

# Lecture 22

## Deep Learning:

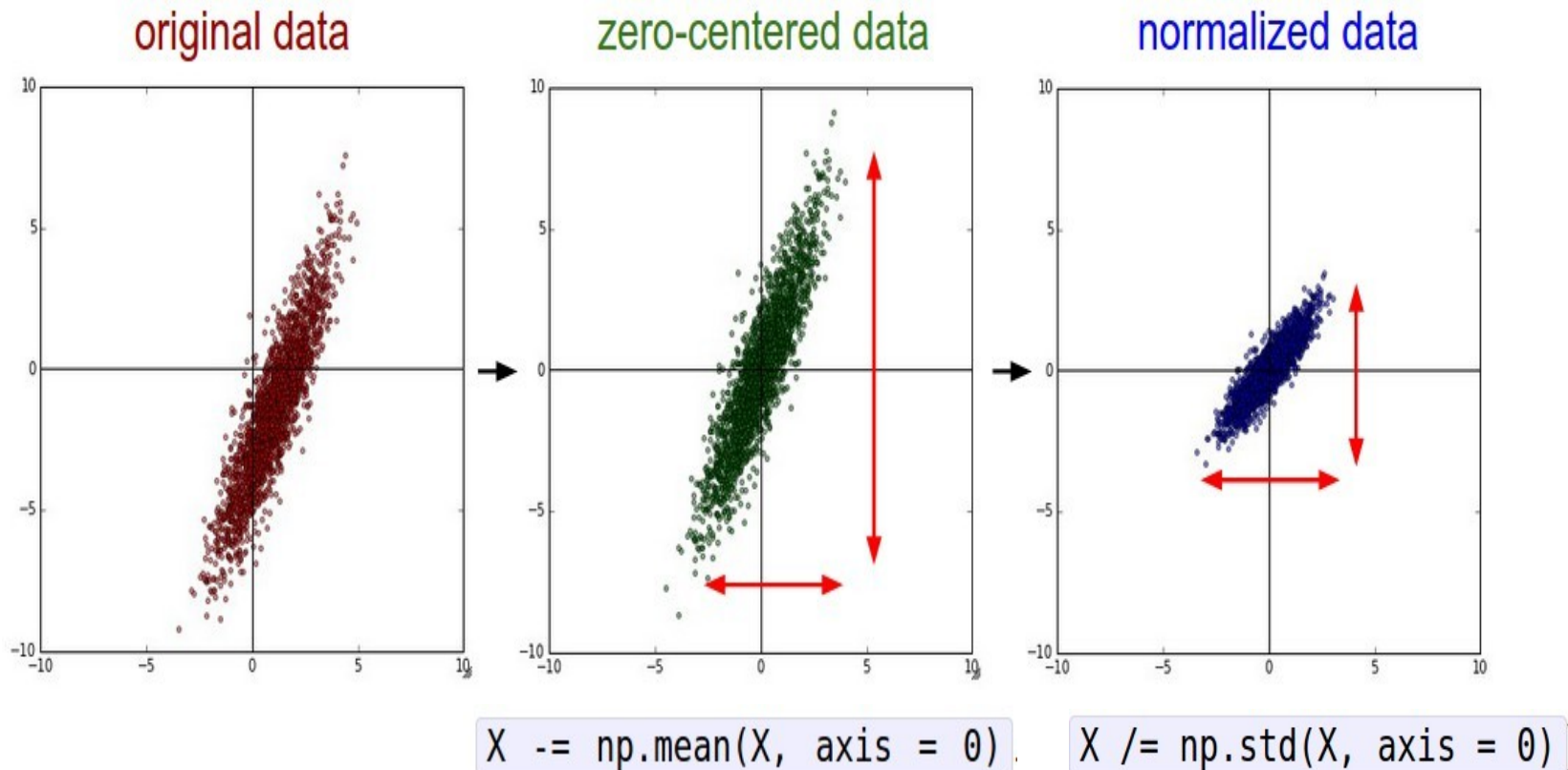
# Object Detection and Segmentation

COS 429: Computer Vision



Thanks: most of these slides shamelessly adapted from  
Stanford CS231n: Convolutional Neural Networks for Visual Recognition  
Fei-Fei Li, Andrej Karpathy, Justin Johnson  
<http://cs231n.stanford.edu/>

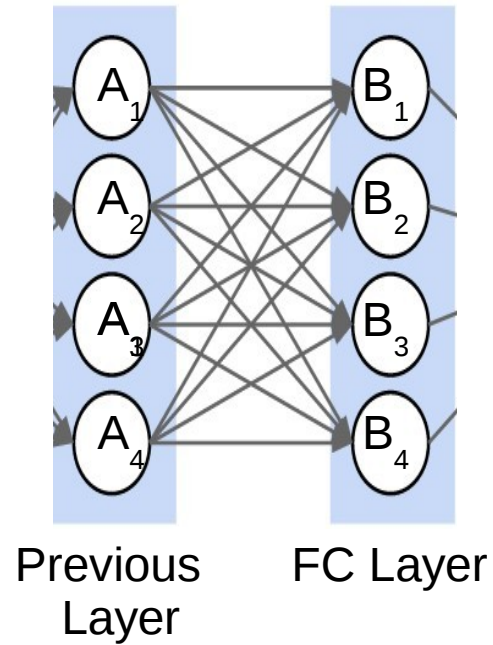
# Step 1: Preprocess the data



(Assume  $X$  [NxD] is data matrix,  
each example in a row)

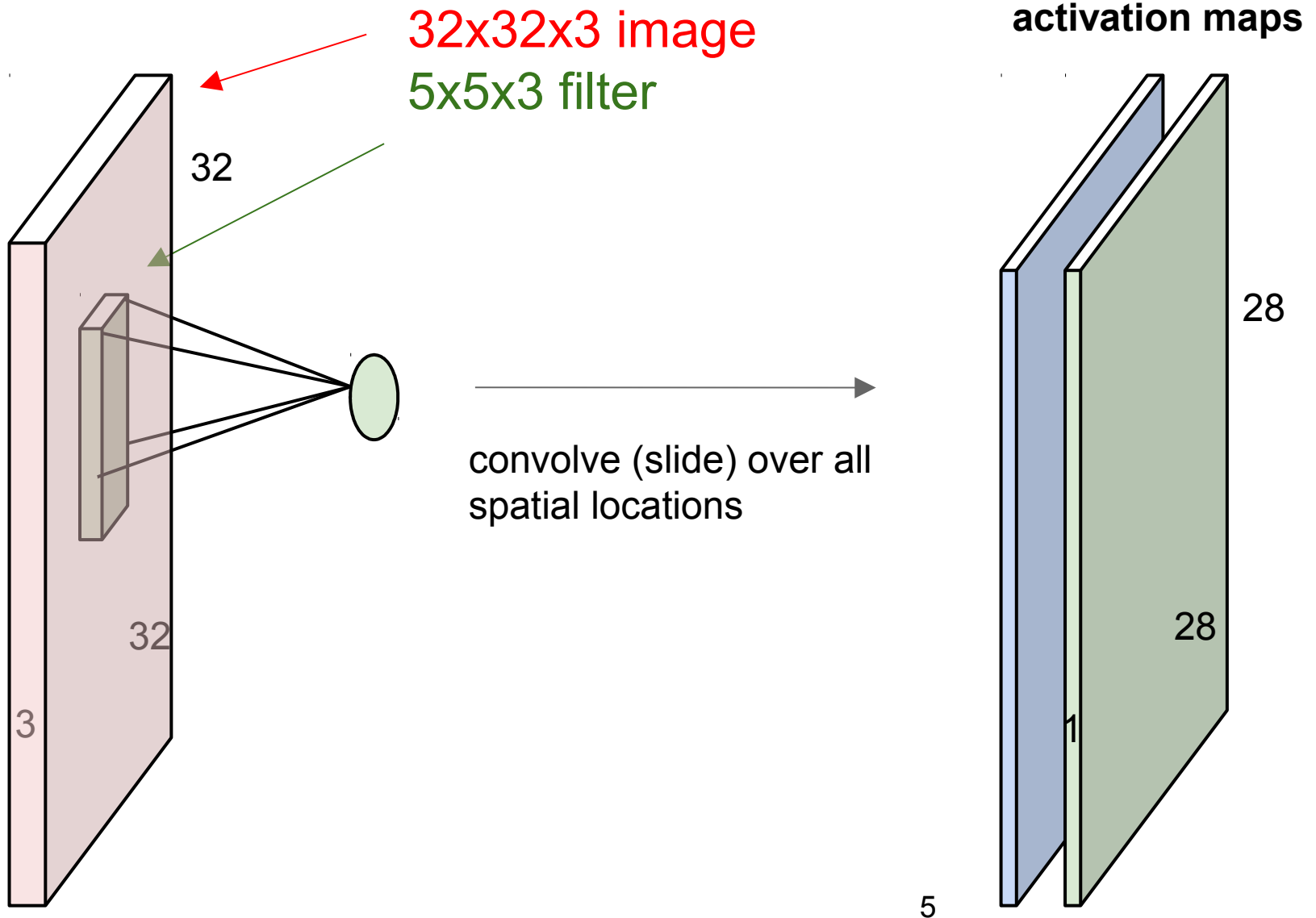


# Fully Connected layer



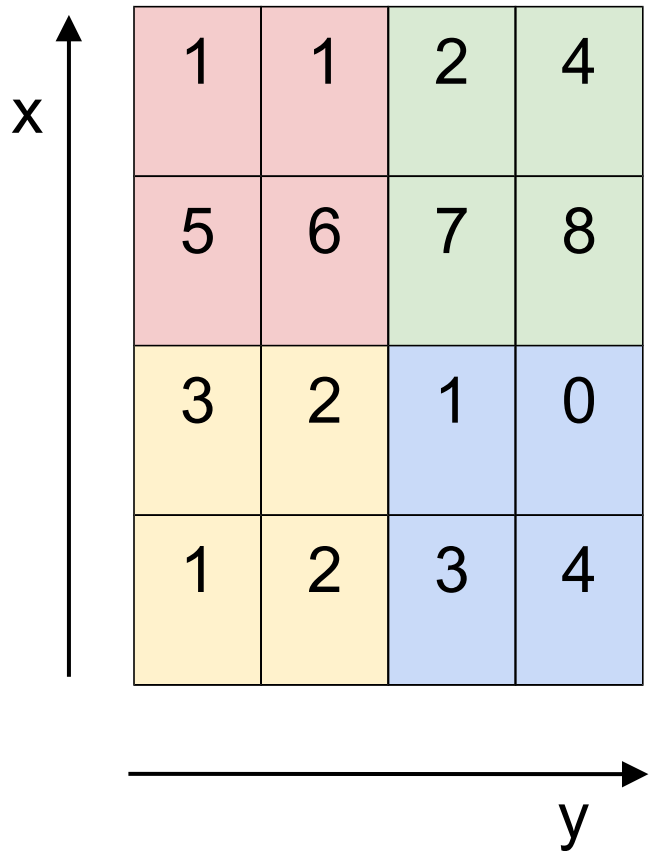
$$B_j = \sum_i (W_{ij} * A_i) + b_j$$

# Convolution Layer

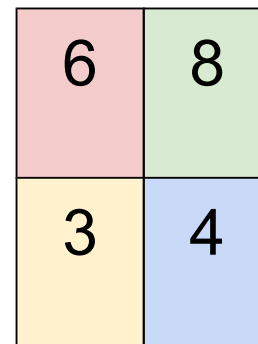


# MAX POOLING

Single depth slice



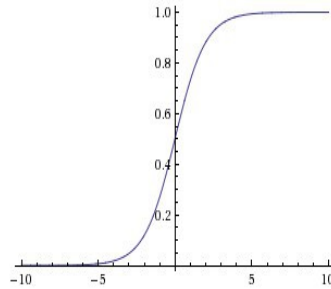
max pool with 2x2 filters  
and stride 2



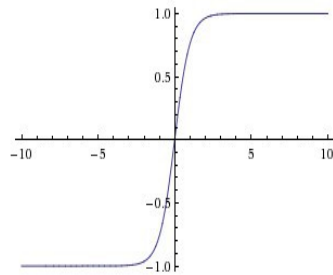
# Activation Layer

## Sigmoid

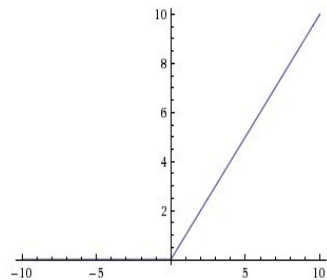
$$\sigma(x) = 1/(1 + e^{-x})$$



## tanh tanh(x)

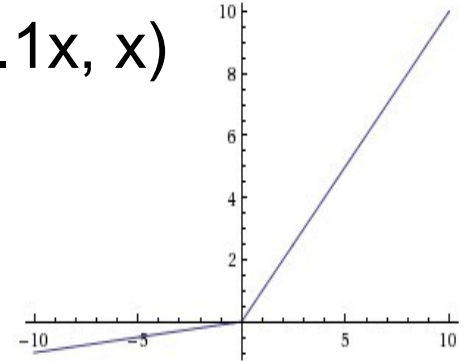


## ReLU max(0,x)



## Leaky ReLU

$$\max(0.1x, x)$$

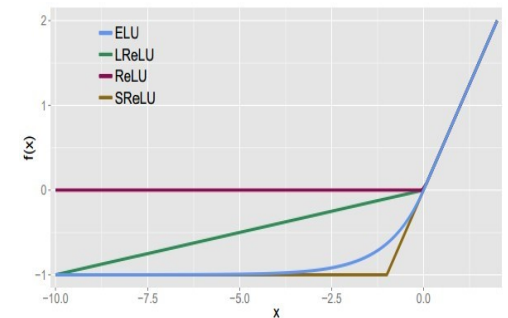


## Maxout

$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

## ELU

$$f(x) = \begin{cases} x & \text{if } x > 0 \\ \alpha (\exp(x) - 1) & \text{if } x \leq 0 \end{cases}$$



**Input:** Values of  $x$  over a mini-batch:  $\mathcal{B} = \{x_{1\dots m}\}$ ;  
Parameters to be learned:  $\gamma, \beta$

**Output:**  $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{ mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{ mini-batch variance}$$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{ normalize}$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{ scale and shift}$$

- Improves gradient flow through the network
- Allows higher learning rates
- Reduces the strong dependence on initialization
- Acts as a form of regularization in a funny way, and slightly reduces the need for dropout, maybe

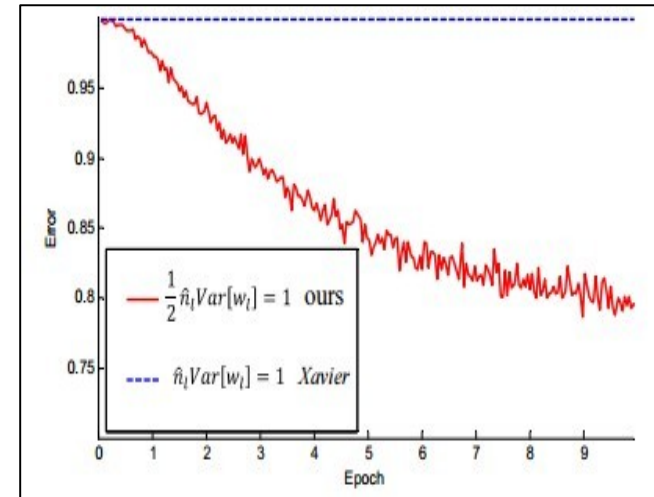
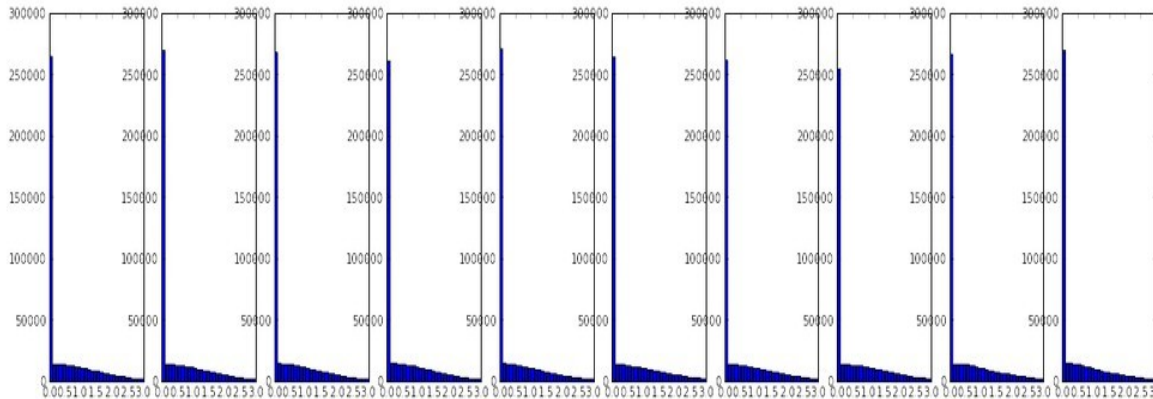
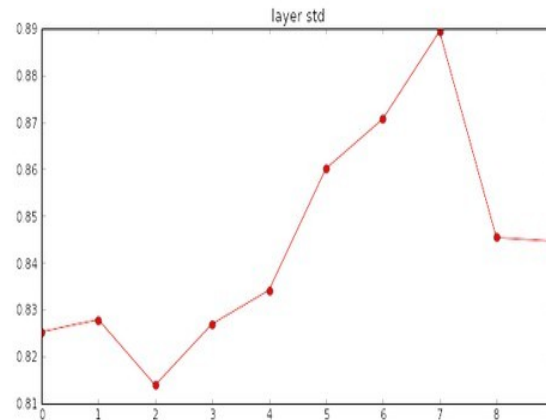
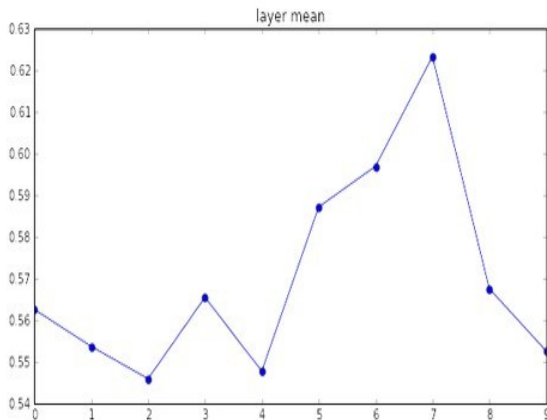


# Weight Initialization

input layer had mean 0.000501 and std 0.999444  
 hidden layer 1 had mean 0.562488 and std 0.825232  
 hidden layer 2 had mean 0.553614 and std 0.827835  
 hidden layer 3 had mean 0.545867 and std 0.813855  
 hidden layer 4 had mean 0.565396 and std 0.826902  
 hidden layer 5 had mean 0.547678 and std 0.834092  
 hidden layer 6 had mean 0.587103 and std 0.860035  
 hidden layer 7 had mean 0.596867 and std 0.870610  
 hidden layer 8 had mean 0.623214 and std 0.889348  
 hidden layer 9 had mean 0.567498 and std 0.845357  
 hidden layer 10 had mean 0.552531 and std 0.844523

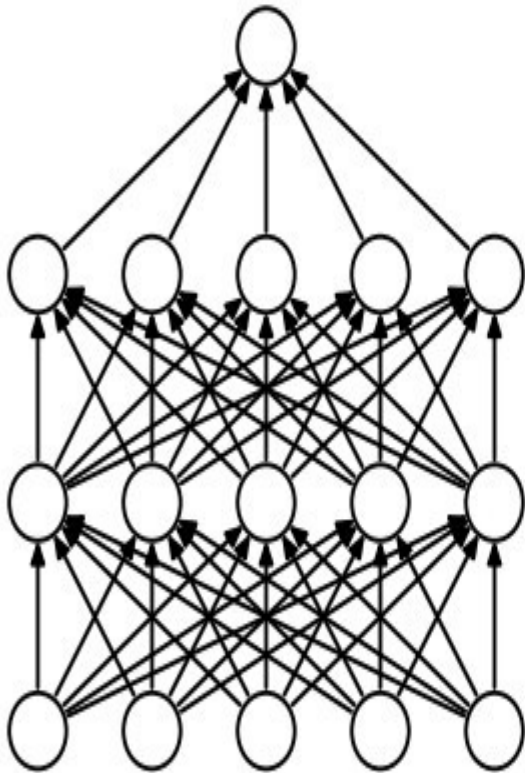
He et al., 2015  
 (note additional /2)

```
W = np.random.randn(fan_in, fan_out) / np.sqrt(fan_in/2) # layer initialization
```

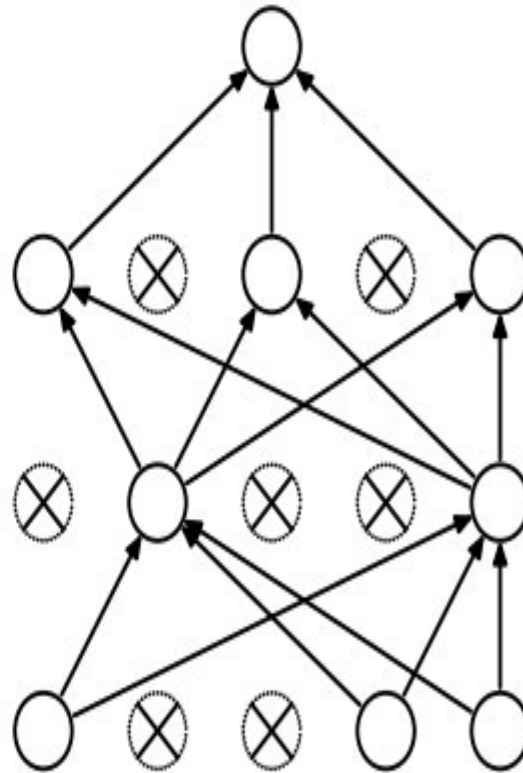


# Regularization: Dropout

“randomly set some neurons to zero in the forward pass”



(a) Standard Neural Net



(b) After applying dropout.

*[Srivastava et al., 2014]*

# Regularization: **DisturbLabel**

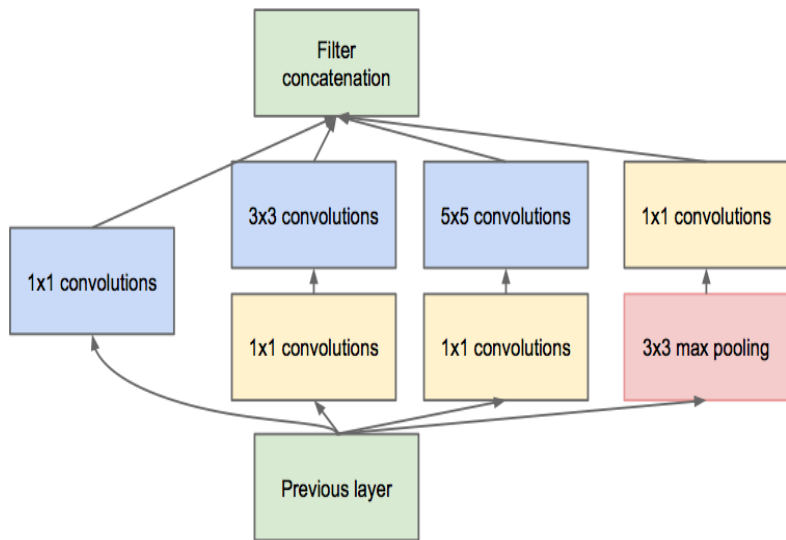
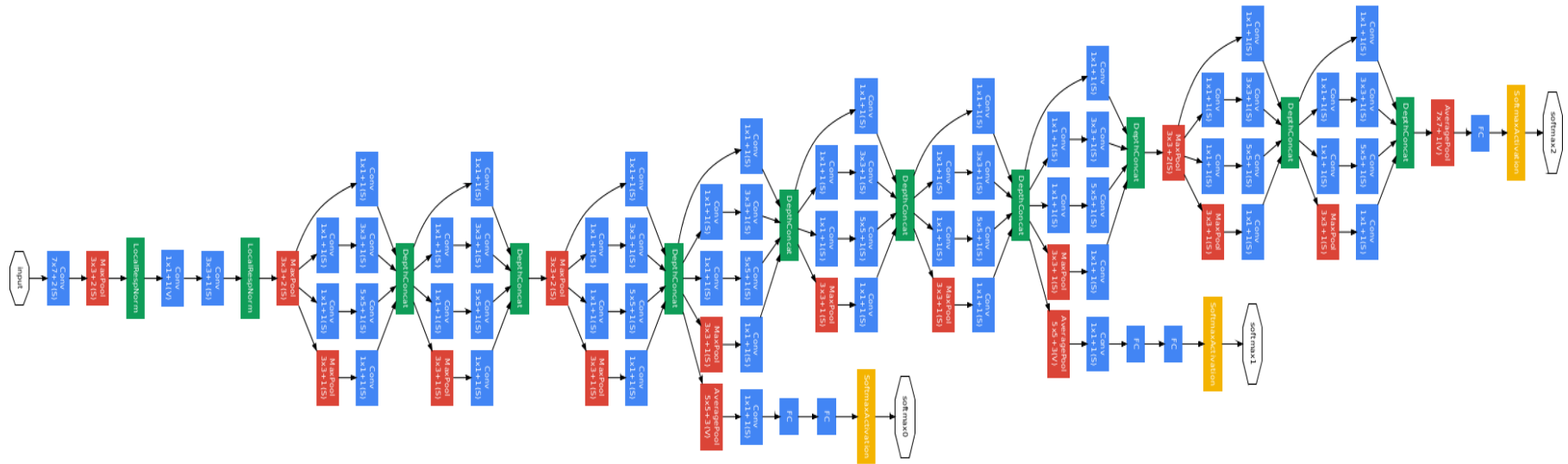
---

“randomly change ground truth label of small % of examples”

=>Improves generalization, reduces need for dropout

# Case Study: GoogLeNet

[Szegedy et al., 2014]



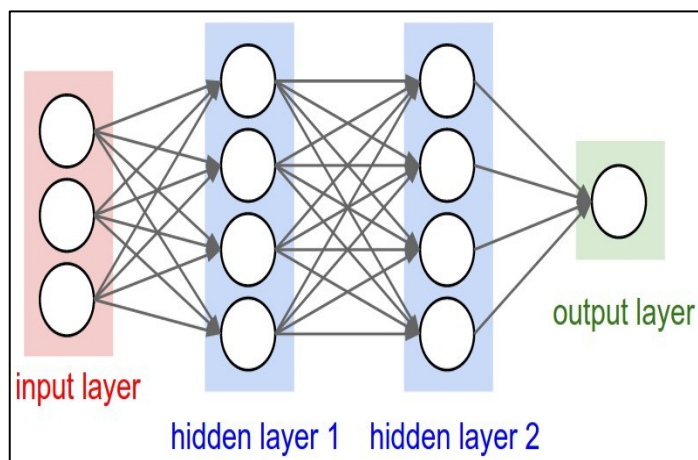
## Inception module

ILSVRC 2014 winner (6.7% top 5 error)

