

CNN Applications and Architectures

Chetan Arora

IIIT Delhi

Adapted from the slides of
Fei Fei Li, Andrej Karpathy, Justin Johnson, A. Torralba,
Svetlana Lazebnik, Ross Girshick, Kaiming He and Rob Fergus

Visual Recognition Tasks

Visual Recognition Tasks



Image Classification

Is it a natural or man made scene

Is it a forest or a beach?

Does this image contain a building?

Visual Recognition Tasks



Object Detection

Does this image contain a building?
[where]

Which objects does this image contain?

Visual Recognition Tasks



**Object Detection
[pixel wise
localization]**

Which pixels are building?

Visual Recognition Tasks



Semantic Segmentation

Classifying each pixel
in the image

Applications: Computational Photography



Face Detection



Dynamic Range Enhancement

Applications: Object Attributes

Building:
42 m height
100 m away



Car:
Police Car
Frontal View
Autonomous
and Assistive
Driving

Applications: Instance Recognition



Does this image contains “India Gate”?

Recognizing landmarks in images

Recognizing products in supermarket

Applications: Assistive Vision



Applications: Security and Surveillance

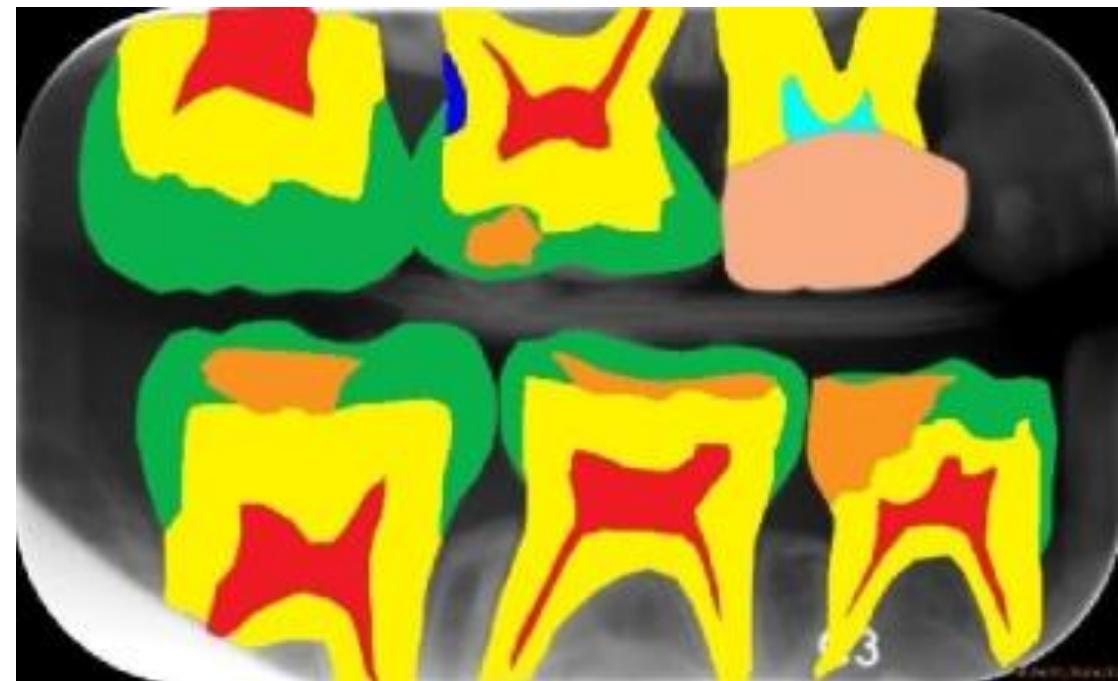
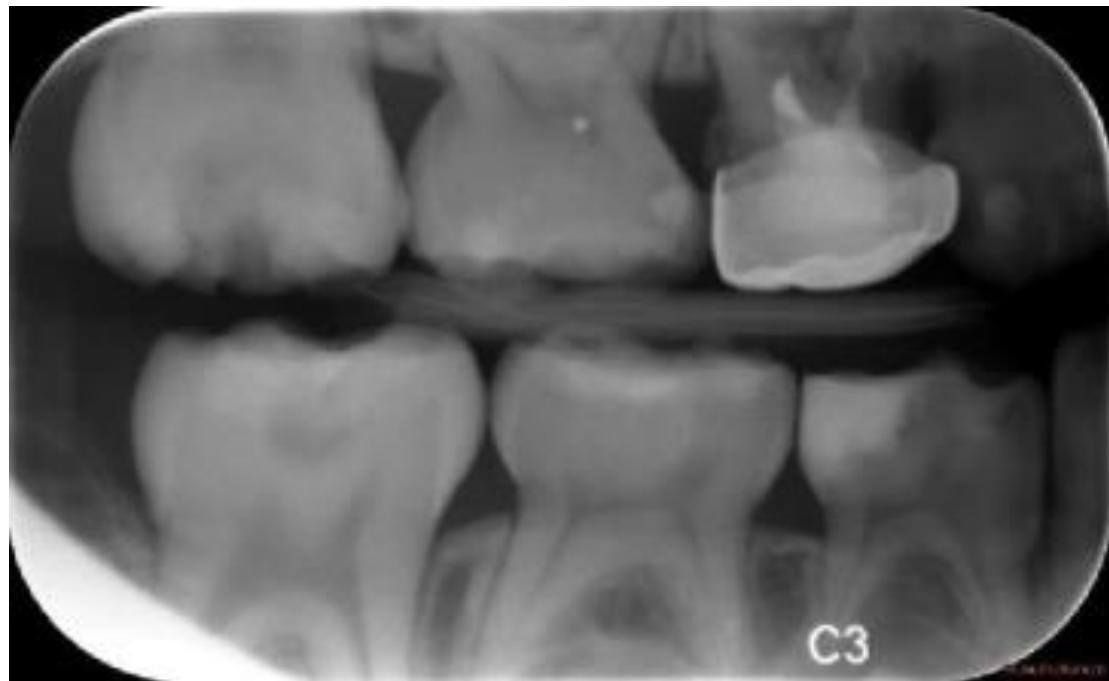


