

Object Detection (Plus some bonuses)

EECS 442 – David Fouhey

Fall 2019, University of Michigan

http://web.eecs.umich.edu/~fouhey/teaching/EECS442_F19/

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

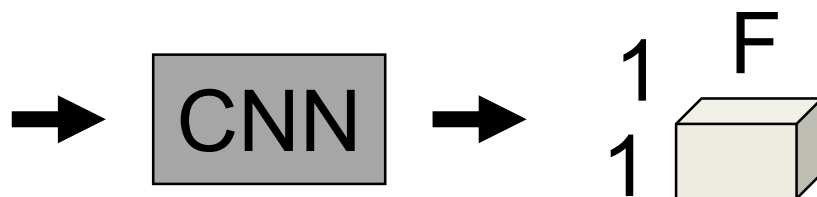
Input



Target



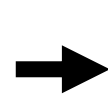
The Wrong Way To Do It



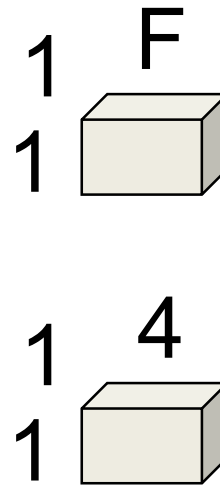
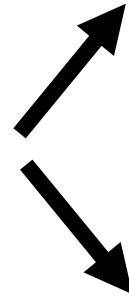
Starting point:

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

The Wrong Way To Do It



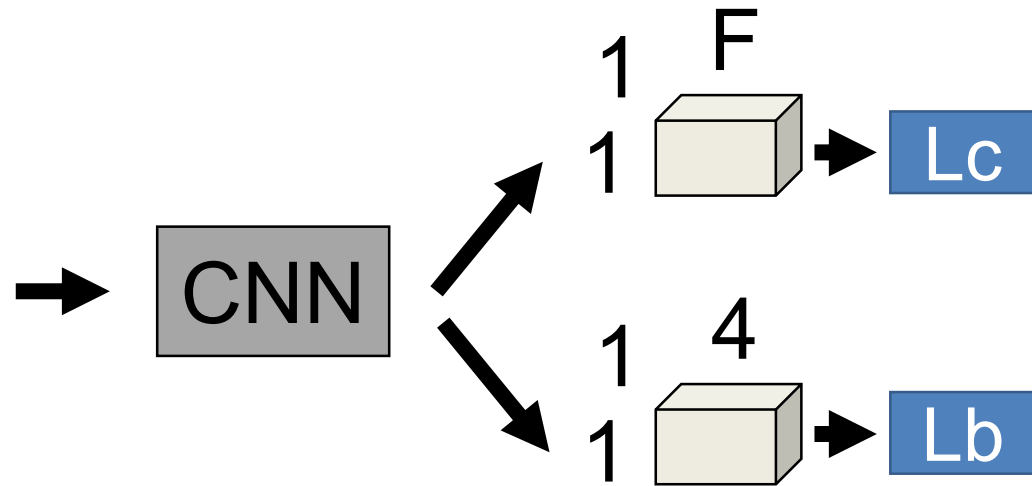
CNN



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



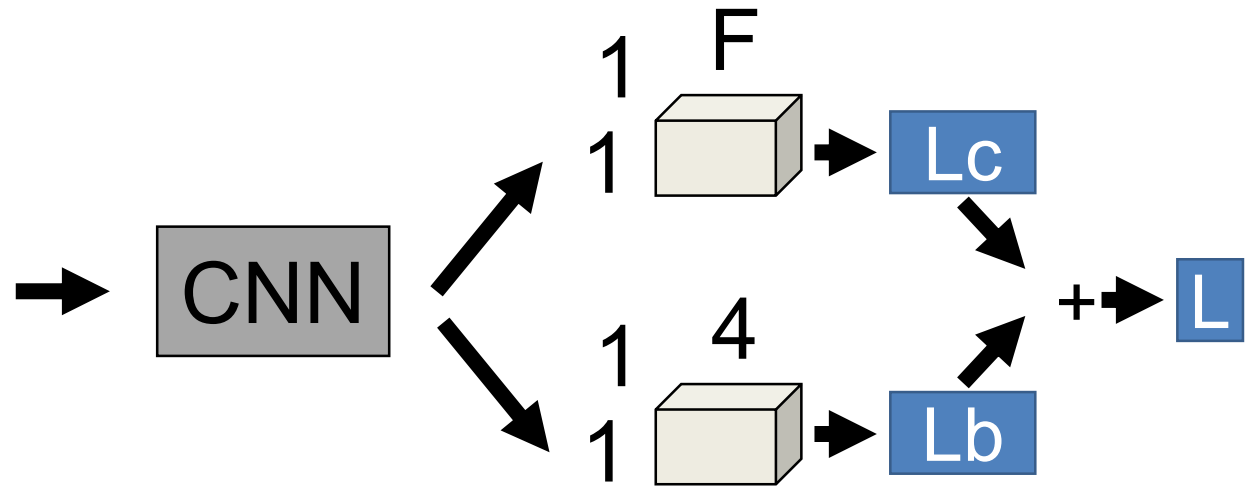
Put a loss on it:

Penalize mistakes on the classes with

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

$L_b = \text{L2 loss}$

The Wrong Way To Do It

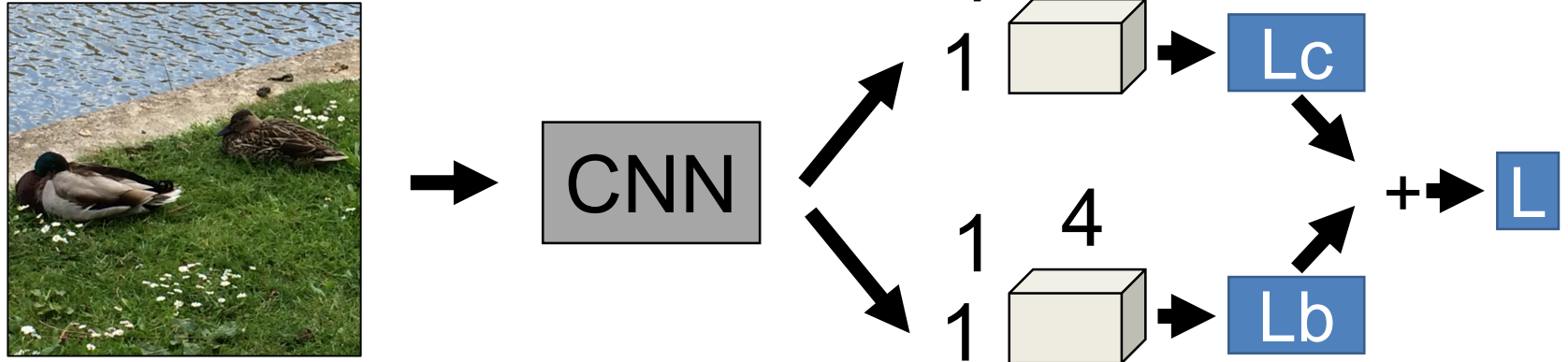


Add losses, backpropagate

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

Why do we need the λ ?

The Wrong Way To Do It

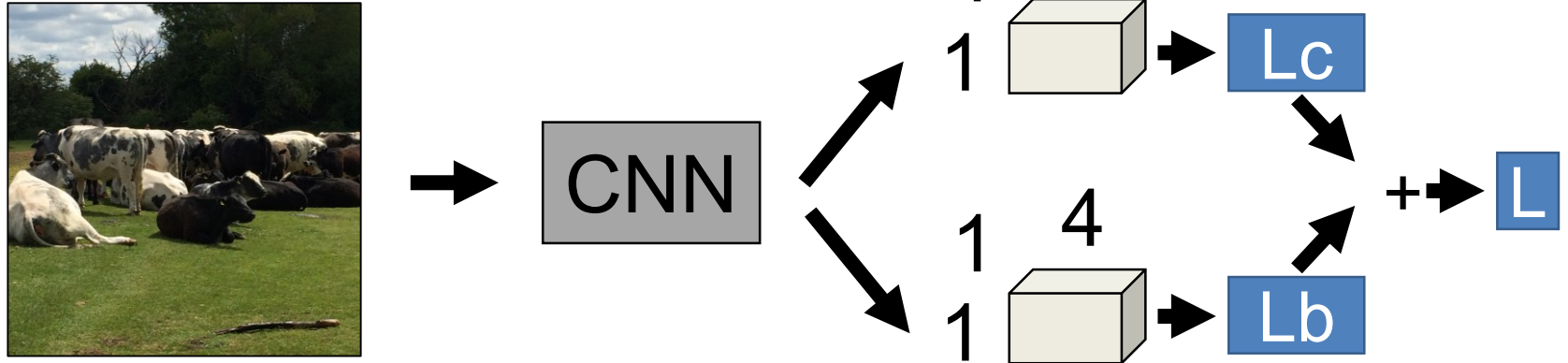


Now there are two ducks.

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.

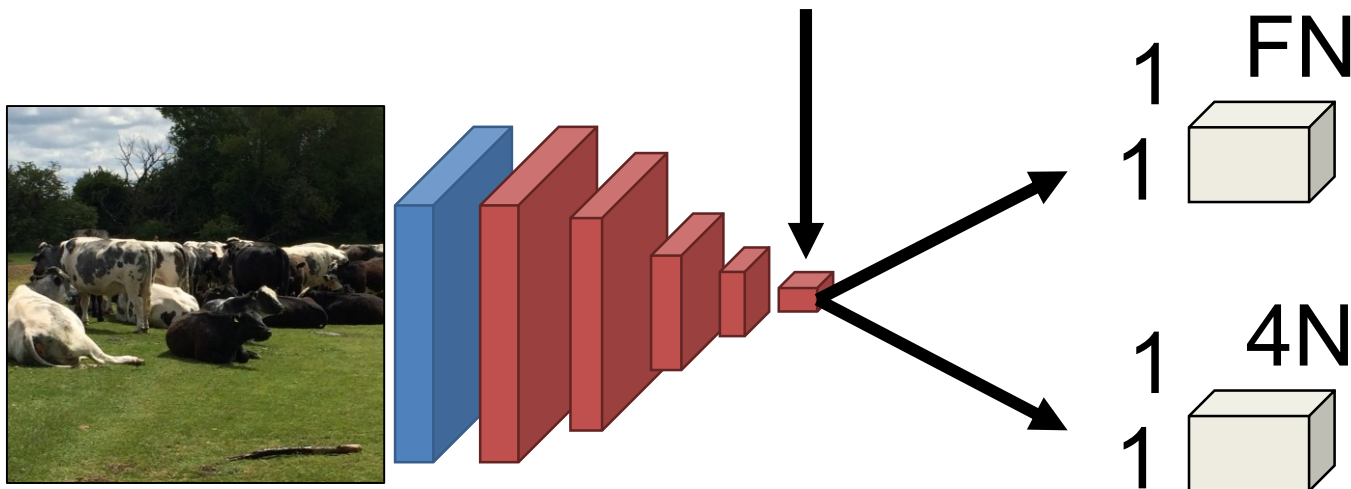
We need *lots* of outputs

(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

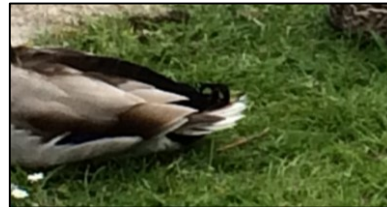


An Alternate Approach

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

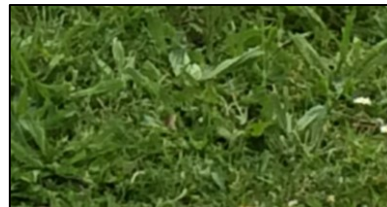


Yes



No?

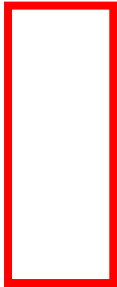
Hold this thought



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



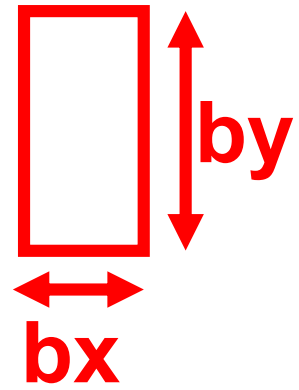
Slide credit: J. Hays

How Many Boxes Are There?

Given a $H \times W$ image and a “template” of size by, bx .

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

A. $(H-by) \times (W-bx)$



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 5th Grader?

Adults compete with
5th 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.

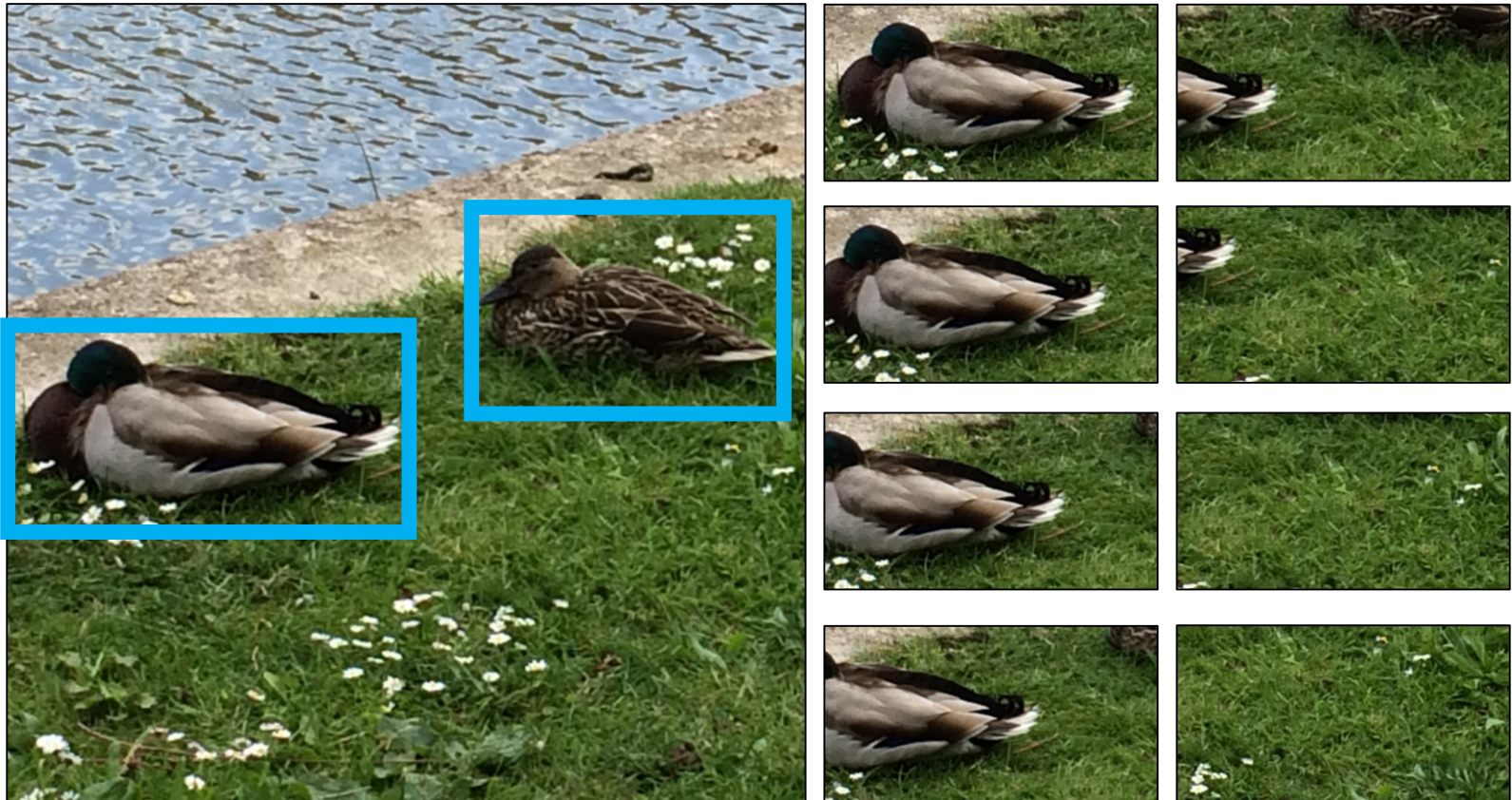
Models often not
better.