# Mini-batch SGD

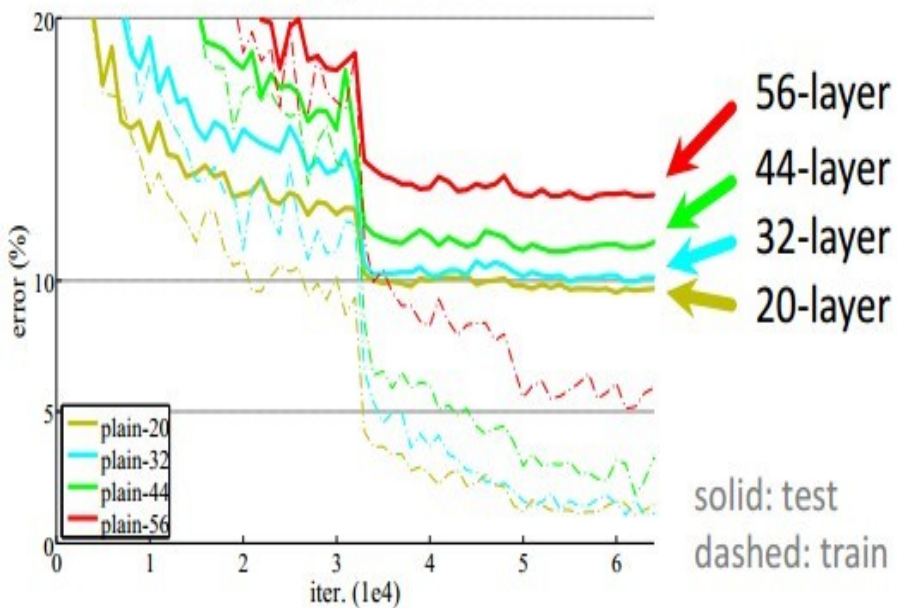
Loop:

1. **Sample** a batch of data
2. **Forward** prop it through the graph, get loss
3. **Backprop** to calculate the gradients
4. **Update** the parameters using the gradient

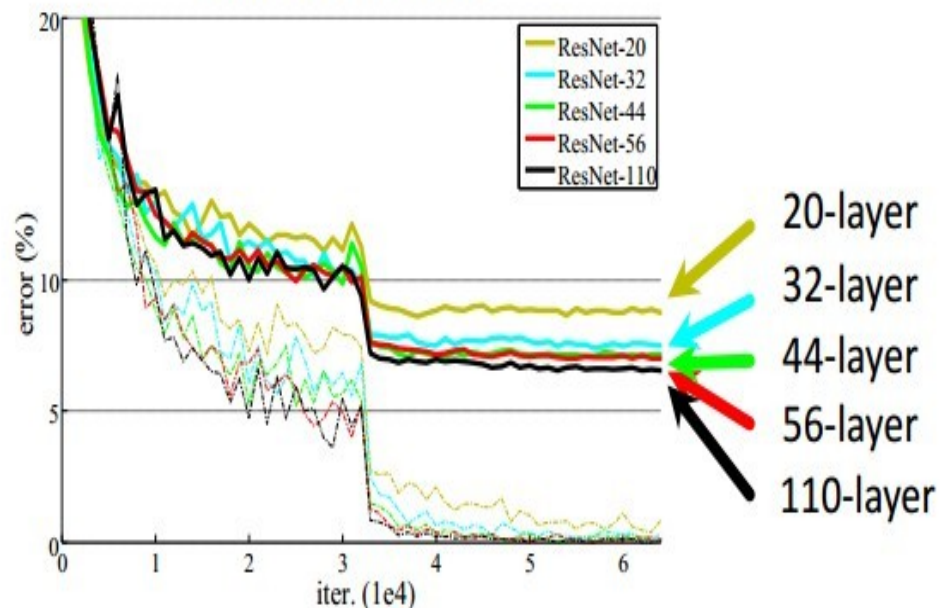


# CIFAR-10 experiments

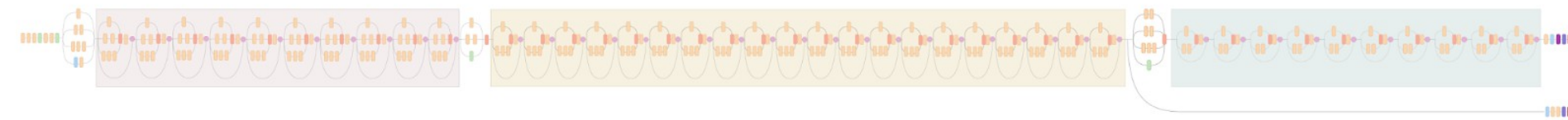
## CIFAR-10 plain nets



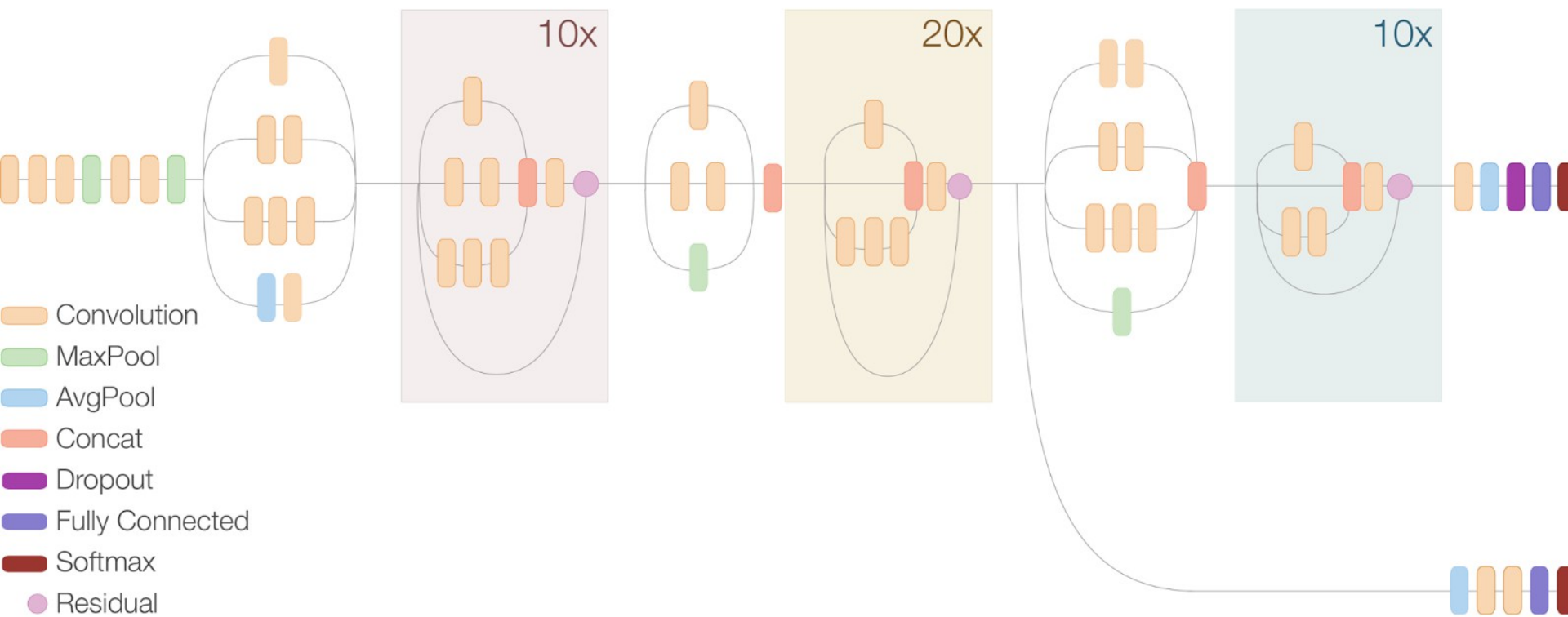
## CIFAR-10 ResNets



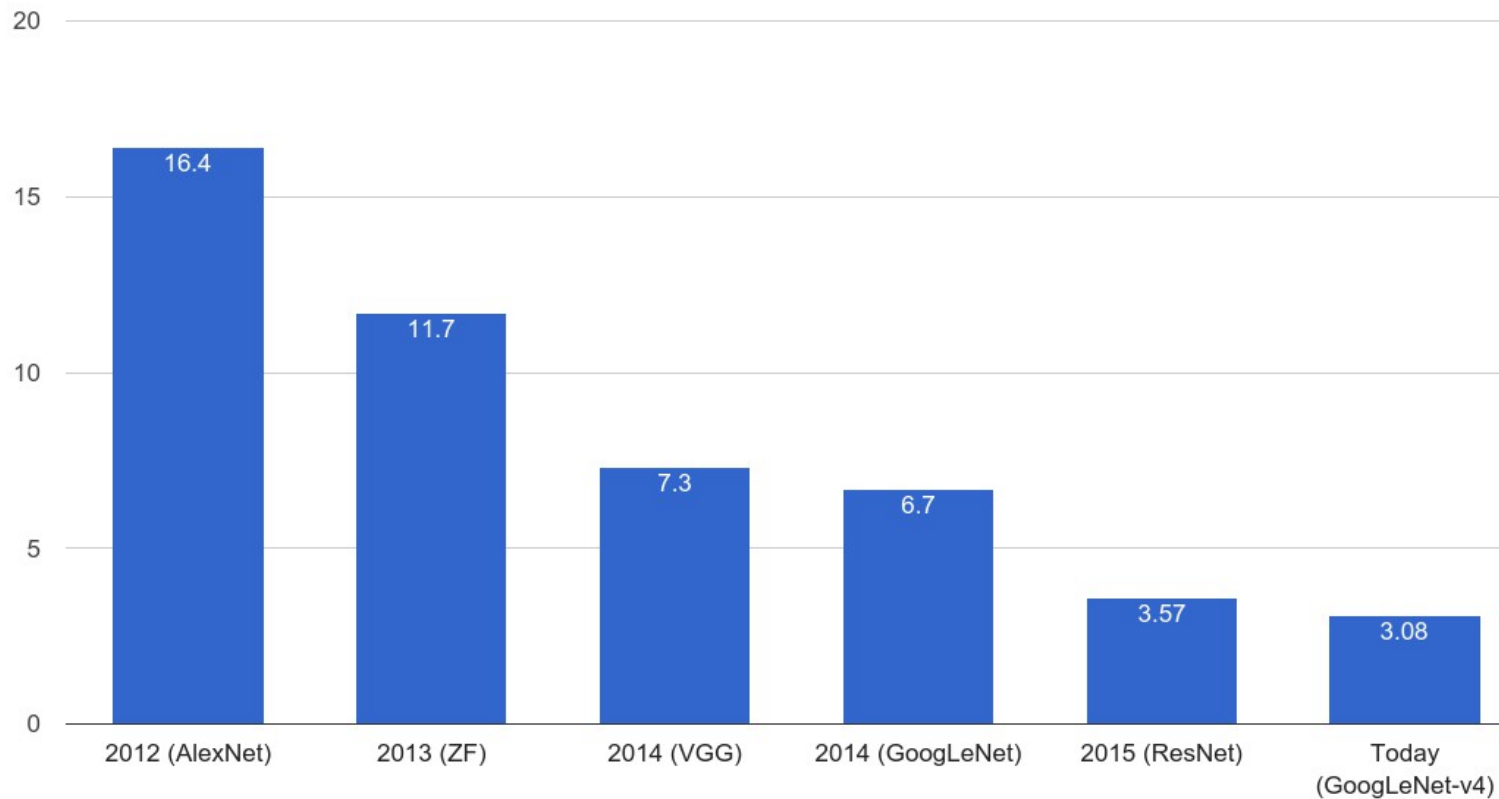
# Inception Resnet V2 Network



## Compressed View



## ImageNet Classification Error (Top 5)



Szegedy et al, Inception-v4, Inception-ResNet and the Impact of Residual Connections on Learning, arXiv 2016



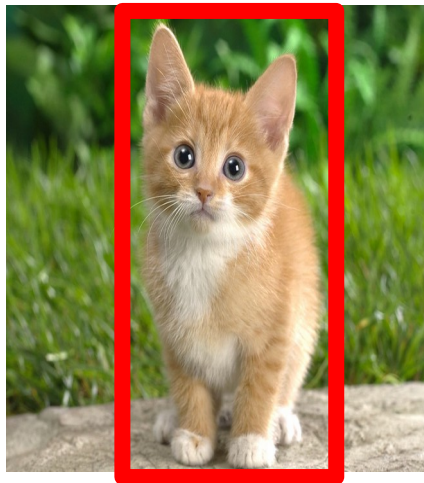
# Computer Vision Tasks

**Classification**



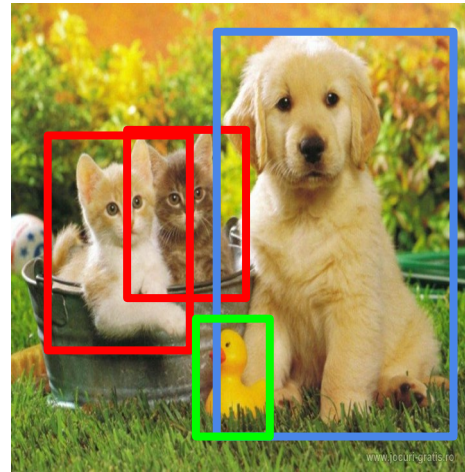
CAT

**Classification  
+ Localization**



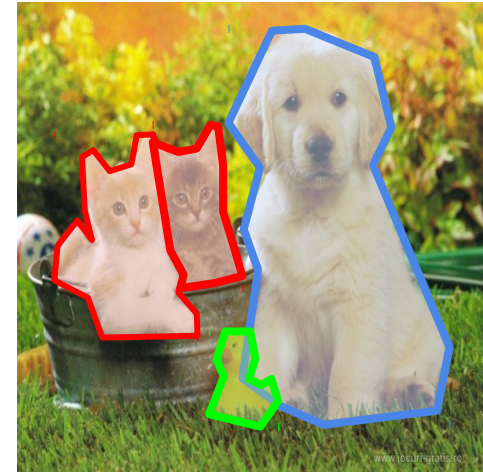
CAT

**Object Detection**



CAT, DOG, DUCK

**Instance  
Segmentation**



CAT, DOG, DUCK

Single object

Multiple objects

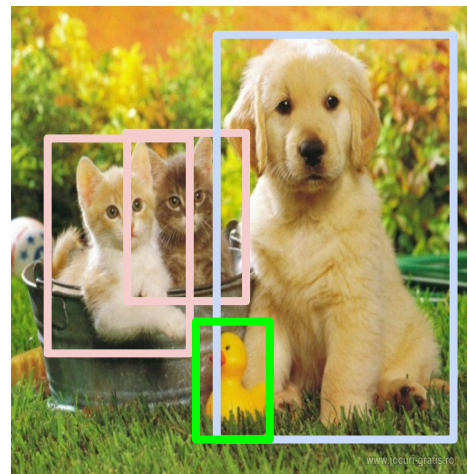
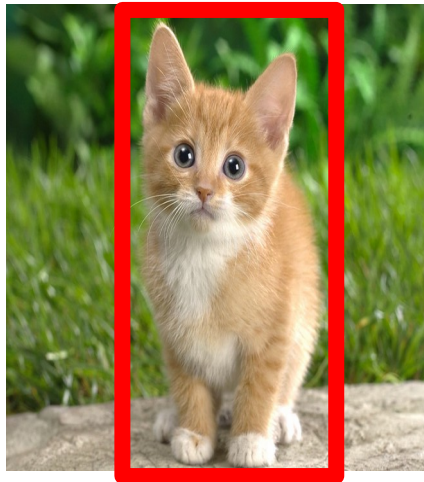
# Computer Vision Tasks

Classification

**Classification  
+ Localization**

Object Detection

Instance  
Segmentation



# Classification + Localization: Task

**Classification:** C classes

**Input:** Image

**Output:** Class label

**Evaluation metric:** Accuracy



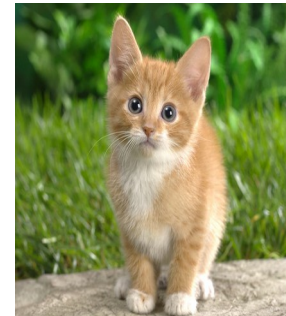
→ CAT

**Localization:**

**Input:** Image

**Output:** Box in the image (x, y, w, h)

**Evaluation metric:** Intersection over Union



→ (x, y, w, h)

**Classification + Localization:** Do both

# Classification + Localization: ImageNet

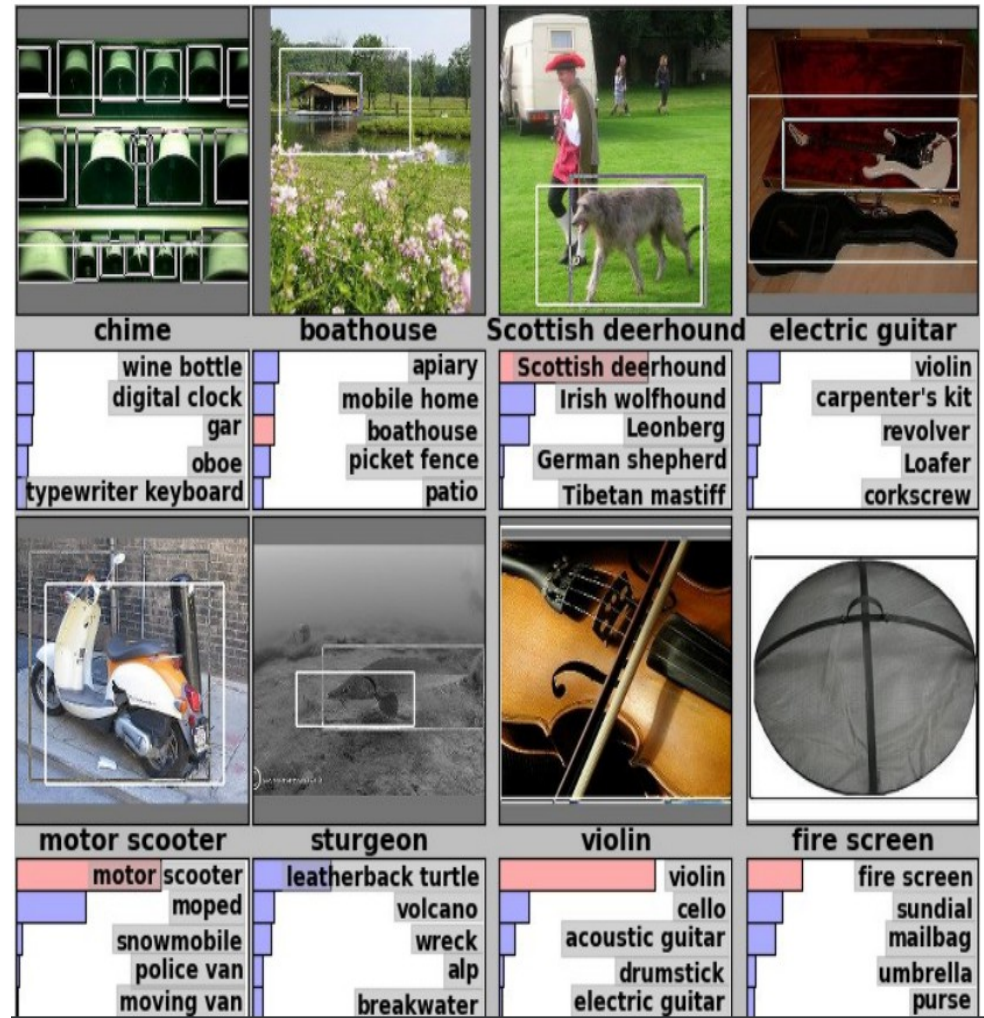
1000 classes (same as classification)

Each image has 1 class, at least one bounding box

~800 training images per class

Algorithm produces 5 (class, box) guesses

Example is correct if at least one one guess has correct class AND bounding box at least 0.5 intersection over union (IoU)



Krizhevsky et. al. 2012

# Idea #1: Localization as Regression

**Input:** image



Neural Net  
→

**Output:**  
Box coordinates  
(4 numbers)

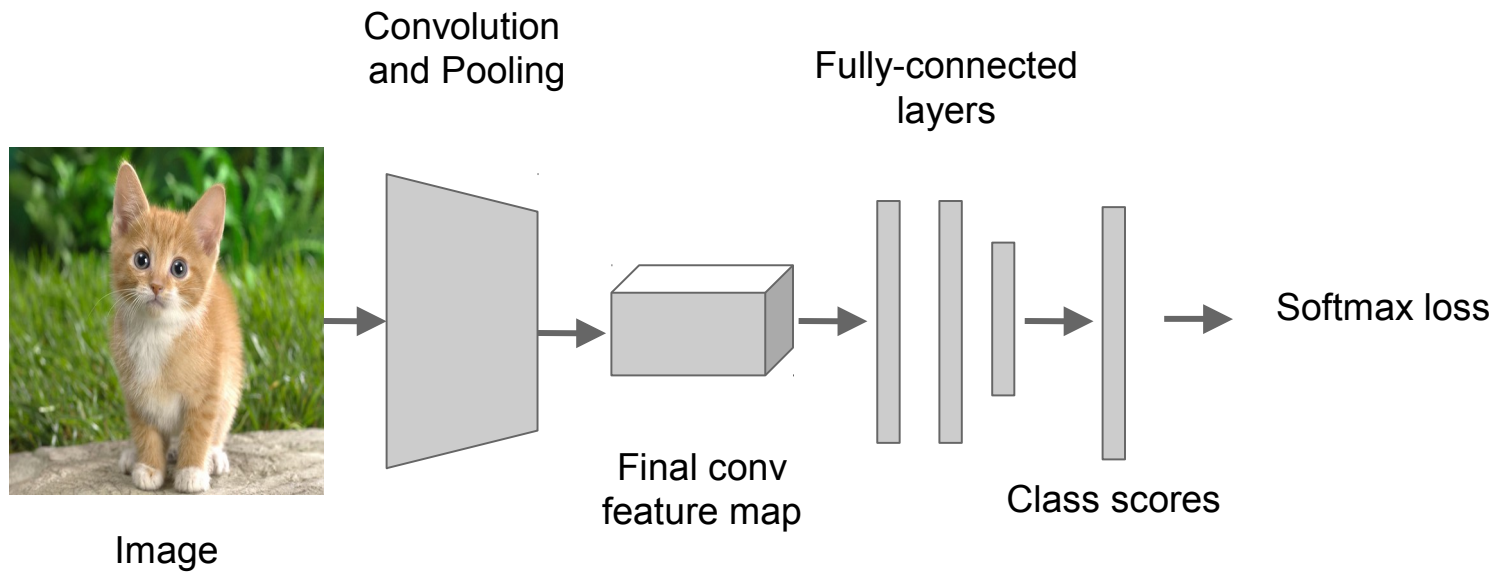
**Correct output:**  
box coordinates  
(4 numbers)

**Loss:**  
L2 distance

Only one object,  
simpler than detection

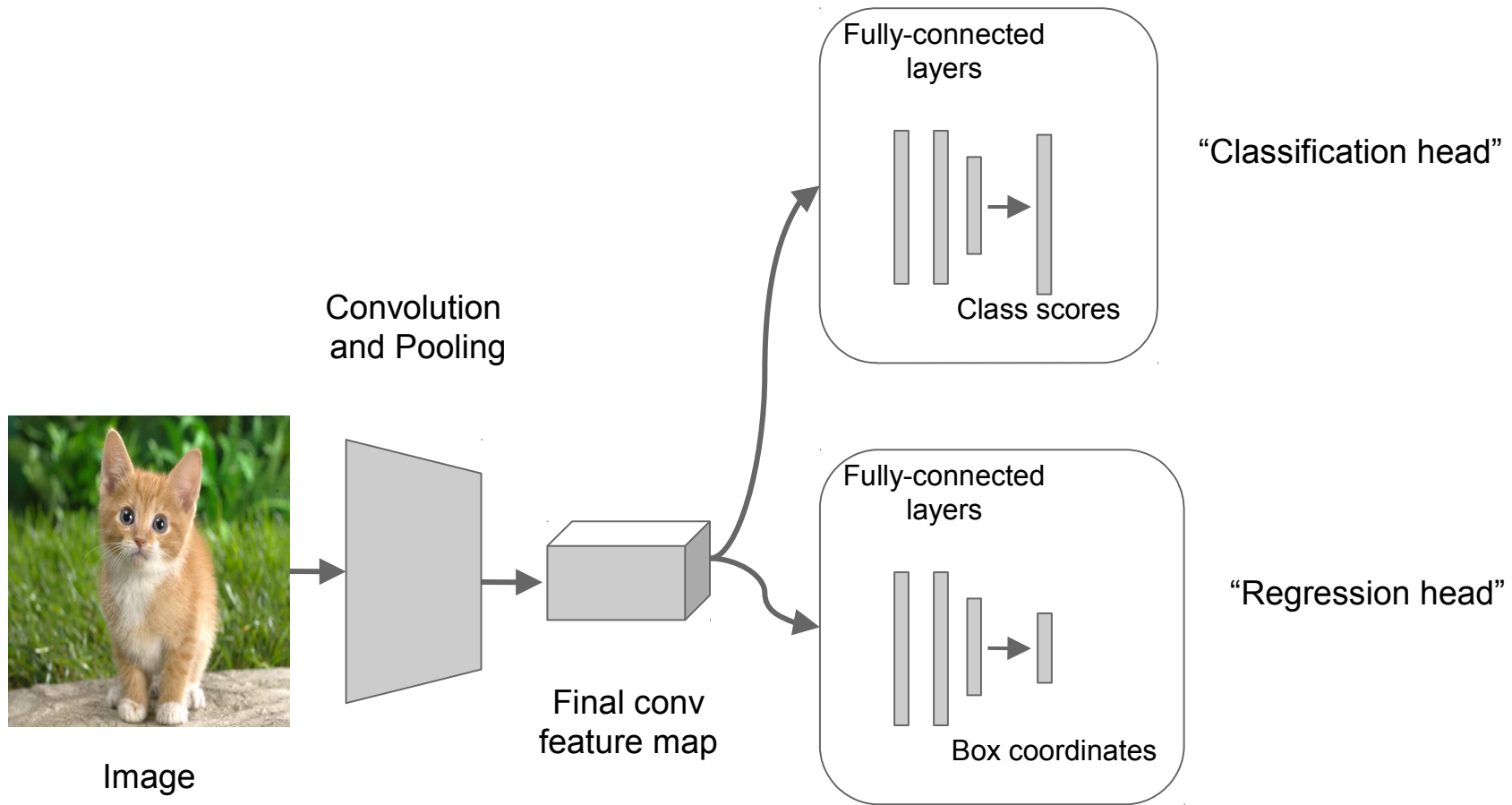
# Simple Recipe for Classification + Localization