Applications: Activity Recognition



What are these
guys doing?

Applications: Medical



Applications: ADAS, Autonomous Driving



The Problem: Semantic Gap

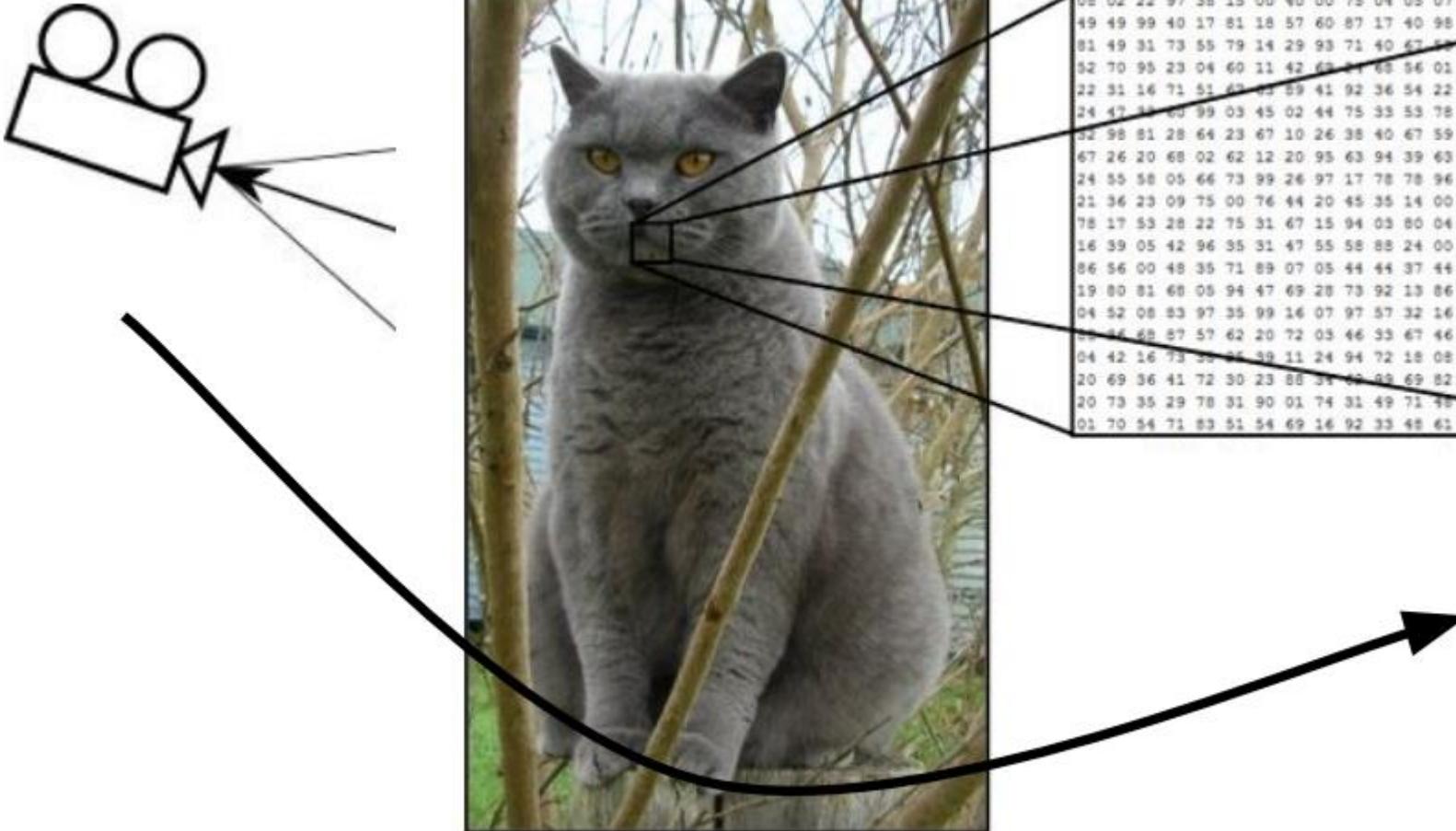
- Images are represented as 3D arrays of numbers, with integers between [0, 255].
- E.g. $300 \times 100 \times 3$ (3 for 3 color channels RGB)