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

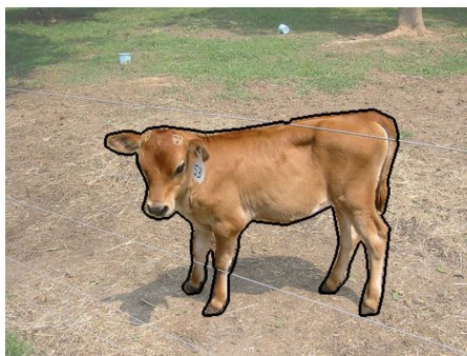
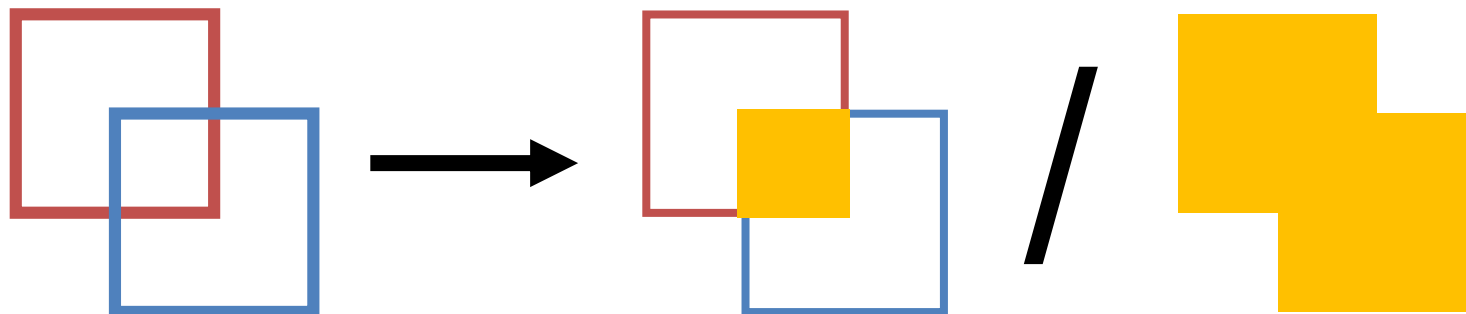
Raise your hand when you think the detection stops being correct.



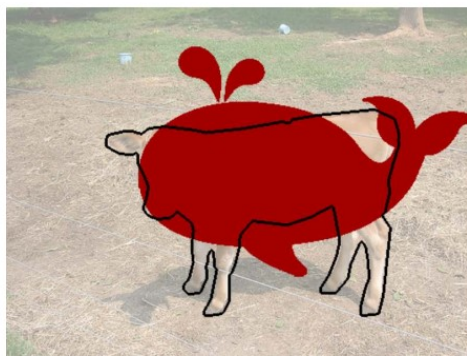
Evaluating – Bounding Boxes

Standard metric for two boxes:

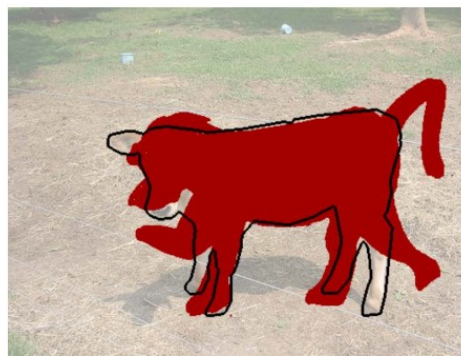
Intersection over union/IoU/Jaccard coefficient



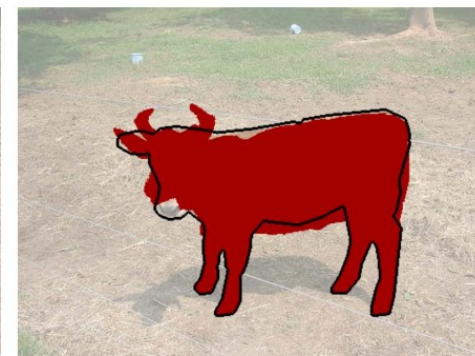
(a) Ground truth



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



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



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

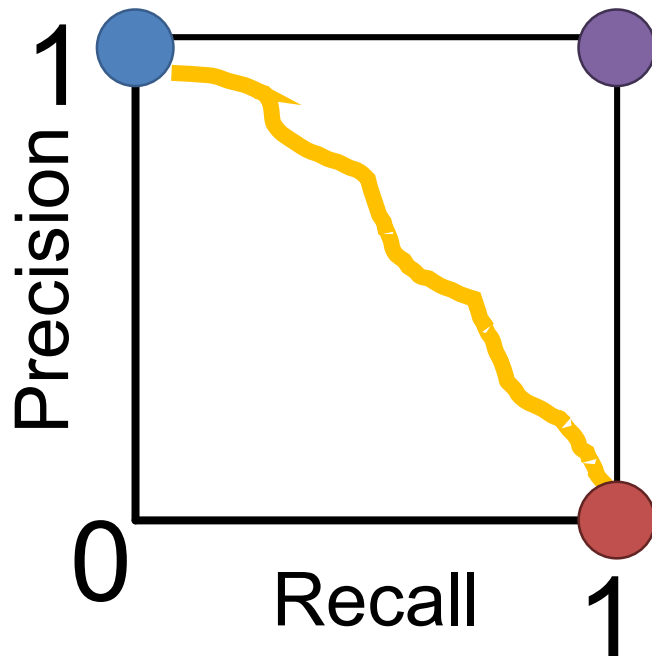
Evaluating Performance

- 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: high intersection over union
- Precision: $\# \text{true detections} / \# \text{detections}$
- Recall: $\# \text{true detections} / \# \text{true positives}$

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

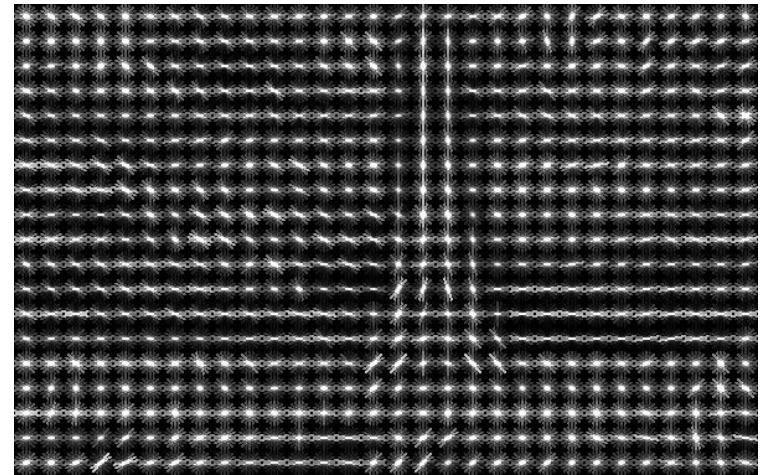


Image credit: N. Snavely

Pedestrian detection with HOG

- Train a pedestrian template using a linear support vector machine

positive training examples



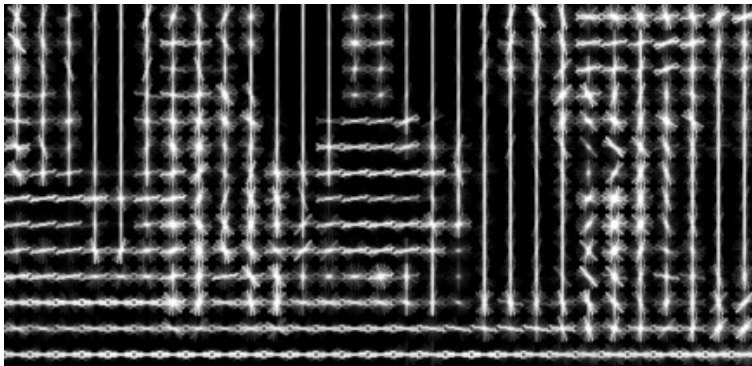
negative training examples



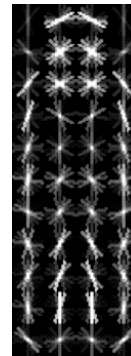
Pedestrian detection with HOG

- 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*

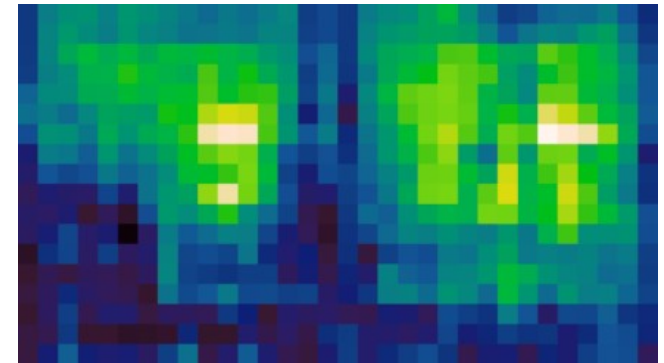
HOG feature map



Template



Detector response map



N. Dalal and B. Triggs, [Histograms of Oriented Gradients for Human Detection](#),

Slide Credit: S. Lazebnik

CVPR 2005

Example detections



[Dalal and Triggs, CVPR 2005]

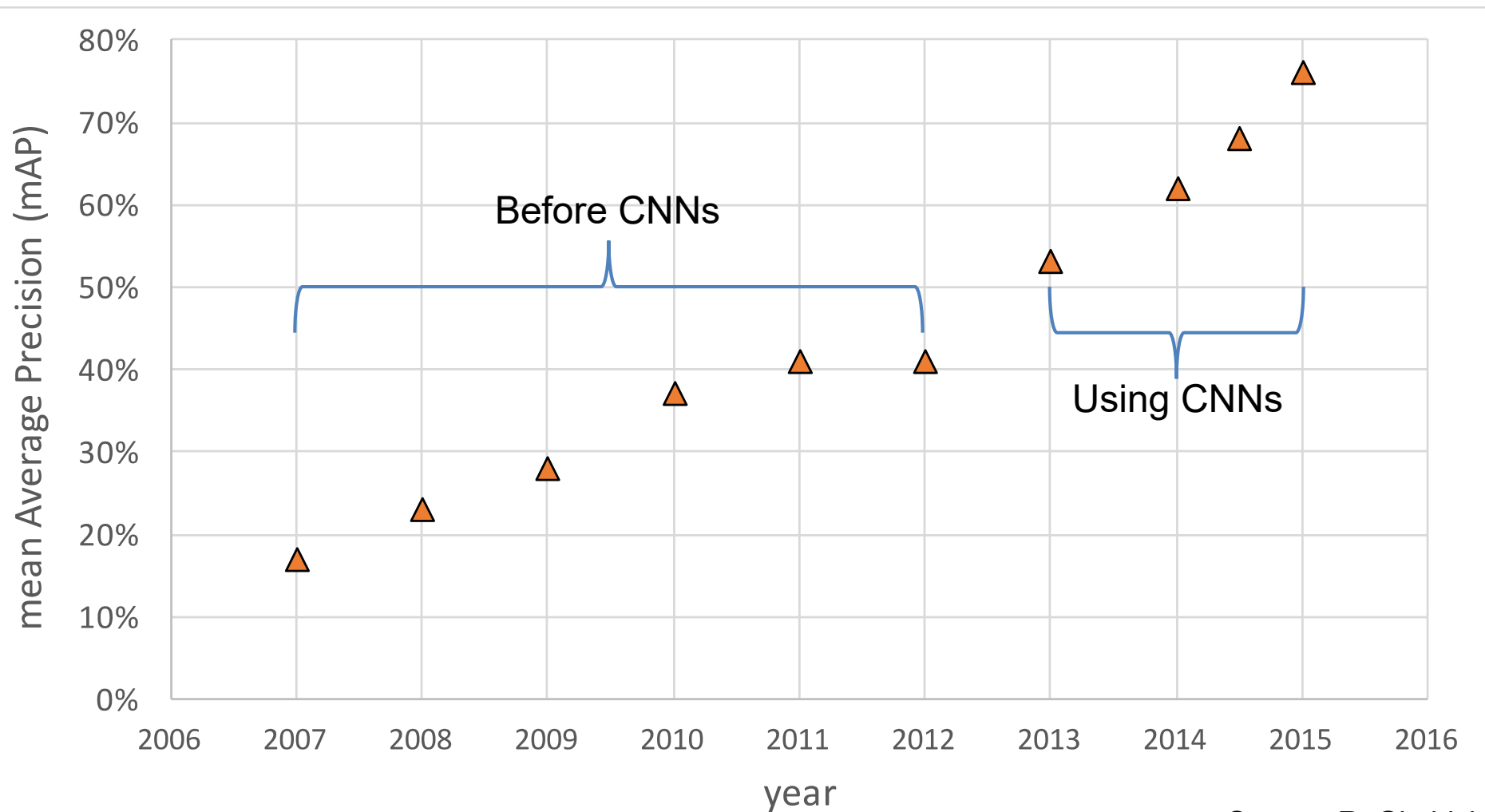
PASCAL VOC Challenge (2005-2012)