**Step 1:** Train (or download) a classification model (AlexNet, VGG, GoogLeNet)



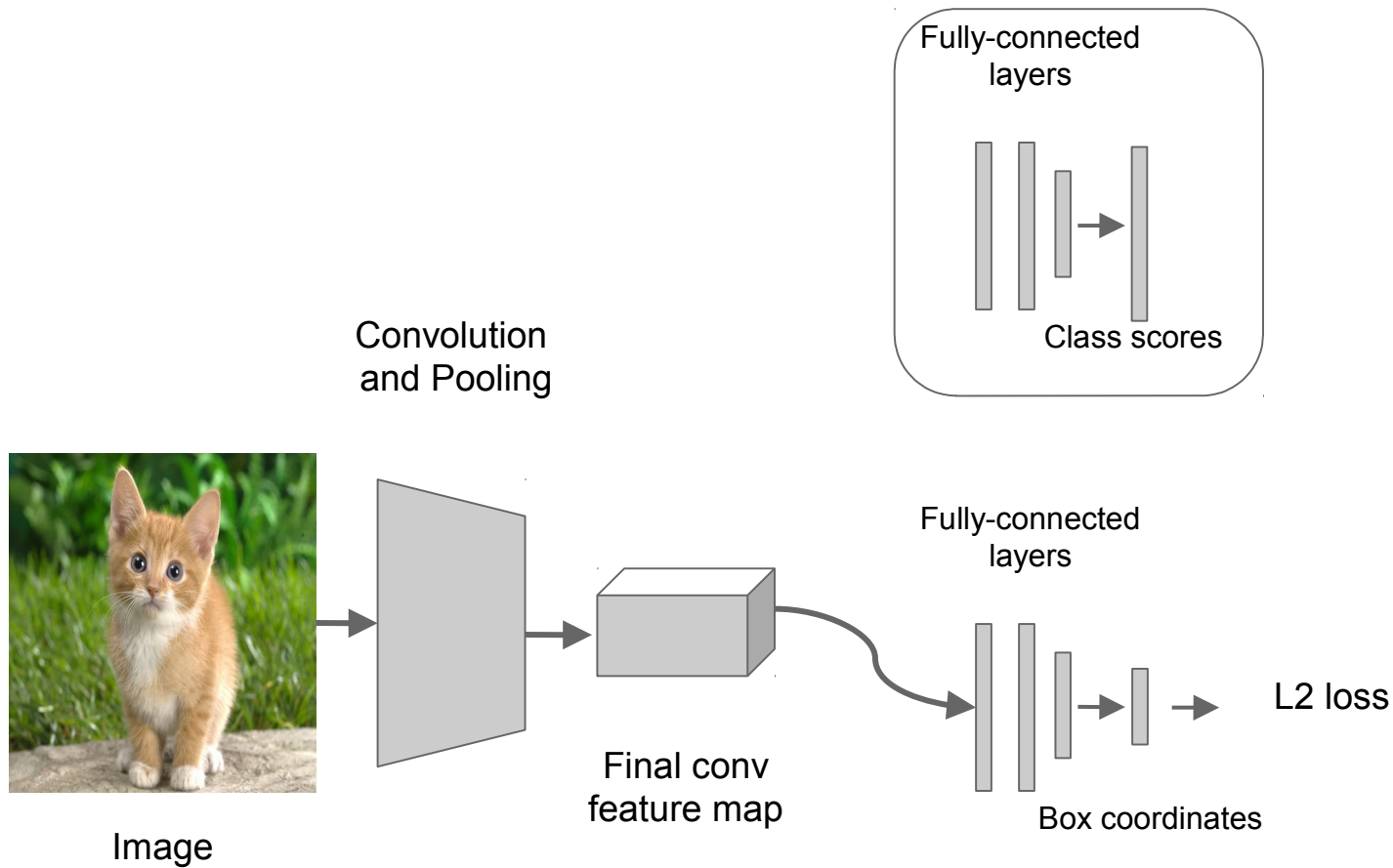
# Simple Recipe for Classification + Localization

**Step 2:** Attach new fully-connected “regression head” to the network



# Simple Recipe for Classification + Localization

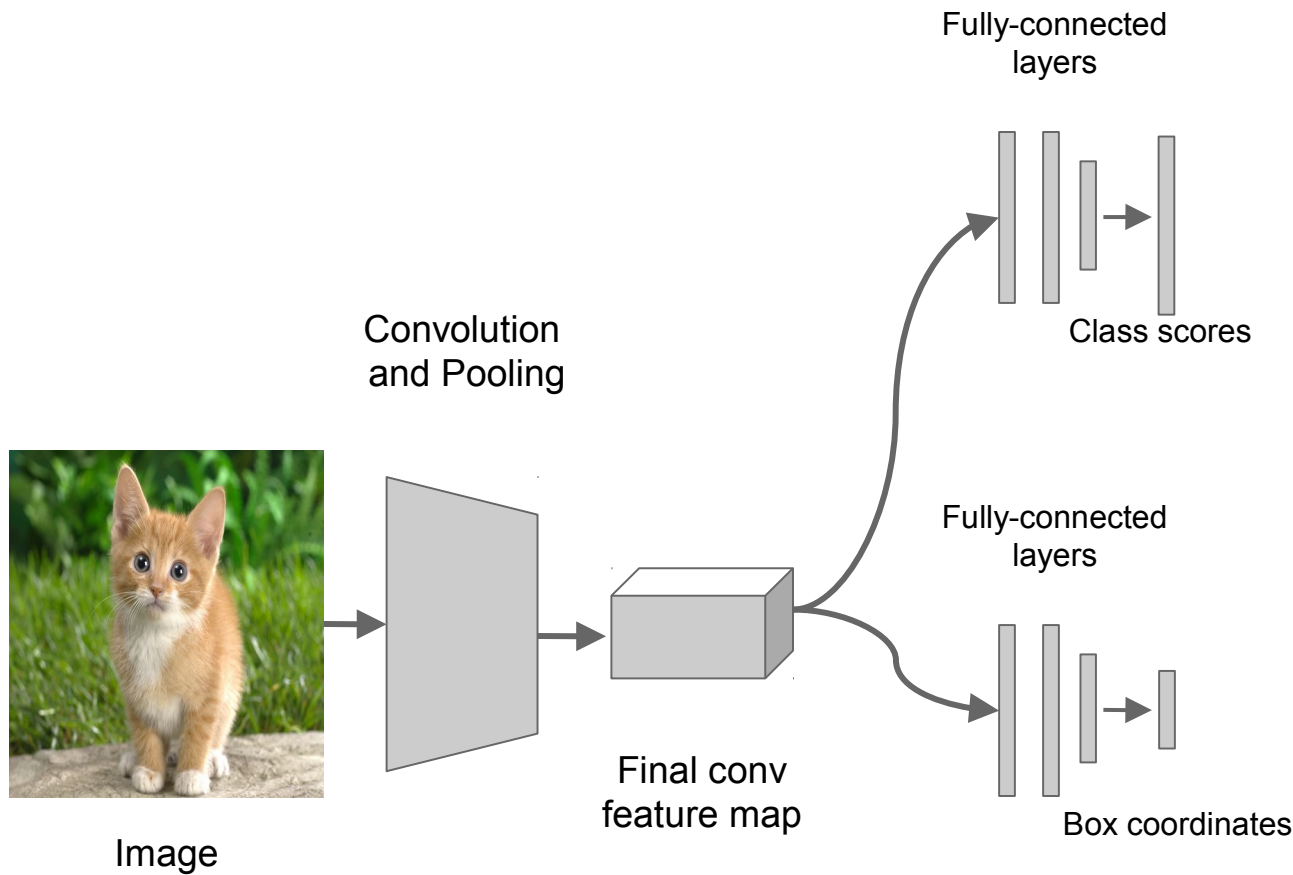
**Step 3:** Train the regression head only with SGD and L2 loss





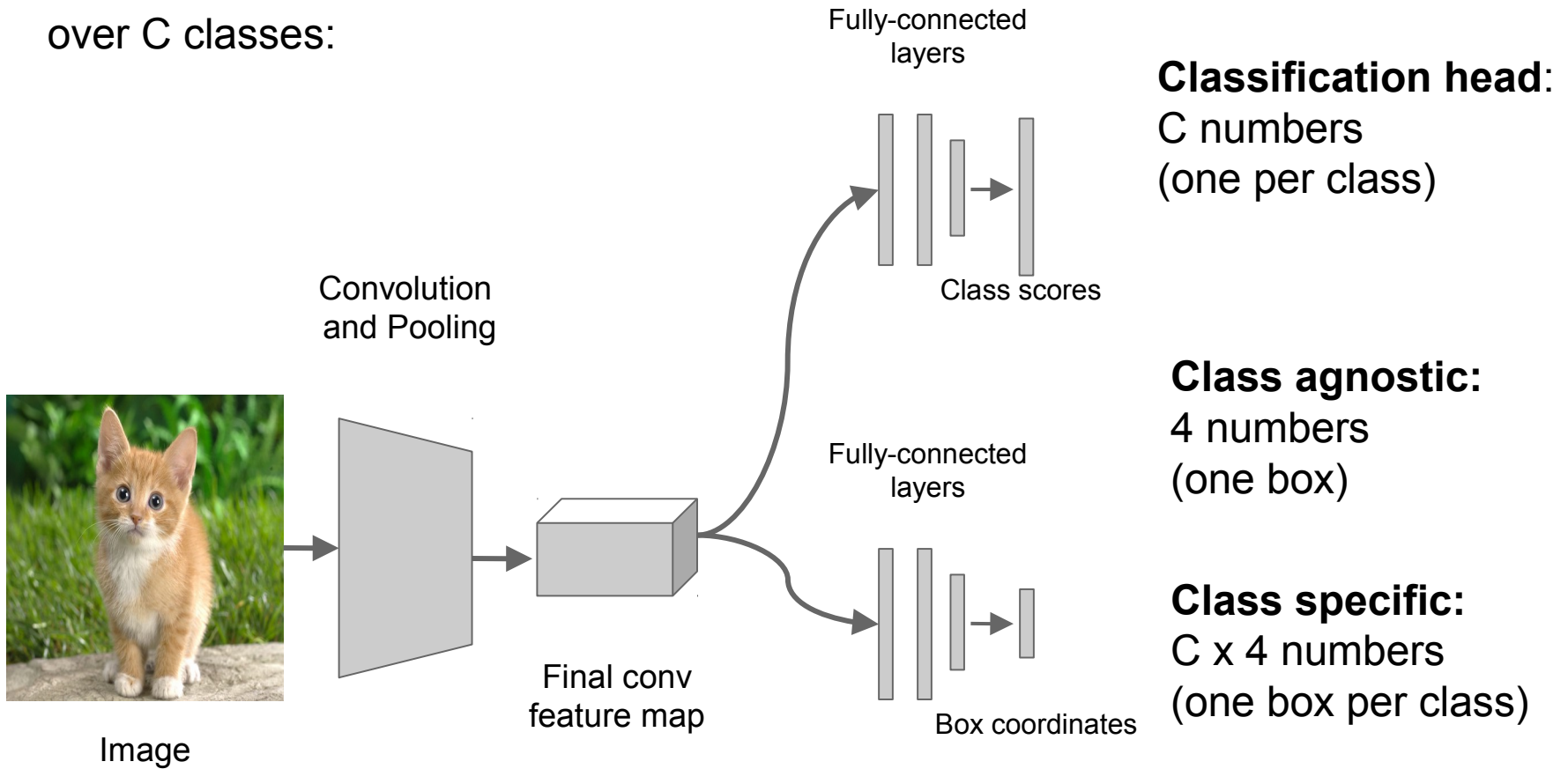
# Simple Recipe for Classification + Localization

**Step 4:** At test time use both heads

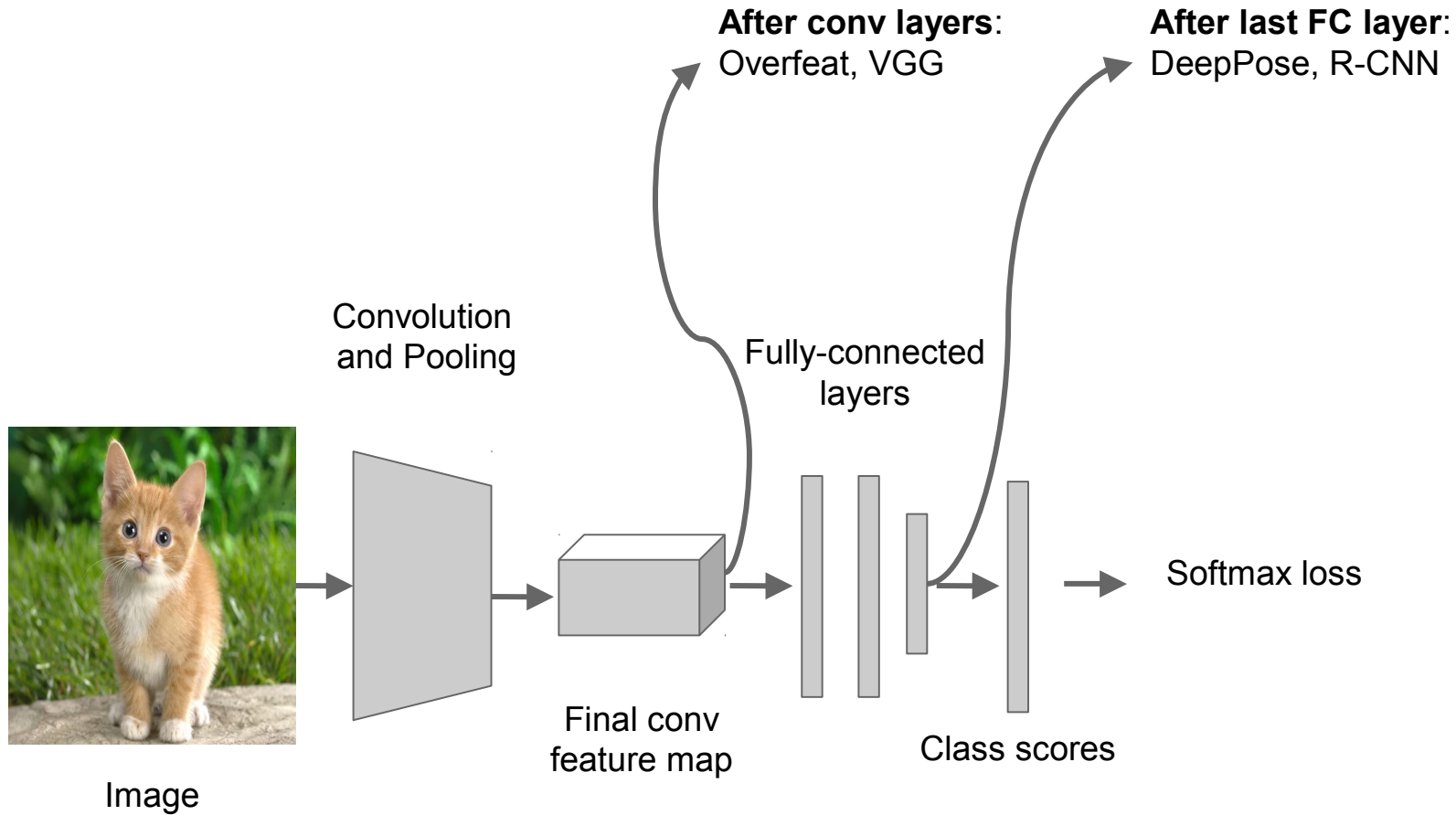


# Per-class vs class agnostic regression

Assume classification  
over  $C$  classes:



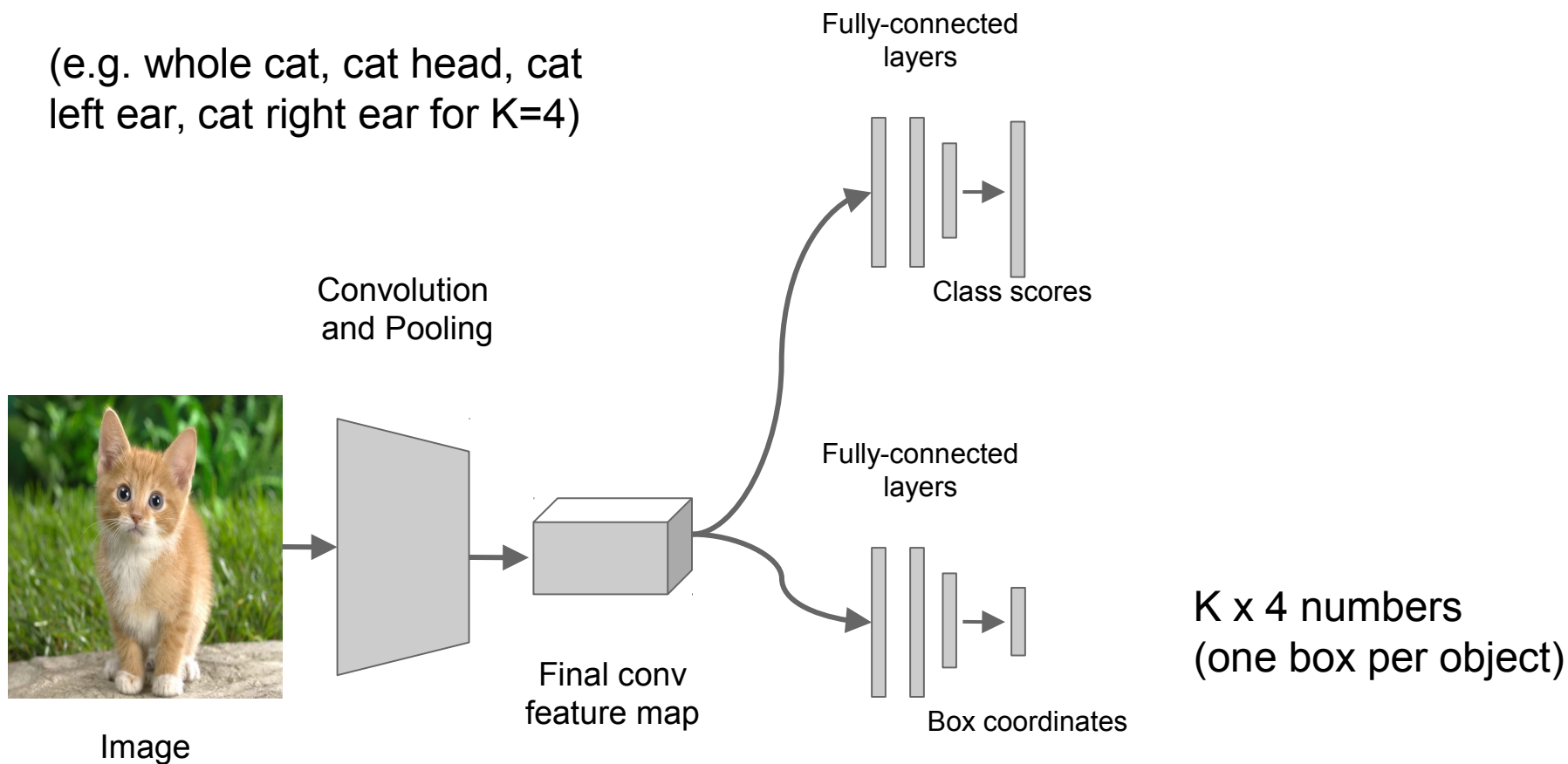
# Where to attach the regression head?



# Aside: Localizing multiple objects

Want to localize **exactly**  $K$  objects in each image

(e.g. whole cat, cat head, cat left ear, cat right ear for  $K=4$ )

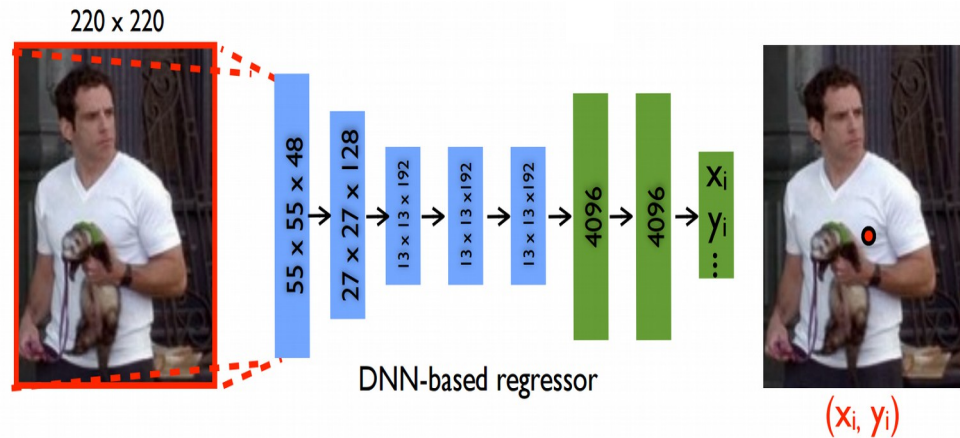


# Aside: Human Pose Estimation

Represent a person by  $K$  joints

Regress  $(x, y)$  for each joint from last fully-connected layer of AlexNet

(Details: Normalized coordinates, iterative refinement)



Toshev and Szegedy, "DeepPose: Human Pose Estimation via Deep Neural Networks", CVPR 2014

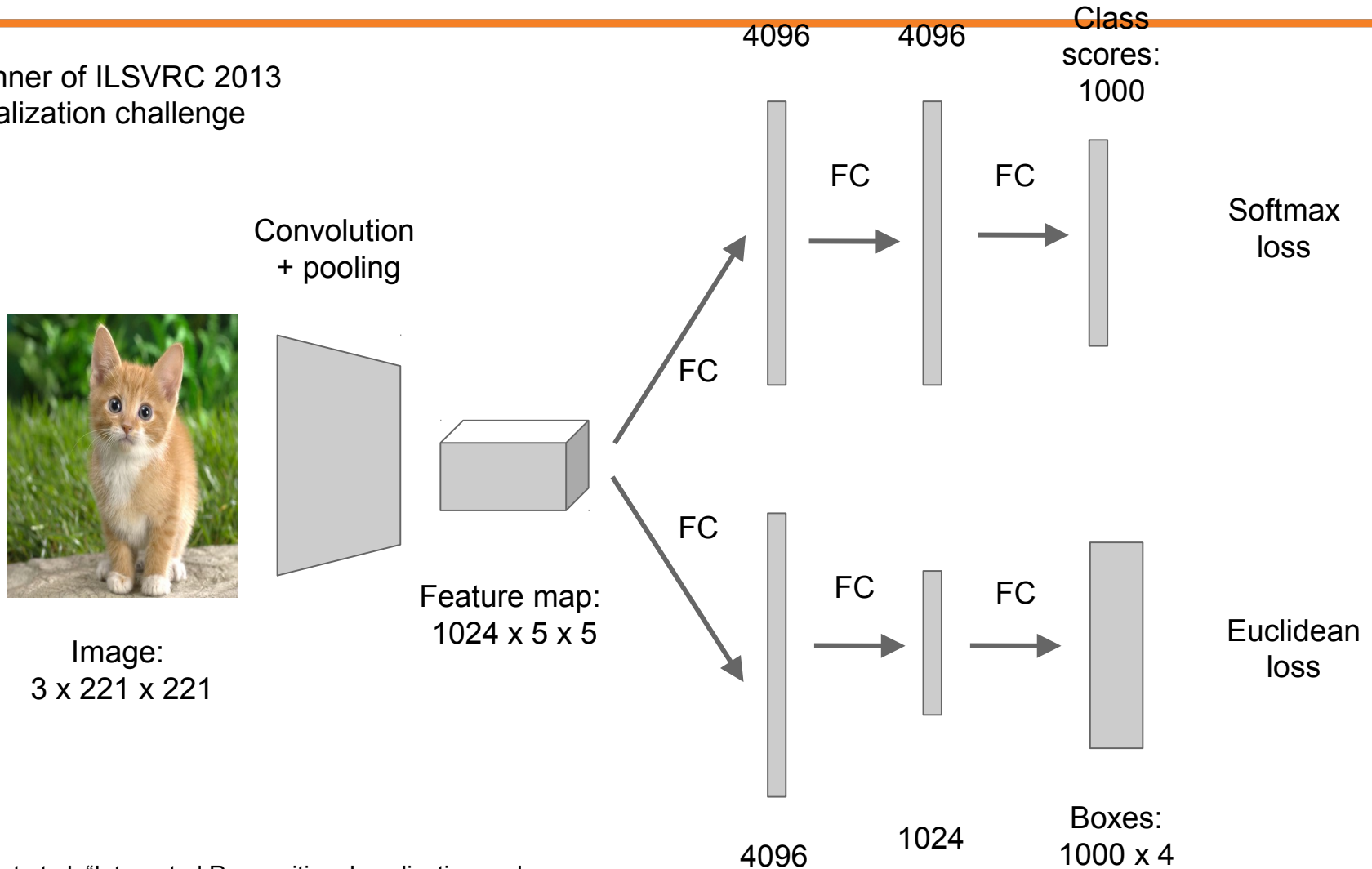
# Idea #2: Sliding Window

---

- Run classification + regression network at multiple locations on a high-resolution image
- Convert fully-connected layers into convolutional layers for efficient computation
- Combine classifier and regressor predictions across all scales for final prediction

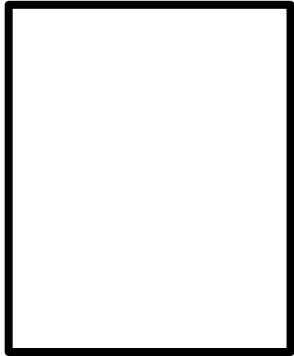
# Sliding Window: Overfeat

Winner of ILSVRC 2013  
localization challenge

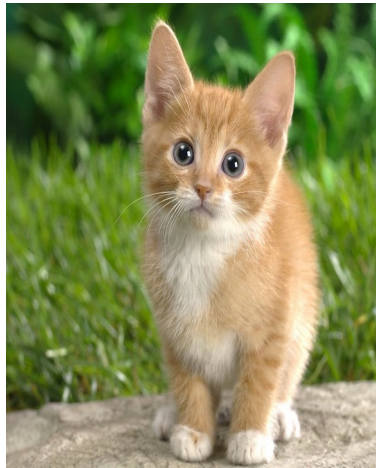


Sermanet et al, "Integrated Recognition, Localization and Detection using Convolutional Networks", ICLR 2014

# Sliding Window: Overfeat



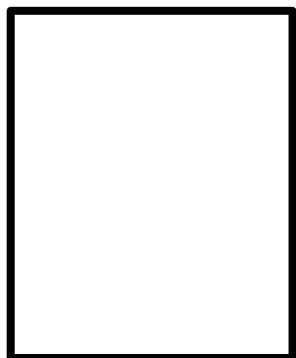
Network input:  
3 x 221 x 221



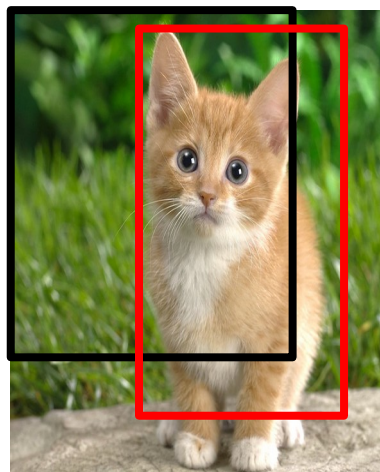
Larger image:  
3 x 257 x 257



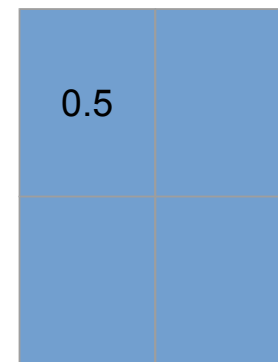
# Sliding Window: Overfeat



Network input:  
3 x 221 x 221

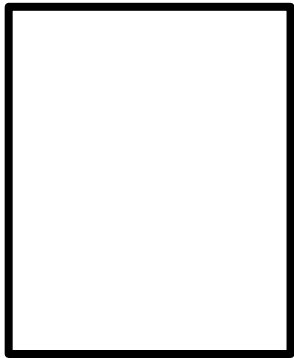


Larger image:  
3 x 257 x 257

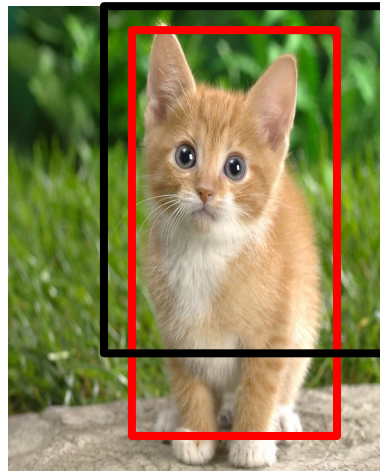


Classification scores:  
P(cat)