08	02	22	97	38	15	00	40	00	75	04	05	07	78	52	12	50	77	81	08
49	49	99	40	17	81	18	57	60	87	17	40	98	43	69	44	81	56	62	00
81	49	31	73	55	79	14	29	93	71	40	67	51	68	30	03	49	13	36	65
52	70	95	23	04	60	21	42	63	21	68	56	01	32	56	71	37	02	36	91
22	31	16	71	51	65	63	89	41	92	36	54	22	40	40	28	66	33	13	80
24	47	34	00	99	03	45	02	44	75	33	53	78	36	84	20	35	17	12	50
32	98	81	28	64	23	67	10	26	38	40	67	59	54	70	66	18	38	64	70
67	26	20	68	02	62	12	20	95	63	94	39	63	08	40	91	66	49	94	21
24	55	58	05	66	73	99	26	97	17	78	78	96	83	14	88	34	89	63	72
21	36	23	09	75	00	76	44	20	45	35	14	00	61	33	97	36	31	33	95
78	17	53	28	22	75	31	67	15	94	03	80	04	62	16	14	09	53	56	92
16	39	05	42	96	35	31	47	55	58	88	24	00	17	54	24	36	29	85	57
86	56	00	48	35	71	89	07	05	44	44	37	44	60	21	58	51	54	17	58
19	80	81	68	05	94	47	69	28	73	92	13	86	52	17	77	04	89	55	40
04	52	08	83	97	35	99	16	07	97	57	32	16	26	26	79	33	27	98	66
04	34	68	87	57	62	20	72	03	46	33	67	46	55	12	32	63	93	53	69
04	42	16	73	55	35	59	11	24	94	72	18	08	46	29	32	40	62	76	36
20	69	36	41	72	30	23	88	31	97	92	69	82	67	59	85	76	04	36	16
20	73	35	29	78	31	90	01	74	31	49	71	48	84	51	16	23	57	05	54
01	70	54	71	63	51	54	69	16	92	33	48	61	43	52	01	89	17	47	48

What the computer sees

Challenges: View Point Variation



Challenges: Illumination Variation



Challenges: Scale



©Warren Photographic



Challenges: Deformation



Challenges: Occlusion



Challenges: Background Clutter



Challenges: Interclass Variation



Image Classification

Data Driven Approach

- Collect a dataset of images and labels
- Use Machine Learning to train an image classifier
- Evaluate the classifier on a withheld set of test images

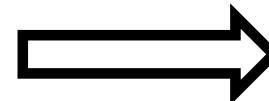
```
def train(train_images, train_labels):  
    # build a model for images -> labels...  
  
    return model  
  
def predict(model, test_images):  
    # predict test_labels using the model...  
  
    return test_labels
```

Example Training Set



First classifier: Nearest Neighbor Classifier

```
def train(train_images, train_labels):  
    # build a model for images -> labels...  
  
    return model  
  
def predict(model, test_images):  
    # predict test_labels using the model...  
  
    return test_labels
```



Remember all training images and their labels



Predict the label of the most similar training image

How do we compare the images?