- 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

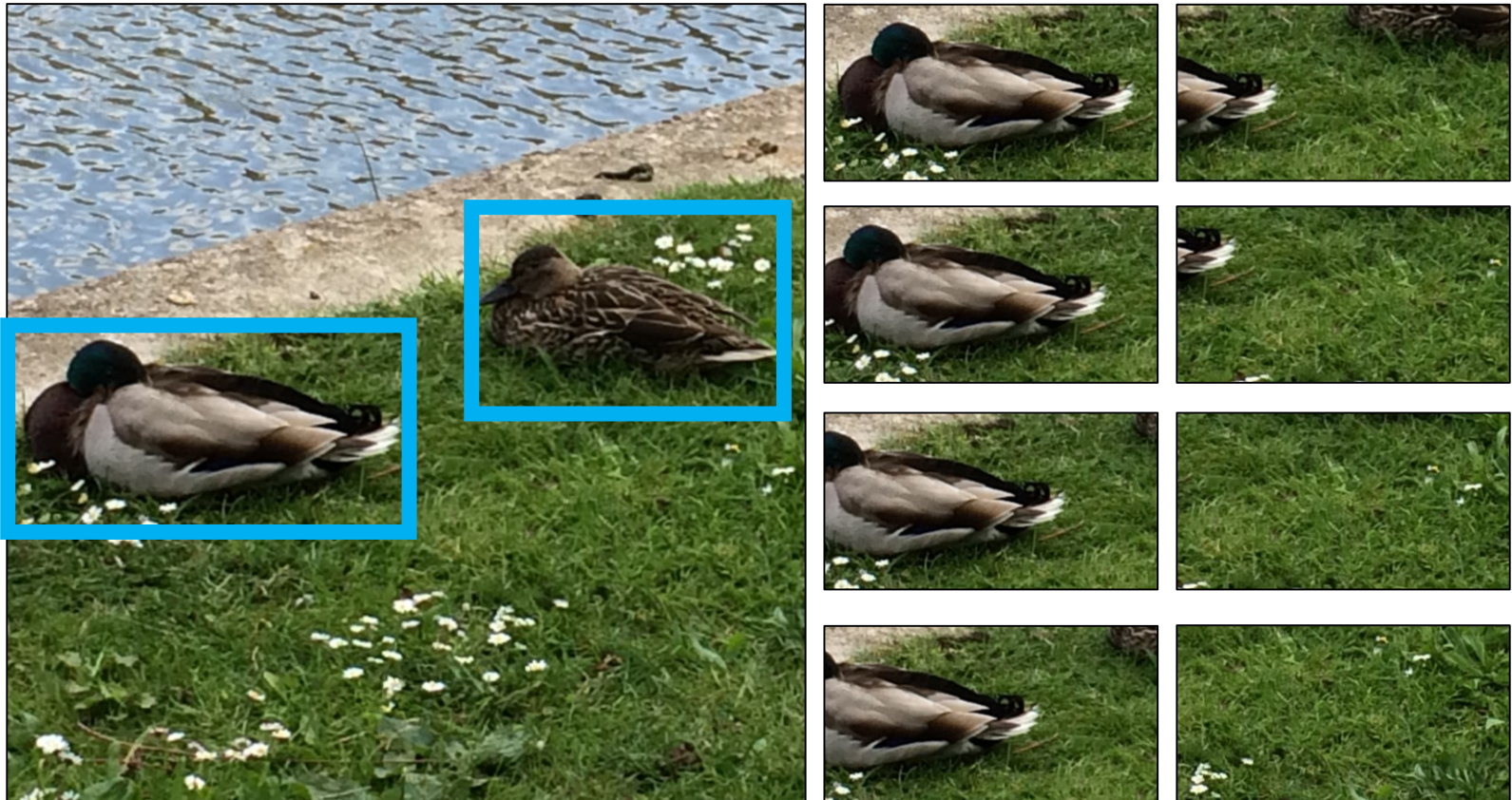
PASCAL VOC



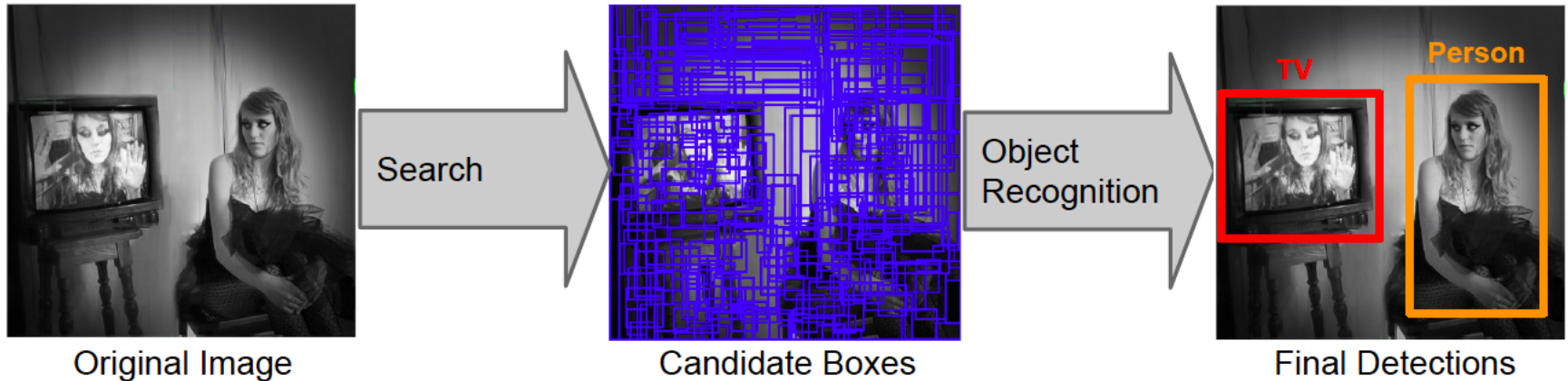
Source: R. Girshick

Region Proposals

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



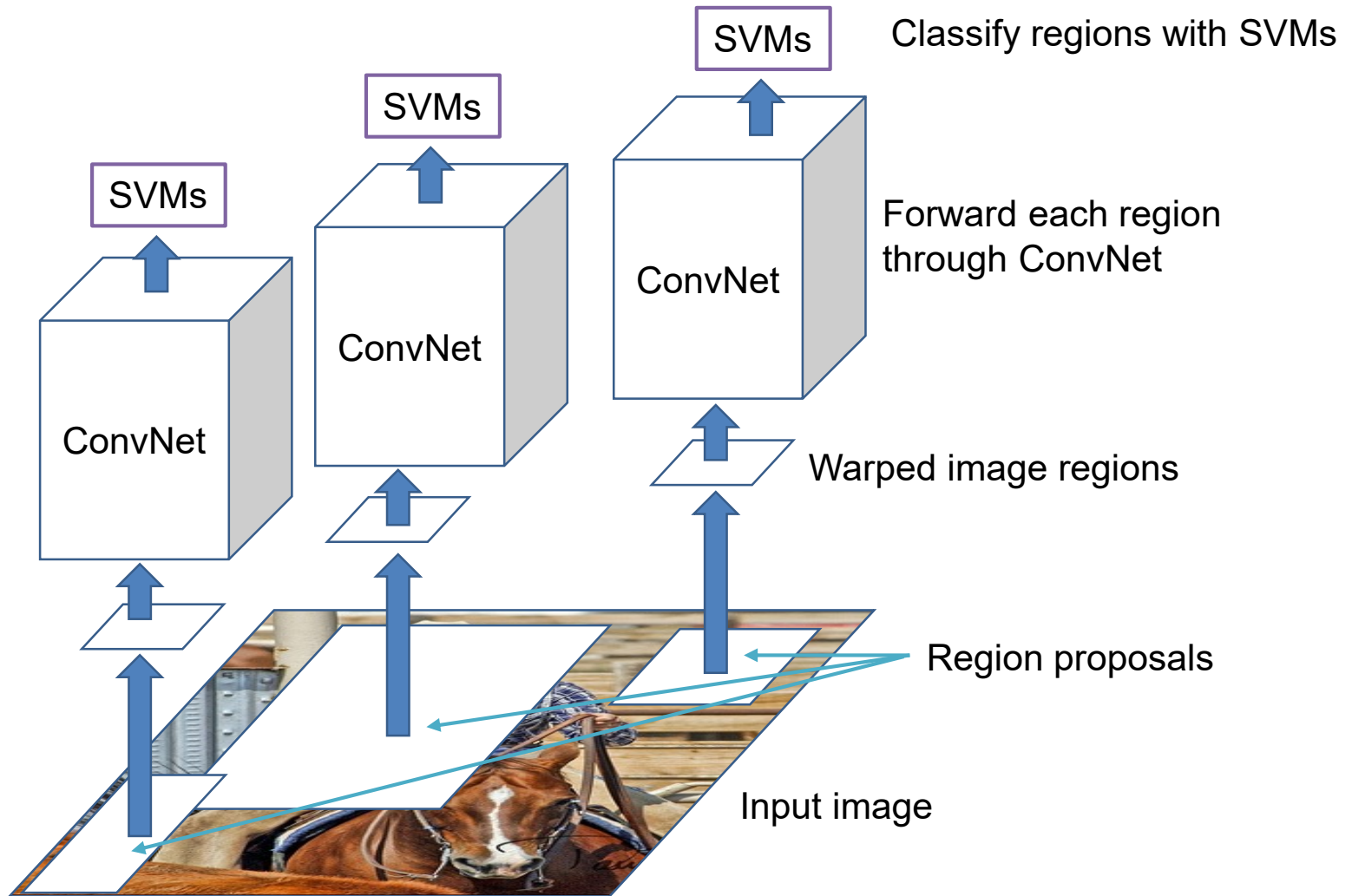
Region proposals



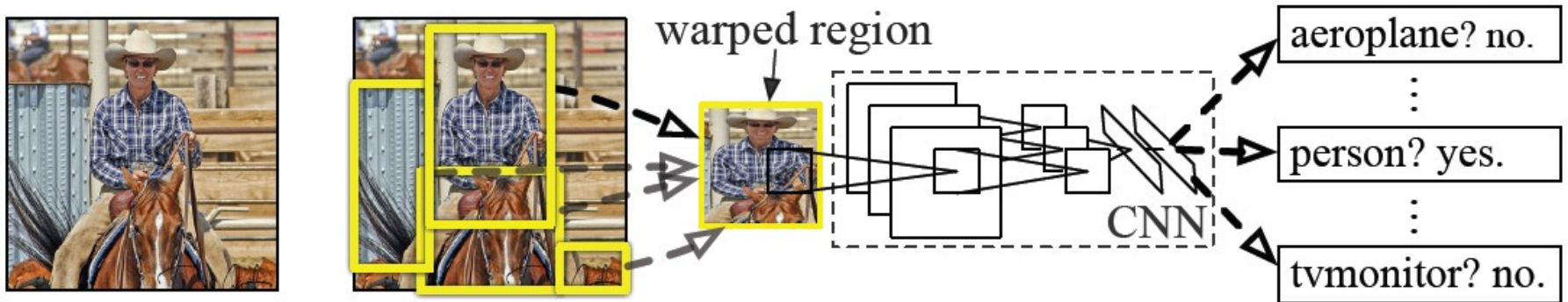
- As an alternative to sliding window search, evaluate a few hundred *region proposals*
 - Can use slower but more powerful features and classifiers
 - Proposal mechanism can be category-independent
 - Proposal mechanism can be trained

R-CNN: Region proposals + CNN features

Source: R. Girshick



R-CNN details

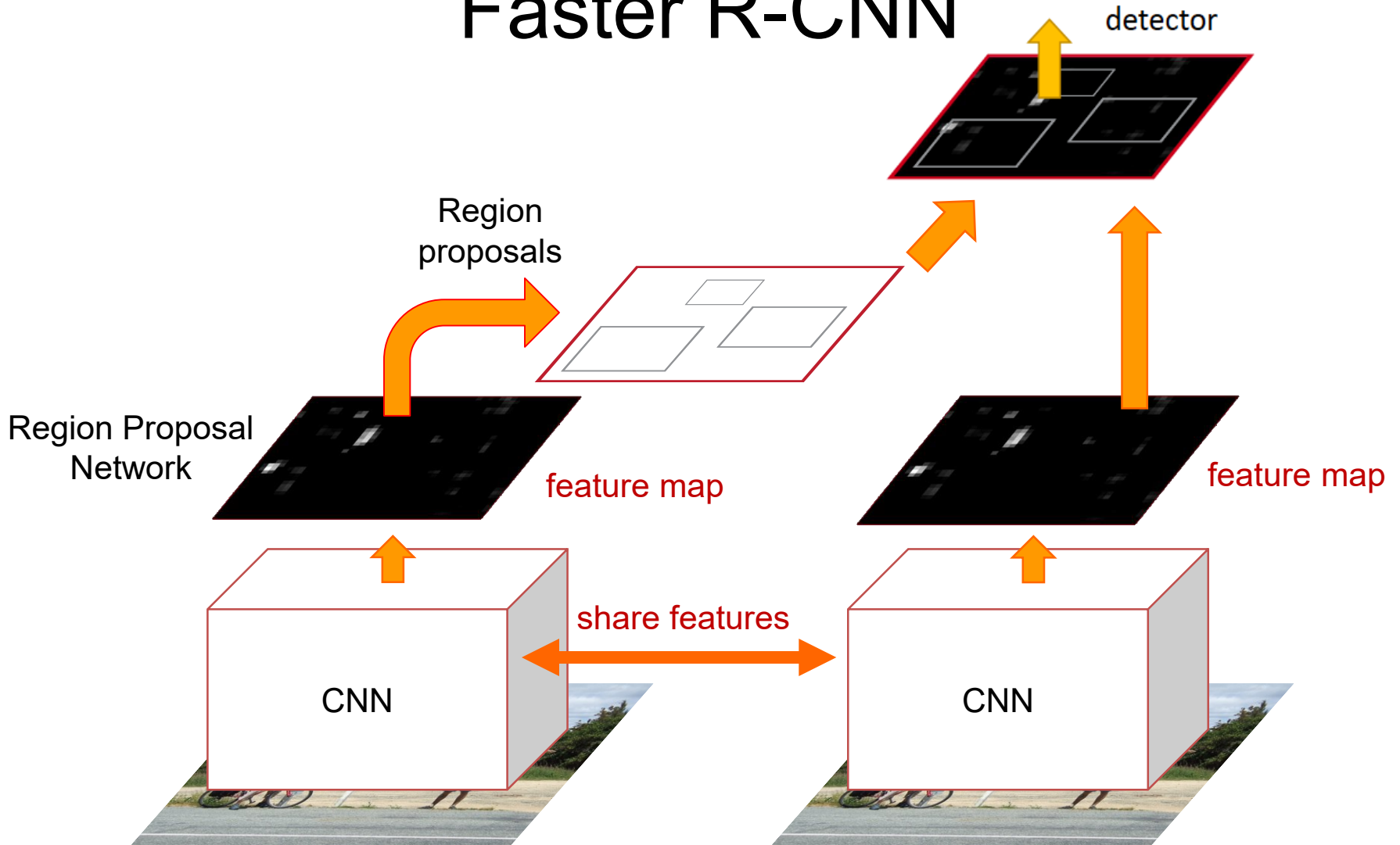


- **Regions:** ~2000 Selective Search proposals
- **Network:** AlexNet *pre-trained* on ImageNet (1000 classes), *fine-tuned* on PASCAL (21 classes)
- **Final detector:** warp proposal regions, extract fc7 network activations (4096 dimensions), classify with linear SVM
- **Bounding box regression** to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

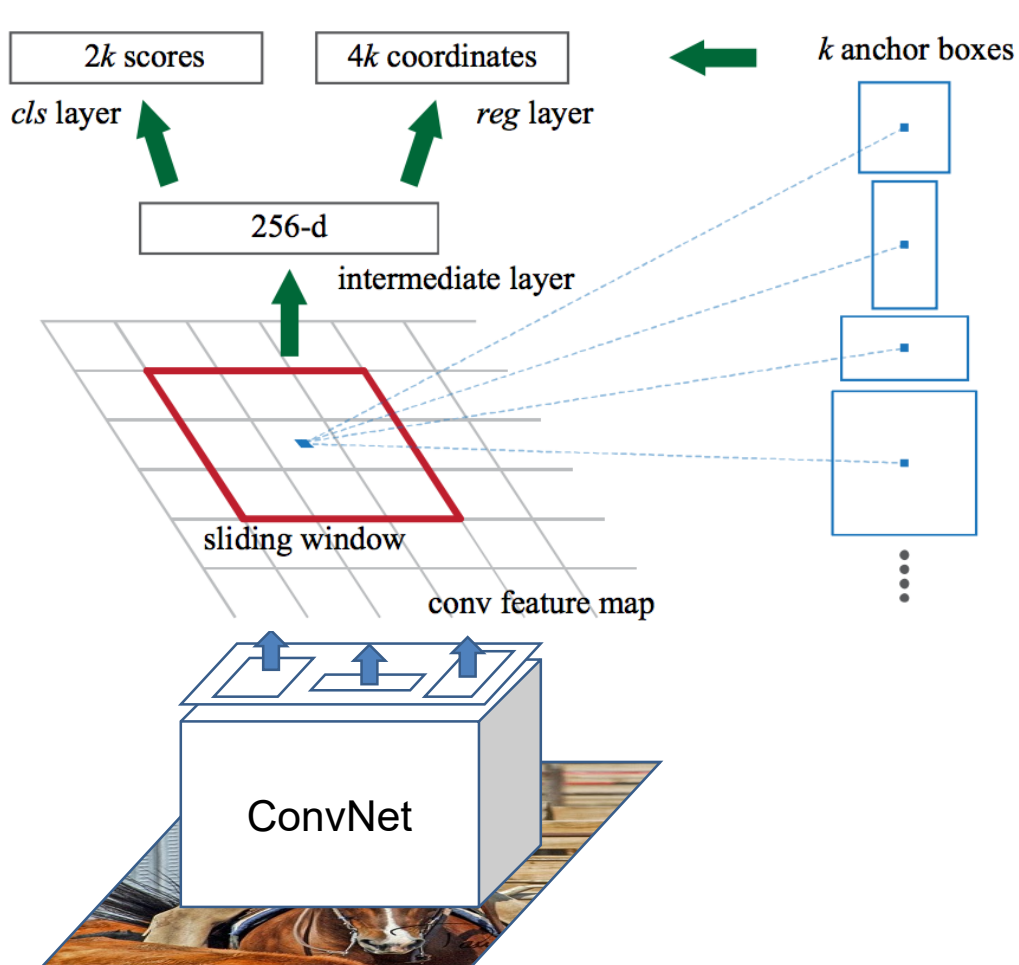
R-CNN pros and cons

- **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)

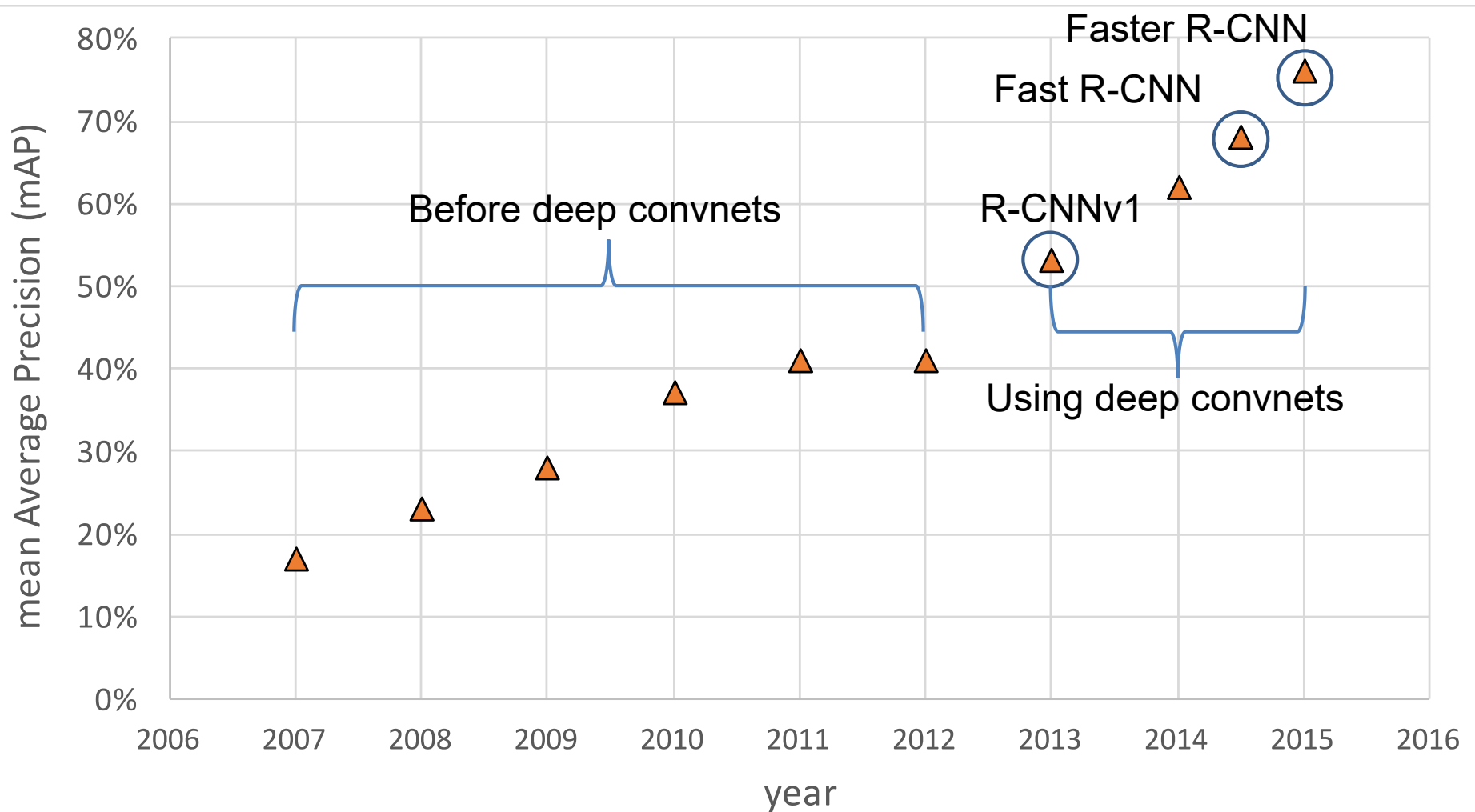
for k “anchors” or boxes relative to the position in feature map.

