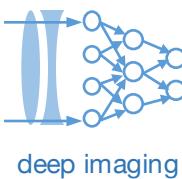


Lecture 14: Beyond classification – object detection and segmentation

Machine Learning and Imaging

BME 590L
Roarke Horstmeyer

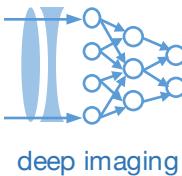


Class project details

- Can work alone or in a group (up to 4 people), required effort will scale with # of people
- Select a “base” dataset (online, or from a list I’ll make)
- Simulate parameters of a physical (imaging) system with base dataset
- Train deep neural net with simulated dataset
- Report results

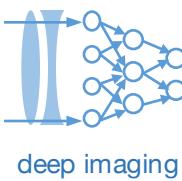
What you'll need to submit:

- 1) The project's source code
- 2) A short research-style paper (3 pages minimum, 5 pages maximum) that includes an introduction, results, a discussion section, references and at least 2 figures
- 3) A completed web template containing the main results from the research paper
- 4) An 8-minute presentation that each student will deliver to the class



Class project – what are the first steps?

1. Think about it!
2. Discuss with your friends/others in the class (feel free to use Slack!)
3. Schedule a short 15 meeting with me:
 - **Wednesday Oct. 30, 10:00am – 1:00pm**
 - **Thursday Oct. 31, 9:00am – 1:00pm**
 - **Friday Oct. 1, 9:00am – 2:00pm**
4. Start to write-up a proposal
 - General aim: 1 paragraph with specification of physical layer
 - Discussion: (a) data source(s), (b) expected simulations, (c) expected CNN, (d) quantitative analysis of physical layer/physical component (comparison, plot, etc).
 - Project proposal due date: **Thursday November 7, 2019**
 - Revised project proposal due date: **Tuesday November 12, 2019**
 - **Final Presentation: December 13, 7 pm – 10 pm**



Example project topics:

Can we design a new lens/transducer/antenna shape to improve classification of X?

What is the tradeoff between image resolution and accuracy for X (classification, segmentation, etc.)? What if we had access to n low-resolution cameras – how might we position them to get the best performance?

Can we determine an optimal set of colors to improve fluorophore distinguishability?

How does classification accuracy change with sensor bit depth, down to the 1-bit level for single-photon detectors?

If we just had a few sensors, how should we arrange them e.g. a mask to be able to predict the position of X?

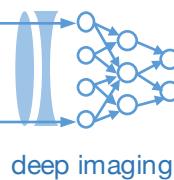
Is there some optimal shift-variant blur that we can use for a particular task?

Or, given a shift-variant blurry image, can we establish a good deconvolution using locally connected layers?

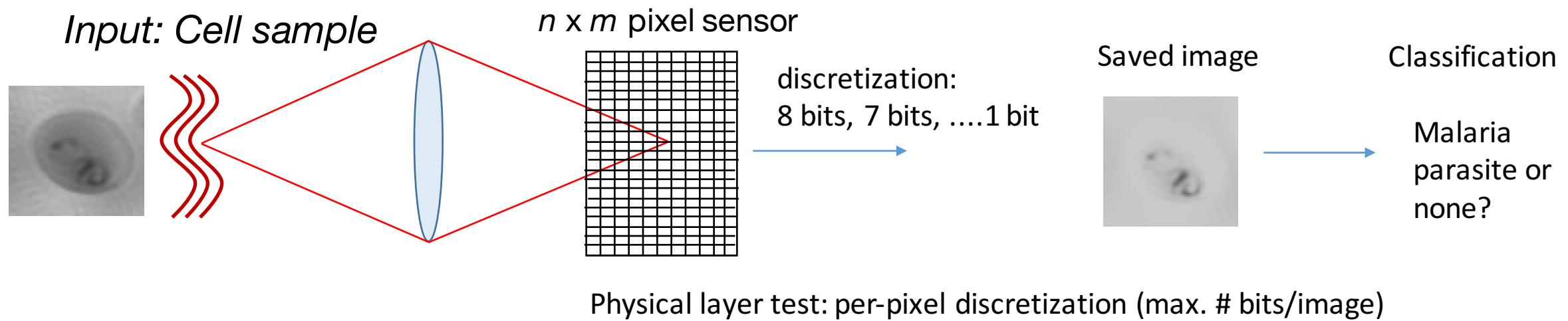
What is the optimal way to layout filters on a sensor to capture a color image for classification? Or an HDR image?

HDR image generation with filters over pixels – what is optimal design?

What if we could make a sensor with different sized pixels – how should they be laid out to achieve the best X?



How does classification accuracy change with sensor bit depth, down to the 1-bit level for single-photon detectors?

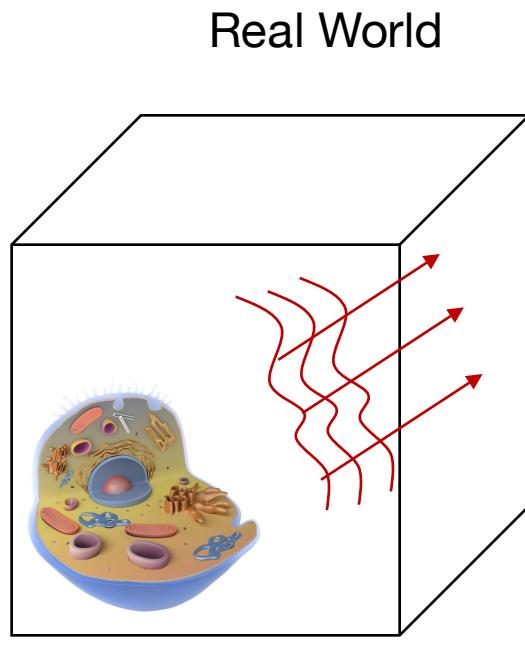
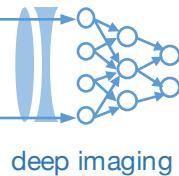


I propose to test the classification performance of a microscope as a function of sensor bit depth (i.e., image discretization). I will plot average classification test accuracy as a function of number of sensor bits from 1 bit to 8 bits. I will additionally test whether the pixel discretization value can be optimized as a physical layer parameter. I will simulate a pixel discretization value, at each pixel, by multiplying the associated raw intensity value at each pixel by a weight, and will then use the `max()` operator to set a threshold. I will examine how classification accuracy varies with this additional constraint, and will attempt to draw insights into where the network prefers to have more bits/pixel.

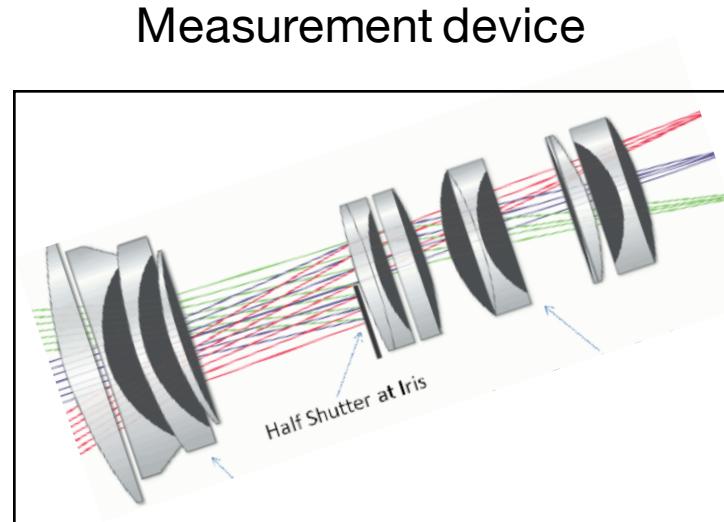
Dataset: 12,500 images of 4 types of blood cell <https://www.kaggle.com/paultimothymooney/blood-cells>

(Specify more details about simulation network, physical layer implementation and quantitative analysis)