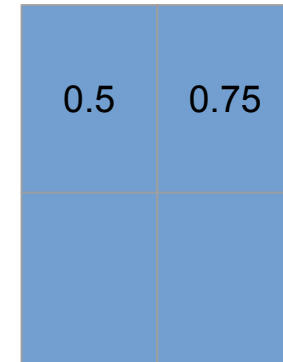
# Sliding Window: Overfeat



Network input:  
 $3 \times 221 \times 221$

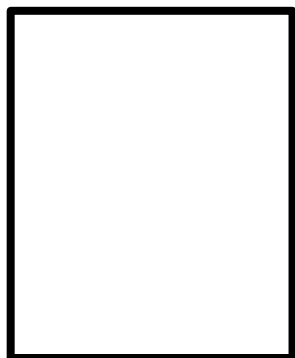


Larger image:  
 $3 \times 257 \times 257$

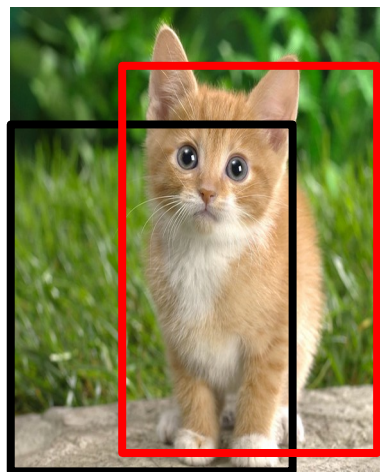


Classification scores:  
 $P(\text{cat})$

# Sliding Window: Overfeat



Network input:  
3 x 221 x 221

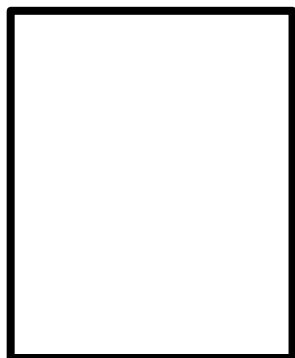


Larger image:  
3 x 257 x 257

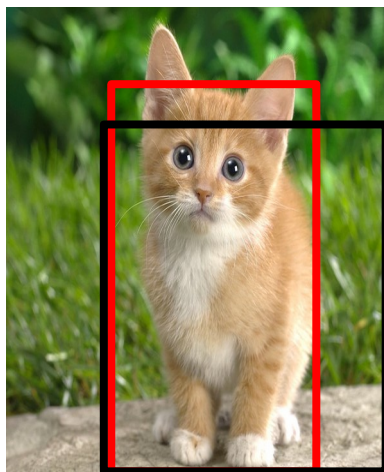
0.5	0.75
0.6	

Classification scores:  
P(cat)

# Sliding Window: Overfeat



Network input:  
3 x 221 x 221

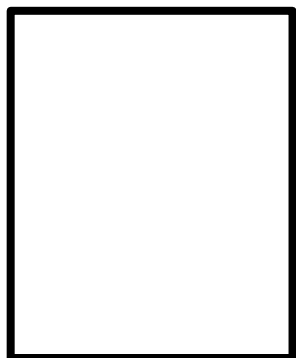


Larger image:  
3 x 257 x 257

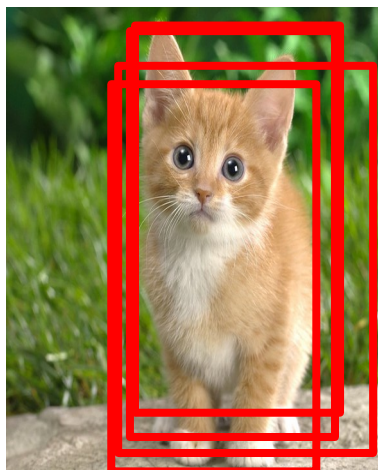
0.5	0.75
0.6	0.8

Classification scores:  
P(cat)

# Sliding Window: Overfeat



Network input:  
3 x 221 x 221



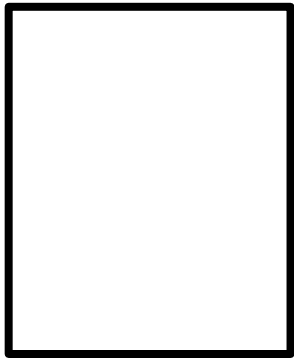
Larger image:  
3 x 257 x 257

0.5	0.75
0.6	0.8

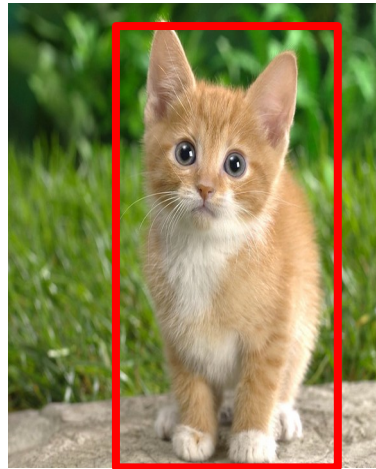
Classification scores:  
P(cat)

# Sliding Window: Overfeat

Greedily merge boxes and scores (details in paper)



Network input:  
3 x 221 x 221



Larger image:  
3 x 257 x 257

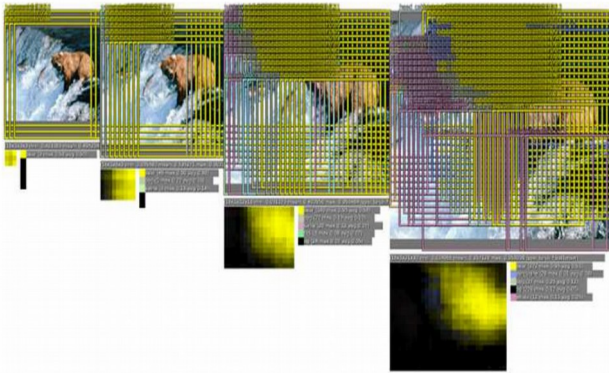
0.8

Classification score:  
P(cat)

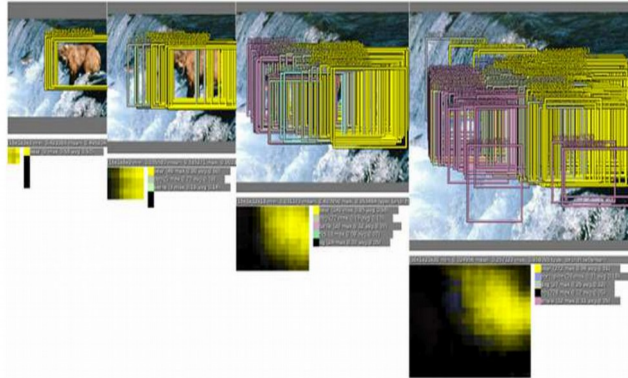
# Sliding Window: Overfeat

In practice use many sliding window locations and multiple scales

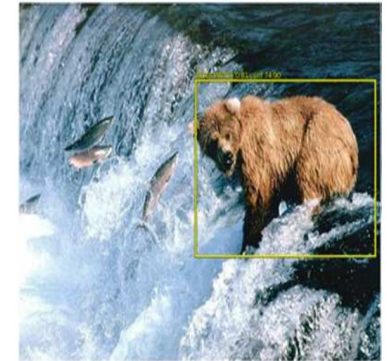
Window positions + score maps



Box regression outputs

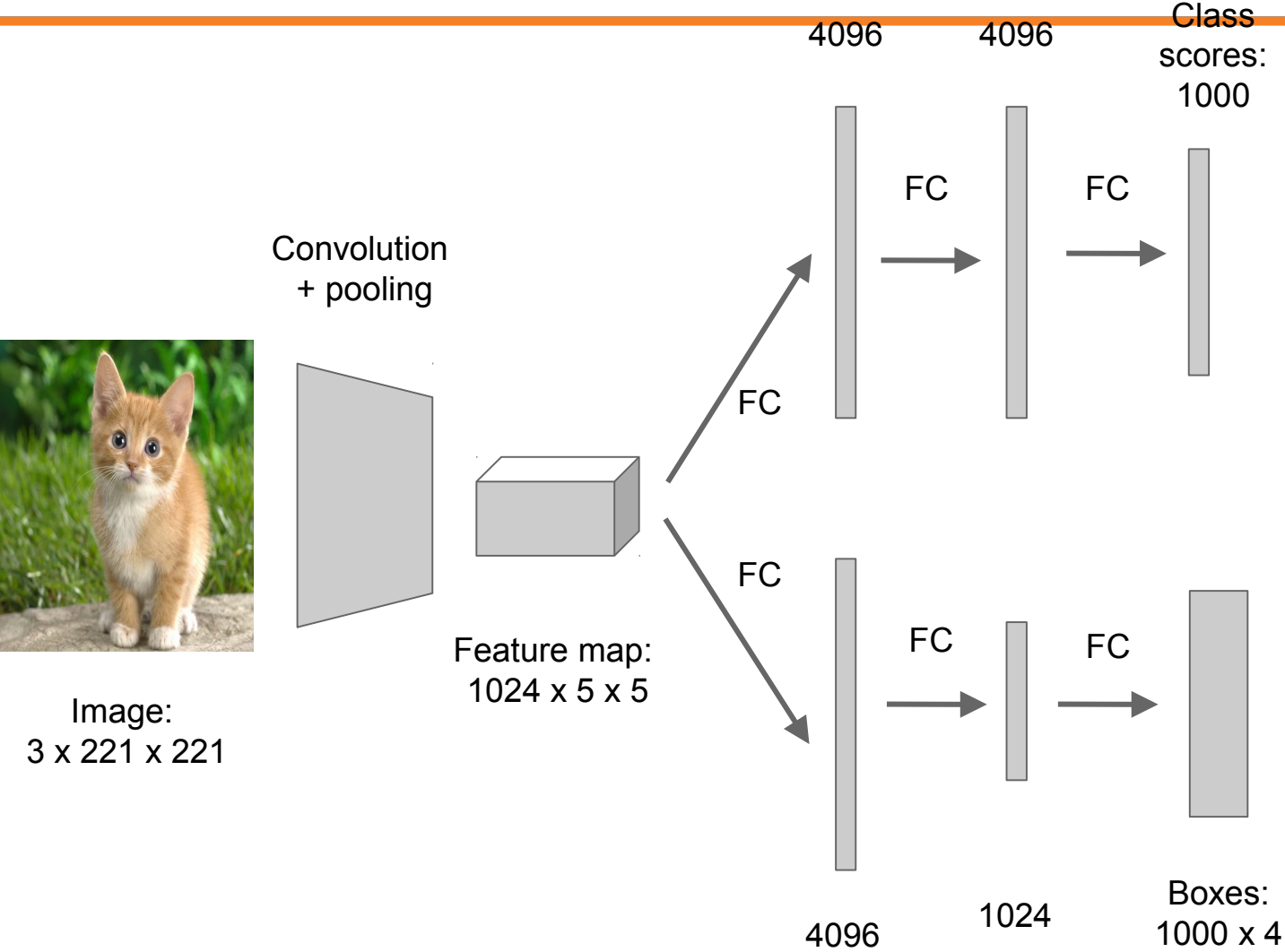


Final Predictions



Sermanet et al, "Integrated Recognition, Localization and Detection using Convolutional Networks", ICLR 2014

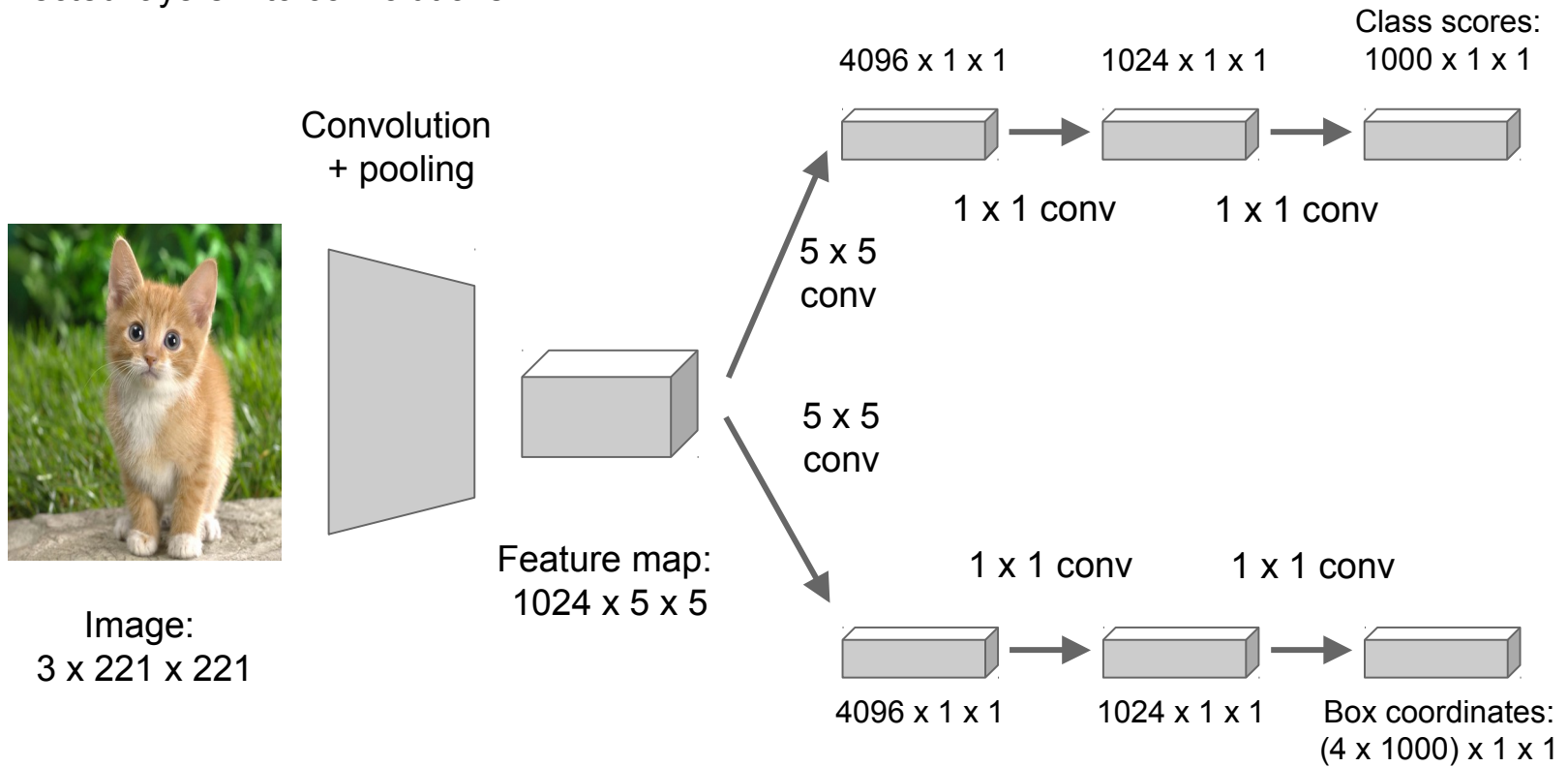
# Efficient Sliding Window: Overfeat





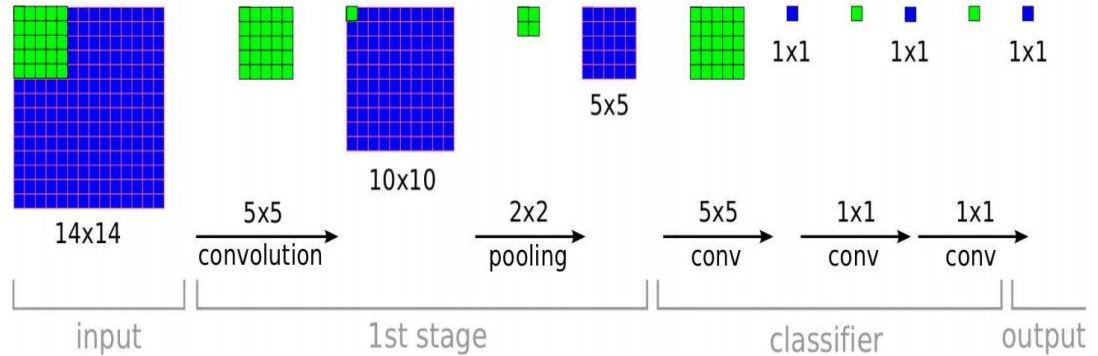
# Efficient Sliding Window: Overfeat

Efficient sliding window by converting fully-connected layers into convolutions

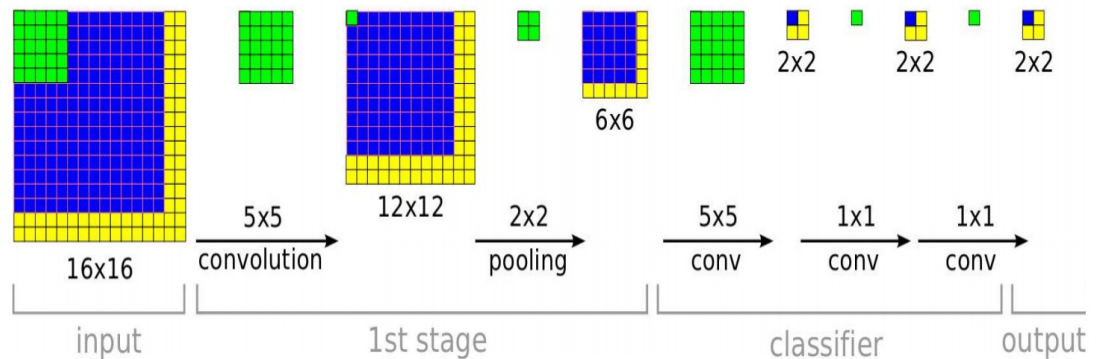


# Efficient Sliding Window: Overfeat

**Training time:** Small image, 1 x 1 classifier output



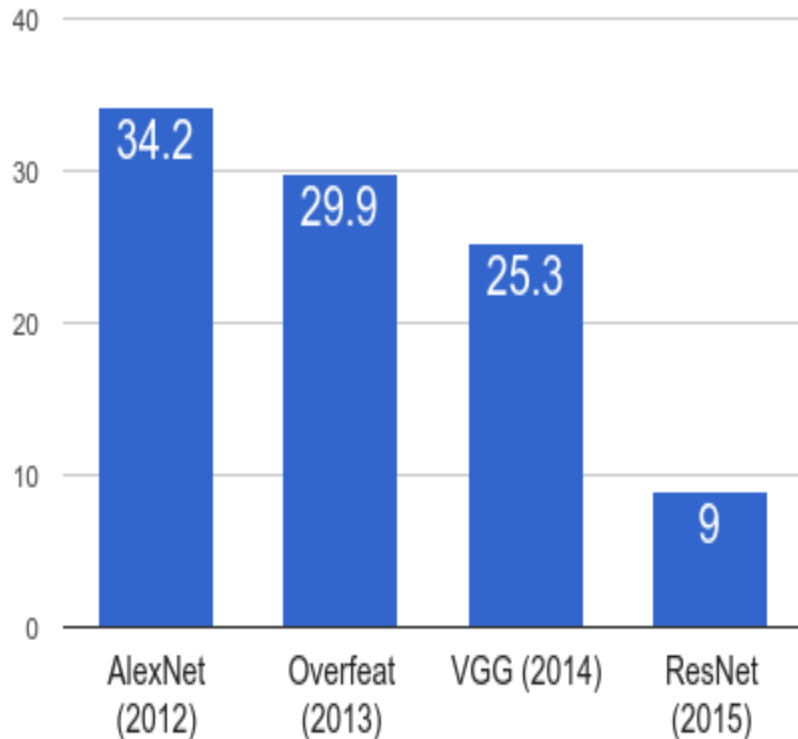
**Test time:** Larger image, 2 x 2 classifier output, only extra compute at yellow regions



Sermanet et al, "Integrated Recognition, Localization and Detection using Convolutional Networks", ICLR 2014

# ImageNet Classification + Localization

Localization Error (Top 5)



**AlexNet:** Localization method not published

**Overfeat:** Multiscale convolutional regression with box merging

**VGG:** Same as Overfeat, but fewer scales and locations; simpler method, gains all due to deeper features

**ResNet:** Different localization method (RPN) and much deeper features

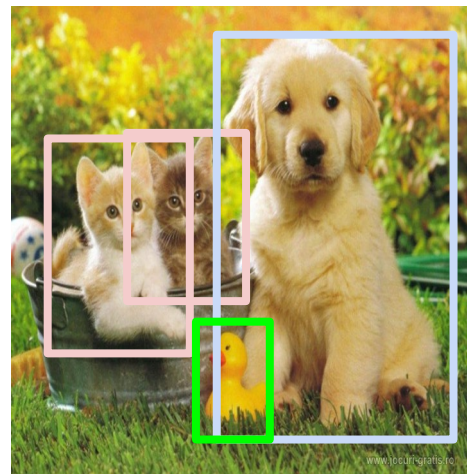
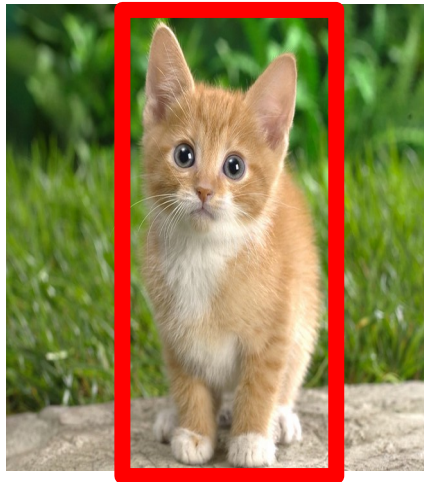
# Computer Vision Tasks

Classification

**Classification  
+ Localization**

Object Detection

Instance  
Segmentation



# Computer Vision Tasks

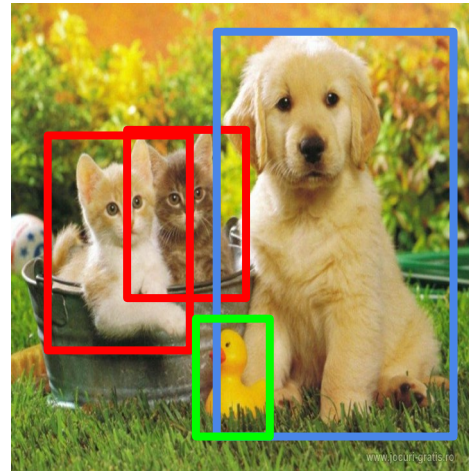
Classification



Classification  
+ Localization



Object Detection



Instance  
Segmentation



# Detection as Regression?



DOG, (x, y, w, h)  
CAT, (x, y, w, h)  
CAT, (x, y, w, h)  
DUCK (x, y, w, h)

= 16 numbers

# Detection as Regression?



DOG,  $(x, y, w, h)$   
CAT,  $(x, y, w, h)$

= 8 numbers

# Detection as Regression?



CAT, (x, y, w, h)

CAT, (x, y, w, h)

....

CAT (x, y, w, h)