Faster R-CNN results

system	time	07 data	07+12 data
R-CNN	~50s	66.0	-
Fast R-CNN	~2s	66.9	70.0
Faster R-CNN	198ms	69.9	73.2

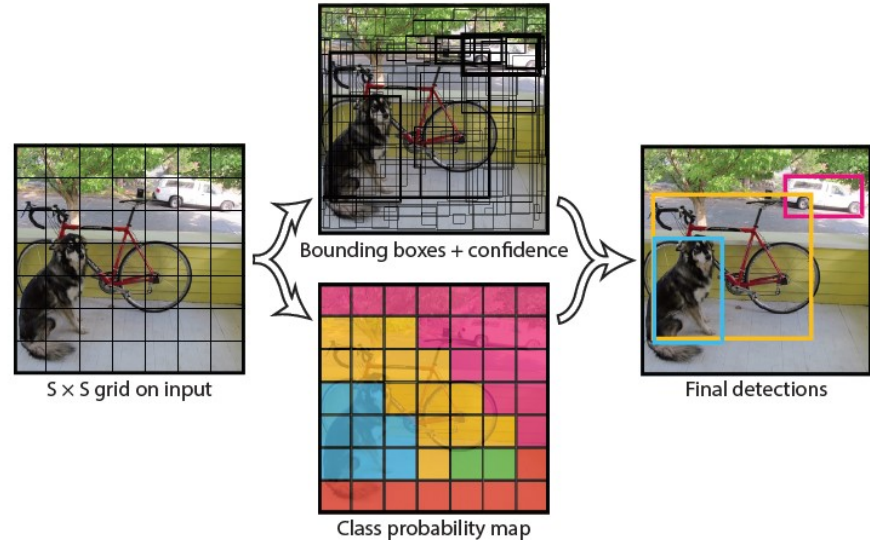
detection mAP on PASCAL VOC 2007, with VGG-16 pre-trained on ImageNet

Object detection progress

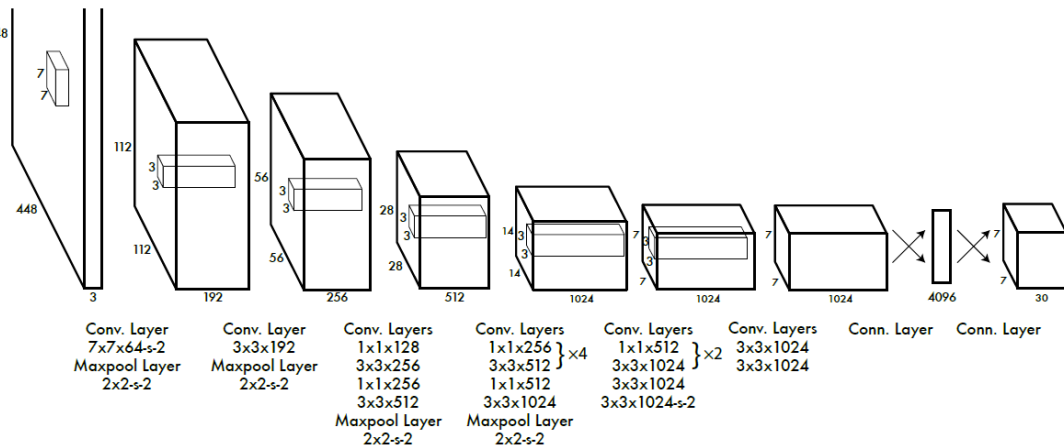


YOLO

1. Take conv feature maps at 7x7 resolution
2. Add two FC layers to predict, at each location, score for each class and 2 bboxes w/ confidences



- 7x speedup over Faster RCNN (45-155 FPS vs. 7-18 FPS)
- Some loss of accuracy due to lower recall, poor localization



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 body
- 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



(a) Person

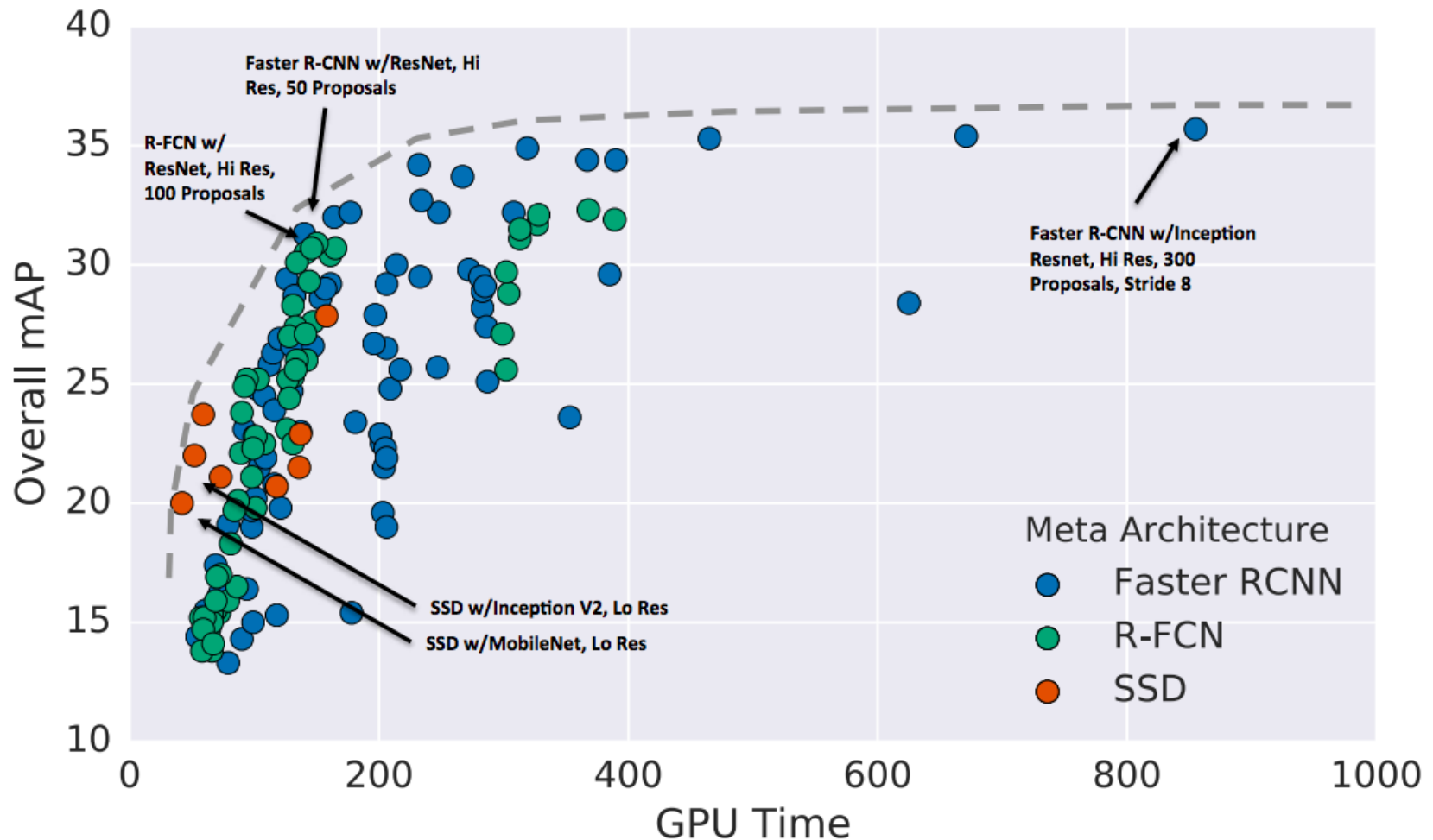
(b) Bicycle

(c) Giraffe

Kitchens from Googling



New detection benchmark: COCO (2014)



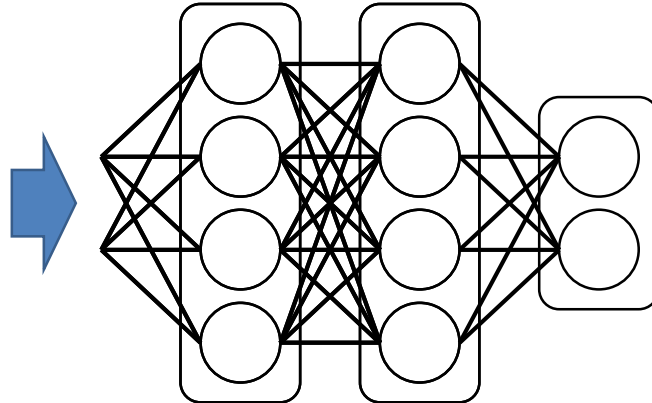
J. Huang et al., [Speed/accuracy trade-offs for modern convolutional object detectors](#), CVPR 2017

Summary: Object detection with CNNs

- R-CNN: region proposals + CNN on cropped, resampled regions
- Fast R-CNN: region proposals + RoI pooling on top of a conv feature map
- Faster R-CNN: RPN + RoI pooling
- Next generation of detectors
 - Direct prediction of BB offsets, class scores on top of conv feature maps
 - Get better context by combining feature maps at multiple resolutions

And Now For Something
Completely Different

ImageNet + Deep Learning

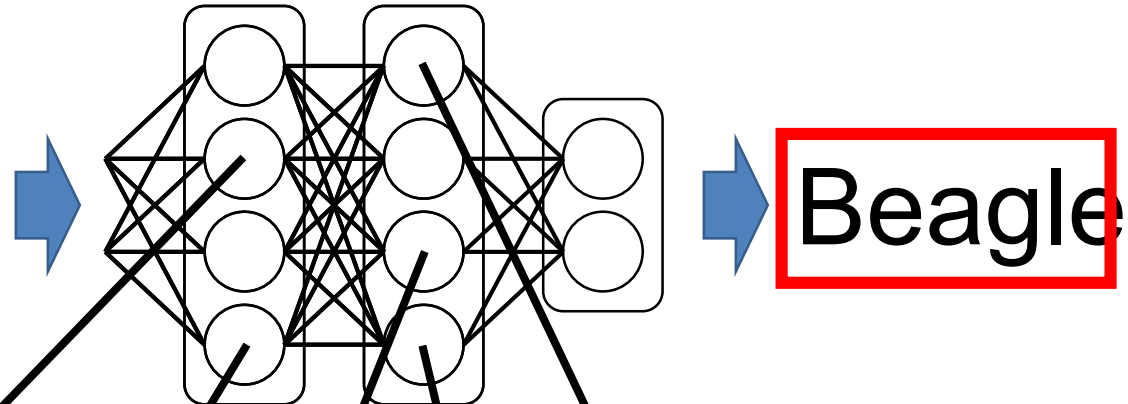


Beagle



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

ImageNet + Deep Learning



Materials?

Pose?

Do we even need semantic labels?

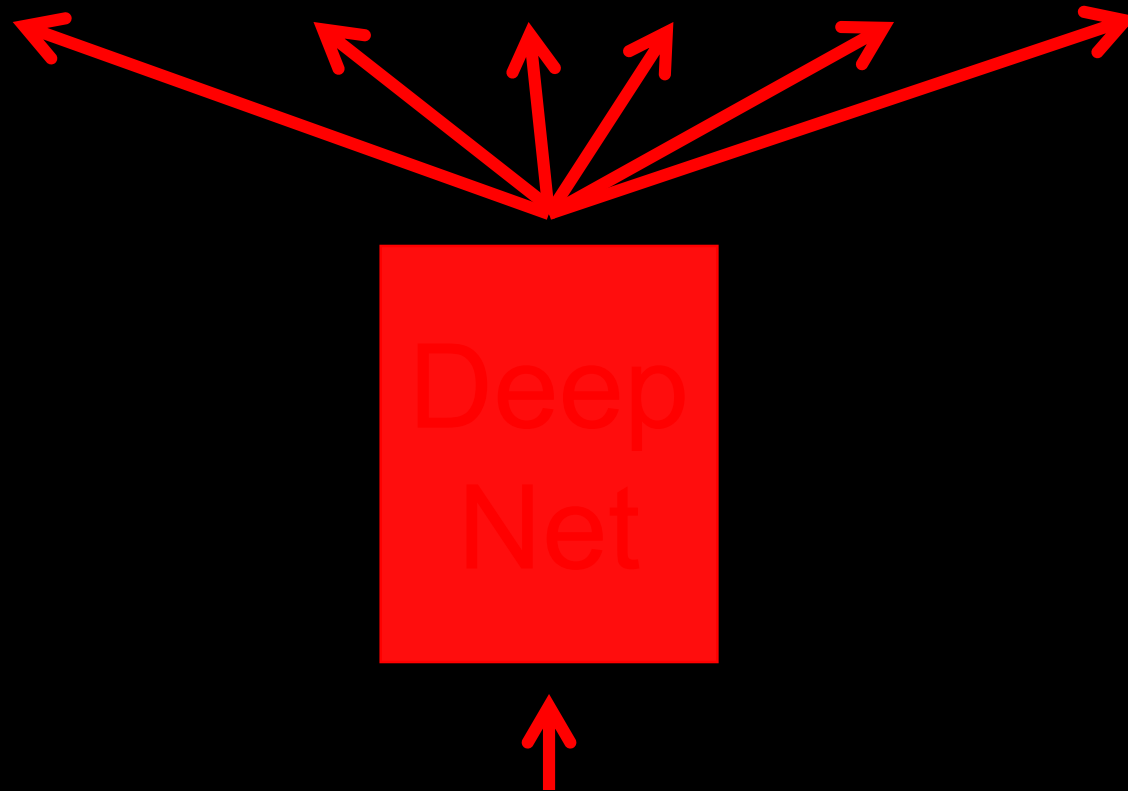
Parts?

Geometry?

Boundaries?

Context as Supervision

[Collobert & Weston 2008; Mikolov et al. 2013]