ML+Imaging pipeline + plan



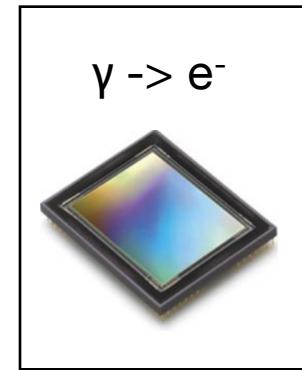
Continuous
complex fields



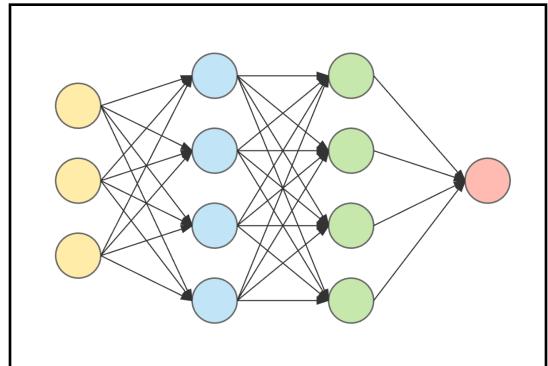
Black box transformations

- Convolution
- Fourier Transform

Digitization



Machine Learning



Linear classification

Logistic classifier

Neural networks

Convolutional NN's

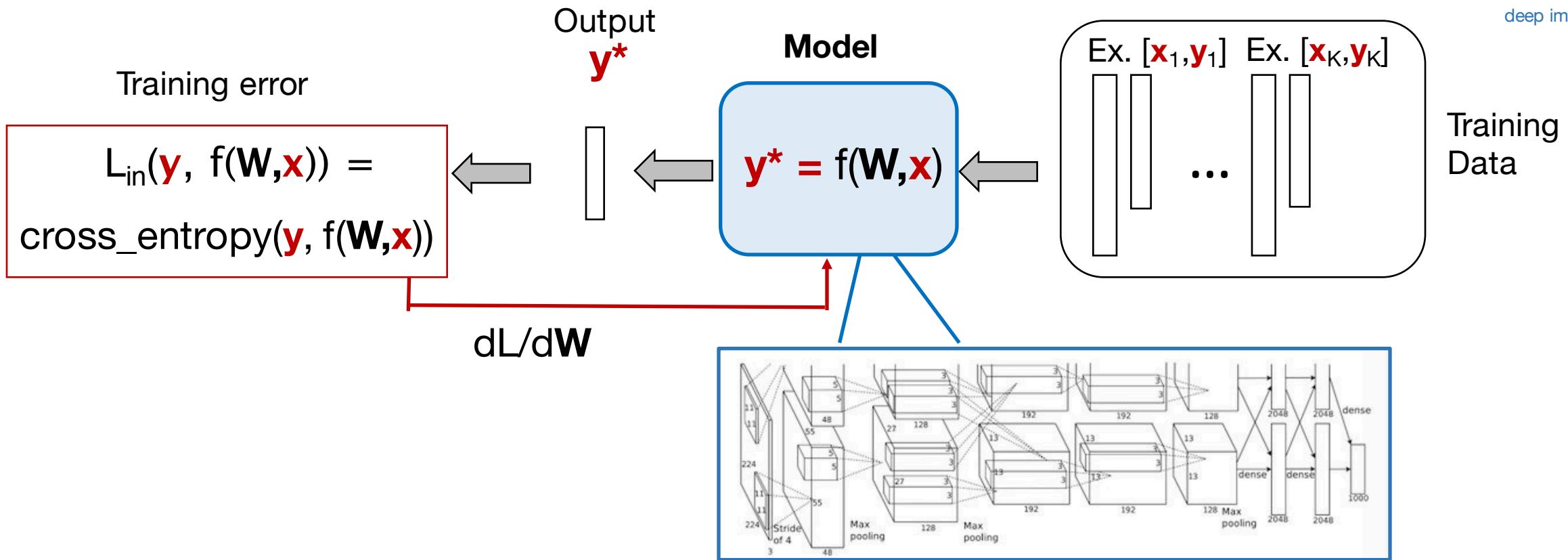
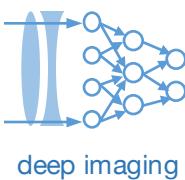
Sampling Theorem

Discrete math &
Linear algebra

← This week

Next week

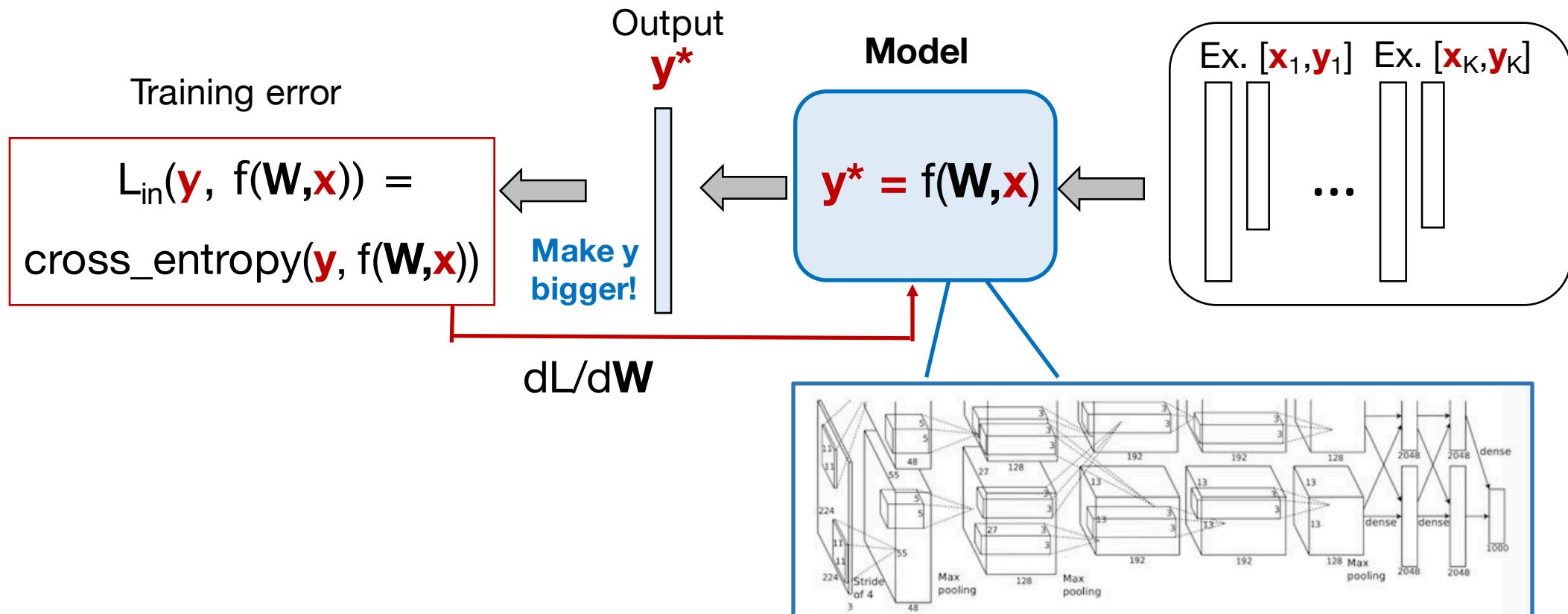
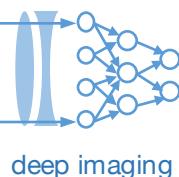
→ End of Class



Dimensional analysis for classification:

Input $x: \sim R^{1000}$

Output $y^*: \sim R^2 - R^{10}$



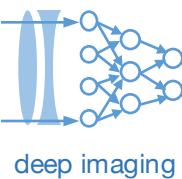
Dimensional analysis for classification:

Input $x: \sim R^{1000}$

Output $y^*: \sim R^2 - R^{10}$

This class – let's make y^* bigger!