L1 Distance: $\text{dist}(I_1, I_2) = \sum_p |I_1^p - I_2^p|$

test image				training image				pixel-wise absolute value differences					
56	32	10	18	-	10	20	24	17	=	46	12	14	1
90	23	128	133		8	10	89	100		82	13	39	33
24	26	178	200		12	16	178	170		12	10	0	30
2	0	255	220		4	32	233	112		2	32	22	108

add → 456

Dataset

- Example dataset: CIFAR-10
- 10 labels
- 50,000 training images
- Each image is tiny: 32x32
- 10,000 test images

airplane



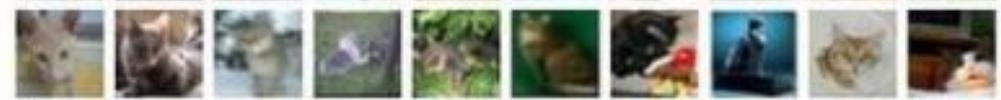
automobile



bird



cat



deer



dog



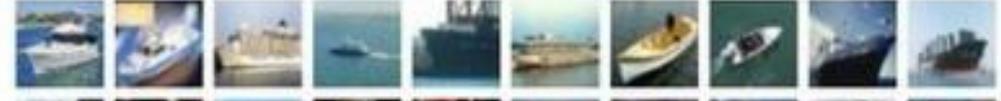
frog



horse



ship

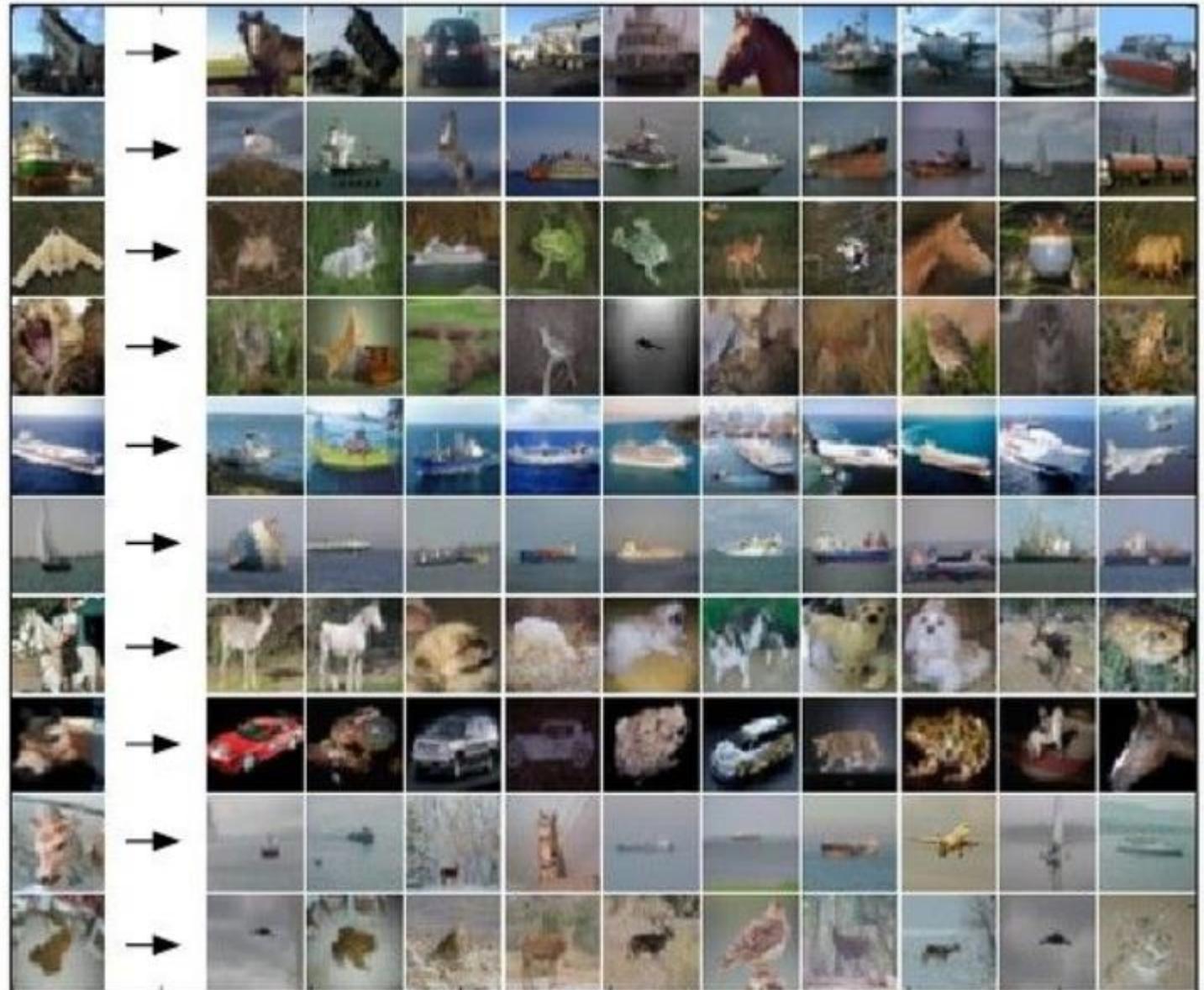


truck



Output: Nearest Neighbor Classifier on CIFAR

For every test image
(first column),
examples of nearest
neighbors in rows