store-bought gimmicks and appliances, the toasters and

Context Prediction for Images

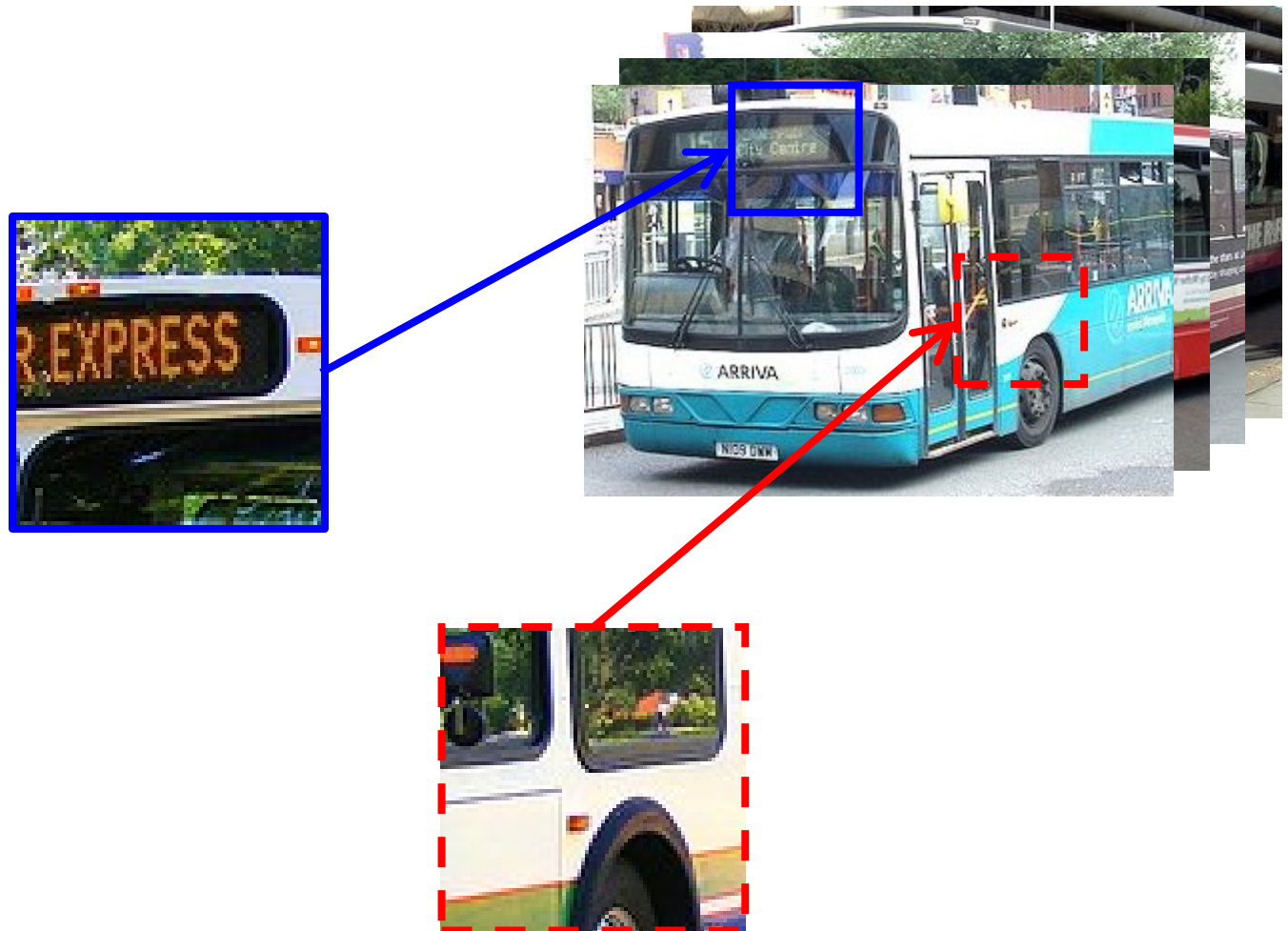


A

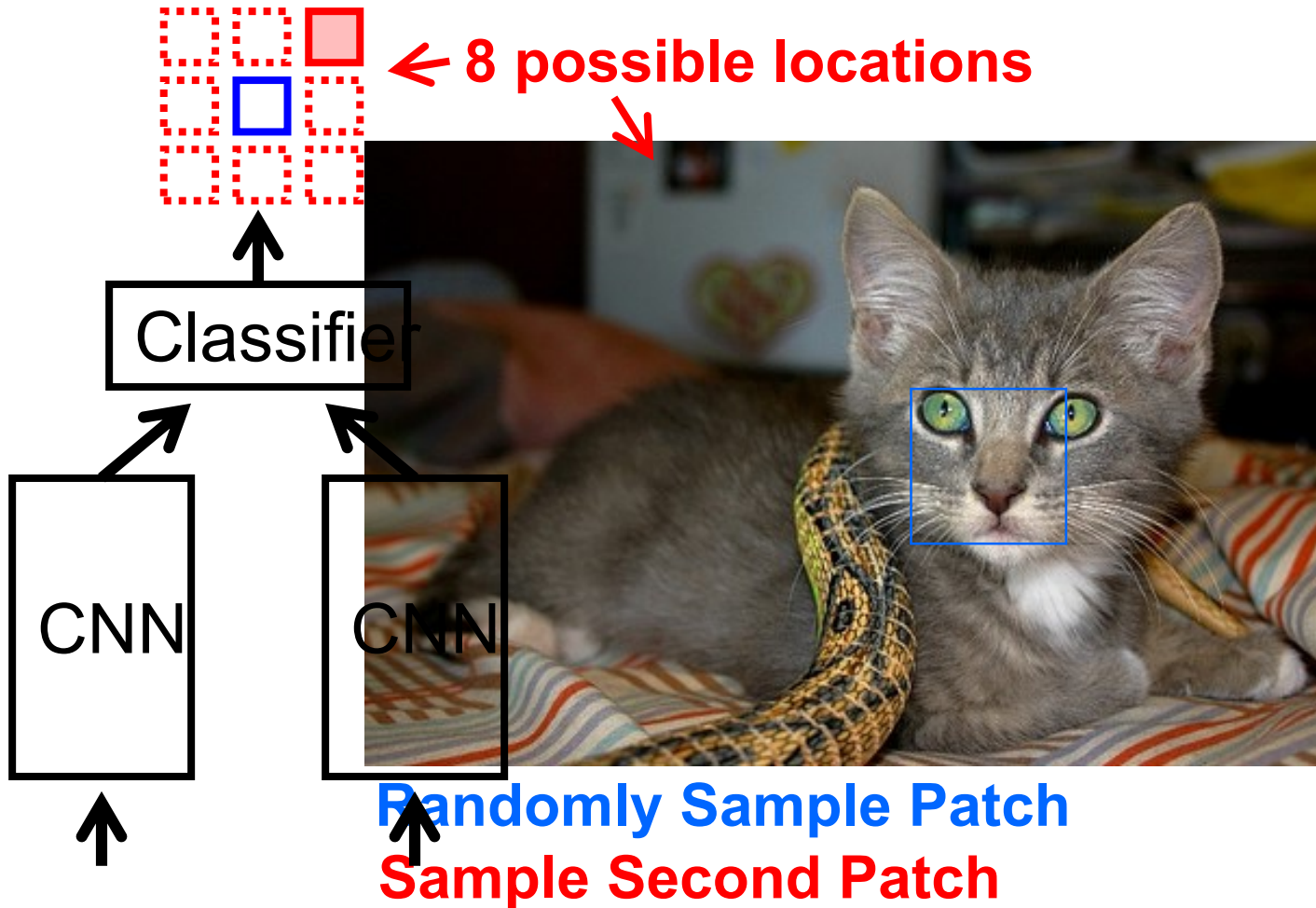
B



Semantics from a non-semantic task



Relative Position Task



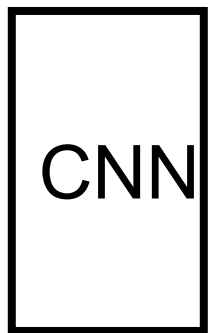
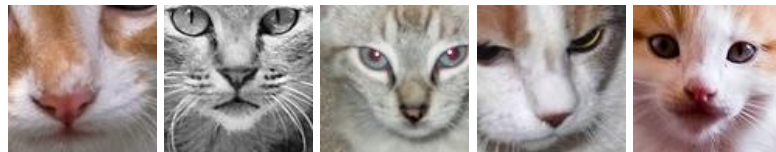
Patch Embedding

Input



!

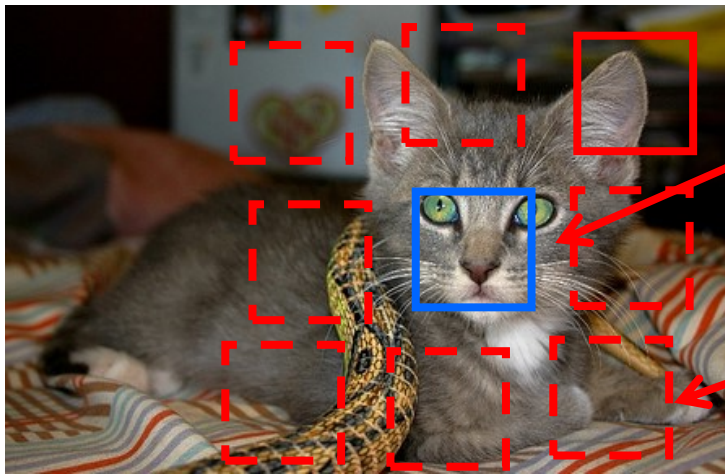
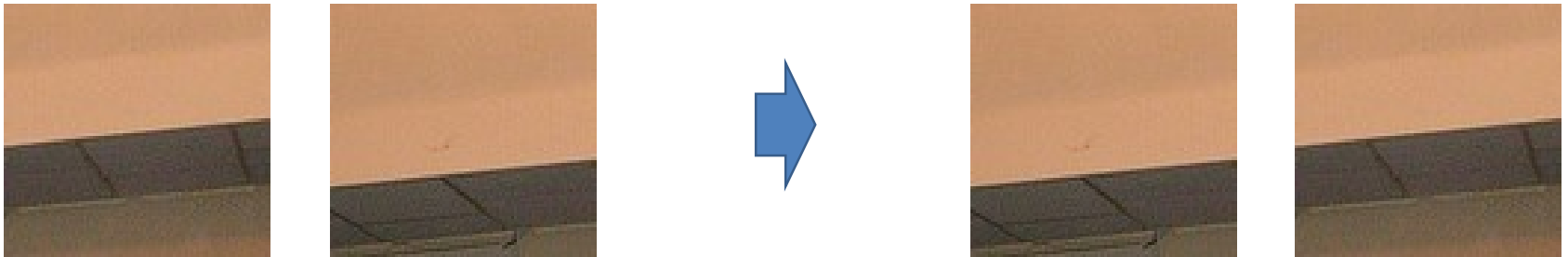
Nearest Neighbors



Note: connects ***across*** instances



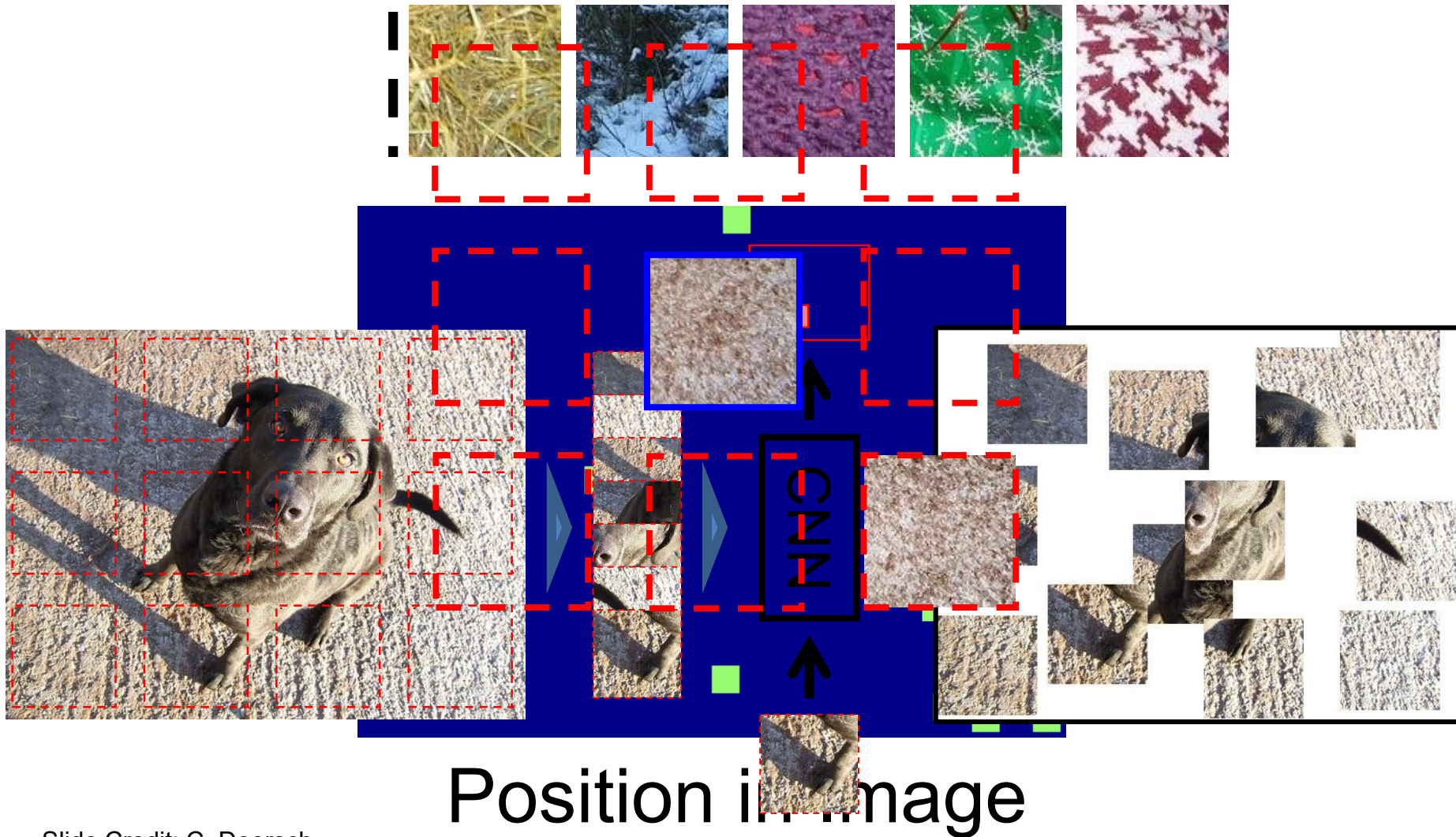
Avoiding Trivial Shortcuts



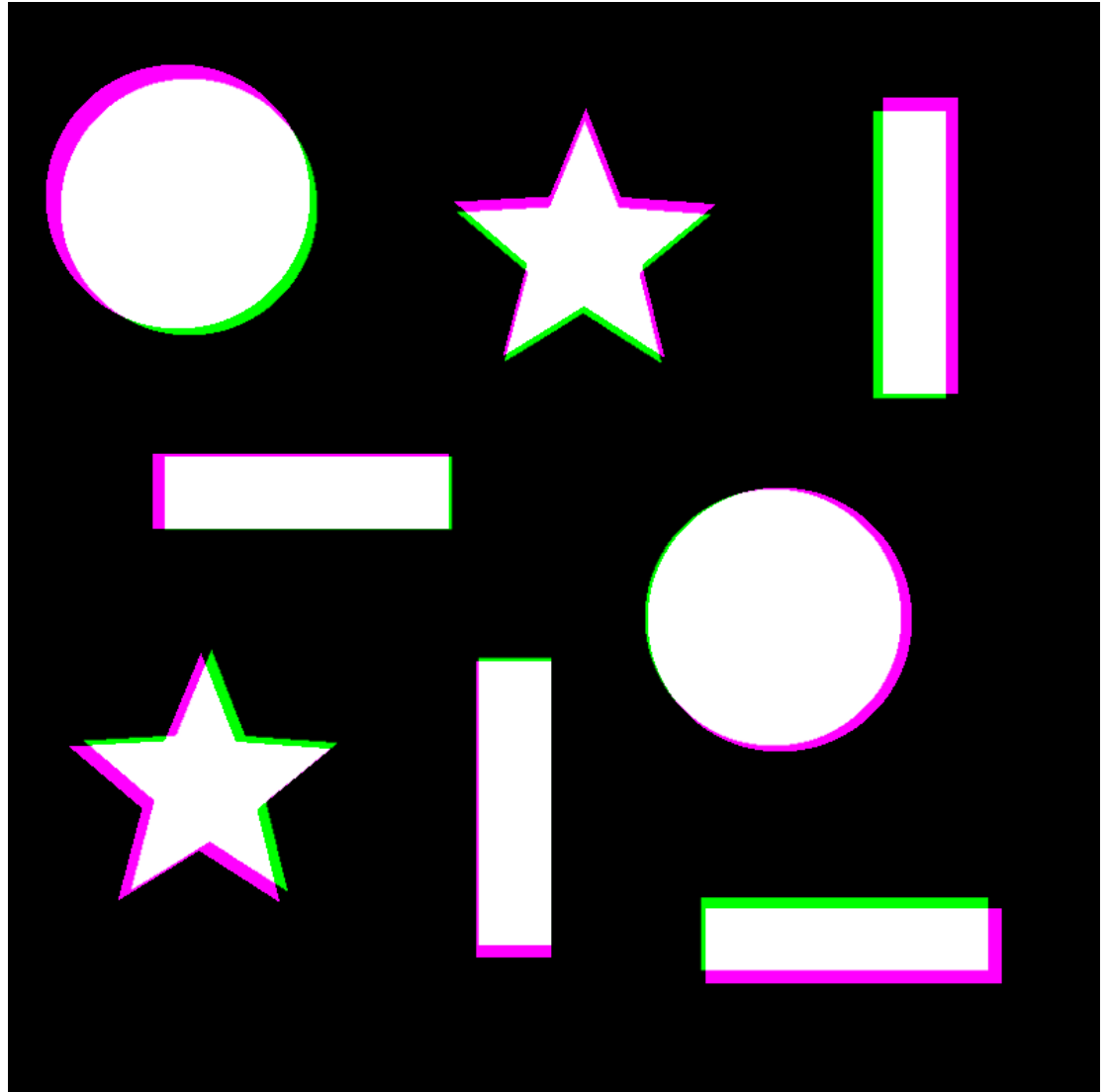
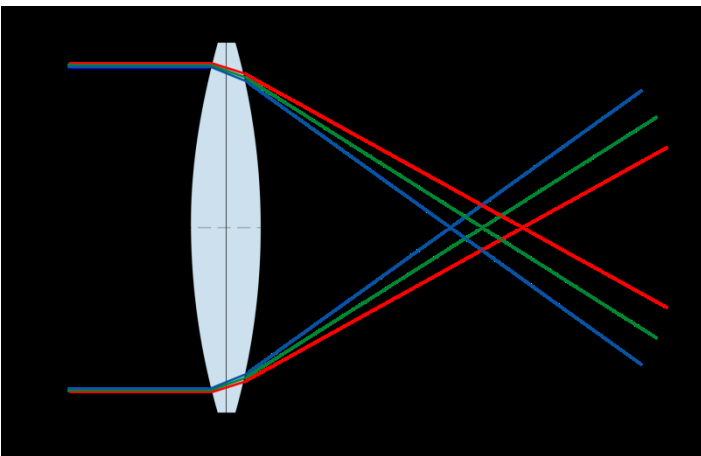
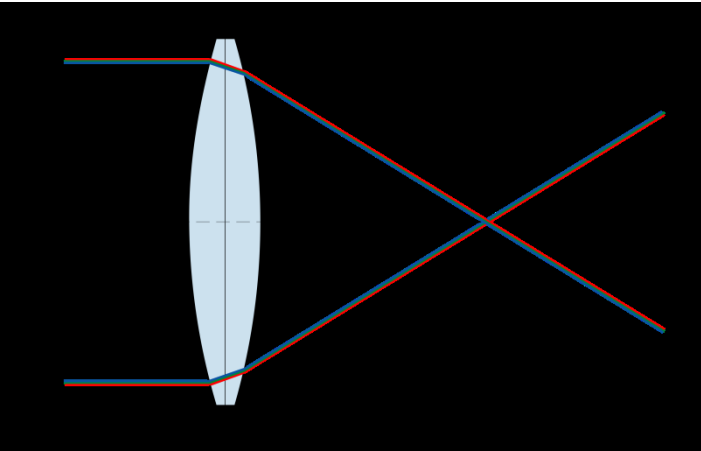
Include a gap