- Object detection
- Segmentation
- Creating 3D volumes
- Better resolution



Other Computer Vision Tasks

Semantic Segmentation



GRASS, CAT,
TREE, SKY

No objects, just pixels

Classification + Localization



CAT

Single Object

Object Detection



DOG, DOG, CAT

Multiple Object

Instance Segmentation



DOG, DOG, CAT

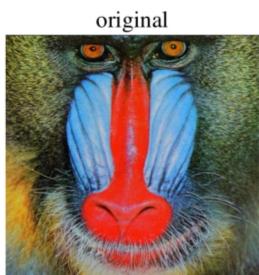
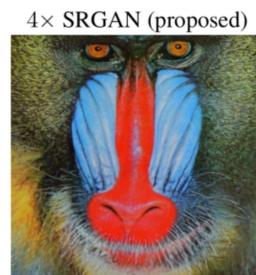
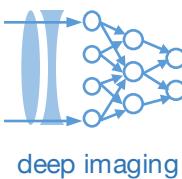
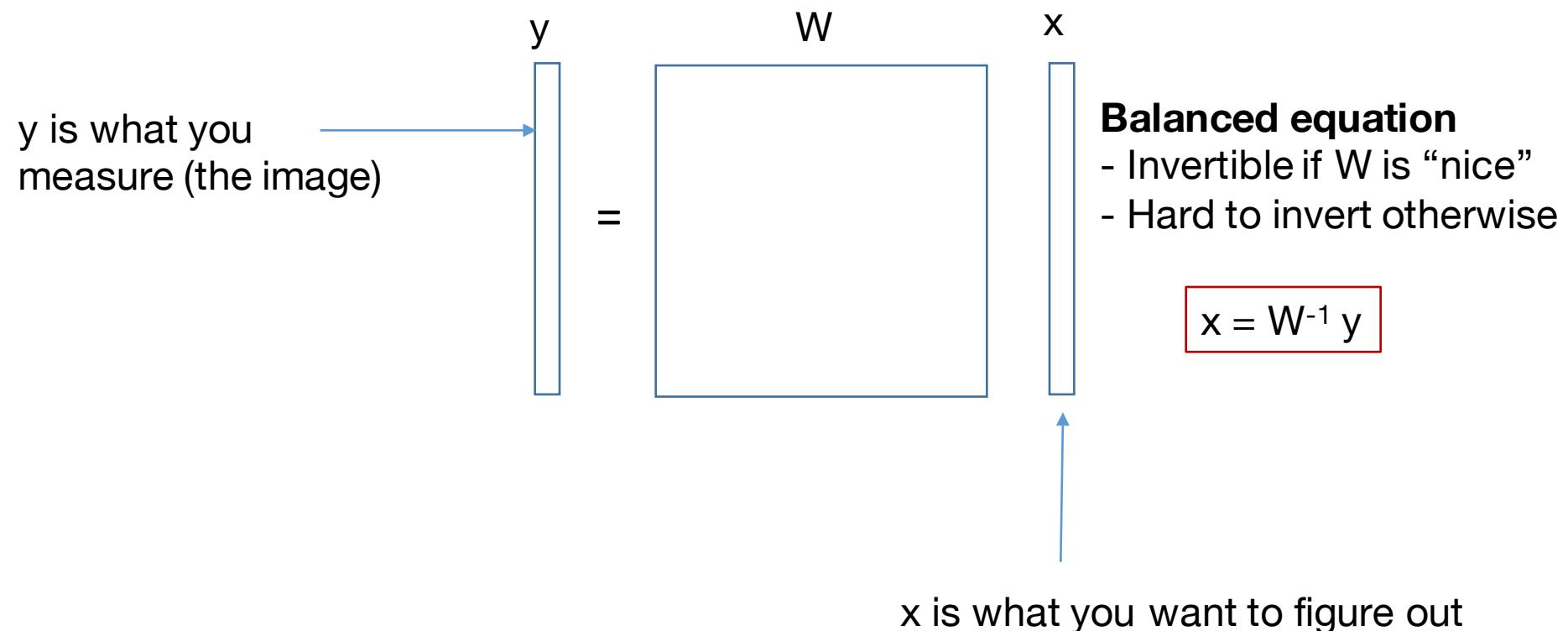


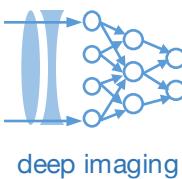
Figure 1: Super-resolved image (left) is almost indistinguishable from original (right). [4× upscaling]

[This image is CC0 public domain](#)



Over-determined, under-determined and balanced inverse equations





Over-determined, under-determined and balanced inverse equations

$$y = Wx$$

$$y = Wx$$

Balanced equation

- Invertible if W is “nice”
- Hard to invert otherwise

$$x = W^{-1}y$$

Over-determined equation

- Unique solution can exist
- If not, it's easy to get close
- Good place – more measurements than unknowns

$$x = W^+y$$

Over-determined, under-determined and balanced inverse equations

$$y = Wx$$

Over-determined equation

- Unique solution can exist
- If not, it's easy to get close
- Good place – more measurements than unknowns

$$y = Wx$$

Balanced equation

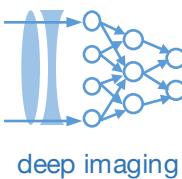
- Invertible if W is “nice”
- Hard to invert otherwise

$$x = W^{-1}y$$

$$y = Wx$$

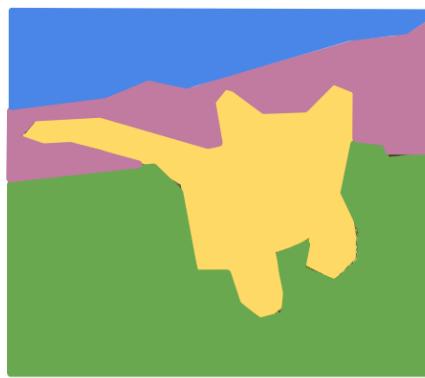
Under-determined equation

- No unique solution for x !
- Hard to invert
- **Not a good place to be**



Other Computer Vision Tasks

Semantic Segmentation



**GRASS, CAT,
TREE, SKY**

No objects, just pixels

Balanced equation

Classification + Localization



CAT

Single Object

Over-determined

Object Detection



DOG, DOG, CAT

Multiple Object

Over-determined

Instance Segmentation



DOG, DOG, CAT

This image is CC0 public domain

Over-determined

Super-resolution

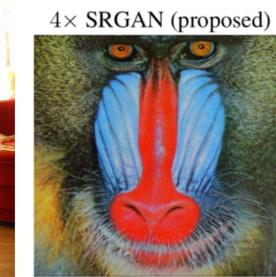
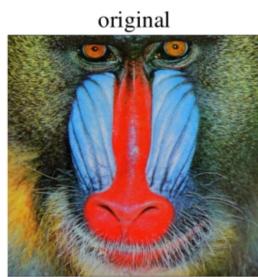
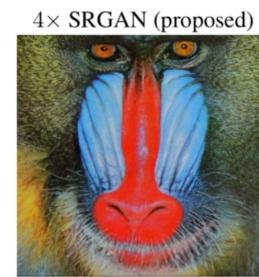
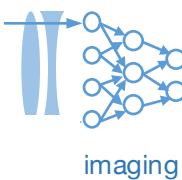


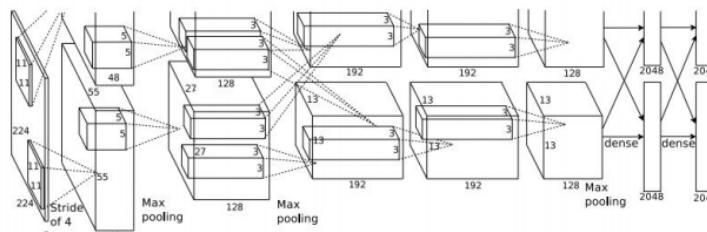
Figure 1: Super-resolved image (left) is almost indistinguishable from original (right). [4× upscaling]



Classification + Localization



This image is CC0 public domain



Treat localization as a regression problem!

Fully Connected:
4096 to 1000

Vector: Fully Connected:
4096 to 4

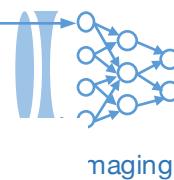
Class Scores
Cat: 0.9
Dog: 0.05
Car: 0.01
...