= many numbers

Need variable sized outputs



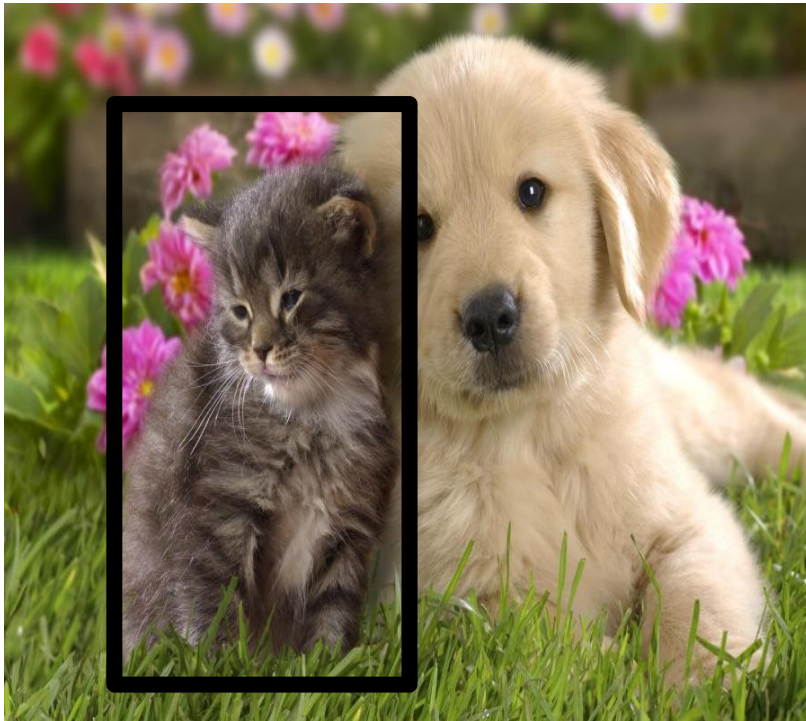
# Detection as Classification



**CAT? NO**

**DOG? NO**

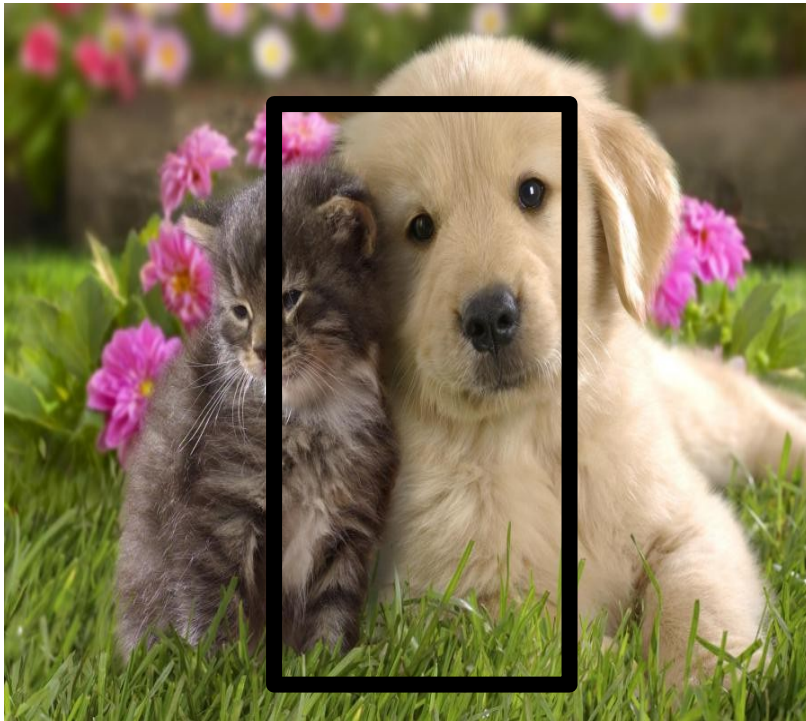
# Detection as Classification



**CAT? YES!**

**DOG? NO**

# Detection as Classification



**CAT? NO**

**DOG? NO**

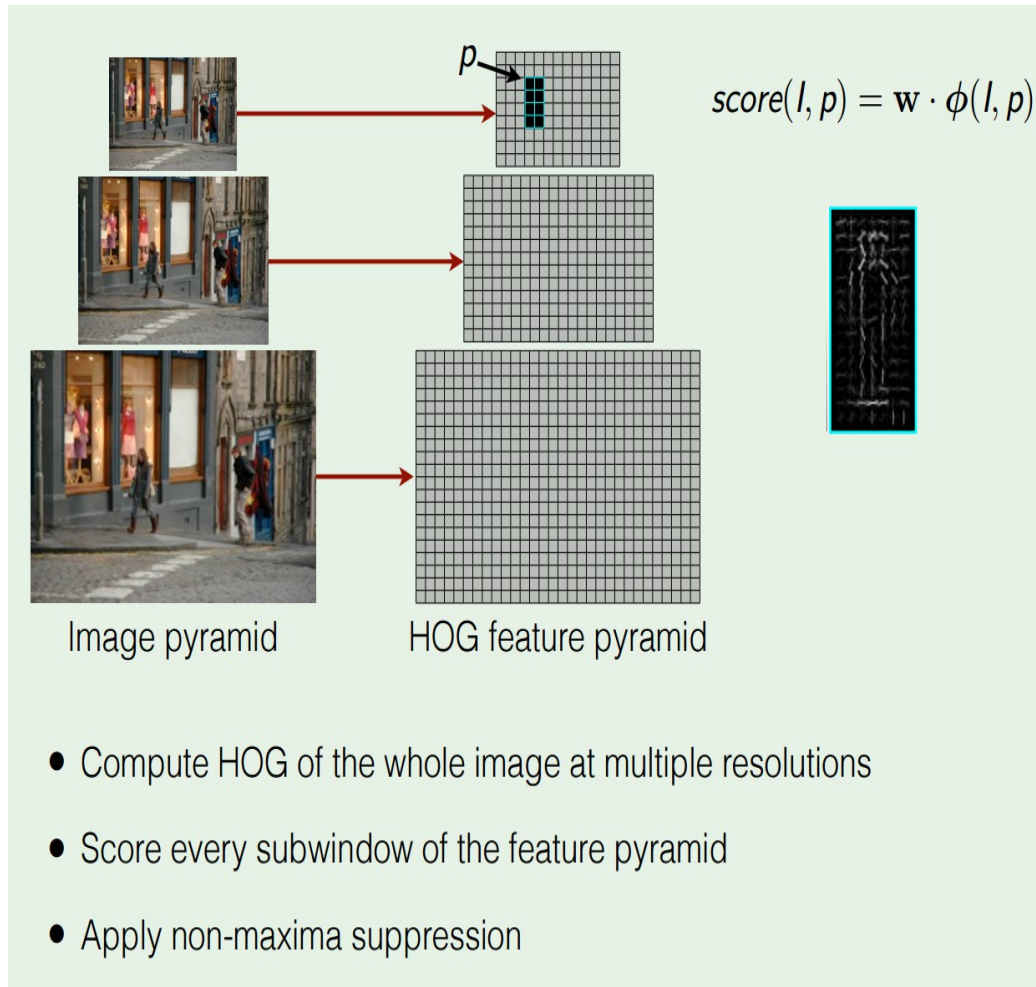
# Detection as Classification

---

**Problem:** Need to test many positions and scales

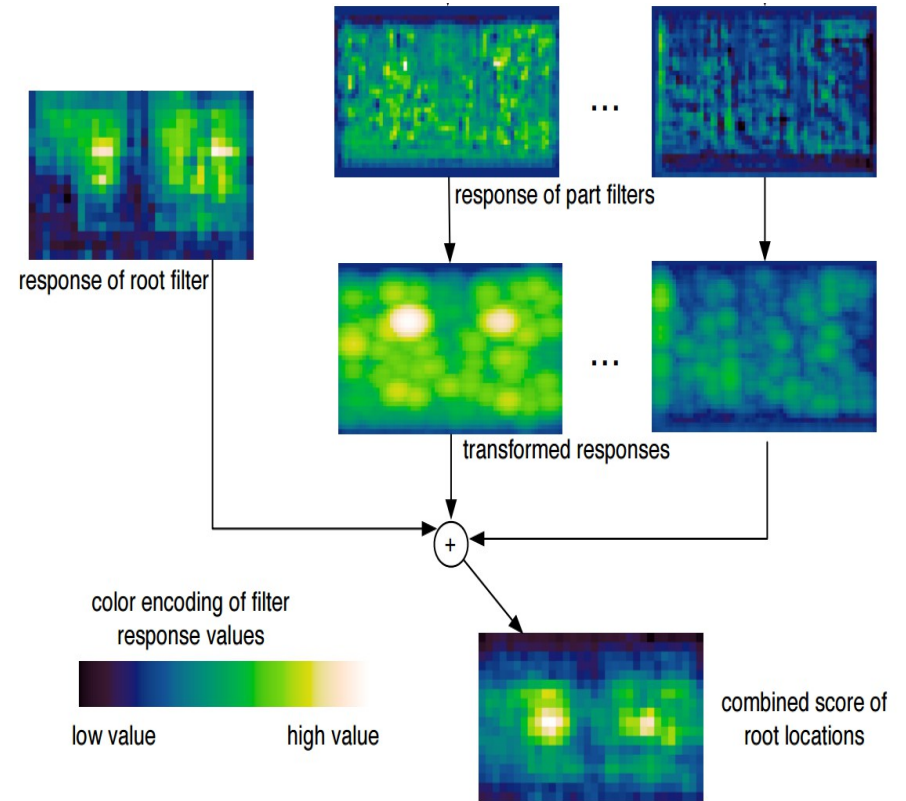
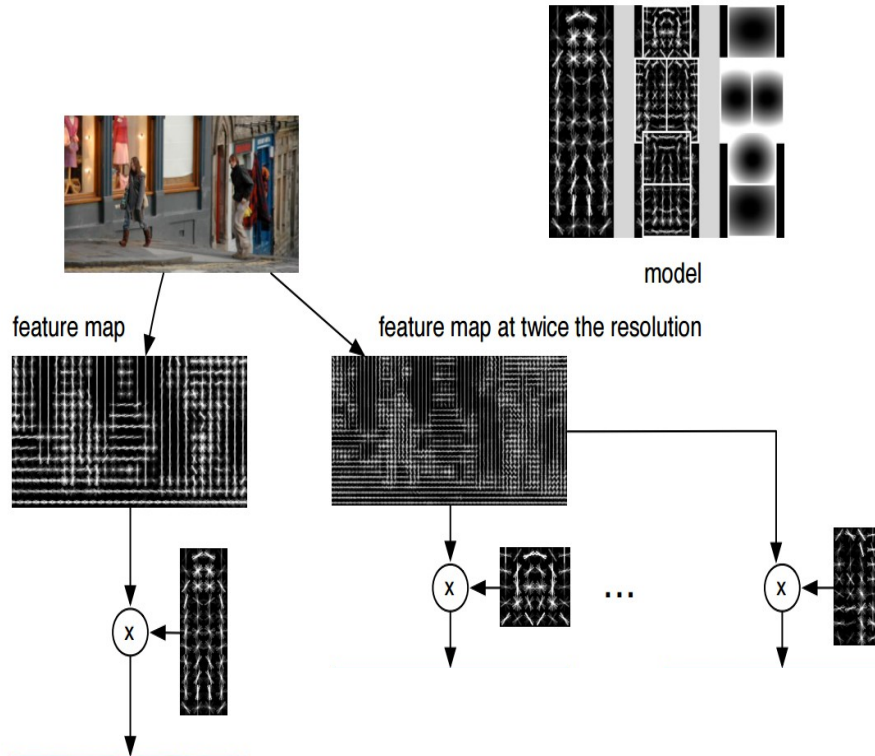
**Solution:** If your classifier is fast enough, just do it

# Histogram of Oriented Gradients



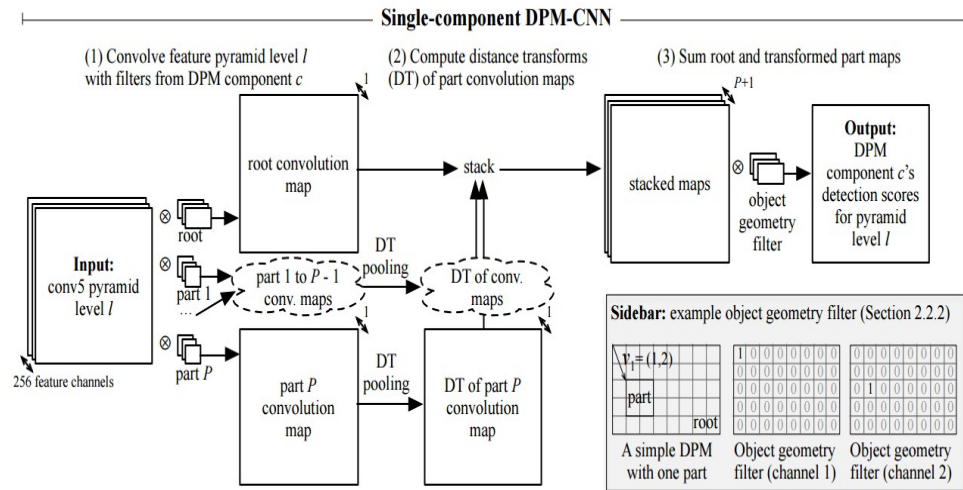
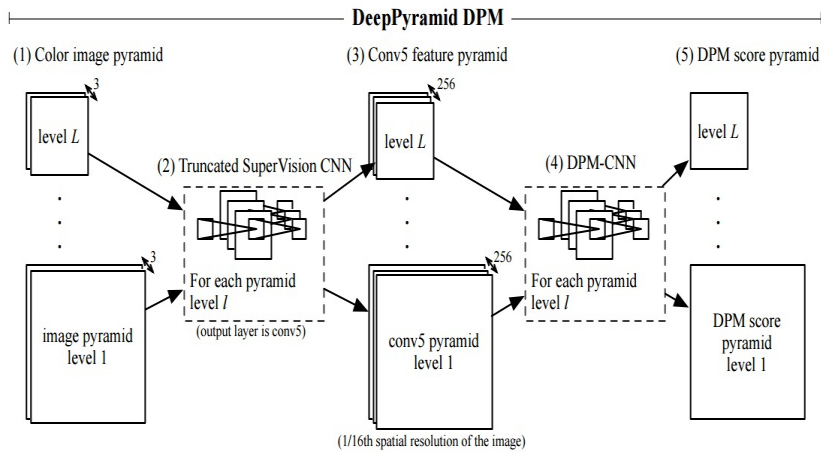
Dalal and Triggs, "Histograms of Oriented Gradients for Human Detection", CVPR 2005  
Slide credit: Ross Girshick

# Deformable Parts Model (DPM)



Felzenszwalb et al, "Object Detection with Discriminatively Trained Part Based Models", PAMI 2010

# Aside: Deformable Parts Models are CNNs?



Girschick et al, "Deformable Part Models are Convolutional Neural Networks", CVPR 2015

# Detection as Classification

---

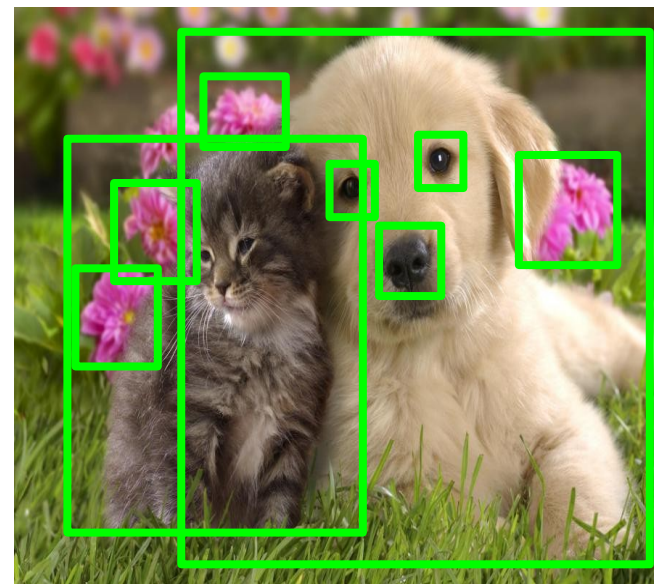
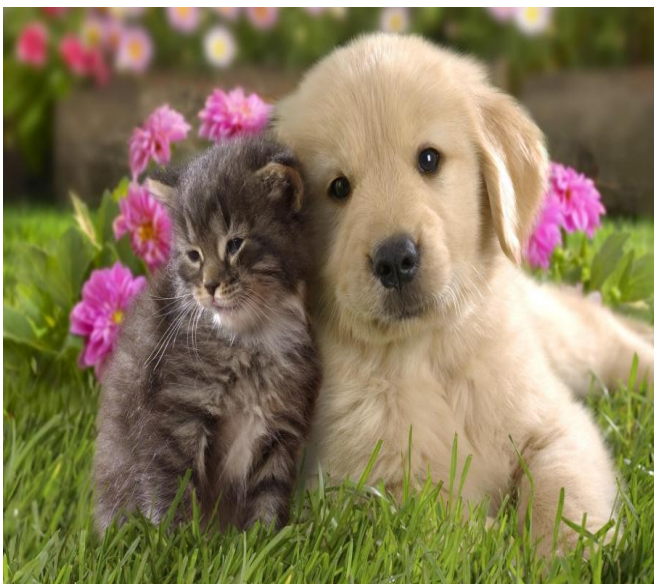
**Problem:** Need to test many positions and scales, and use a computationally demanding classifier (CNN)

**Solution:** Only look at a tiny subset of possible positions



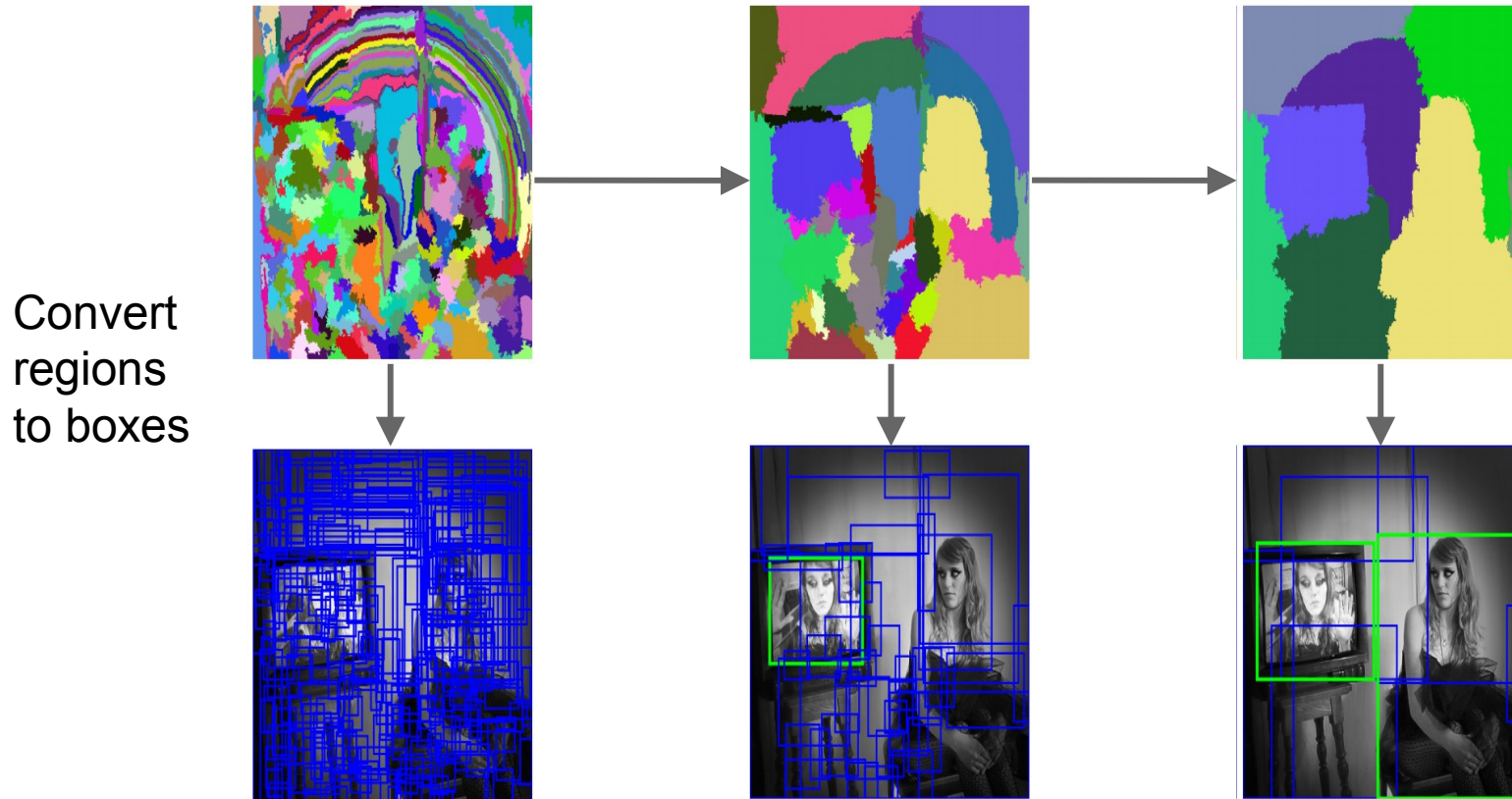
# Region Proposals

- Find “blobby” image regions that are likely to contain objects
- “Class-agnostic” object detector
- Look for “blob-like” regions



# Region Proposals: Selective Search

Bottom-up segmentation, merging regions at multiple scales



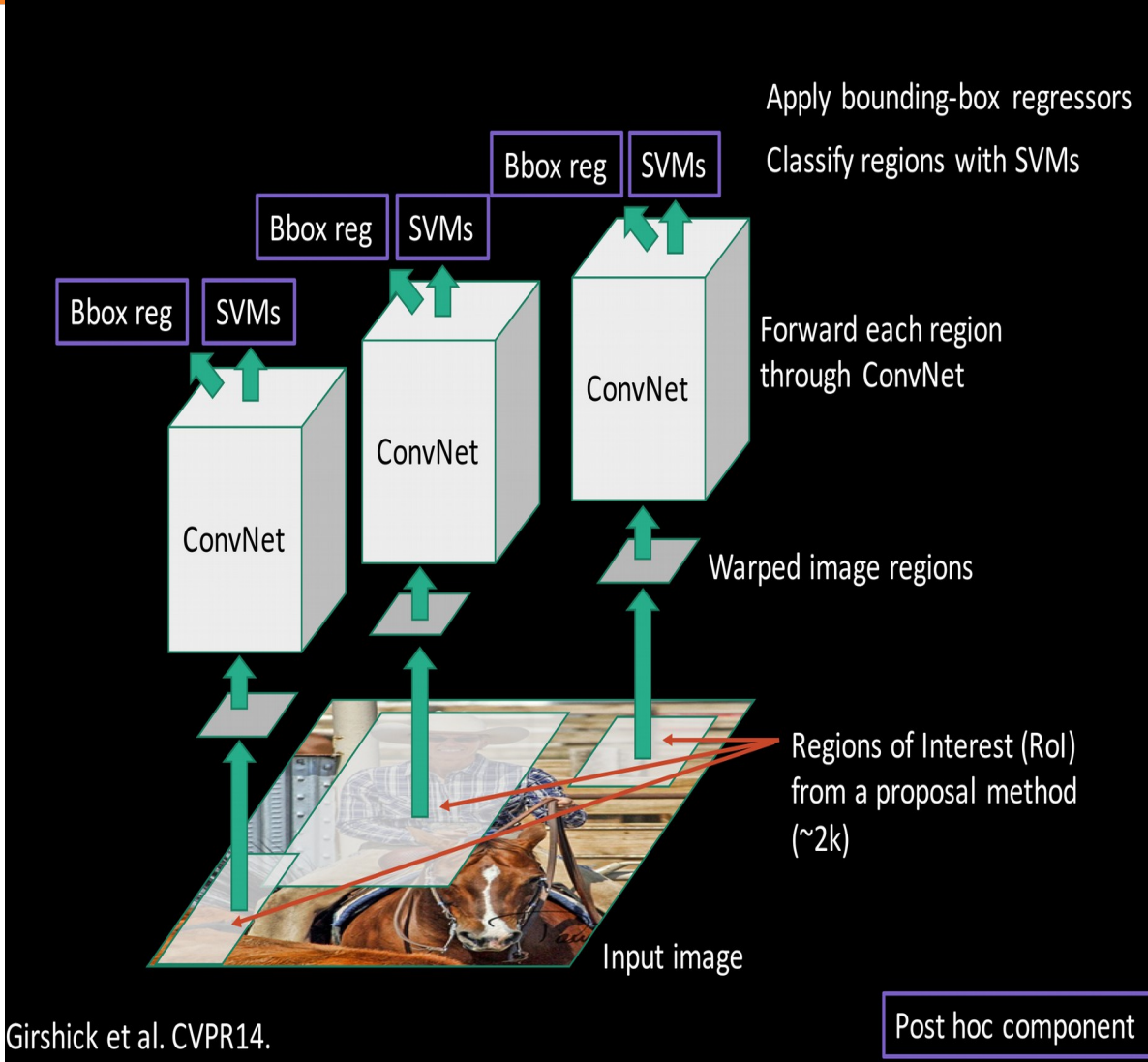
Uijlings et al, "Selective Search for Object Recognition", IJCV 2013

# Region Proposals: Many other choices

Method	Approach	Outputs Segments	Outputs Score	Control #proposals	Time (sec.)	Repeatability	Recall Results	Detection Results
Bing [18]	Window scoring		✓	✓	0.2	***	*	.
CPMC [19]	Grouping	✓	✓	✓	250	-	**	*
EdgeBoxes [20]	Window scoring		✓	✓	0.3	**	***	***
Endres [21]	Grouping	✓	✓	✓	100	-	***	**
Geodesic [22]	Grouping	✓		✓	1	*	***	**
MCG [23]	Grouping	✓	✓	✓	30	*	***	***
Objectness [24]	Window scoring		✓	✓	3	.	*	.
Rahtu [25]	Window scoring		✓	✓	3	.	.	*
RandomizedPrim's [26]	Grouping	✓		✓	1	*	*	**
Rantalankila [27]	Grouping	✓		✓	10	**	.	**
Rigor [28]	Grouping	✓		✓	10	*	**	**
SelectiveSearch [29]	Grouping	✓	✓	✓	10	**	***	***
Gaussian				✓	0	.	.	*
SlidingWindow				✓	0	***	.	.
Superpixels		✓			1	*	.	.
Uniform				✓	0	.	.	.

Hosang et al, "What makes for effective detection proposals?", PAMI 2015

# Putting it together: R-CNN

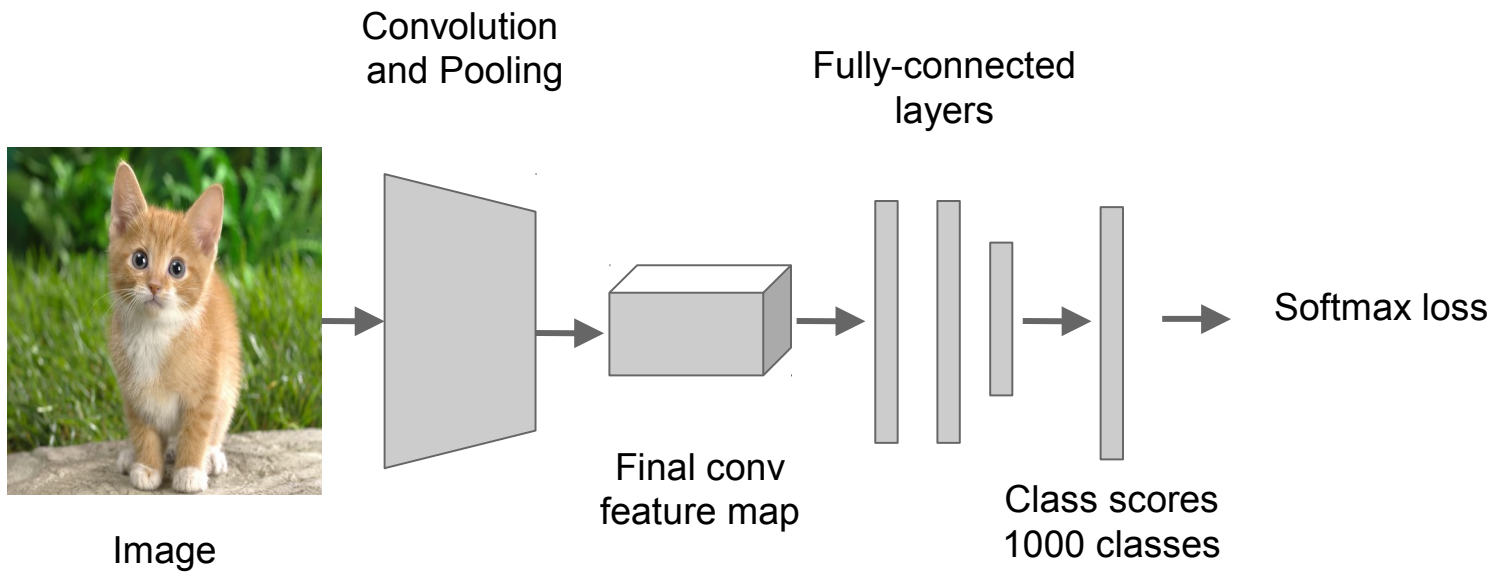


Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014

Slide credit: Ross Girshick

# R-CNN Training

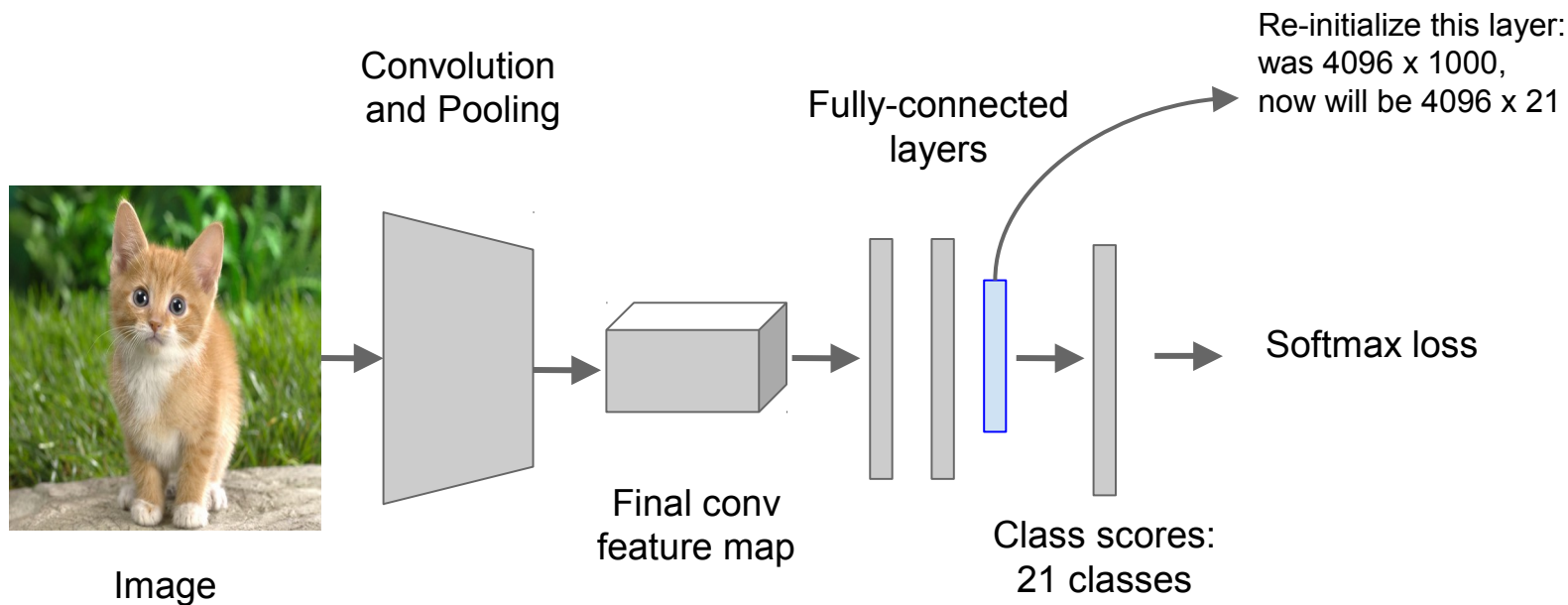
**Step 1:** Train (or download) a classification model for ImageNet (AlexNet)



# R-CNN Training

## Step 2: Fine-tune model for detection

- Instead of 1000 ImageNet classes, want 20 object classes + background
- Throw away final fully-connected layer, reinitialize from scratch
- Keep training model using positive / negative regions from detection images



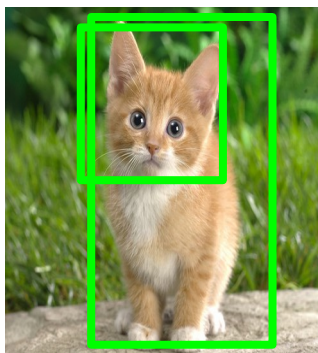
# R-CNN Training

## Step 3: Extract features

- Extract region proposals for all images
- For each region: warp to CNN input size, run forward through CNN, save pool5 features to disk
- Have a big hard drive: features are ~200GB for PASCAL dataset!