Jitter the patch locations

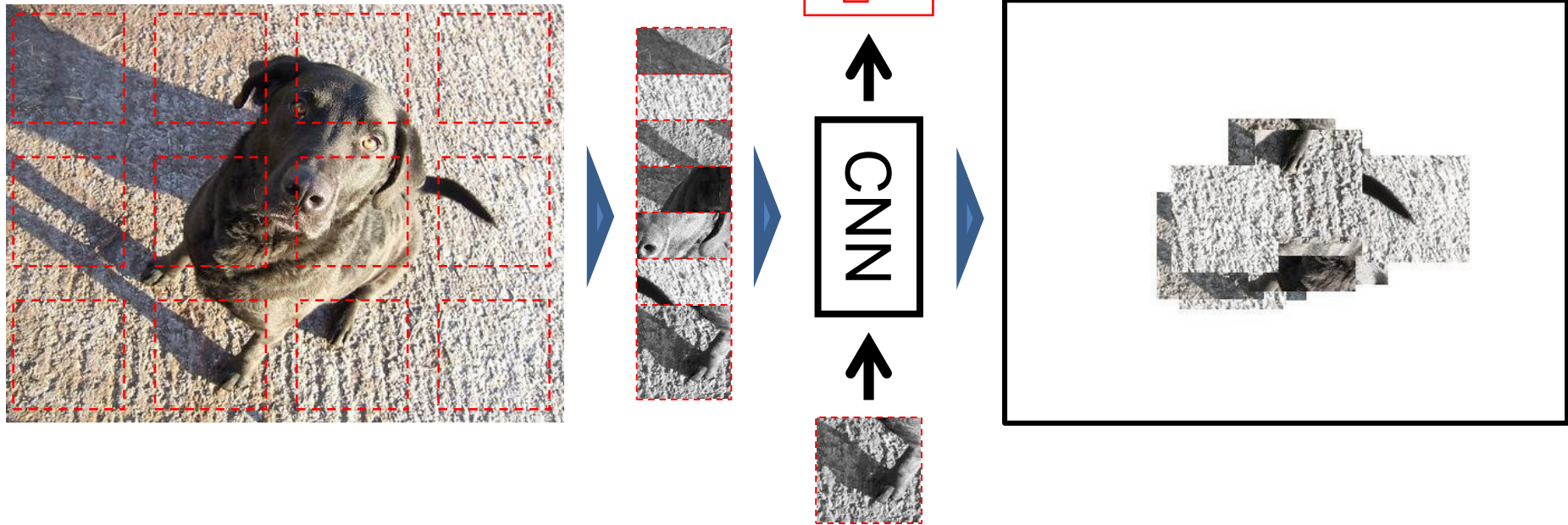
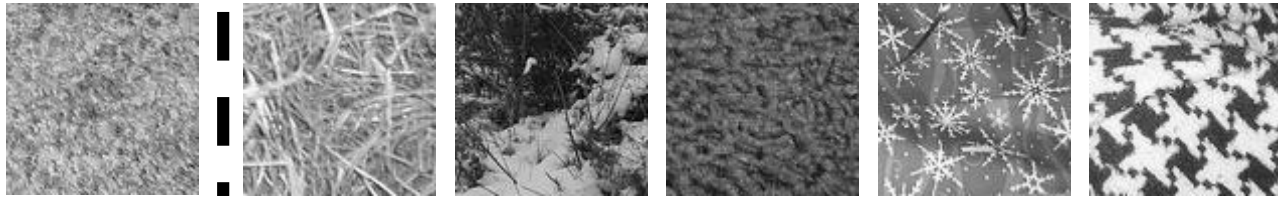
A Not-So “Trivial” Shortcut



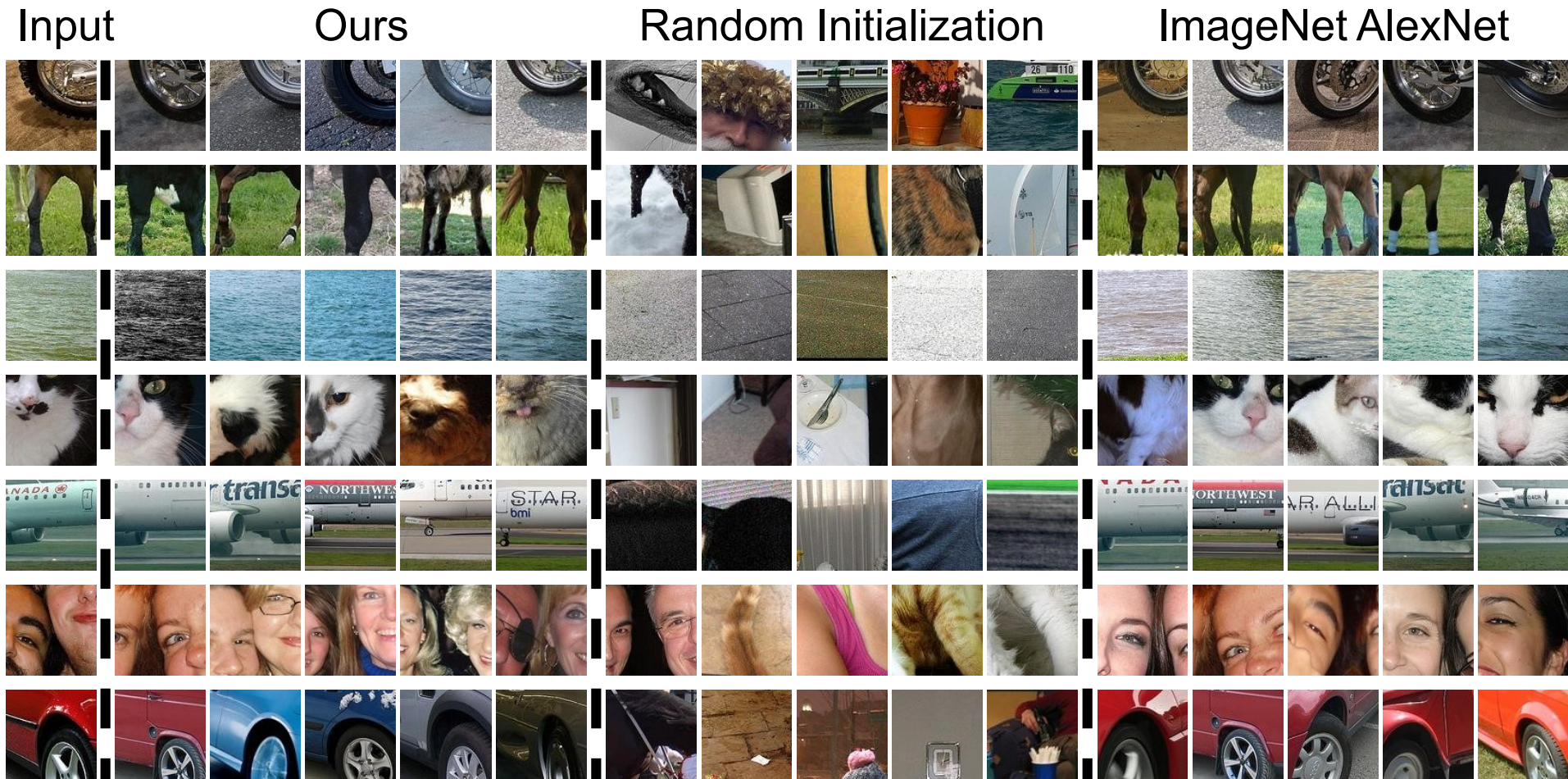
Chromatic Aberration



Chromatic Aberration



What is learned?



Pre-Training for R-CNN

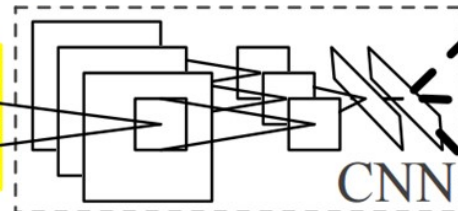


1. Input image



2. Extract region proposals (~2k)

warped region



CNN

3. Compute CNN features

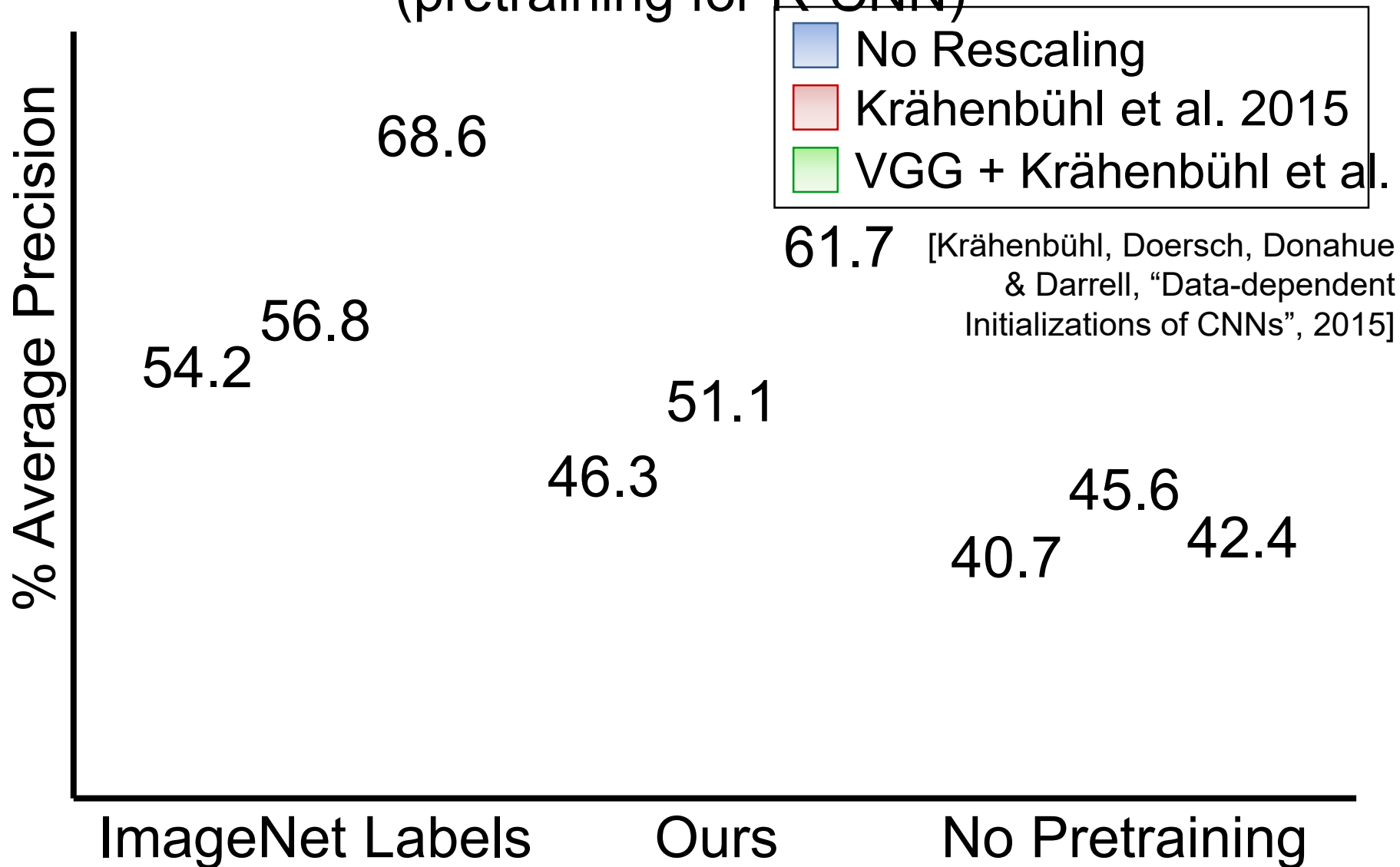
- aeroplane? no.
- ⋮
- person? yes.
- ⋮
- tvmonitor? no.

4. Classify regions

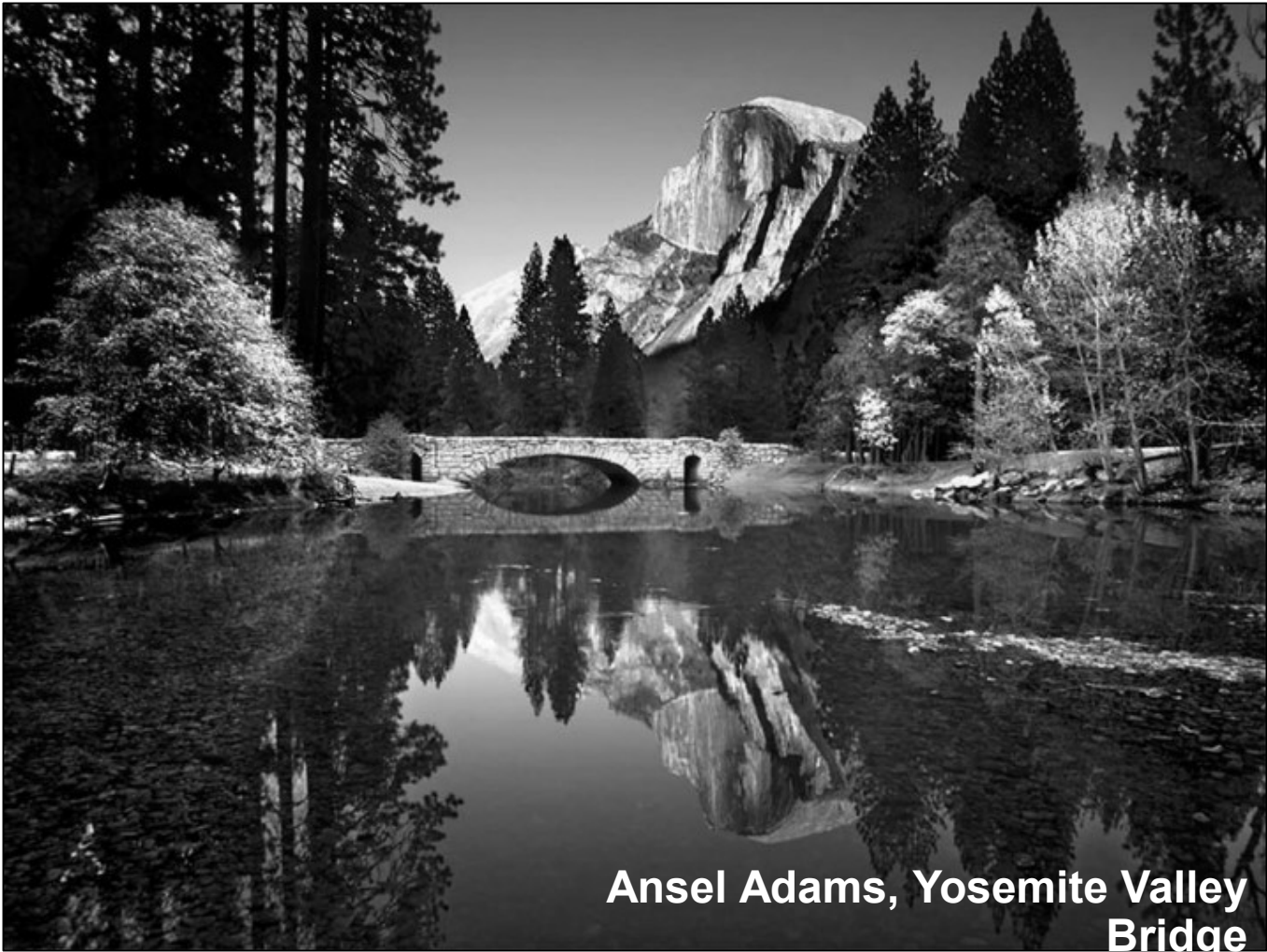
Pre-train on relative-position task, w/o labels