Generalization

- Terrible performance in real life
- Distance metrics on level of whole images can be very unintuitive

Original



Shifted



Occluded



Darkened



ImageNet Challenge

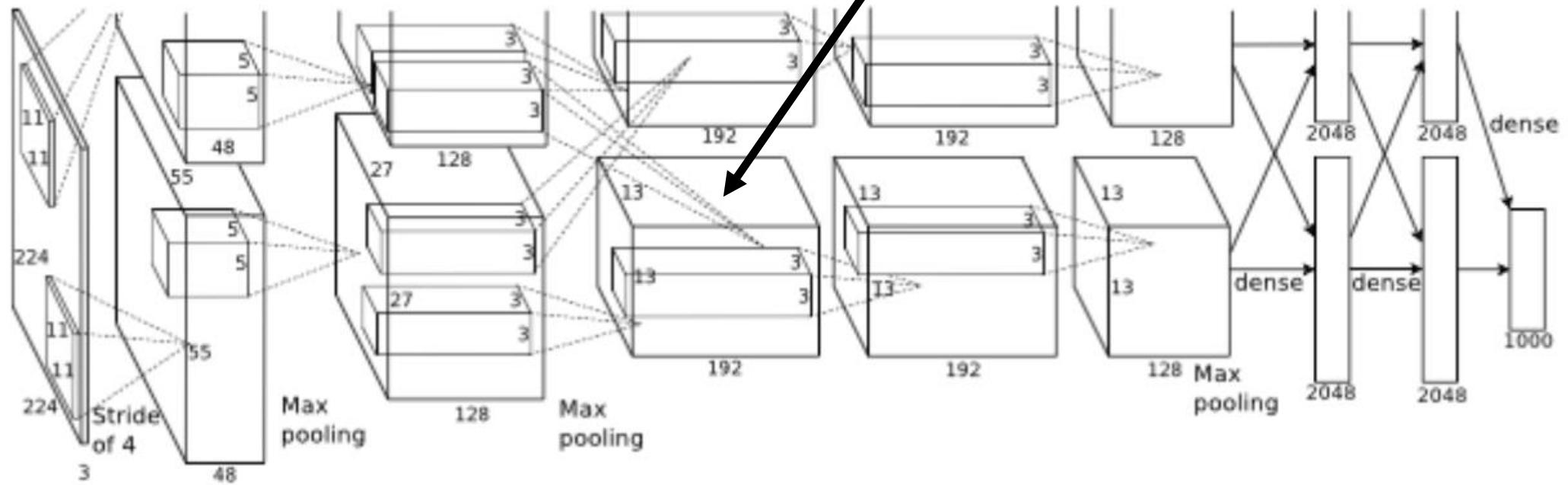
- ImageNet (image-net.org)
 - 15M high resolution images
 - over 22K categories
 - labeled by Mechanical Turk workers
- ImageNet Large-Scale Visual Recognition Challenge (ILSVRC)
- 2010-2015
- 2010 competition
 - 1.2 M training images, 1000 categories (general and specific)
 - 200K test images
 - output a list of 5 object categories in descending order of confidence
 - two error rates: top-1 and top-5

AlexNet (2012)

- 5 convolutional layers
- 3 fully connected layers
- 1000-way softmax output layer

Trained on GTX 580
with 3 GB of RAM.