Image

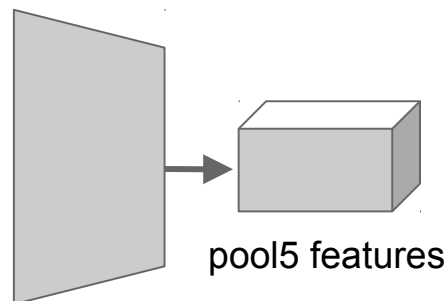


Region Proposals



Crop + Warp

Convolution  
and Pooling



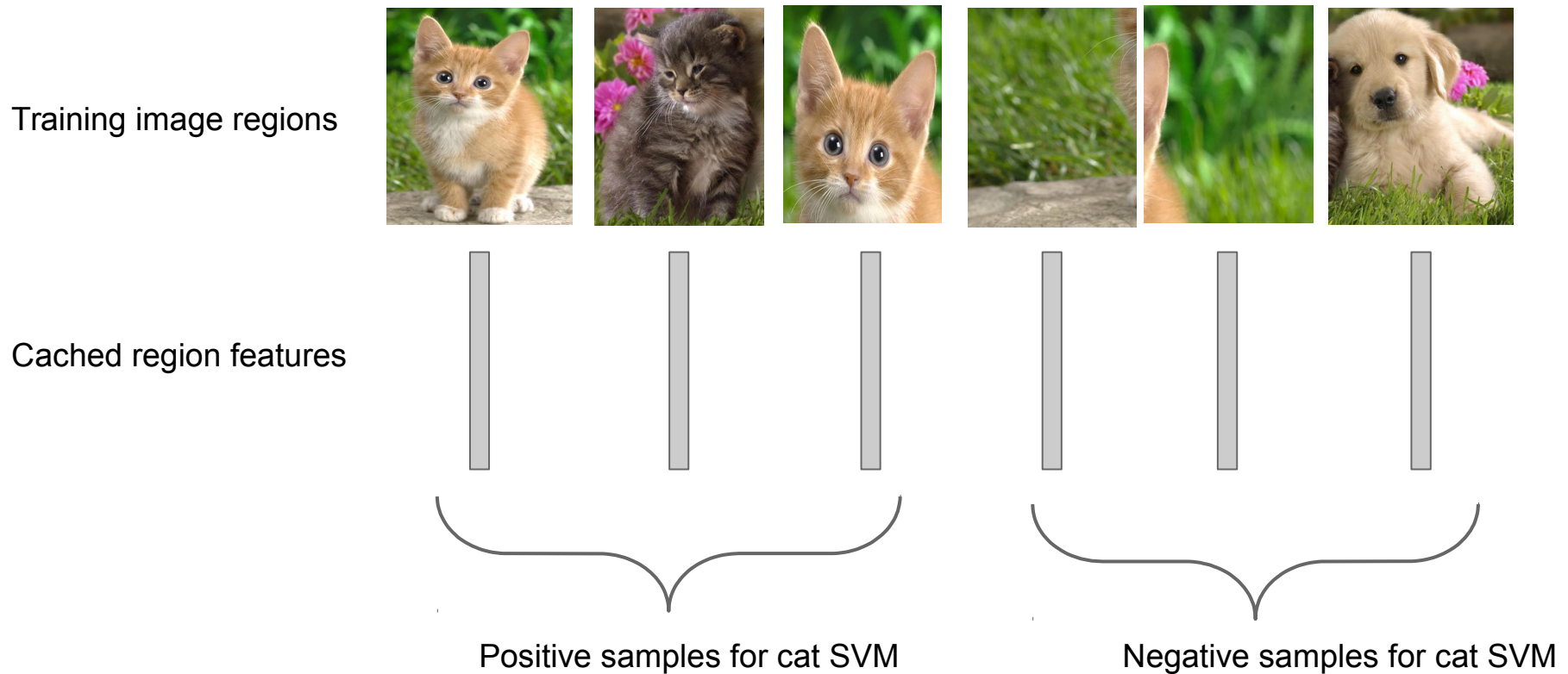
Forward pass



Save to disk

# R-CNN Training

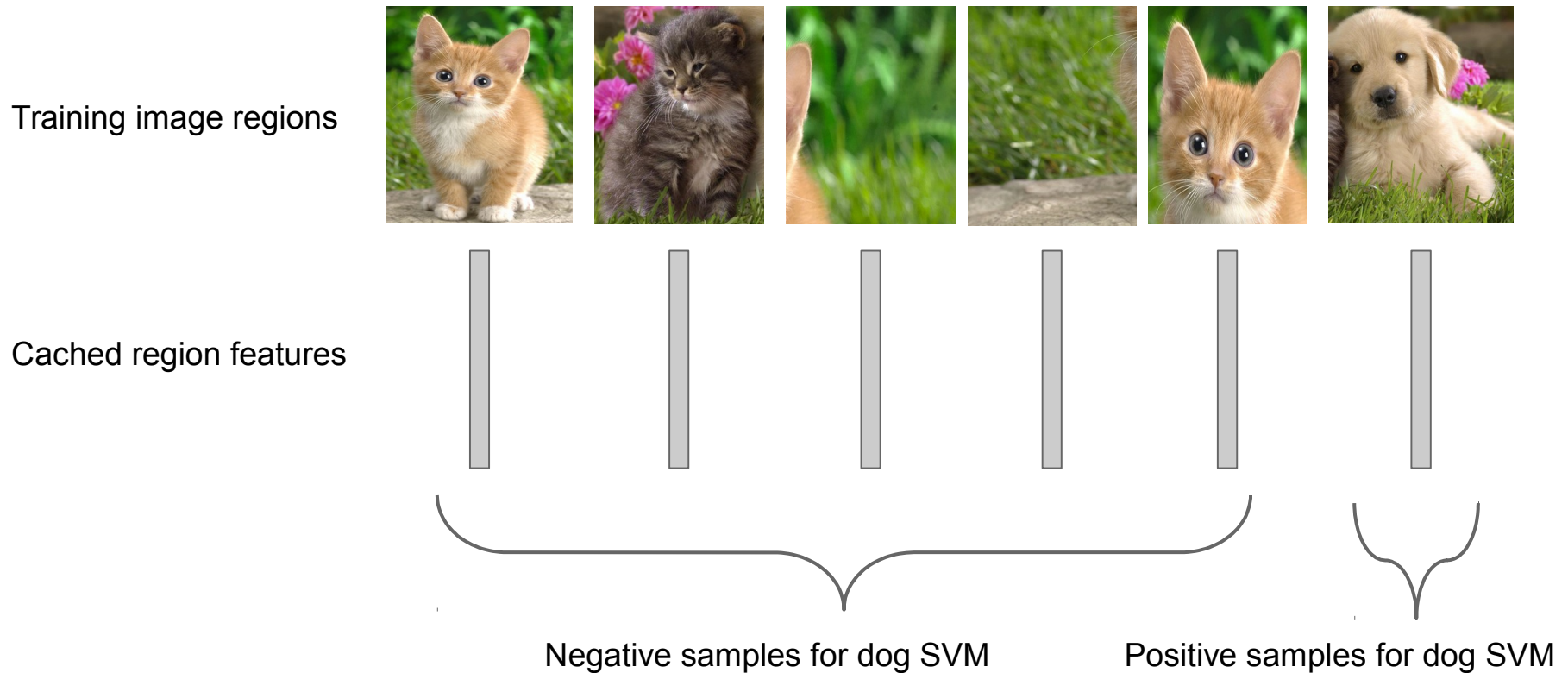
**Step 4:** Train one binary SVM per class to classify region features





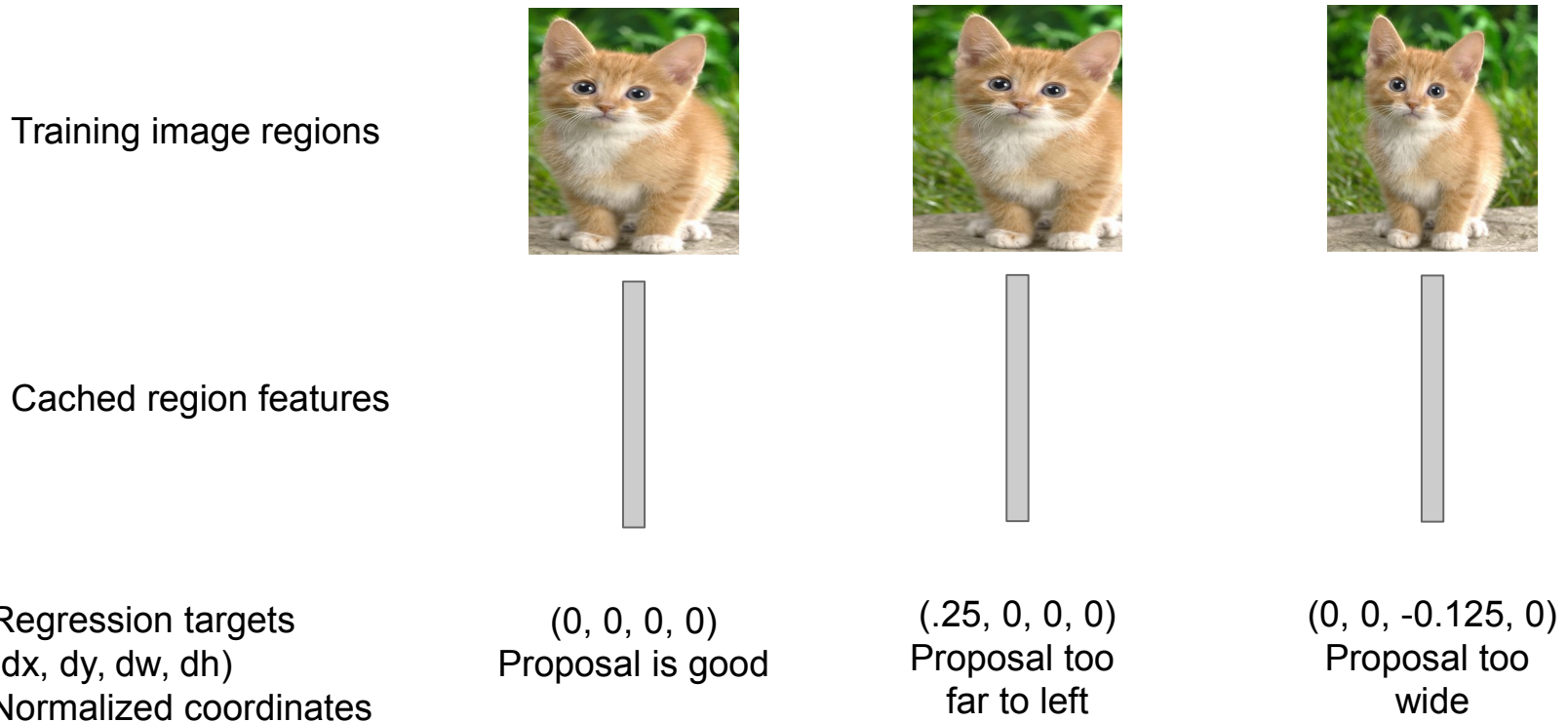
# R-CNN Training

**Step 4:** Train one binary SVM per class to classify region features



# R-CNN Training

**Step 5** (bbox regression): For each class, train a linear regression model to map from cached features to offsets to GT boxes to make up for “slightly wrong” proposals



# Object Detection: Datasets

	PASCAL VOC (2010)	ImageNet Detection (ILSVRC 2014)	MS-COCO (2014)
Number of classes	20	<b>200</b>	80
Number of images (train + val)	~20k	<b>~470k</b>	~120k
Mean objects per image	2.4	1.1	<b>7.2</b>

# Object Detection: Evaluation

We use a metric called “mean average precision” (mAP)

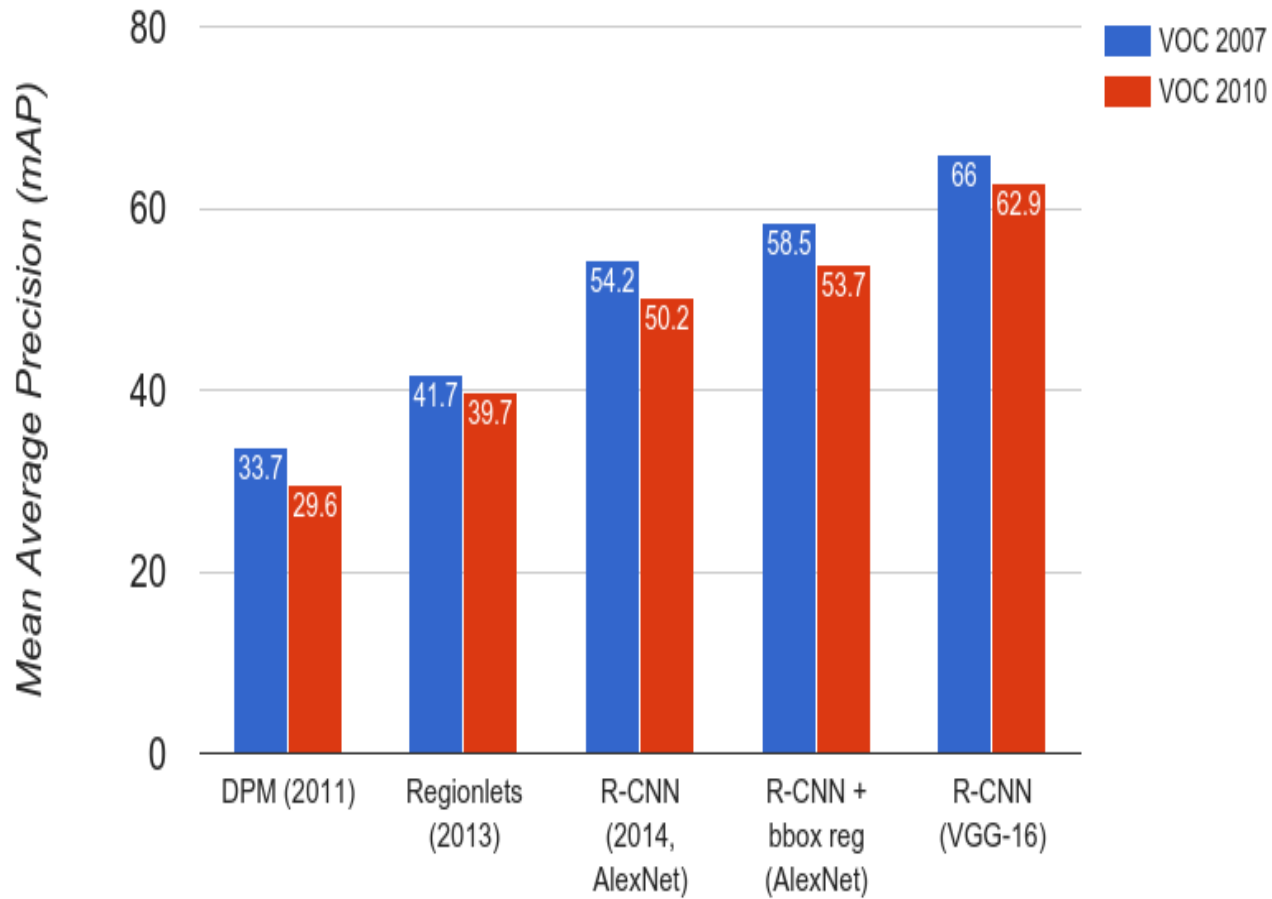
Compute average precision (AP) separately for each class, then average over classes

A detection is a true positive if it has IoU with a ground-truth box greater than some threshold (usually 0.5) (mAP@0.5)

Combine all detections from all test images to draw a precision / recall curve for each class; AP is area under the curve

TL;DR mAP is a number from 0 to 100; high is good

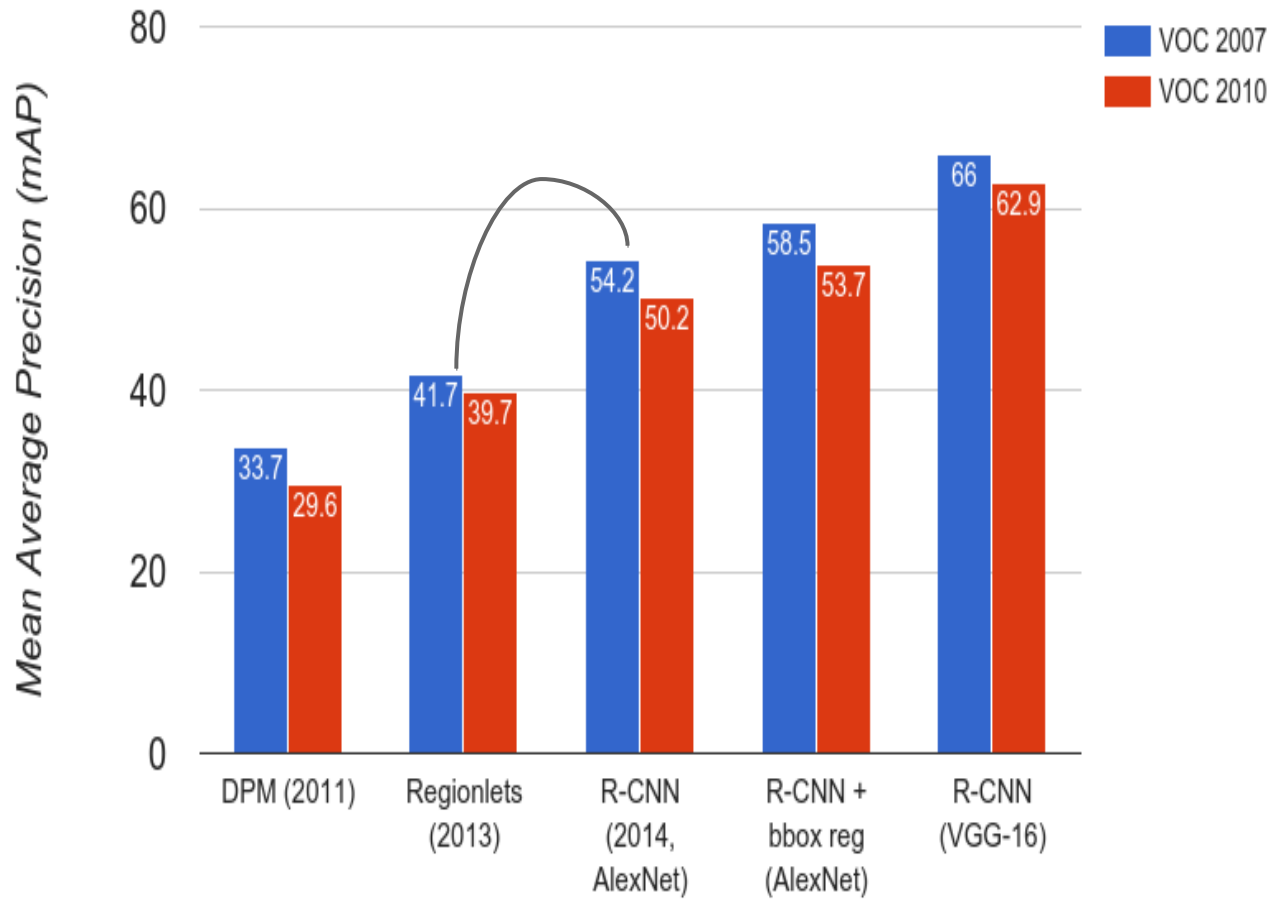
# R-CNN Results



Wang et al, "Regionlets for Generic Object Detection", ICCV 2013

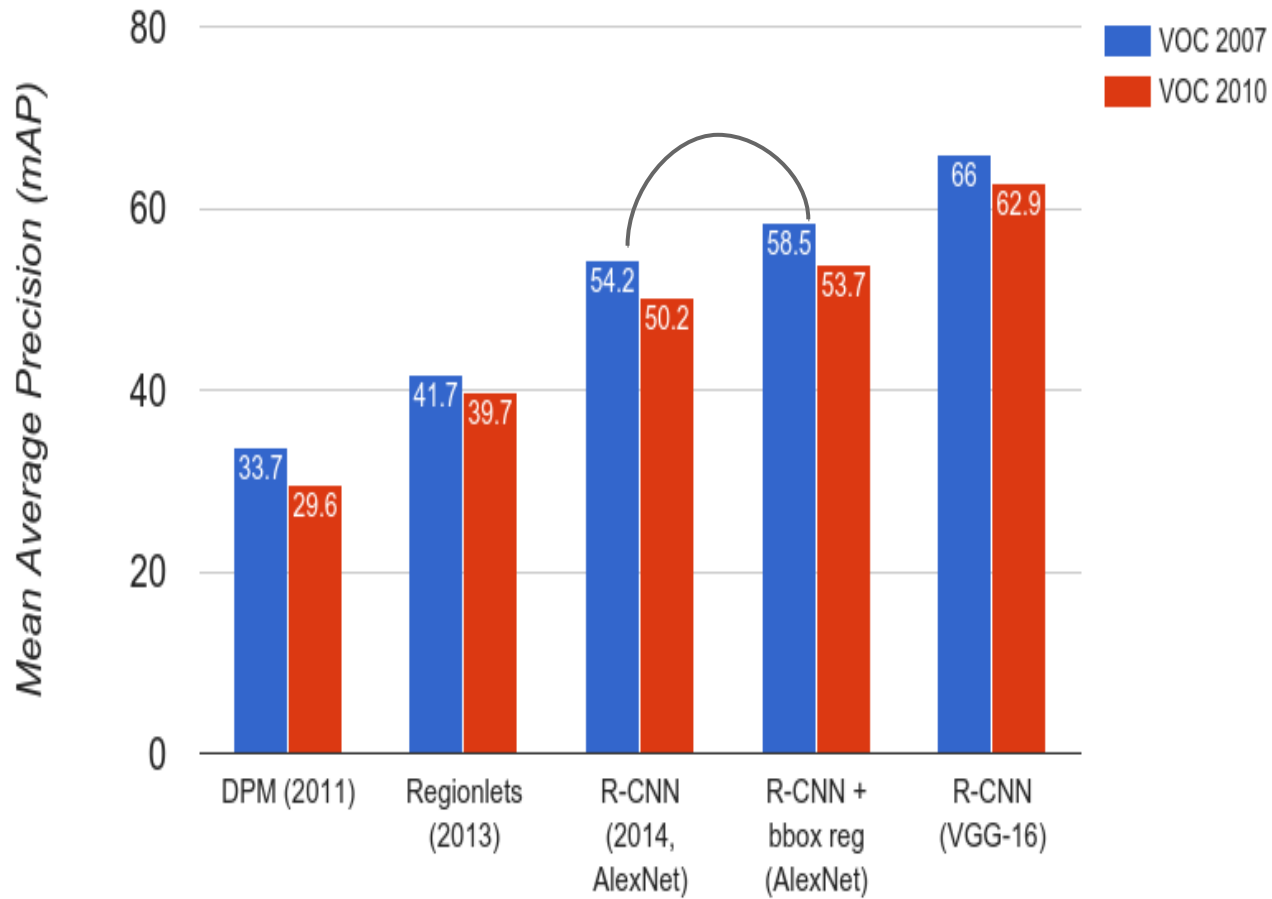
# R-CNN Results

Big improvement compared to pre-CNN methods



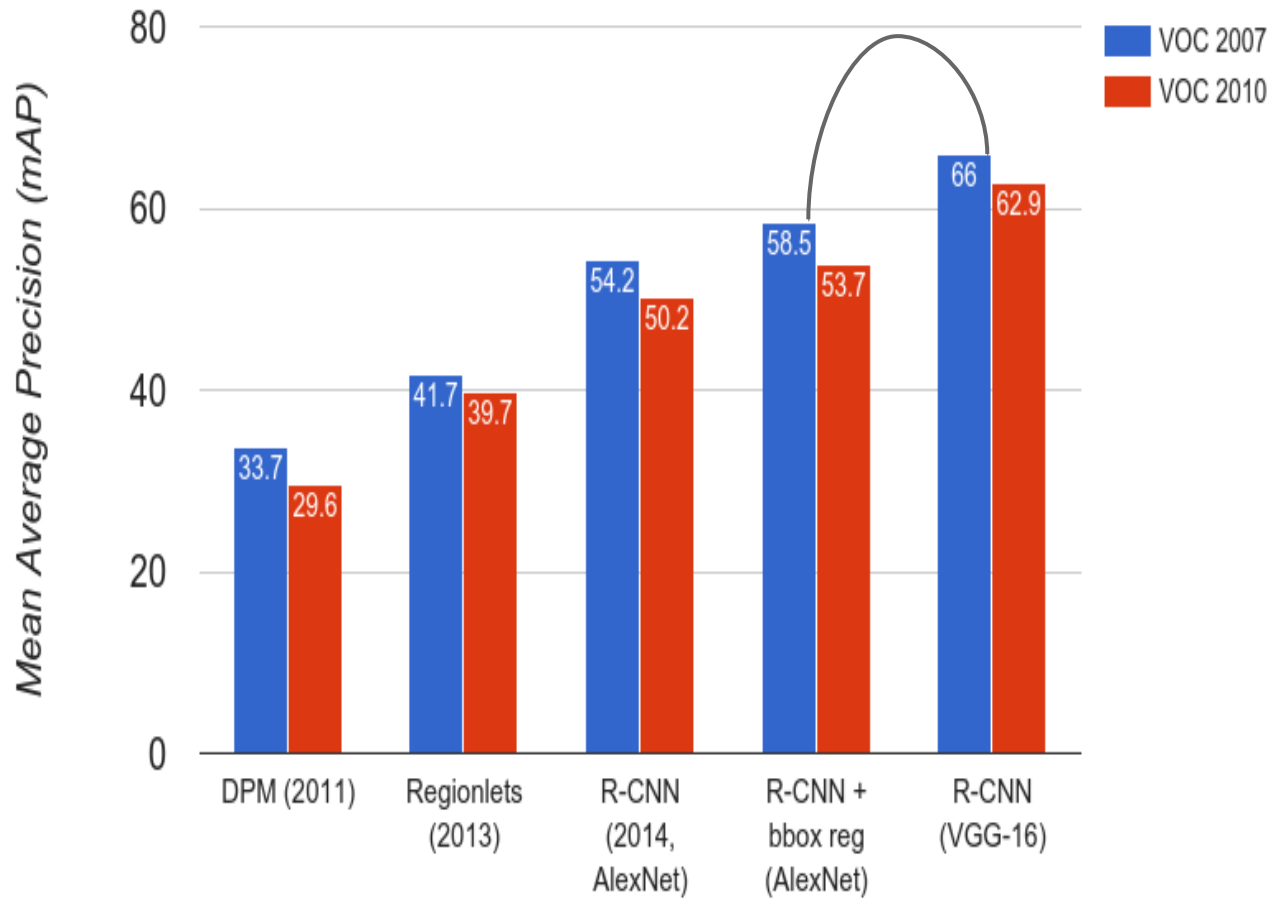
# R-CNN Results

Bounding box regression helps a bit



# R-CNN Results

Features from a deeper network help a lot



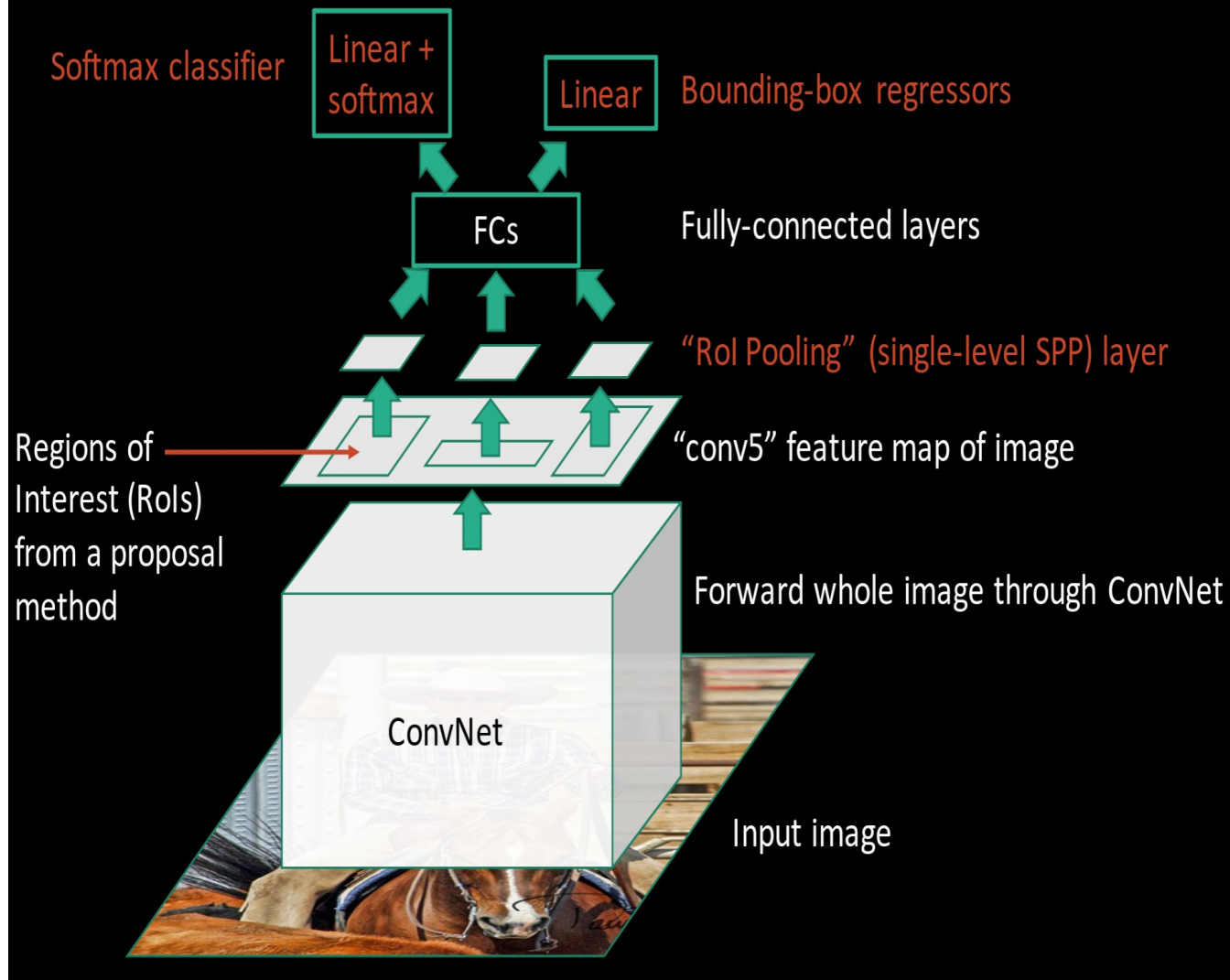


# R-CNN Problems

---

1. Slow at test-time: need to run full forward pass of CNN for each region proposal
2. SVMs and regressors are post-hoc: CNN features not updated in response to SVMs and regressors
3. Complex multistage training pipeline