VOC 2007 Performance

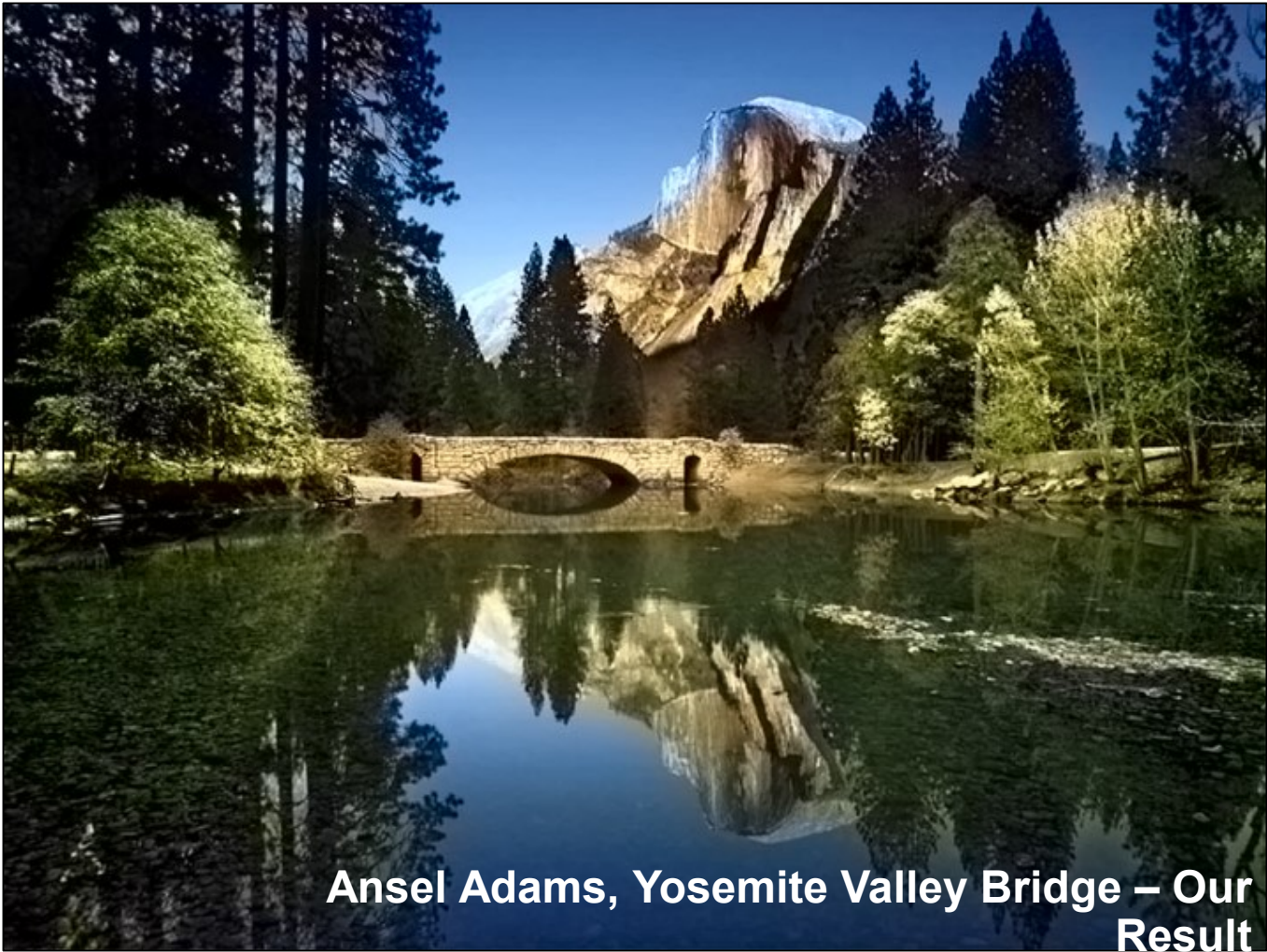
(pretraining for R-CNN)

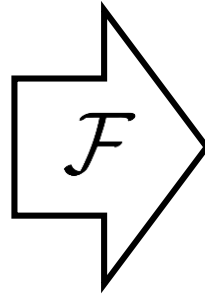


Other Sources Of Signal



**Ansel Adams, Yosemite Valley
Bridge**



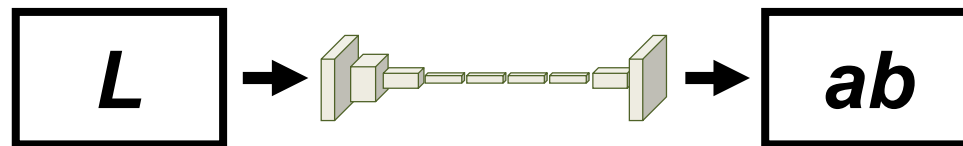


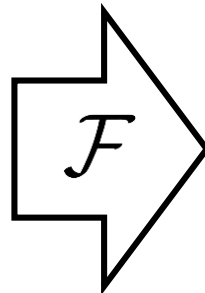
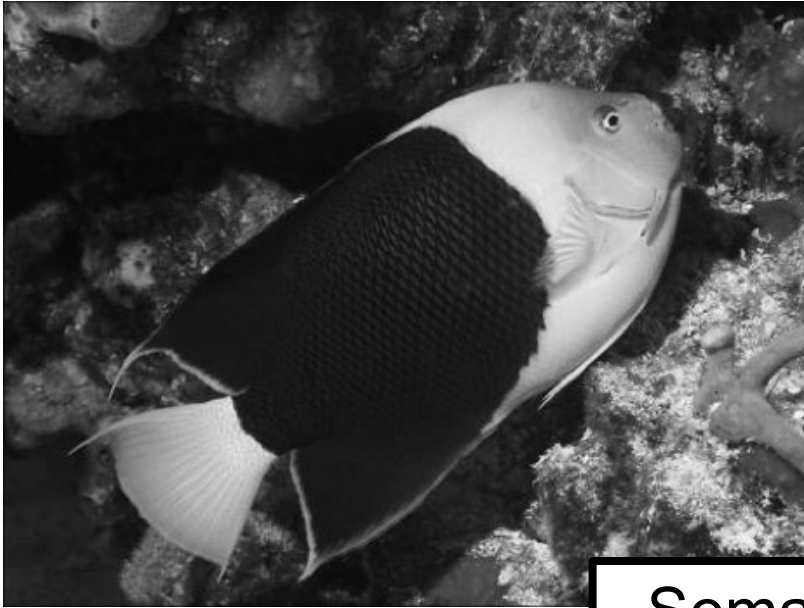
Grayscale image: L channel

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

Color information: ab channels

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

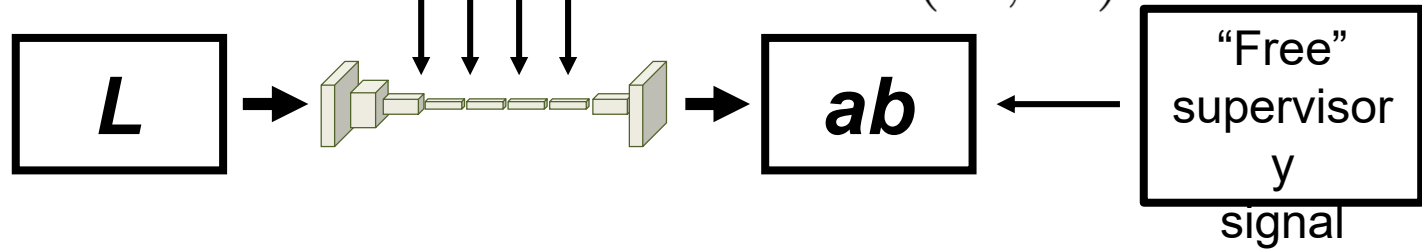




Semantics? Higher-level abstraction?

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

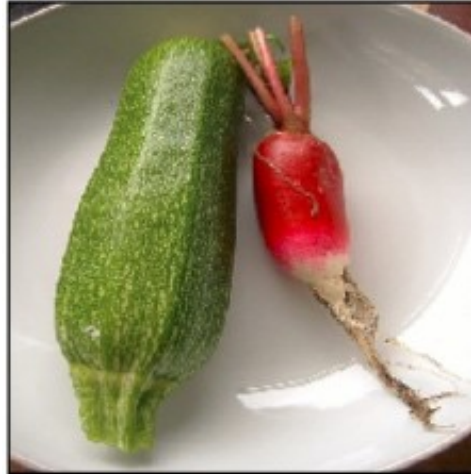
Concatenate (L, ab)
 $(\mathbf{X}, \hat{\mathbf{Y}})$



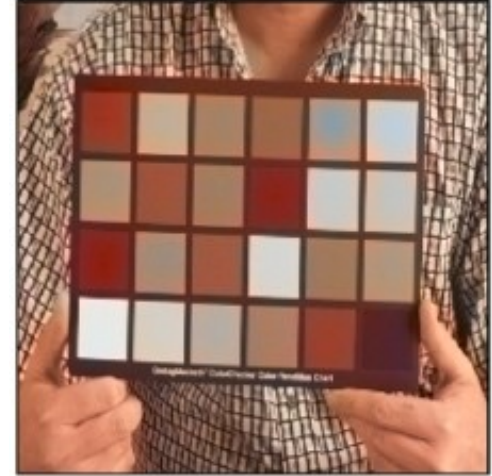
Input



Ground Truth



Output





Visually Indicated Sounds

Andrew Owens

Phillip Isola

Josh McDermott

Antonio Torralba

Edward Adelson

William Freeman

Face Recognition

- Some goals for face/person/recognition:
 - Should be able to recognize lots of people
 - Shouldn't require lots of data per-person
 - Should be able to recognize a new person post-training
- Classification doesn't handle any of these well.

Face Recognition

Goal: given images, would like to be able to compute distances between the images.

Face recognition is then: given reference image, is this person the same as the reference image (here: < 1.1)



1.22

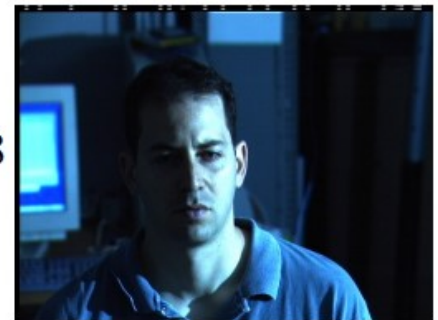


1.04

1.33

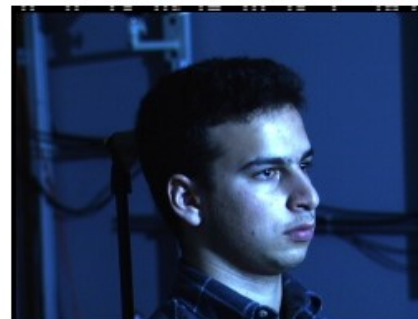


1.33

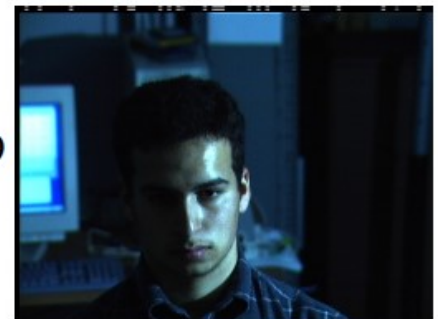


0.78

1.26



0.99

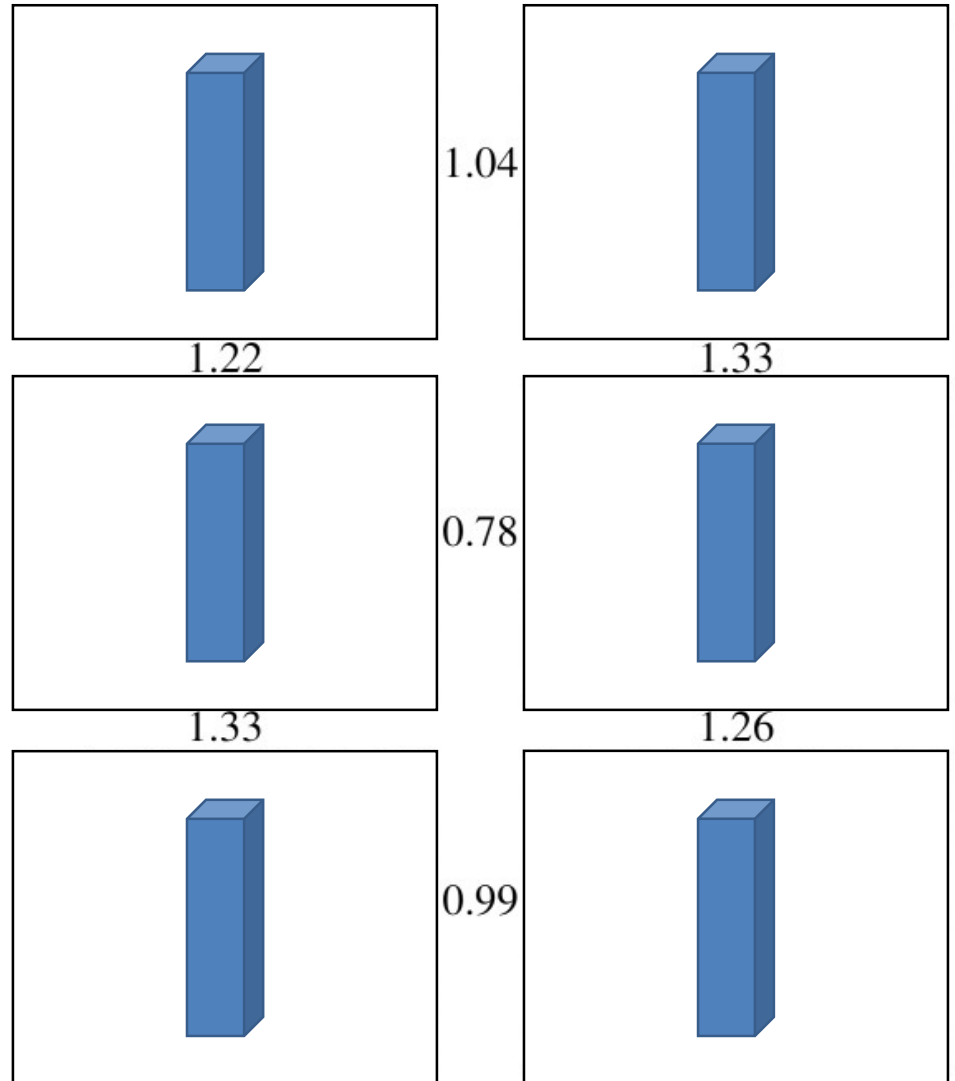


Face Recognition

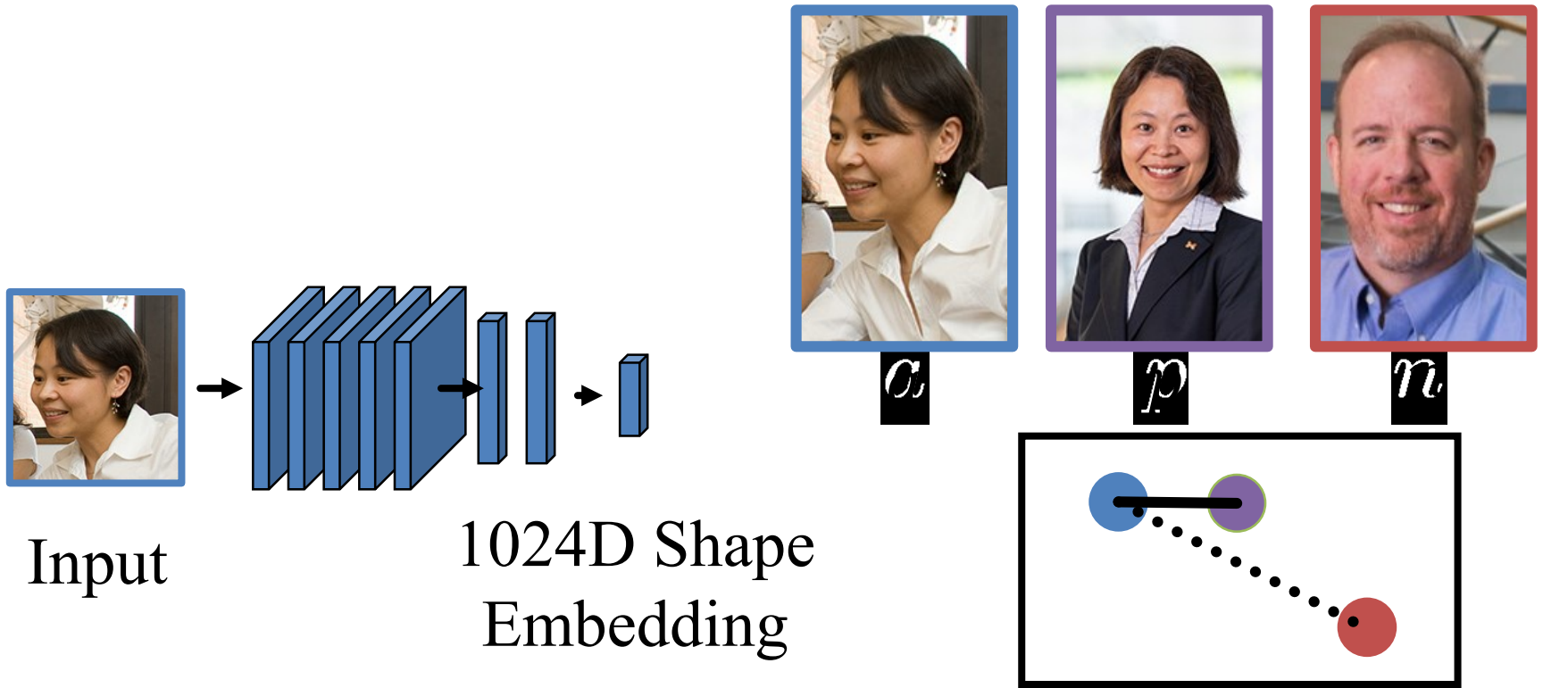
Distances between images are not meaningful.

Need to convert image into vector (typically normalized to unit sphere)

How?