Network (feature
maps spread across
two GPUs)



AlexNet (2012): Tricks of the Trade

- Downsampled images
 - shorter dimension 256 pixels, longer dimension cropped about center to 256 pixels
 - R, G, B channels
- Mean subtraction from inputs
- ReLU instead of logistic or tanh units
- DropOut

AlexNet (2012): Tricks of the Trade

- Data set augmentation
 - Generate image translations by selecting random 224 x 224 sub-images
 - Horizontal reflections (standard trick in computer vision)
 - When testing, extract 10 distinct 224x224 sub-images and average predictions
- More data set augmentation
 - Vary intensity and color of the illumination from epoch to epoch

Visualization

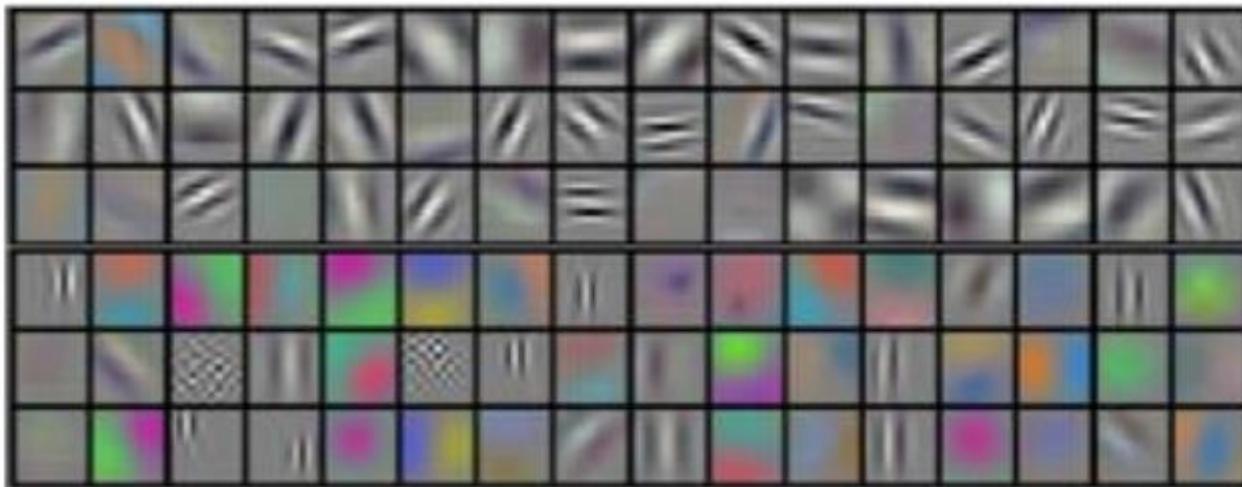
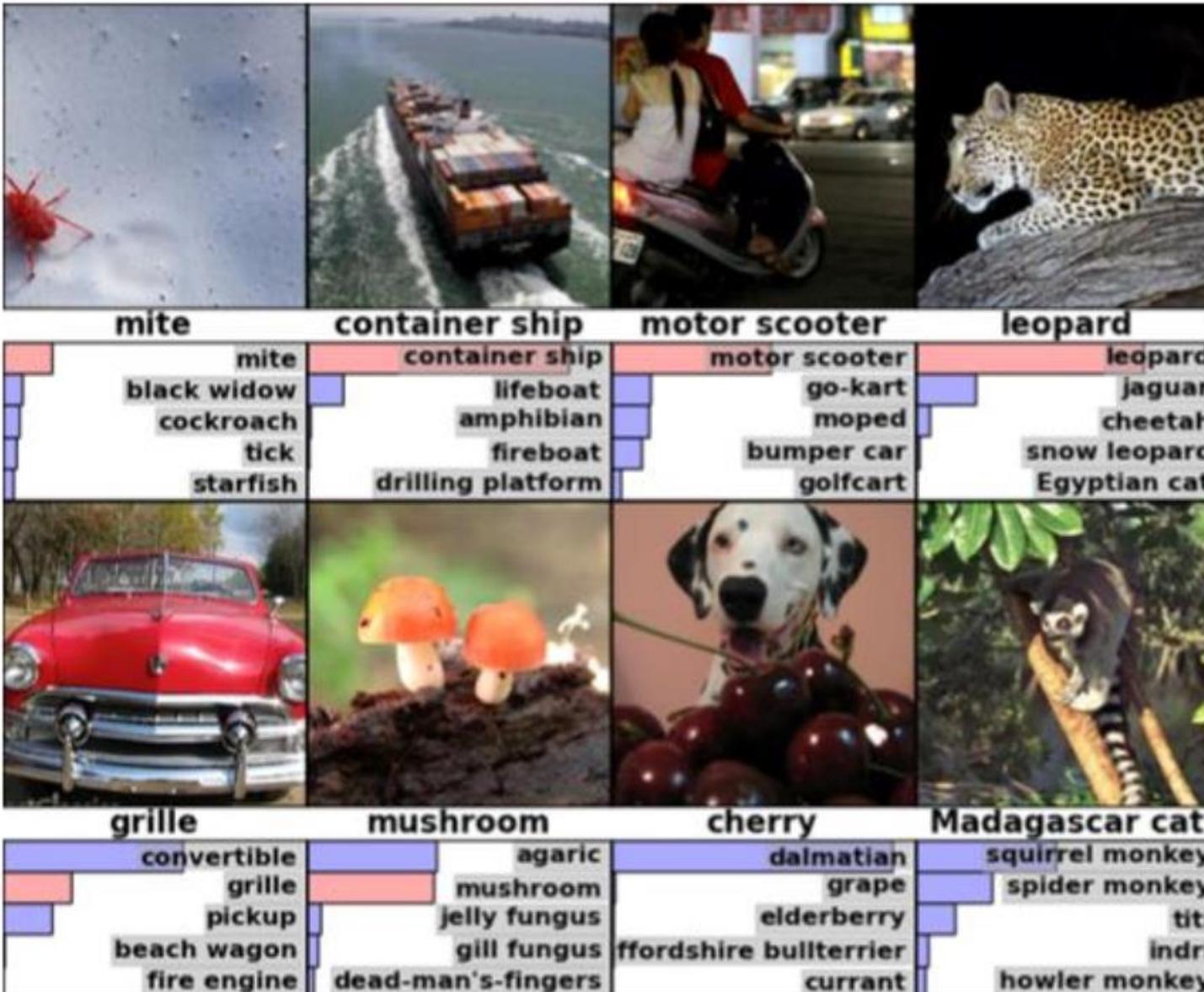