Box Coordinates → L2 Loss
(x, y, w, h)

Correct label:
Cat

Softmax Loss

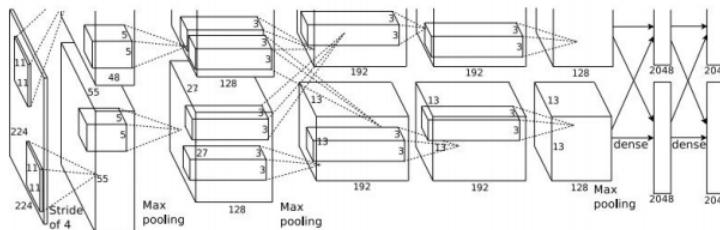
Correct box:
(x', y', w', h')



Classification + Localization



This image is CC0 public domain.



Treat localization as a regression problem!

Fully
Connected:
4096 to 1000

Vector:
4096 Fully
Connected:
4096 to 4

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

Multitask Loss

Box
Coordinates → L2 Loss
(x, y, w, h)

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

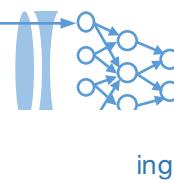
Correct label:
Cat

Softmax
Loss

+

Loss

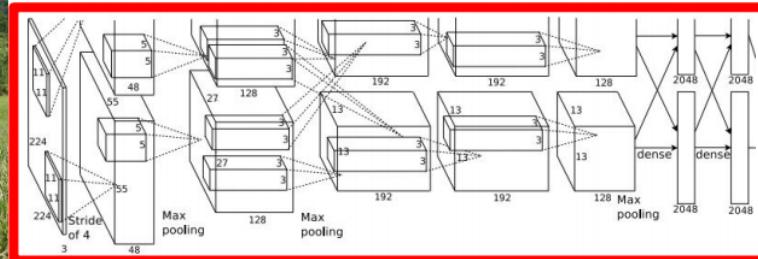
Stanford CS231n - <http://cs231n.stanford.edu>



Classification + Localization



This image is CC0 public domain



Often pretrained on ImageNet
(Transfer learning)

Treat localization as a
regression problem!

Fully Connected: 4096 to 1000

Vector: 4096
Fully Connected: 4096 to 4

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

Box Coordinates \rightarrow L2 Loss
(x, y, w, h)

Correct label:
Cat

Softmax Loss

+

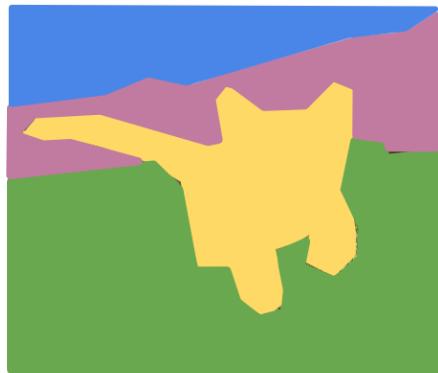
↑

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

Stanford CS231n - <http://cs231n.stanford.edu>

Other Computer Vision Tasks

Semantic Segmentation



**GRASS, CAT,
TREE, SKY**

No objects, just pixels

Classification + Localization



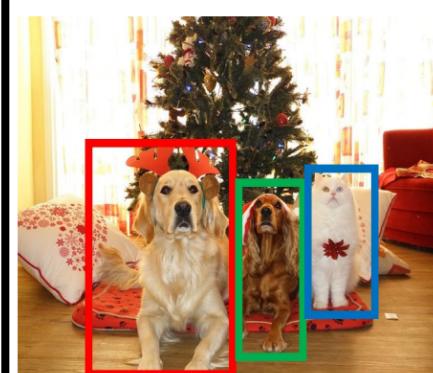
CAT

Single Object

**Balanced
equation**

**Over-
determined**

Object Detection



DOG, DOG, CAT

Multiple Object

**Over-
determined**

Instance Segmentation

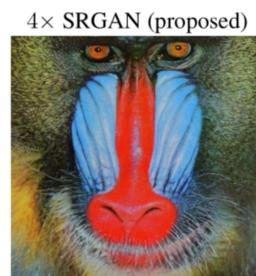


DOG, DOG, CAT

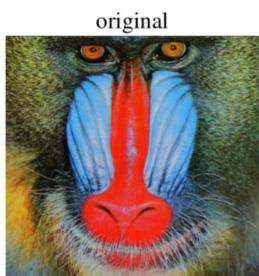
This image is CC0 public domain

**Over-
determined**

Super-resolution



4× SRGAN (proposed)



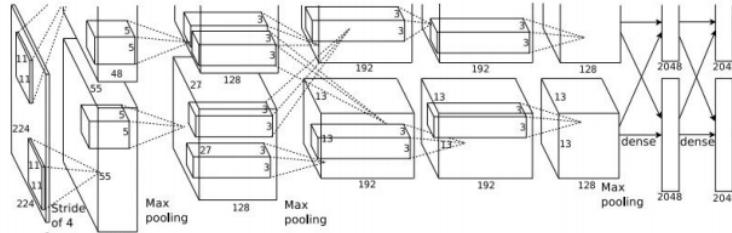
original

Figure 1: Super-resolved image (left) is almost indistinguishable from original (right). [4× upscaling]

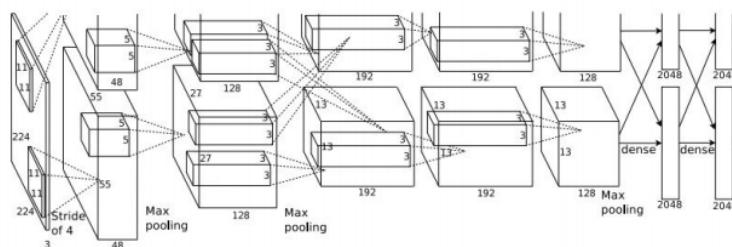


Object Detection as Regression?

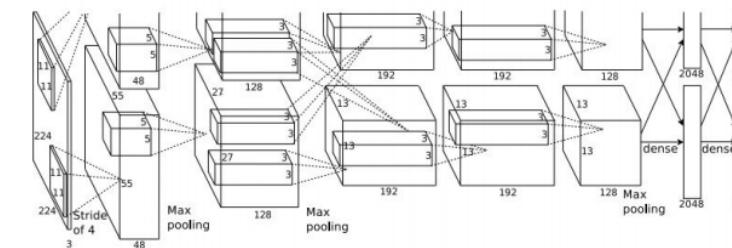
Each image needs a
different number of outputs!



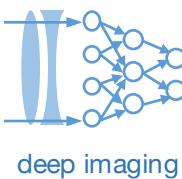
CAT: (x, y, w, h) 4 numbers



DOG: (x, y, w, h)
DOG: (x, y, w, h) 16 numbers
CAT: (x, y, w, h)

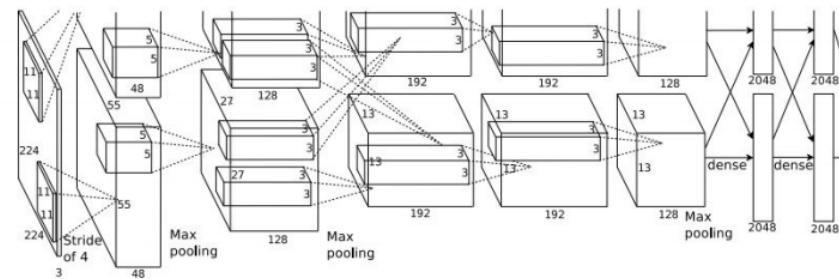


DUCK: (x, y, w, h) Many
DUCK: (x, y, w, h) numbers!
....