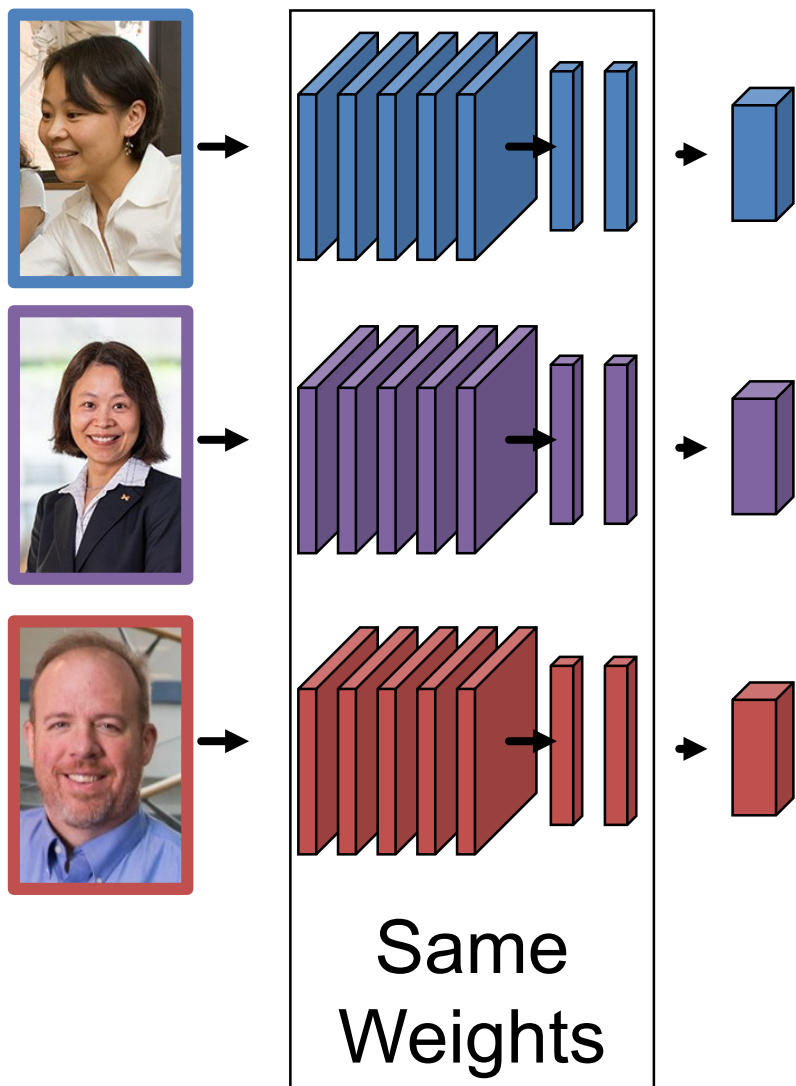


Triplet Network



$$L(a, p, n) = \max(\overline{D(a, p)} - \overline{D(a, n)} + \alpha, 0)$$

Training



$$\max(D(\text{blue}, \text{purple}) - D(\text{blue}, \text{red}) + \alpha, 0)$$

Idea: Pass all three images through same network (e.g., in same batch).

In Practice

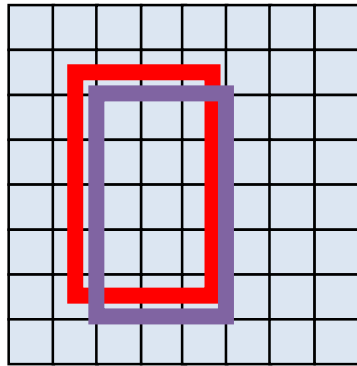
- Picking triplets crucial:
 - Lots of easy triplets (e.g., me and Shaquille O'Neal) – don't learn anything
 - Only hard triplets (e.g., me and my doppelganger) – training fails since it's too hard
 - During training, some of the worst mistakes in practice (triplets with highest loss) are often just mislabeling mistakes
- More generally called “Metric Learning”
- Lots of applications beyond face recognition

Next Time

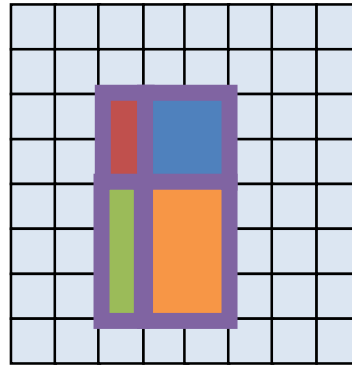
- Video

Extra Stuff

Fast R-CNN – ROI-Pool



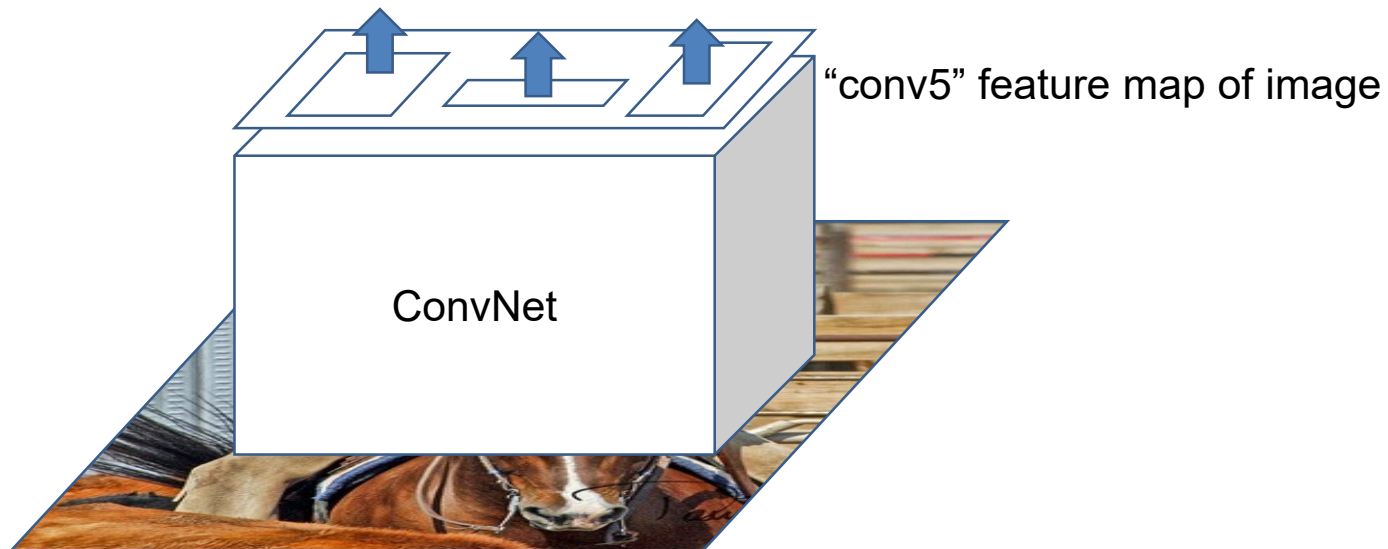
Line up



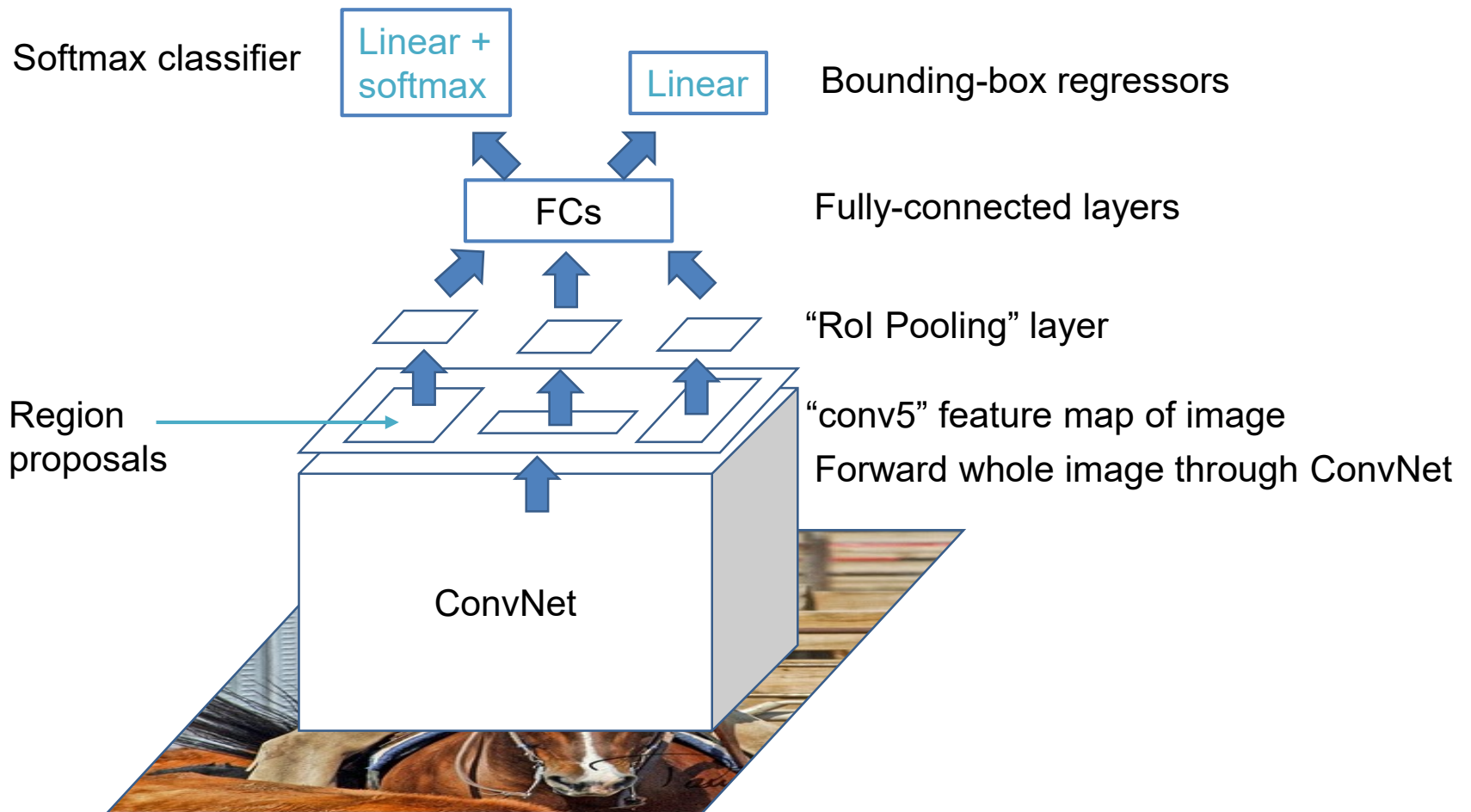
Divide



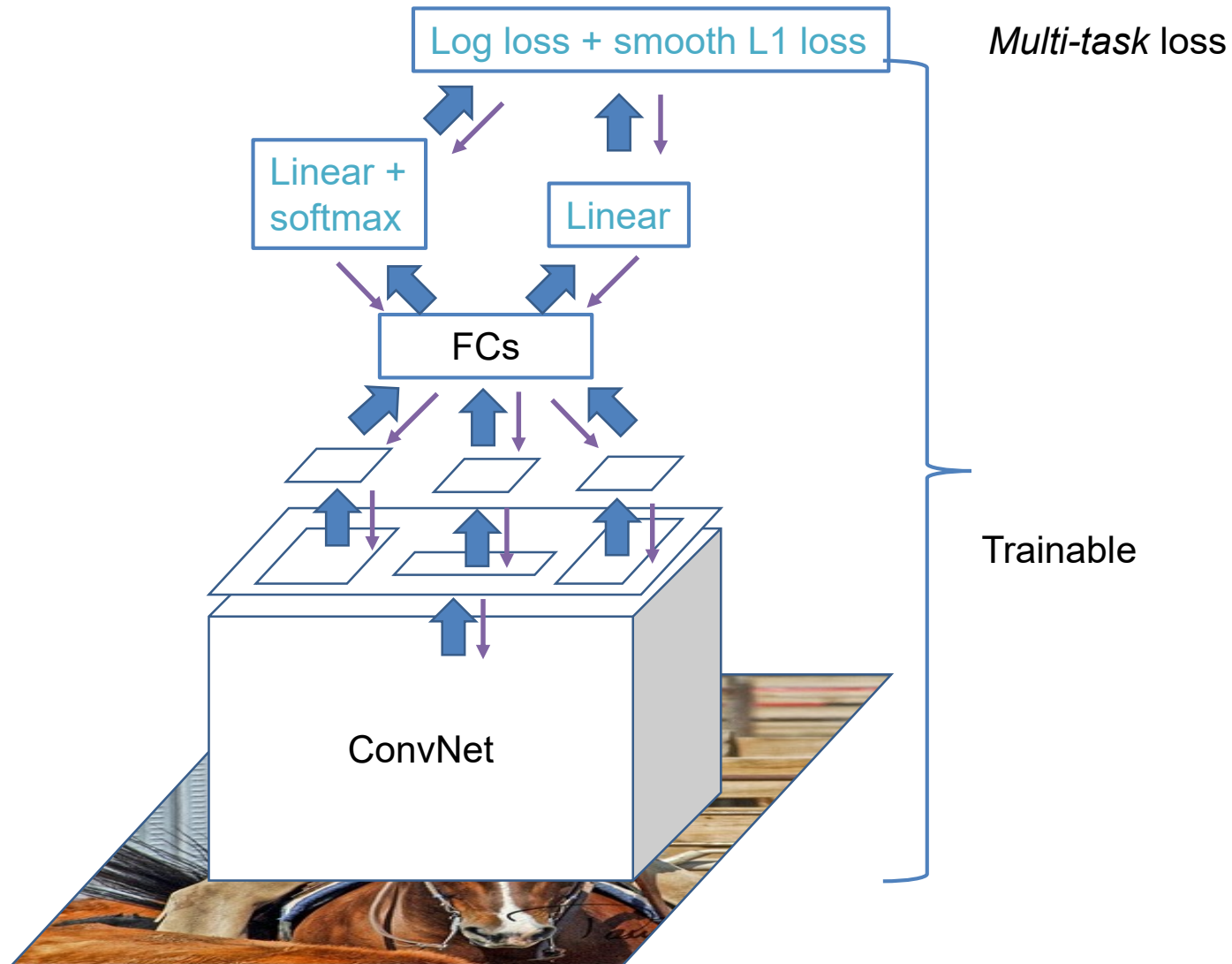
Pool



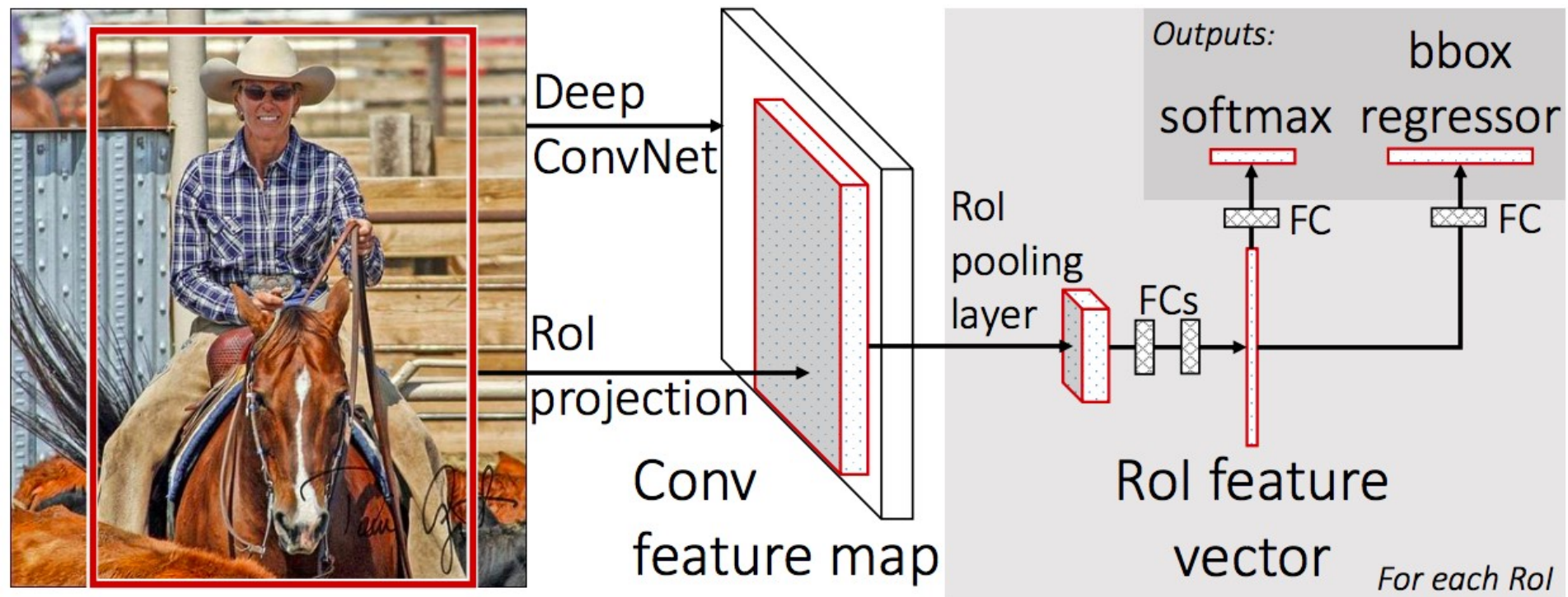
Fast R-CNN



Fast R-CNN training



Fast R-CNN: Another view



Fast R-CNN results

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

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