Figure 3: 96 convolutional kernels of size $11 \times 11 \times 3$ learned by the first convolutional layer on the $224 \times 224 \times 3$ input images. The top 48 kernels were learned on GPU 1 while the bottom 48 kernels were learned on GPU 2. See Section 6.1 for details.

Results



Results

Model	Top-1	Top-5
<i>Sparse coding [2]</i>	47.1%	28.2%
<i>SIFT + FVs [24]</i>	45.7%	25.7%
CNN	37.5%	17.0%

Table 1: Comparison of results on ILSVRC-2010 test set. In *italics* are best results achieved by others.

Model	Top-1 (val)	Top-5 (val)	Top-5 (test)
<i>SIFT + FVs [7]</i>	—	—	26.2%
1 CNN	40.7%	18.2%	—
5 CNNs	38.1%	16.4%	16.4%
1 CNN*	39.0%	16.6%	—
7 CNNs*	36.7%	15.4%	15.3%

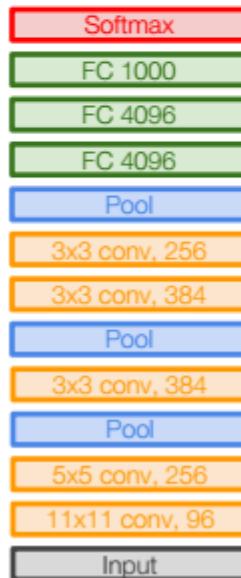
Table 2: Comparison of error rates on ILSVRC-2012 validation and test sets. In *italics* are best results achieved by others. Models with an asterisk* were “pre-trained” to classify the entire ImageNet 2011 Fall release. See Section 6 for details.

VGG Net (2014)

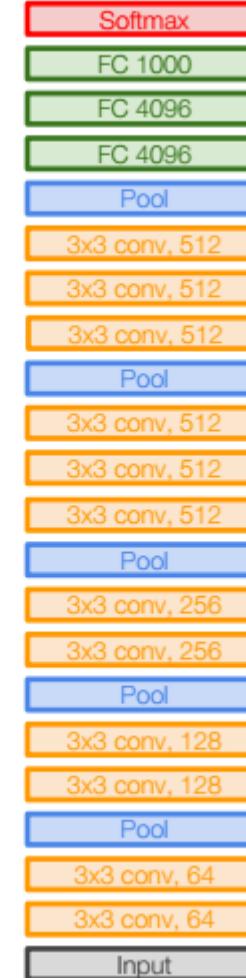
- Smaller filters, deeper networks.

Why?