Object Detection as Classification: Sliding Window

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

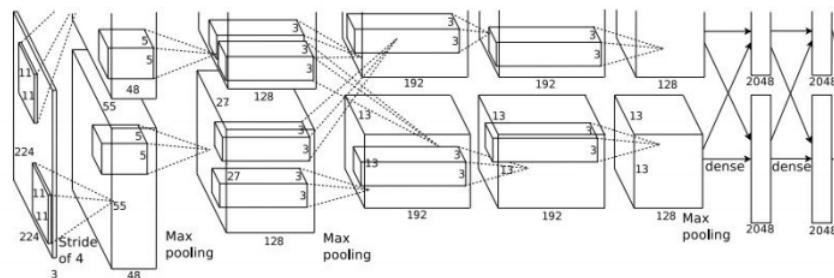


Dog? NO
Cat? NO
Background? YES

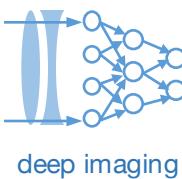
Stanford CS231n - <http://cs231n.stanford.edu>

Object Detection as Classification: Sliding Window

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

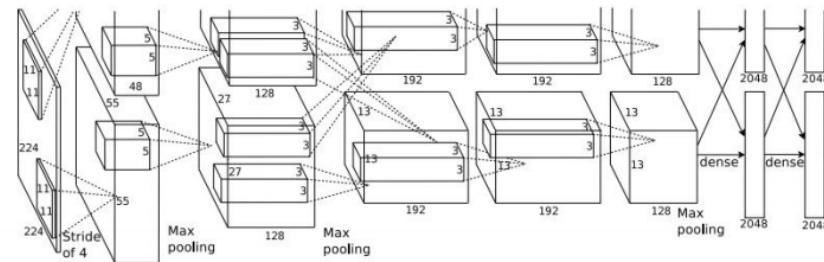
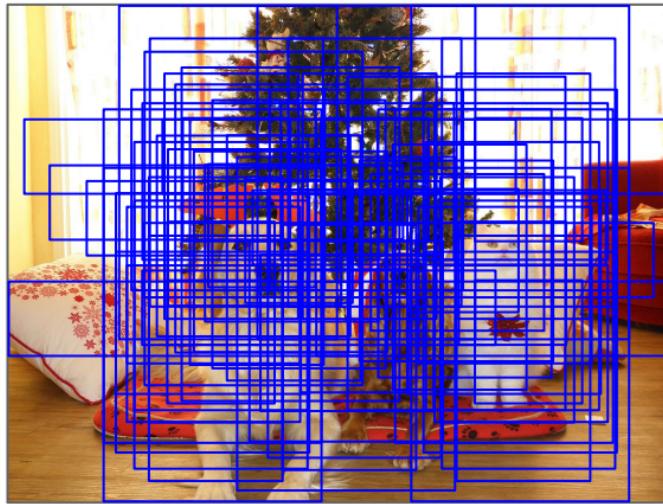


Dog? YES
Cat? NO
Background? NO



Object Detection as Classification: 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

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

Stanford CS231n - <http://cs231n.stanford.edu>

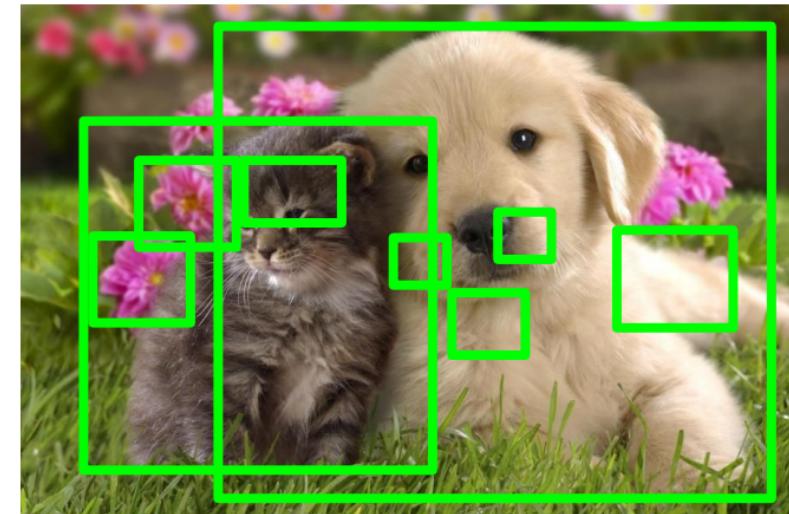
Solution: First apply a fixed ROI scheme to pull out “blobs” of interest



(Image source: van de Sande et al. ICCV'11)

Region Proposals / Selective Search

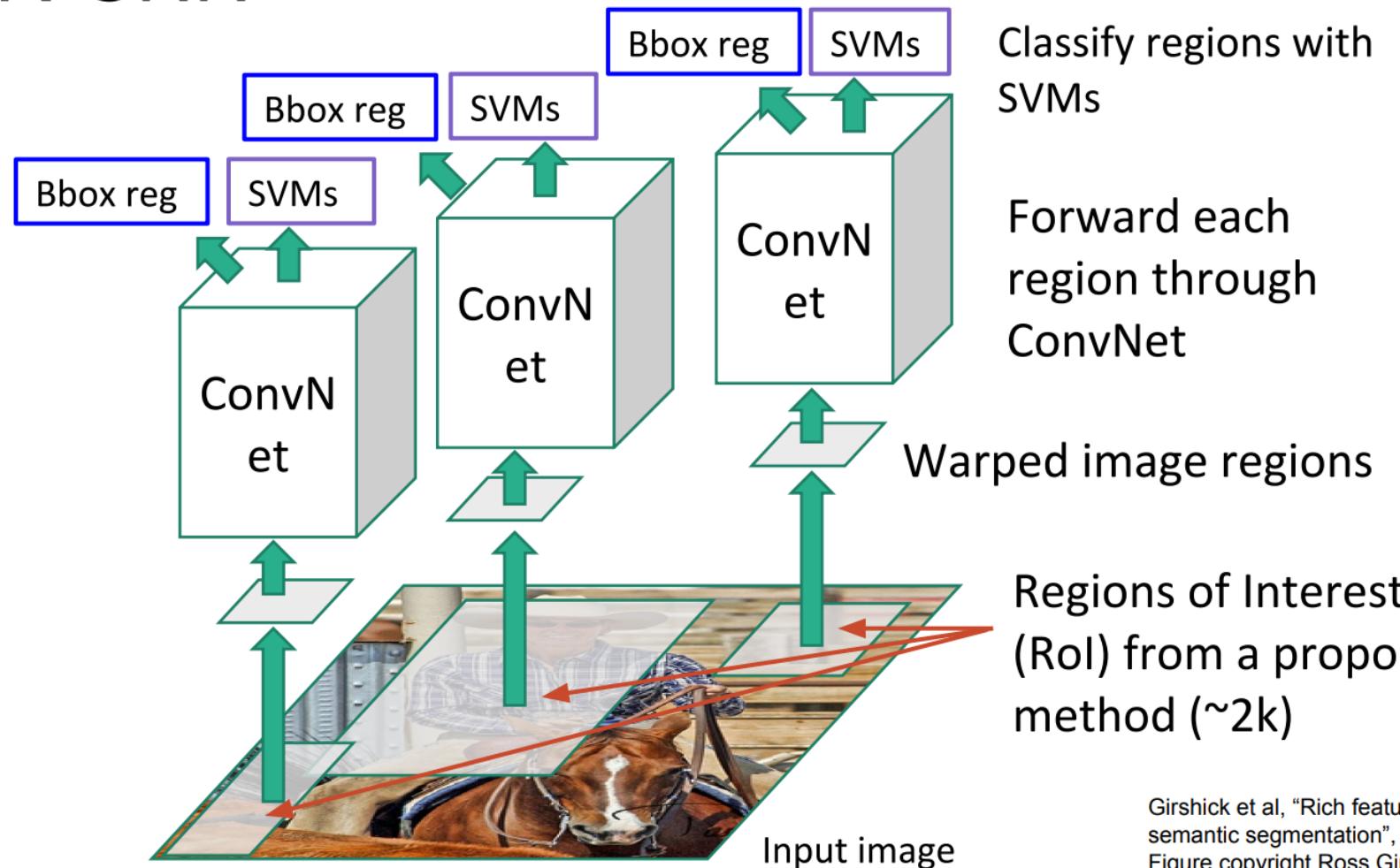
- Find “blobby” image regions that are likely to contain objects
- Relatively fast to run; e.g. Selective Search gives 2000 region proposals in a few seconds on CPU



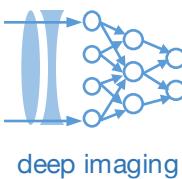
Alexe et al, "Measuring the objectness of image windows", TPAMI 2012
Uijlings et al, "Selective Search for Object Recognition", IJCV 2013
Cheng et al, "BING: Binarized normed gradients for objectness estimation at 300fps", CVPR 2014
Zitnick and Dollar, "Edge boxes: Locating object proposals from edges", ECCV 2014