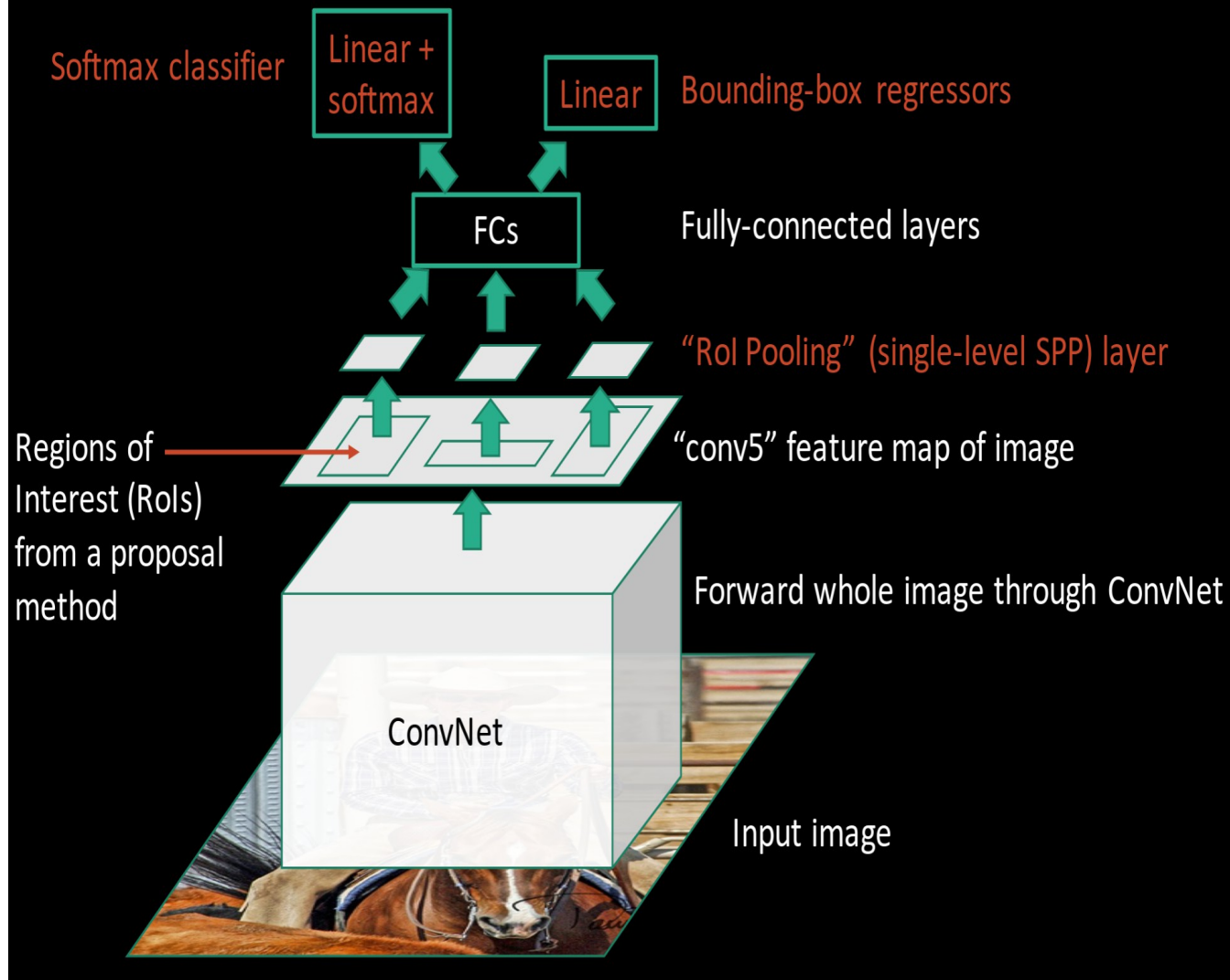
# Fast R-CNN (test time)



Girshick, "Fast R-CNN", ICCV 2015

Slide credit: Ross Girshick

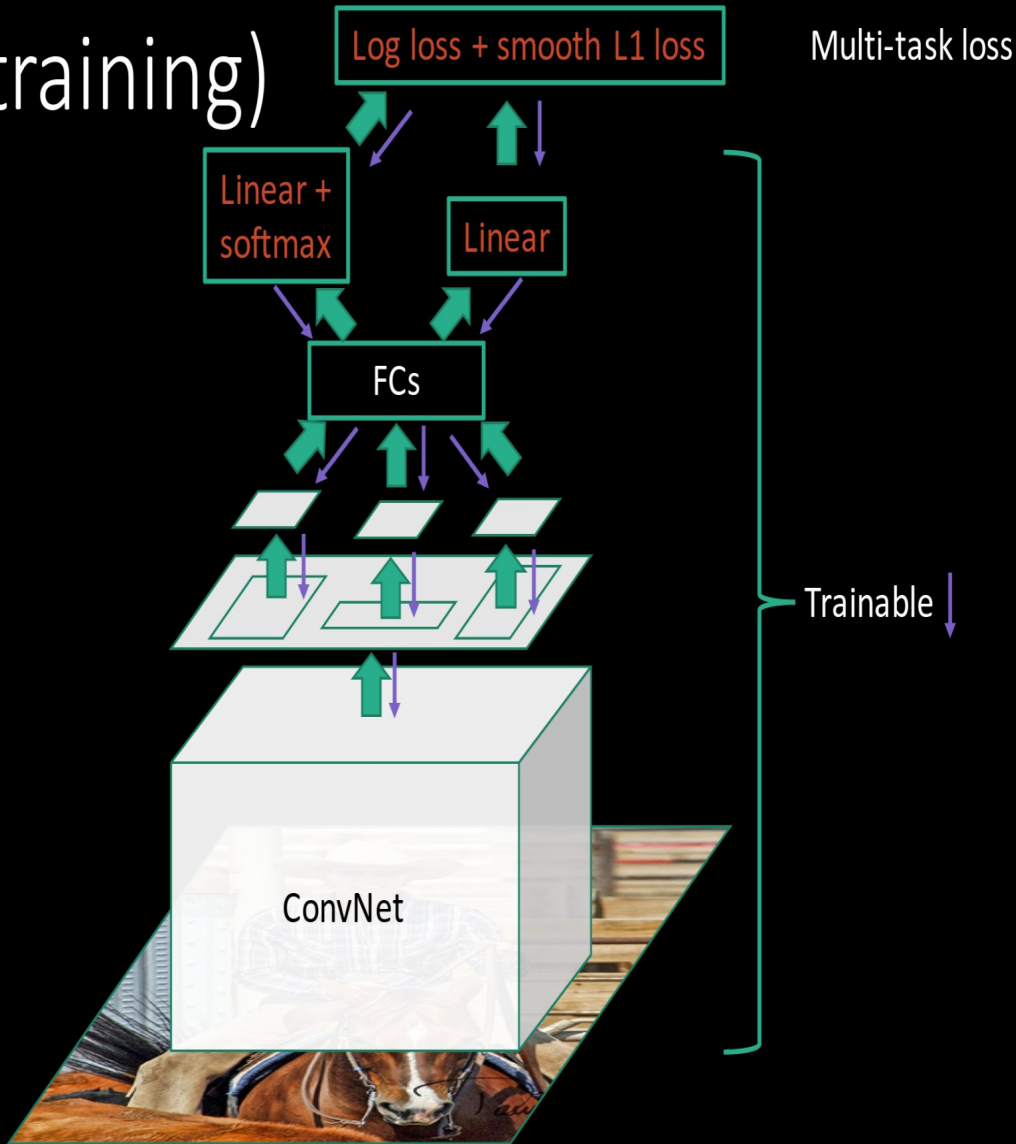
# Fast R-CNN (test time)



**R-CNN Problem #1:**  
Slow at test-time due to independent forward passes of the CNN

**Solution:**  
Share computation of convolutional layers between proposals for an image

# Fast R-CNN (training)



## R-CNN Problem #2:

Post-hoc training: CNN not updated in response to final classifiers and regressors

## R-CNN Problem #3:

Complex training pipeline

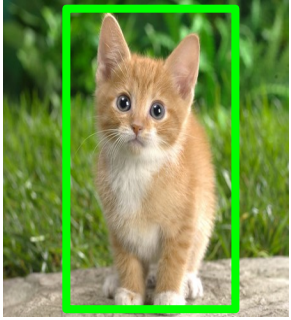
## Solution:

Just train the whole system end-to-end all at once!

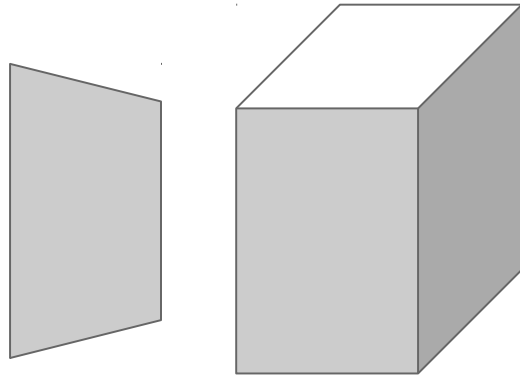
Slide credit: Ross Girschick

# Fast R-CNN: Region of Interest Pooling

Convolution  
and Pooling

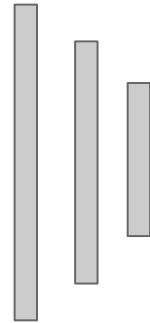


Hi-res input  
image:  
 $3 \times 800 \times 600$   
with region  
proposal



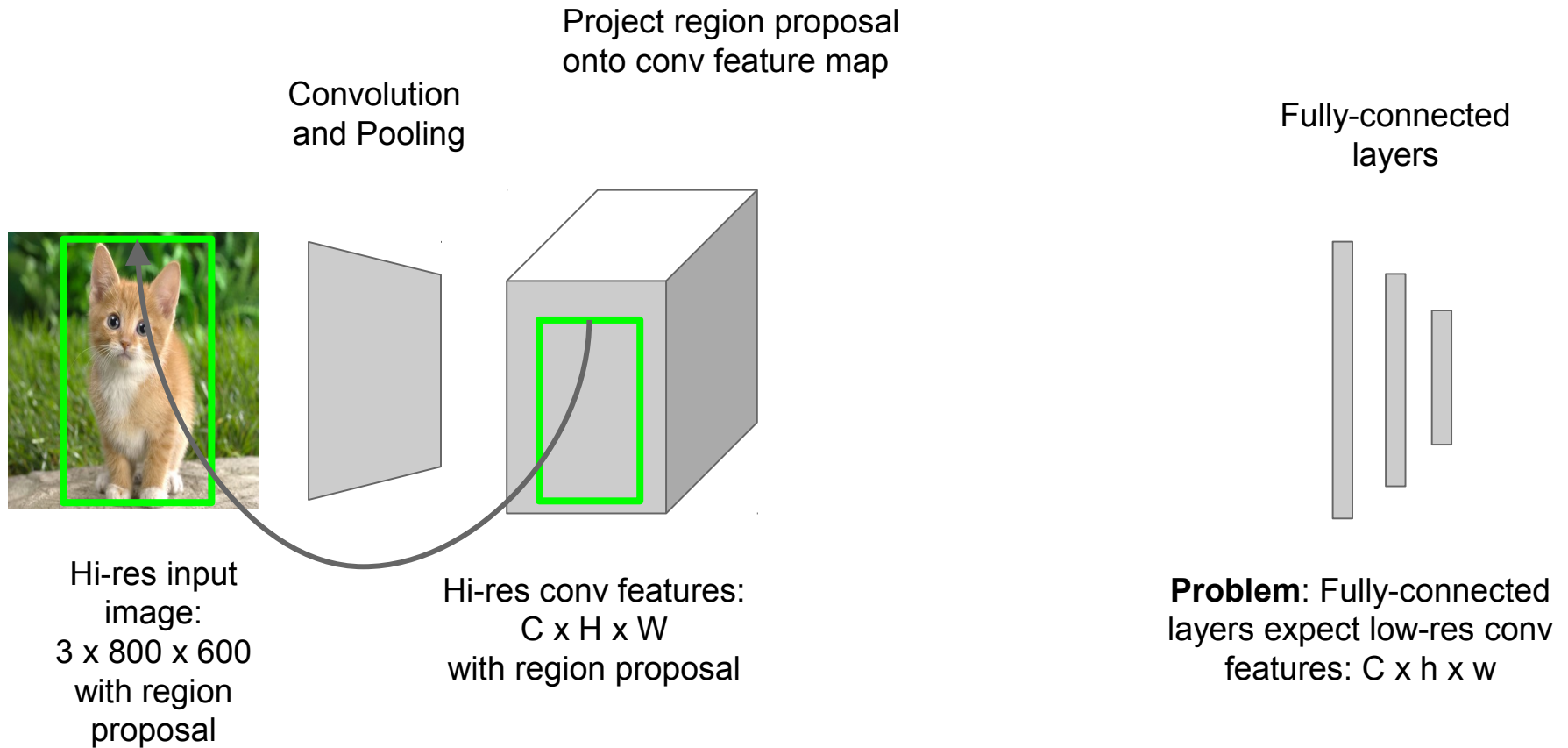
Hi-res conv features:  
 $C \times H \times W$   
with region proposal

Fully-connected  
layers

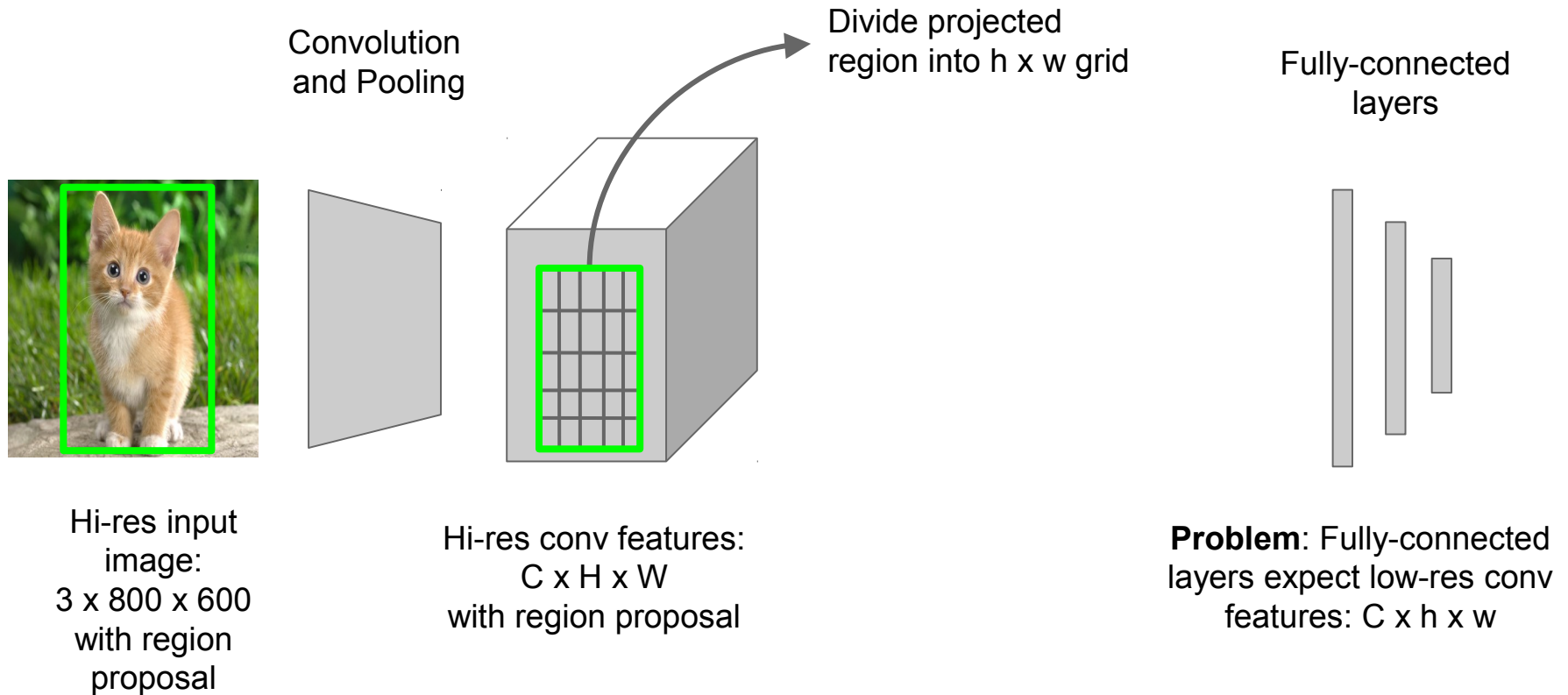


**Problem:** Fully-connected  
layers expect low-res conv  
features:  $C \times h \times w$

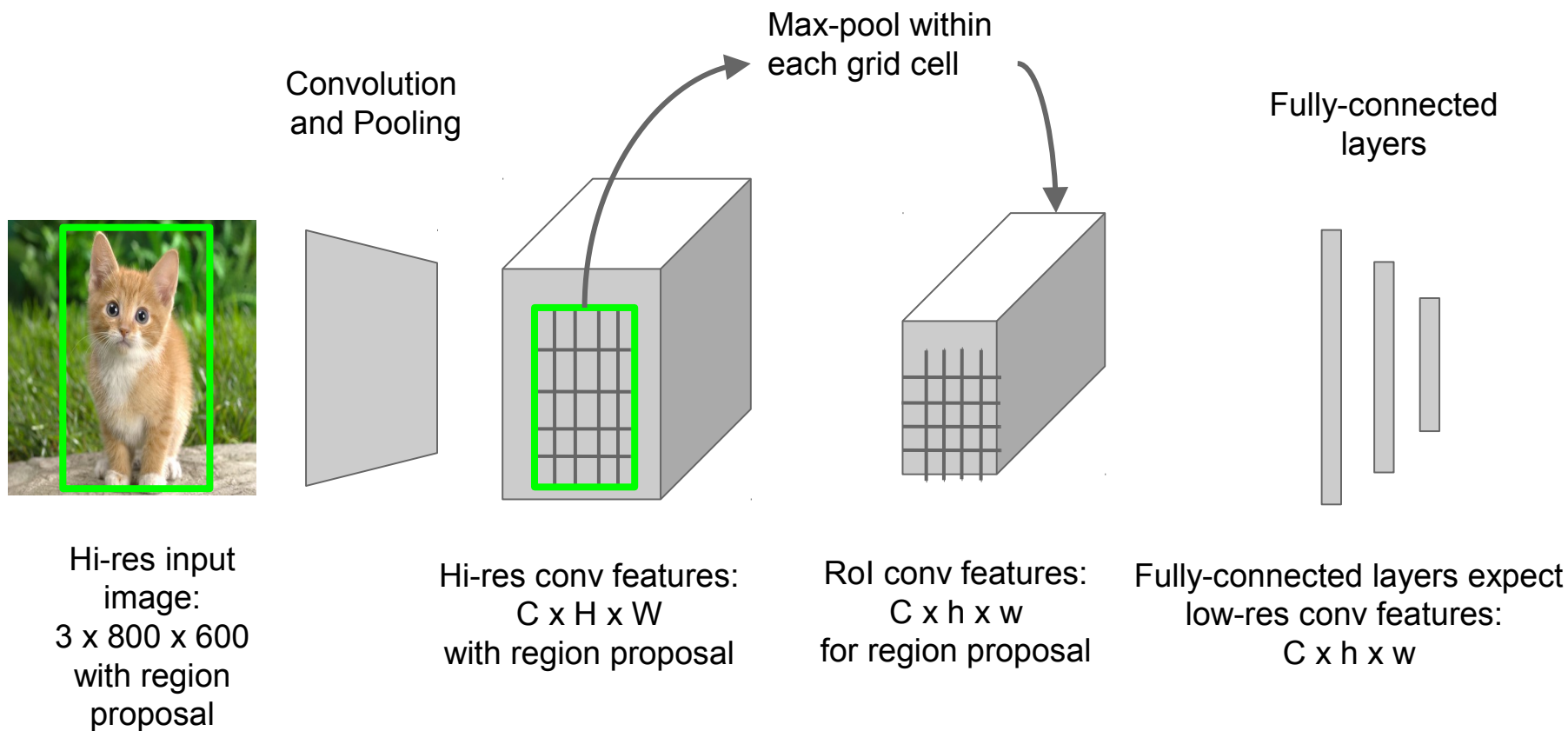
# Fast R-CNN: Region of Interest Pooling



# Fast R-CNN: Region of Interest Pooling

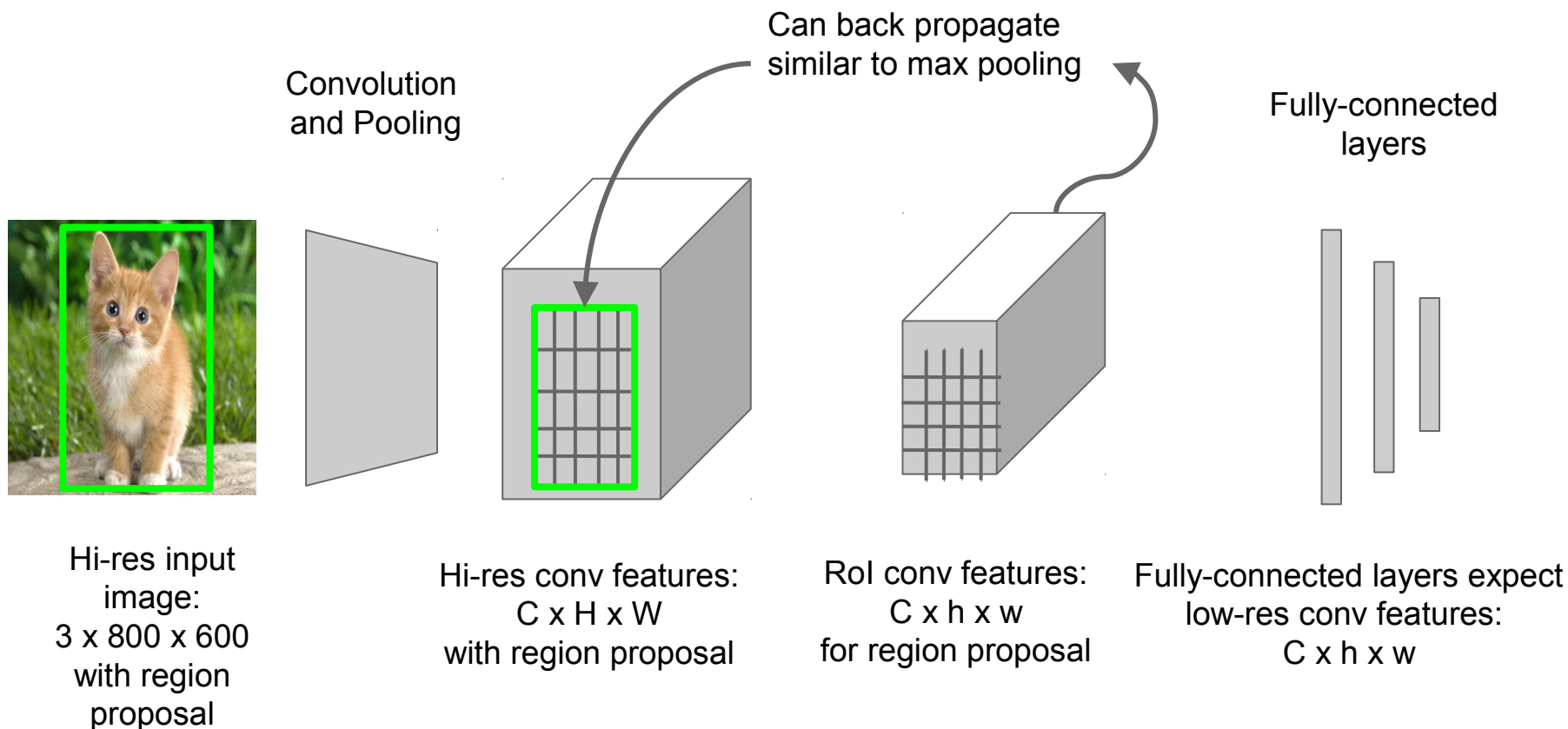


# Fast R-CNN: Region of Interest Pooling





# Fast R-CNN: Region of Interest Pooling



# Fast R-CNN Results

Faster!

