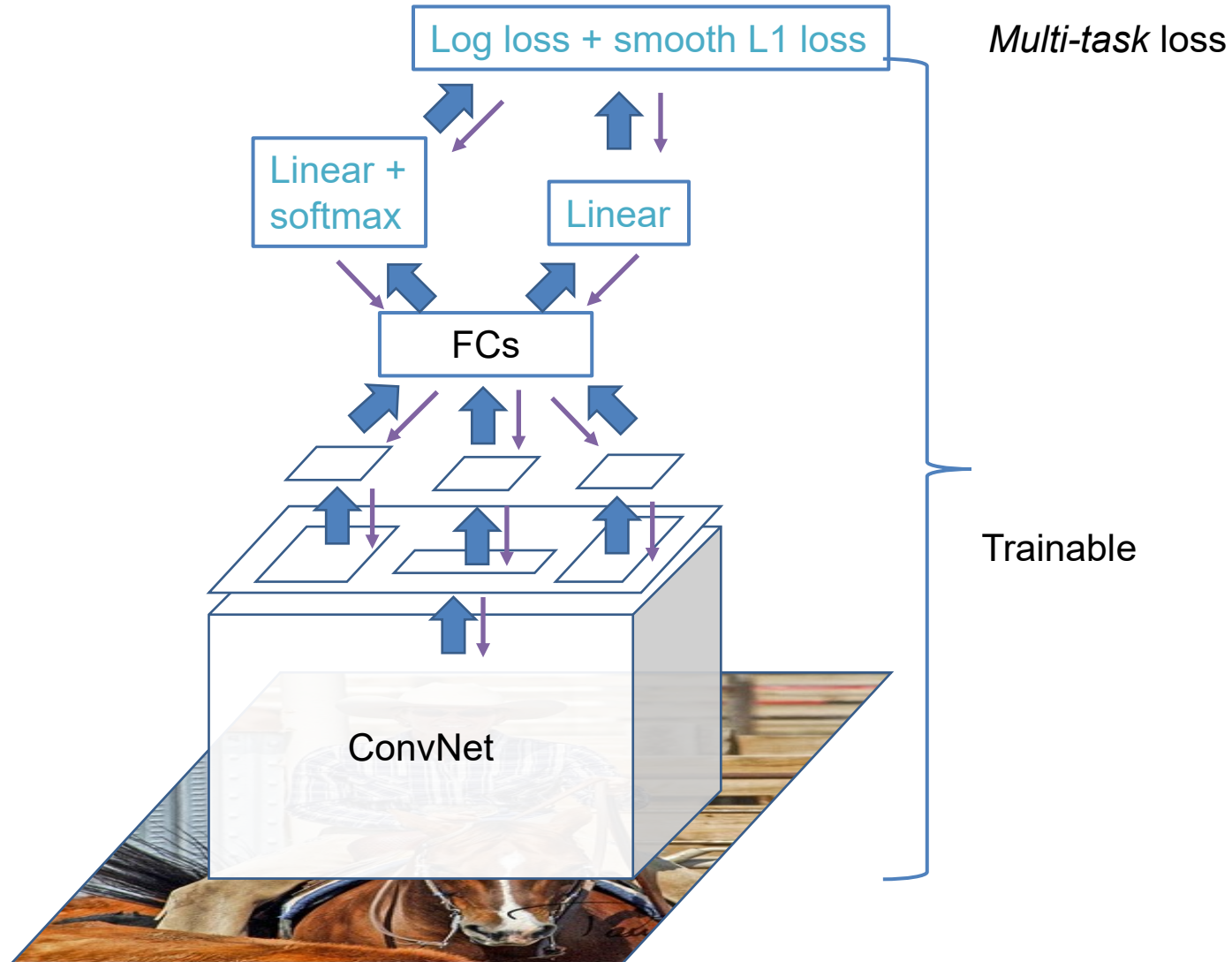
Pool



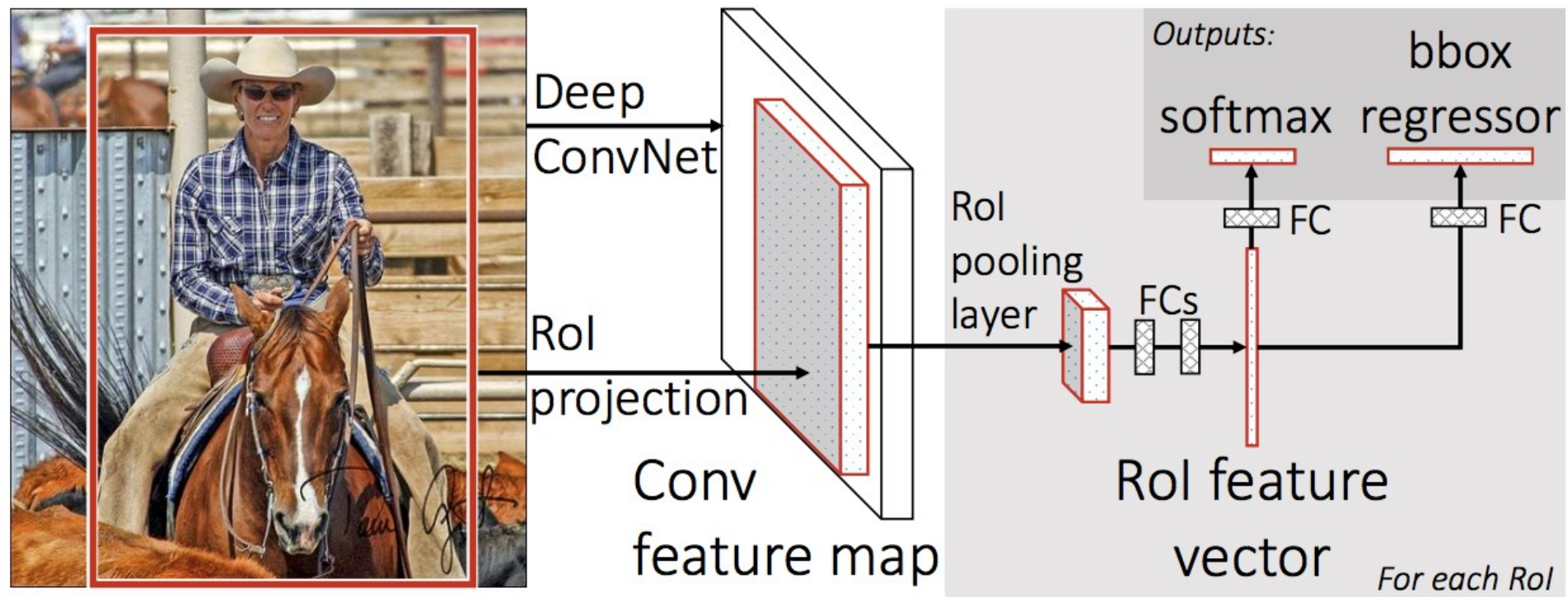
# Fast R-CNN



# Fast R-CNN training



# Fast R-CNN: Another view



# Fast R-CNN results

	Fast R-CNN	R-CNN
Train time (h)	<b>9.5</b>	84
- Speedup	<b>8.8x</b>	1x
Test time / image	<b>0.32s</b>	47.0s
Test speedup	<b>146x</b>	1x
mAP	<b>66.9%</b>	66.0%

Timings exclude object proposal time, which is equal for all methods.  
All methods use VGG16 from Simonyan and Zisserman.