Note: Training dataset has marked boxes, so don't necessarily need to do selective search for training, just evaluation/testing

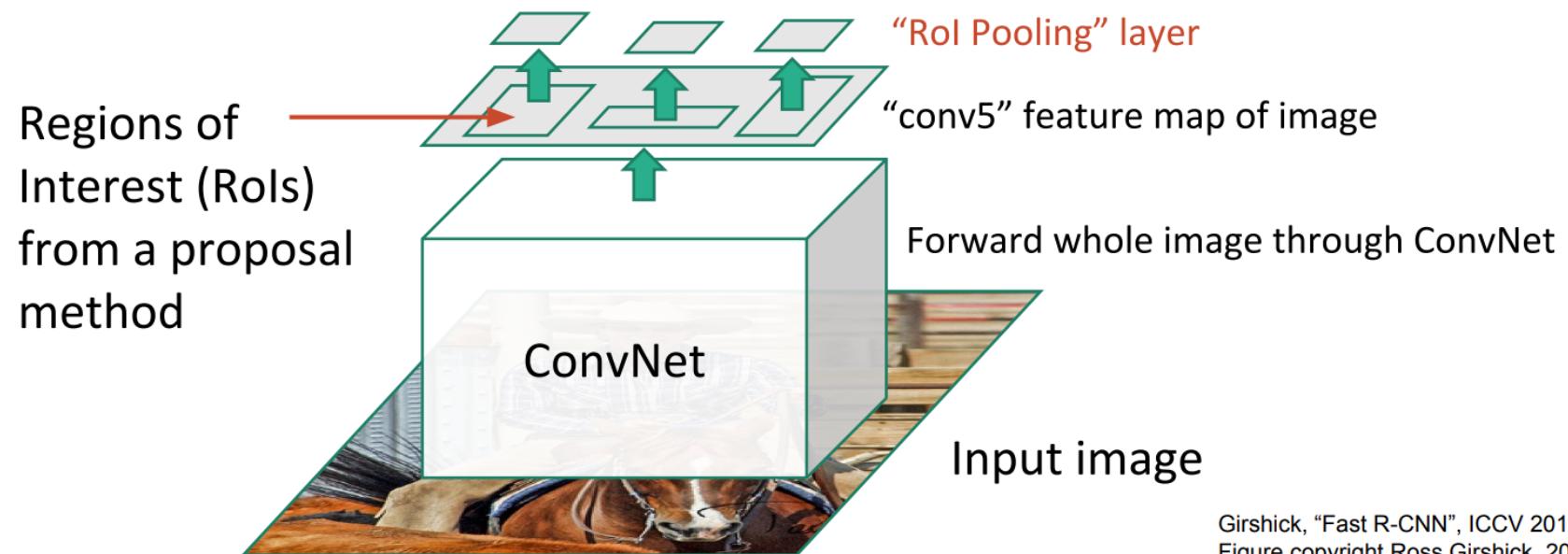
R-CNN



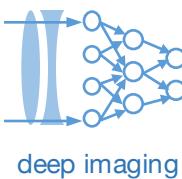
Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.
 Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.



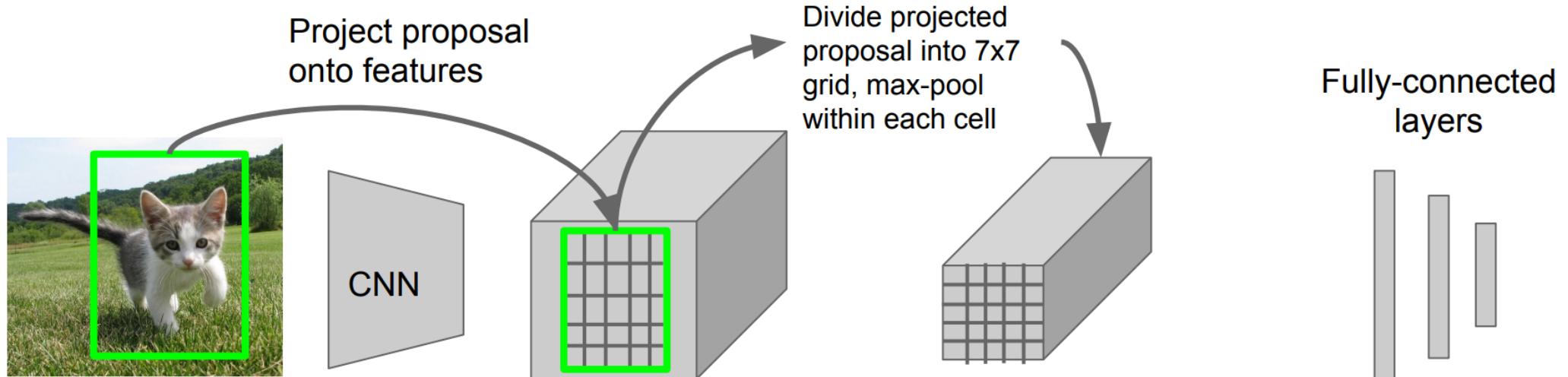
Fast R-CNN



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



Fast R-CNN: ROI Pooling



Hi-res input image:
 $3 \times 640 \times 480$
with region proposal

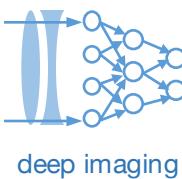
Hi-res conv features:
 $512 \times 20 \times 15$;
Projected region
proposal is e.g.
 $512 \times 18 \times 8$
(varies per proposal)

ROI conv features:
 $512 \times 7 \times 7$
for region proposal

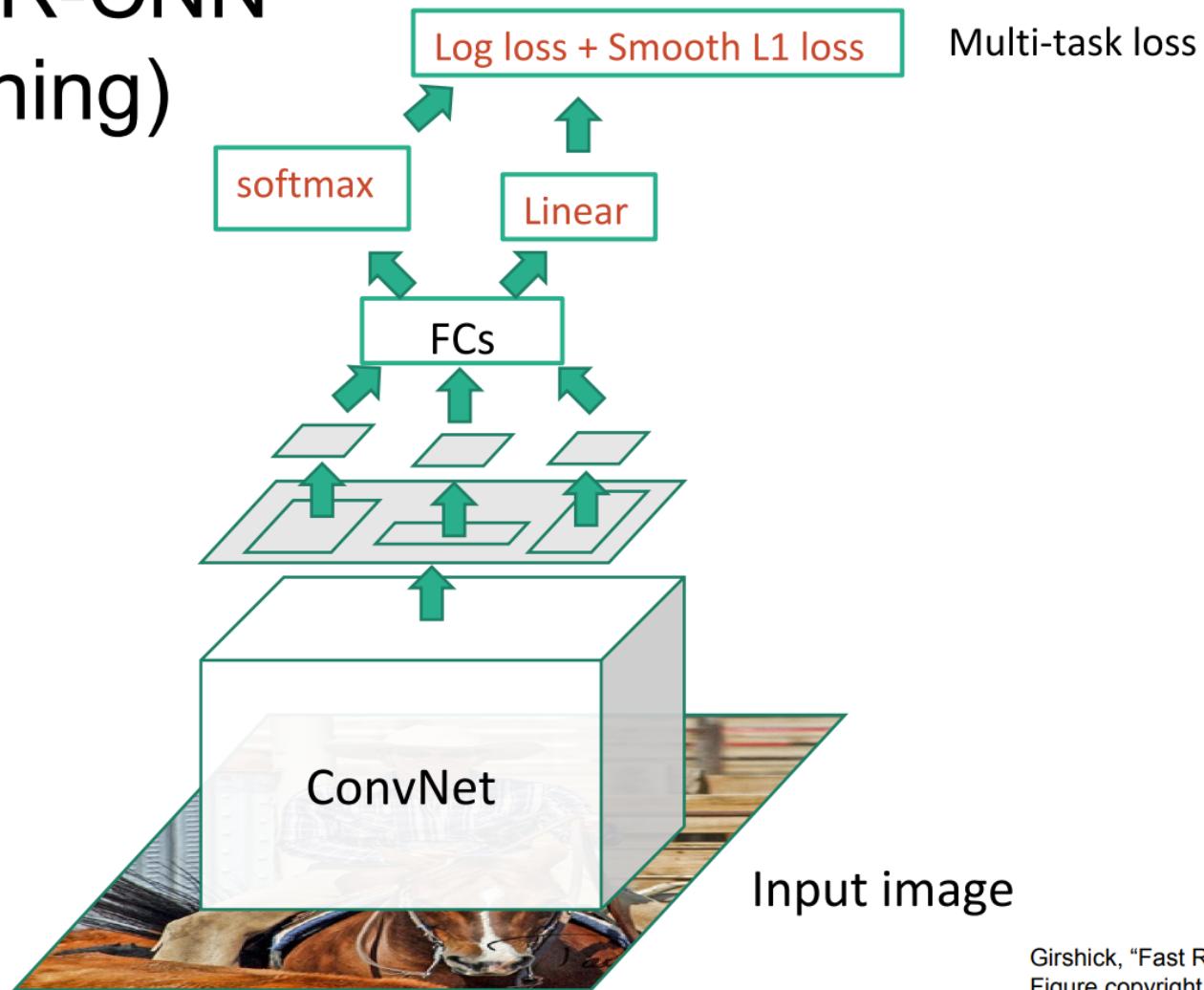
Fully-connected layers expect
low-res conv features:
 $512 \times 7 \times 7$