	<b>R-CNN</b>	<b>Fast R-CNN</b>
Training Time:	84 hours	<b>9.5 hours</b>
(Speedup)	1x	<b>8.8x</b>

Using VGG-16 CNN on Pascal VOC 2007 dataset

# Fast R-CNN Results

	R-CNN	Fast R-CNN
Faster!		
Training Time:	84 hours	<b>9.5 hours</b>
(Speedup)	1x	<b>8.8x</b>
FASTER!		
Test time per image	47 seconds	<b>0.32 seconds</b>
(Speedup)	1x	<b>146x</b>

Using VGG-16 CNN on Pascal VOC 2007 dataset

# Fast R-CNN Results

	R-CNN	Fast R-CNN
Faster!		
Training Time:	84 hours	<b>9.5 hours</b>
(Speedup)	1x	<b>8.8x</b>
FASTER!		
Test time per image	47 seconds	<b>0.32 seconds</b>
(Speedup)	1x	<b>146x</b>
Better!		
mAP (VOC 2007)	66.0	<b>66.9</b>

Using VGG-16 CNN on Pascal VOC 2007 dataset

# Fast R-CNN Problem:

Test-time speeds don't include region proposals

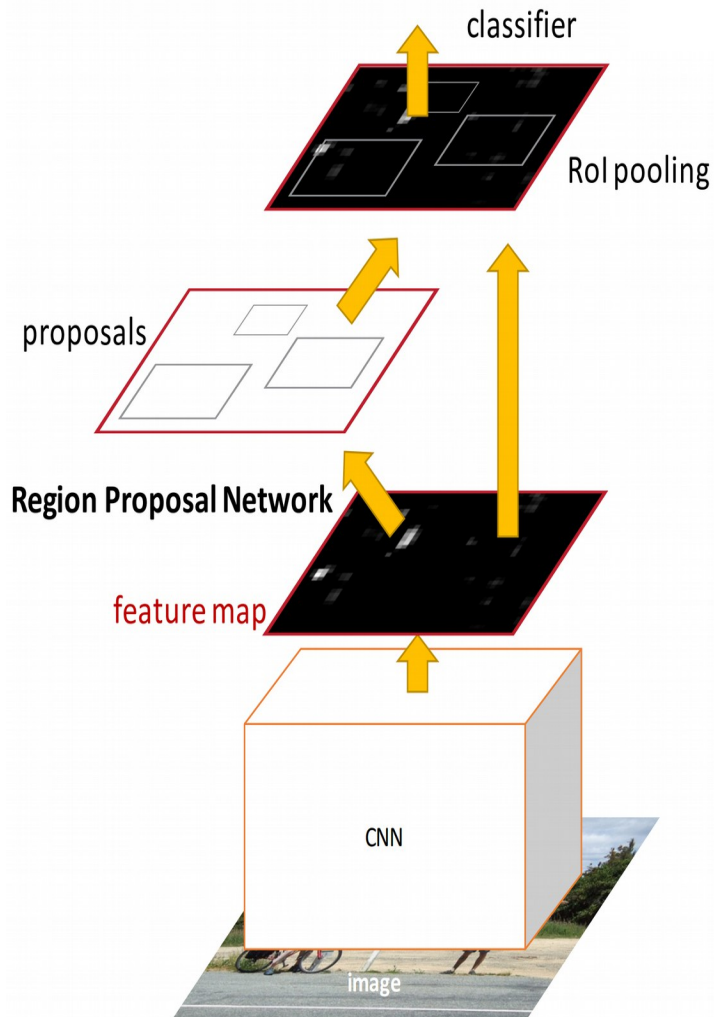
	<b>R-CNN</b>	<b>Fast R-CNN</b>
Test time per image	47 seconds	<b>0.32 seconds</b>
(Speedup)	1x	<b>146x</b>
Test time per image with Selective Search	50 seconds	<b>2 seconds</b>
(Speedup)	1x	<b>25x</b>

# Fast R-CNN Problem Solution:

Test-time speeds don't include region proposals  
Just make the CNN do region proposals too!

	R-CNN	Fast R-CNN
Test time per image	47 seconds	<b>0.32 seconds</b>
(Speedup)	1x	<b>146x</b>
Test time per image with Selective Search	50 seconds	<b>2 seconds</b>
(Speedup)	1x	<b>25x</b>

# Faster R-CNN:



## Insert a **Region Proposal Network (RPN)** after the last

convolutional layer

RPN trained to produce region proposals directly; no need for external region proposals!

After RPN, use RoI Pooling and an upstream classifier and bbox regressor just like Fast R-CNN

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

Slide credit: Ross Girshick

# Faster R-CNN: Region Proposal Network

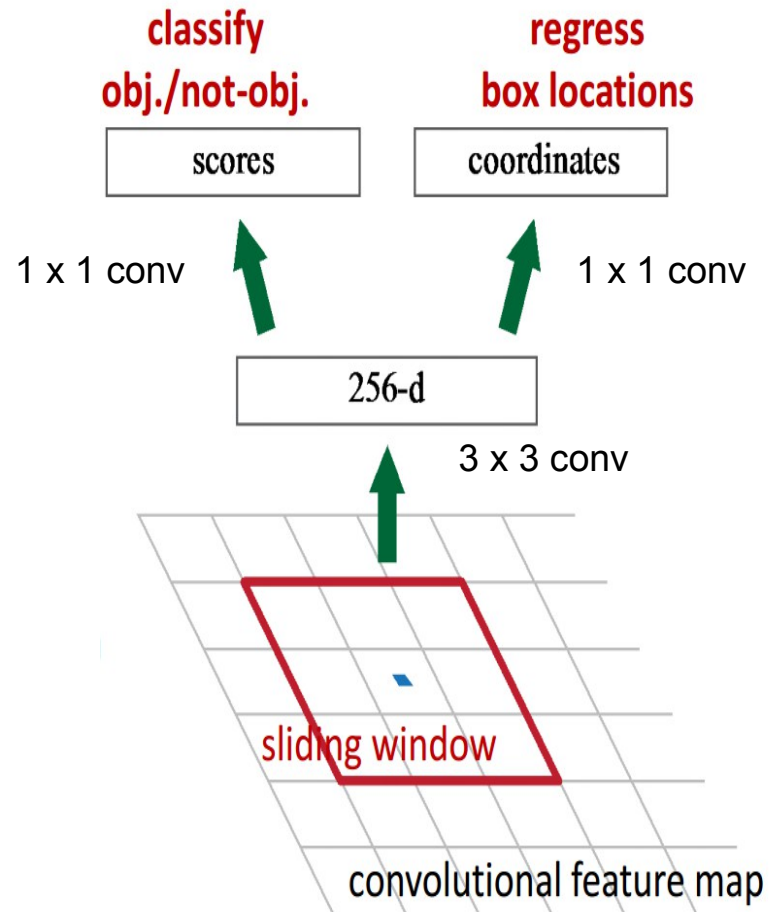
Slide a small window on the feature map

Build a small network for:

- classifying object or not-object, and
- regressing bbox locations

Position of the sliding window provides localization information with reference to the image

Box regression provides finer localization information with reference to this sliding window



Slide credit: Kaiming He



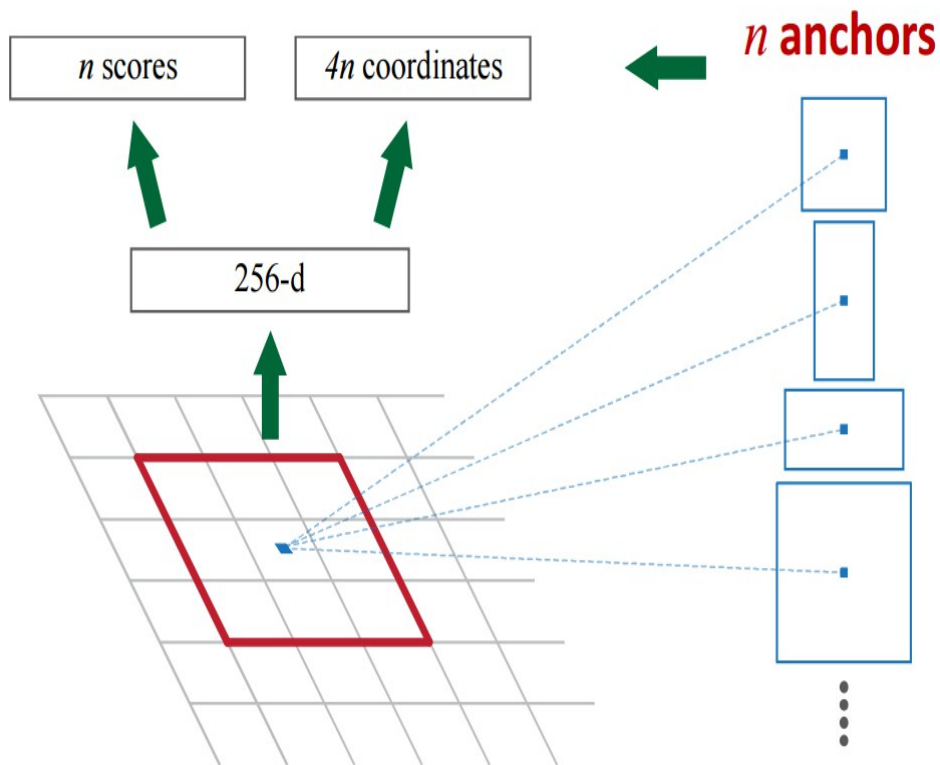
# Faster R-CNN: Region Proposal Network

Use  $N$  anchor boxes at each location

Anchors are **translation invariant**: use the same ones at every location

Regression gives offsets from anchor boxes

Classification gives the probability that each (regressed) anchor shows an object



# Faster R-CNN: Training

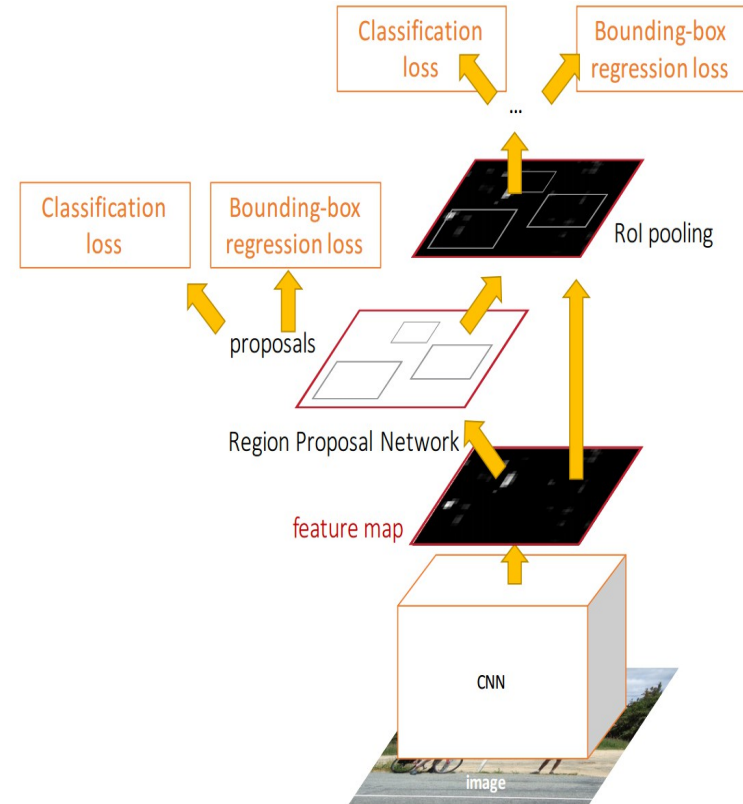
In the paper: Ugly pipeline

- Use alternating optimization to train RPN, then Fast R-CNN with RPN proposals, etc.
- More complex than it has to be

Since publication: Joint training!

One network, four losses

- RPN classification (anchor good / bad)
- RPN regression (anchor  $\rightarrow$  proposal)
- Fast R-CNN classification (over classes)
- Fast R-CNN regression (proposal  $\rightarrow$  box)



Slide credit: Ross Girschick

# Faster R-CNN: Results

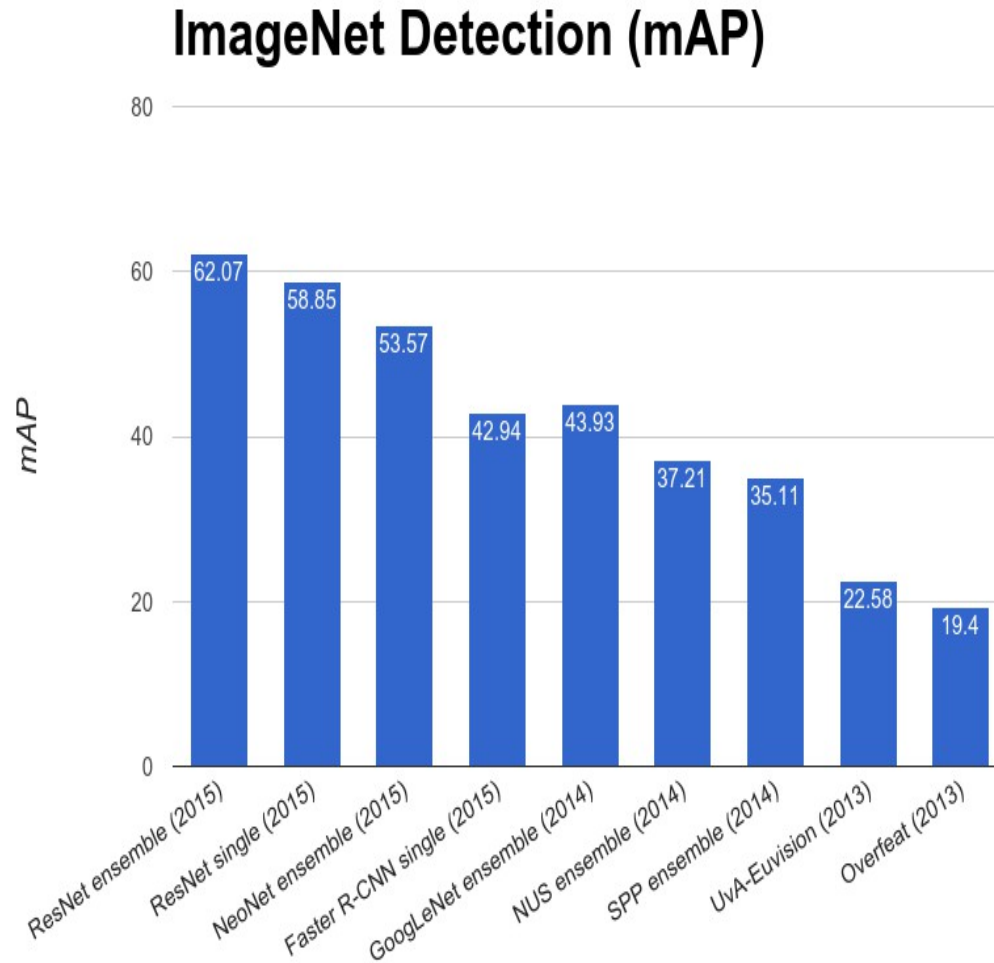
	<b>R-CNN</b>	<b>Fast R-CNN</b>	<b>Faster R-CNN</b>
Test time per image (with proposals)	50 seconds	2 seconds	<b>0.2 seconds</b>
(Speedup)	1x	25x	<b>250x</b>
mAP (VOC 2007)	66.0	<b>66.9</b>	<b>66.9</b>

# Object Detection State-of-the-art: ResNet 101 + Faster R-CNN + some extras

training data	COCO train		COCO trainval	
test data	COCO val		COCO test-dev	
mAP	@.5	@[.5, .95]	@.5	@[.5, .95]
baseline Faster R-CNN (VGG-16)	41.5	21.2		
baseline Faster R-CNN (ResNet-101)	48.4	27.2		
+box refinement	49.9	29.9		
+context	51.1	30.0	53.3	32.2
+multi-scale testing	53.8	32.5	<b>55.7</b>	<b>34.9</b>
ensemble			<b>59.0</b>	<b>37.4</b>

He et. al, "Deep Residual Learning for Image Recognition", arXiv 2015

# ImageNet Detection 2013 - 2015



# YOLO: You Only Look Once Detection as Regression

Divide image into  $S \times S$  grid

Within each grid cell predict:

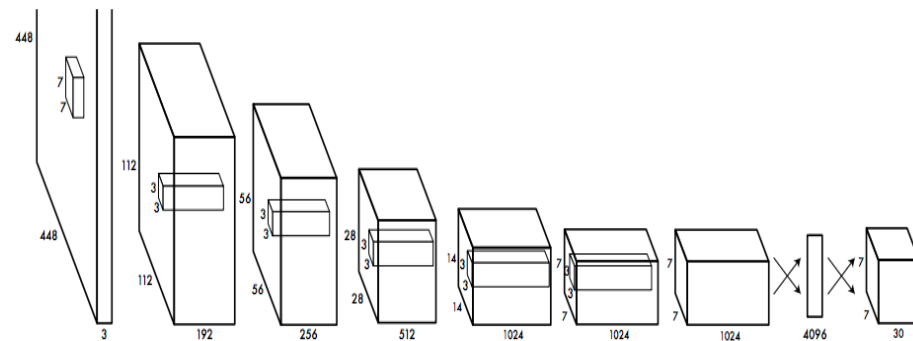
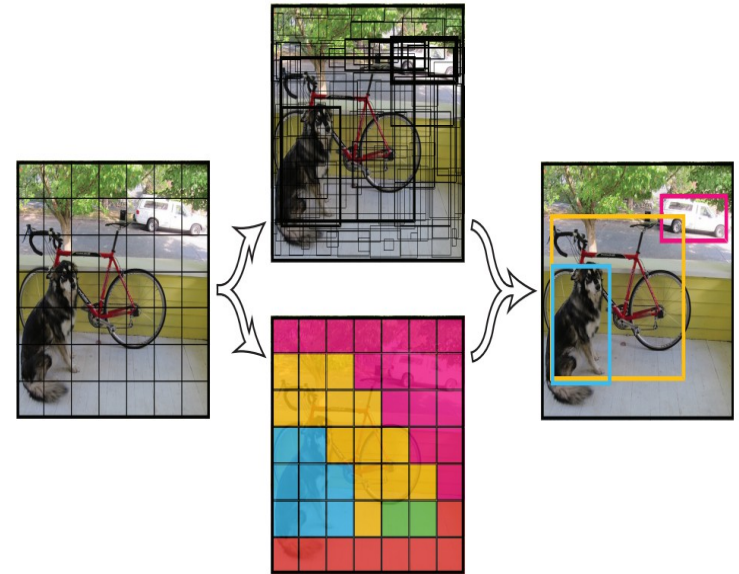
B Boxes: 4 coordinates + confidence

Class scores: C numbers

Regression from image to  
 $7 \times 7 \times (5 * B + C)$  tensor

Direct prediction using a CNN

Redmon et al, "You Only Look Once:  
Unified, Real-Time Object Detection", arXiv 2015



# YOLO: You Only Look Once

## Detection as Regression

Faster than Faster R-CNN, but not as good

Redmon et al, "You Only Look Once: Unified, Real-Time Object Detection", arXiv 2015

Real-Time Detectors	Train	mAP	FPS
100Hz DPM [30]	2007	16.0	100
30Hz DPM [30]	2007	26.1	30
Fast YOLO	2007+2012	52.7	<b>155</b>
YOLO	2007+2012	<b>63.4</b>	45
<hr/>			
Less Than Real-Time			
Fastest DPM [37]	2007	30.4	15
R-CNN Minus R [20]	2007	53.5	6
Fast R-CNN [14]	2007+2012	70.0	0.5
Faster R-CNN VGG-16[27]	2007+2012	73.2	7
Faster R-CNN ZF [27]	2007+2012	62.1	18

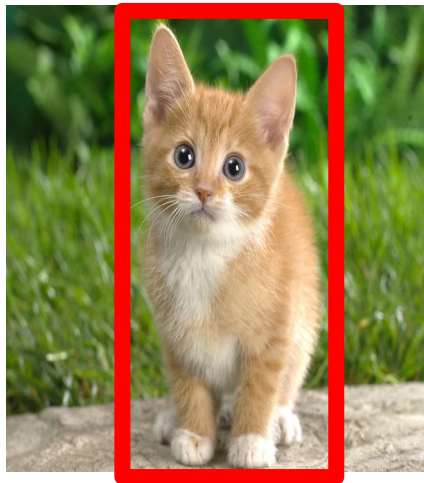
# Computer Vision Tasks

**Classification**

**Classification  
+ Localization**

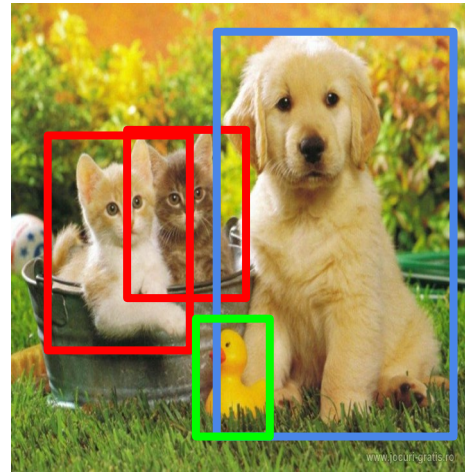
**Object Detection**

**Segmentation**



CAT

CAT



CAT, DOG, DUCK



CAT, DOG, DUCK

Single object

Multiple objects



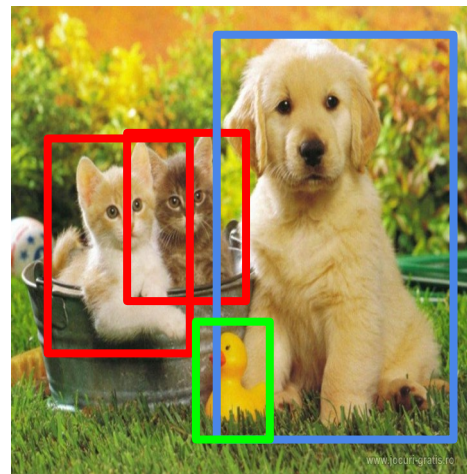
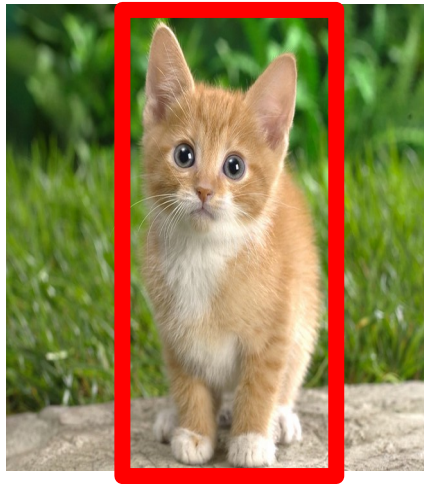
# Computer Vision Tasks

Classification

**Classification  
+ Localization**

**Object Detection**

Segmentation



Lecture 8

# Computer Vision Tasks

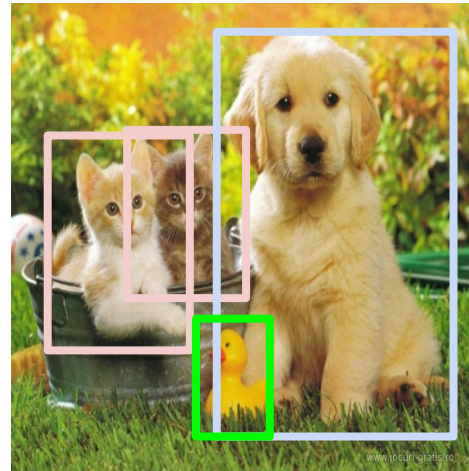
Classification



Classification  
+ Localization



Object Detection



Segmentation



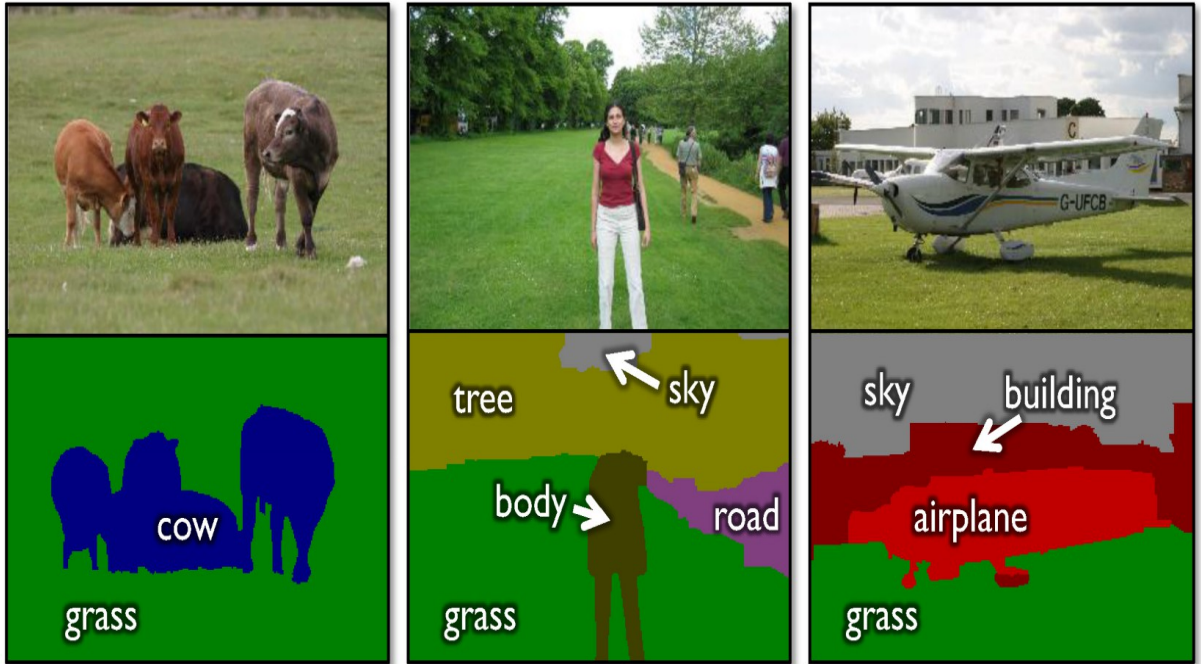
Today

# Semantic Segmentation

Label every pixel!

Don't differentiate instances (cows)

Classic computer vision problem



object classes	building	grass	tree	cow	sheep	sky	airplane	water	face	car	
	bicycle	flower	sign	bird	book	chair	road	cat	dog	body	boat

Figure credit: Shotton et al, "TextonBoost for Image Understanding: Multi-Class Object Recognition and Segmentation by Jointly Modeling Texture, Layout, and Context", IJCV 2007

# Instance Segmentation

Detect instances,  
give category, label  
pixels

“simultaneous  
detection and  
segmentation” (SDS)

Lots of recent work  
(MS-COCO)

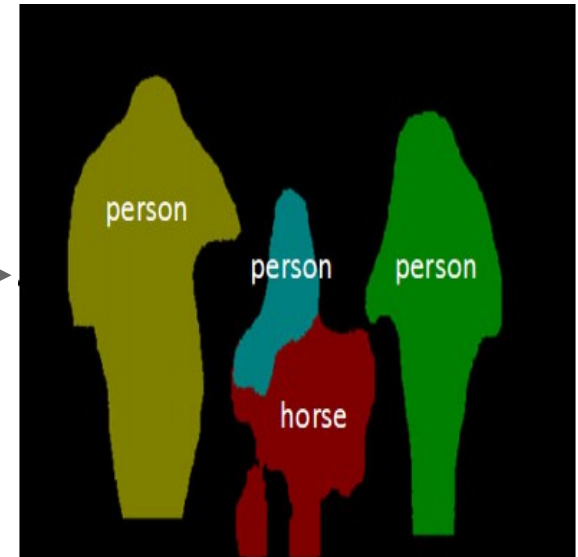


Figure credit: Dai et al, “Instance-aware Semantic Segmentation via Multi-task Network Cascades”, arXiv 2015