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

Stanford CS231n - <http://cs231n.stanford.edu>



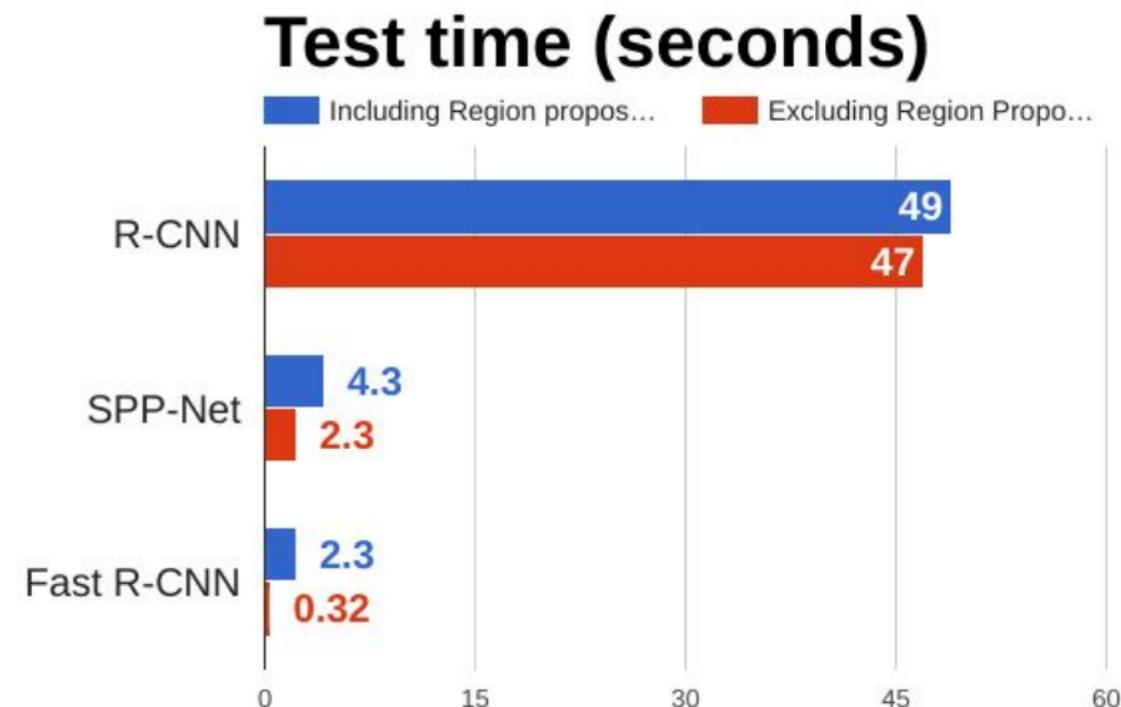
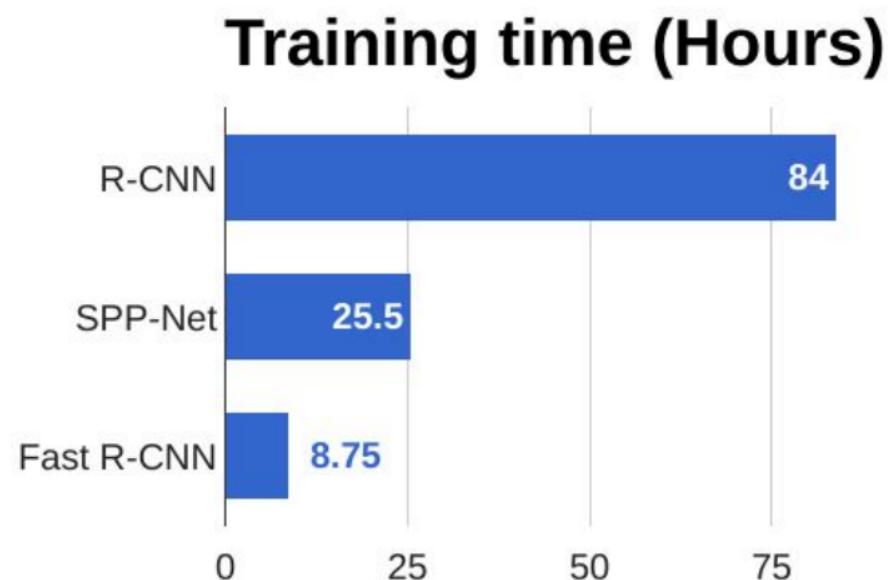
Fast R-CNN (Training)



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

Stanford CS231n - <http://cs231n.stanford.edu>

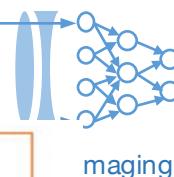
R-CNN vs SPP vs Fast R-CNN



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



Faster R-CNN:

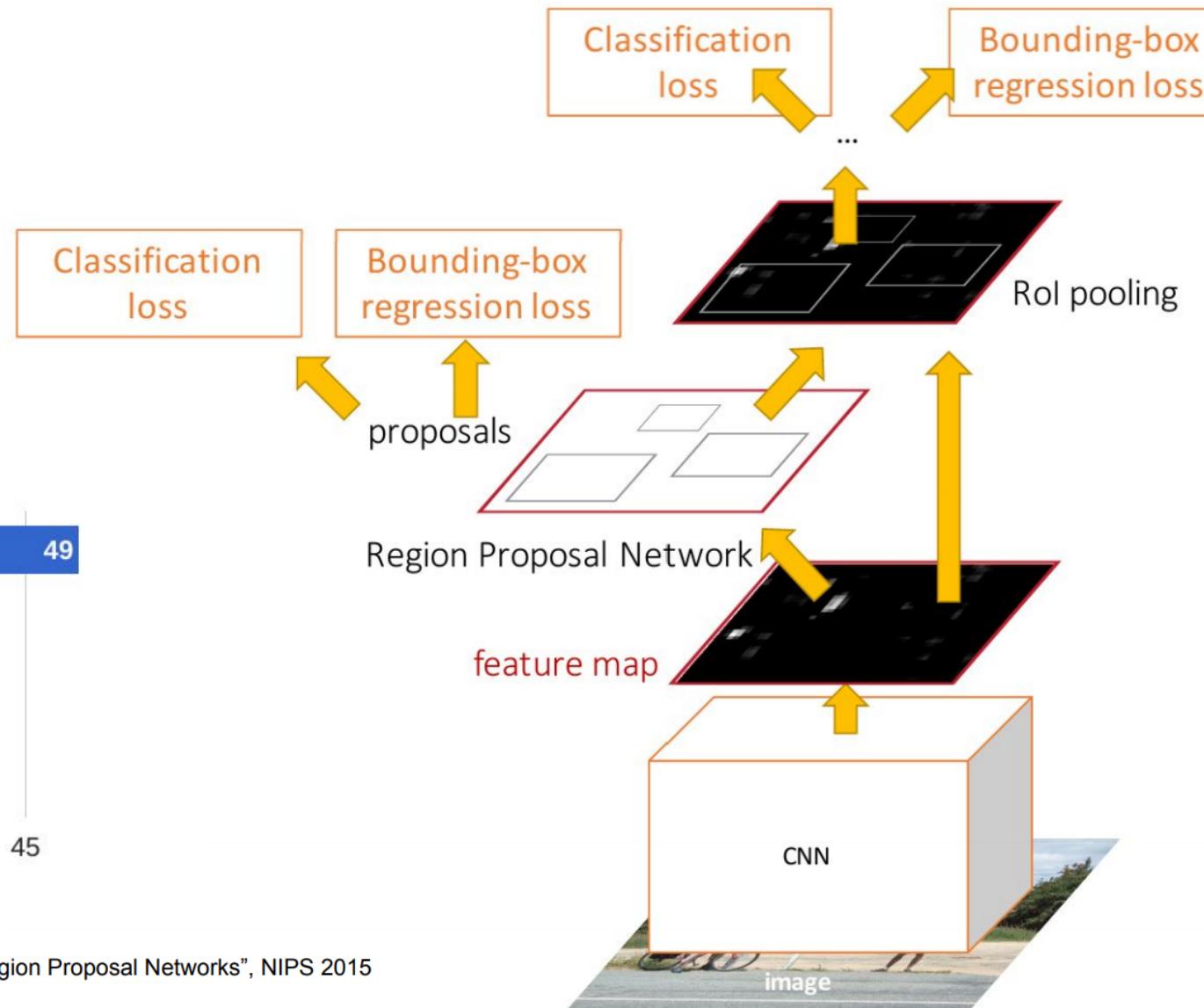
Make CNN do proposals!

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

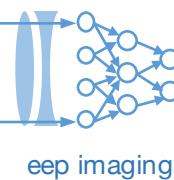
R-CNN Test-Time Speed



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



Stanford CS231n - <http://cs231n.stanford.edu>

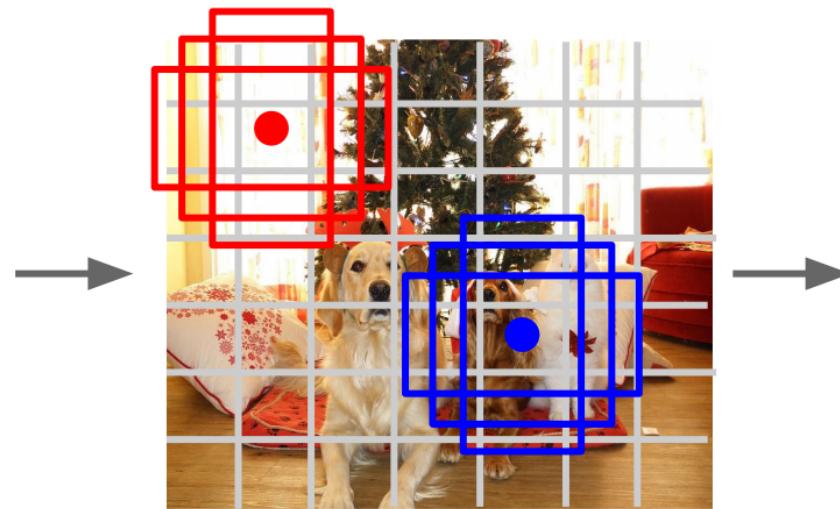


Detection without Proposals: YOLO / SSD

Go from input image to tensor of scores with one big convolutional network!



Input image
 $3 \times H \times W$



Divide image into grid
 7×7

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

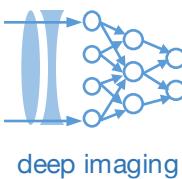
Within each grid cell:

- 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)

Output:
 $7 \times 7 \times (5 * B + C)$

Redmon et al, "You Only Look Once:
Unified, Real-Time Object Detection", CVPR 2016
Liu et al, "SSD: Single-Shot MultiBox Detector", ECCV 2016

Stanford CS231n - <http://cs231n.stanford.edu>



Object Detection: Impact of Deep Learning

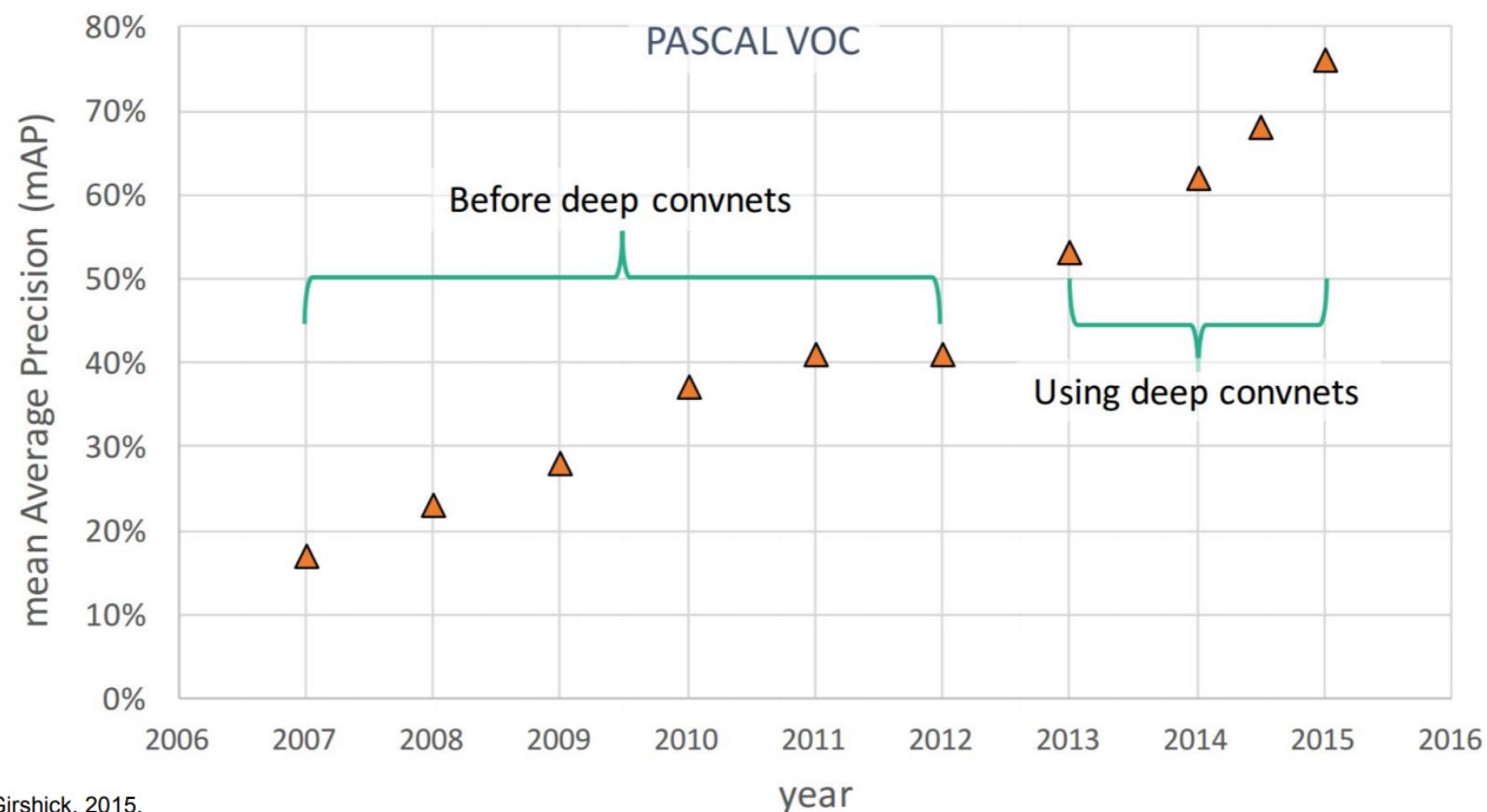


Figure copyright Ross Girshick, 2015.
Reproduced with permission.

Stanford CS231n - <http://cs231n.stanford.edu>