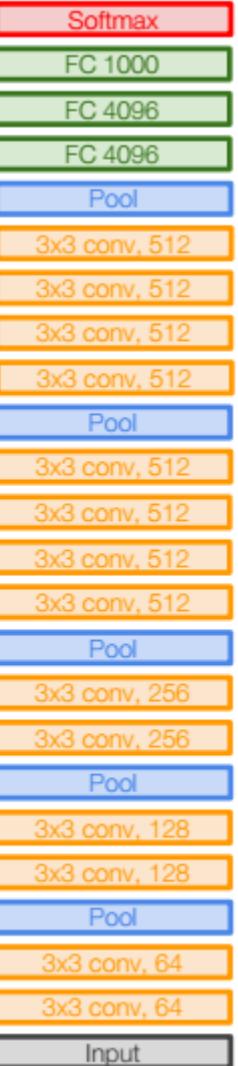
- Stack of three 3×3 conv (stride 1) layers has same effective receptive field as one 7×7 conv layer
- But deeper, more non-linearities
- And fewer parameters: $3 * 3^2 C^2$ vs. $7^2 C^2$ for C channels per layer



AlexNet

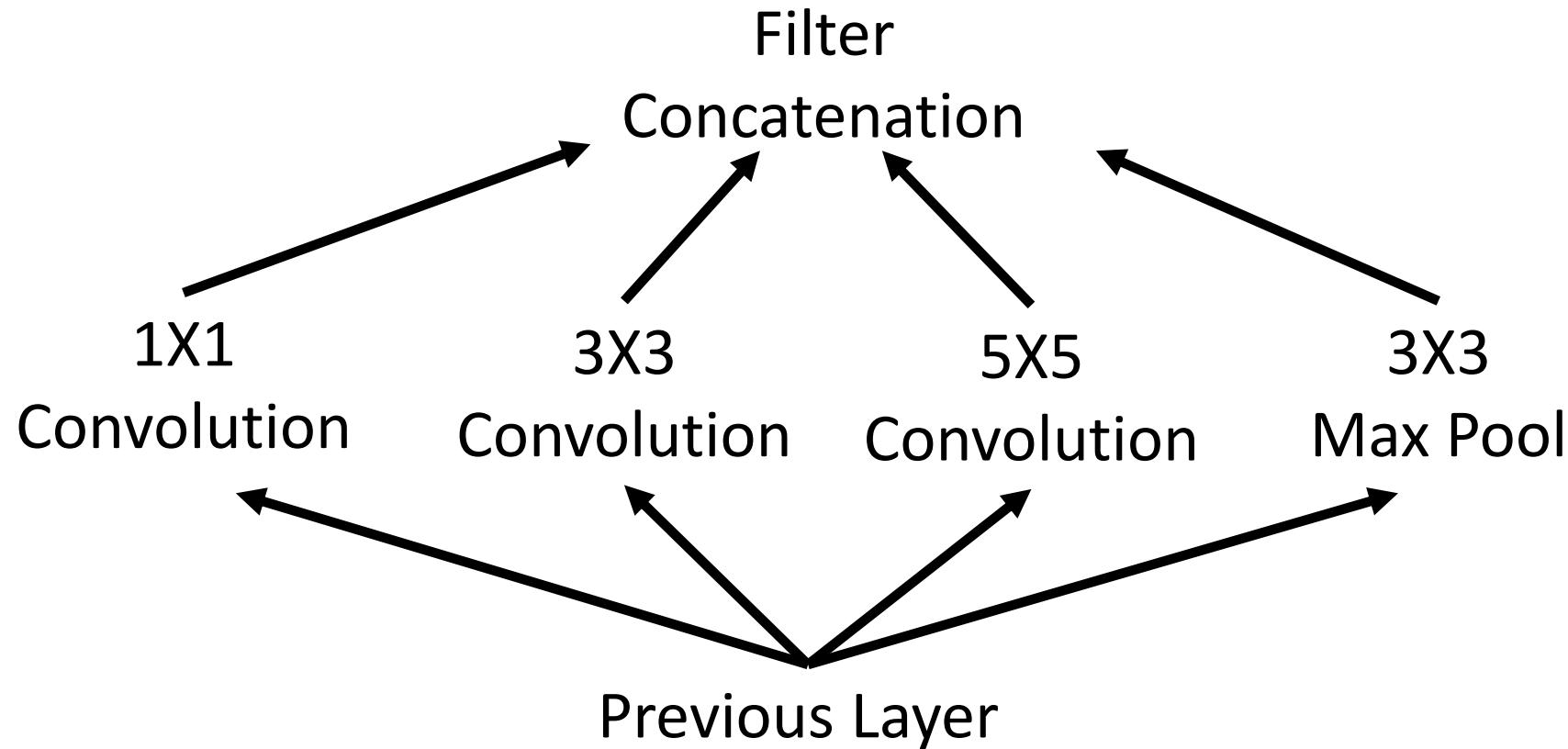


VGG16



VGG19

GoogLeNet (2014): Inception Module



GoogLeNet

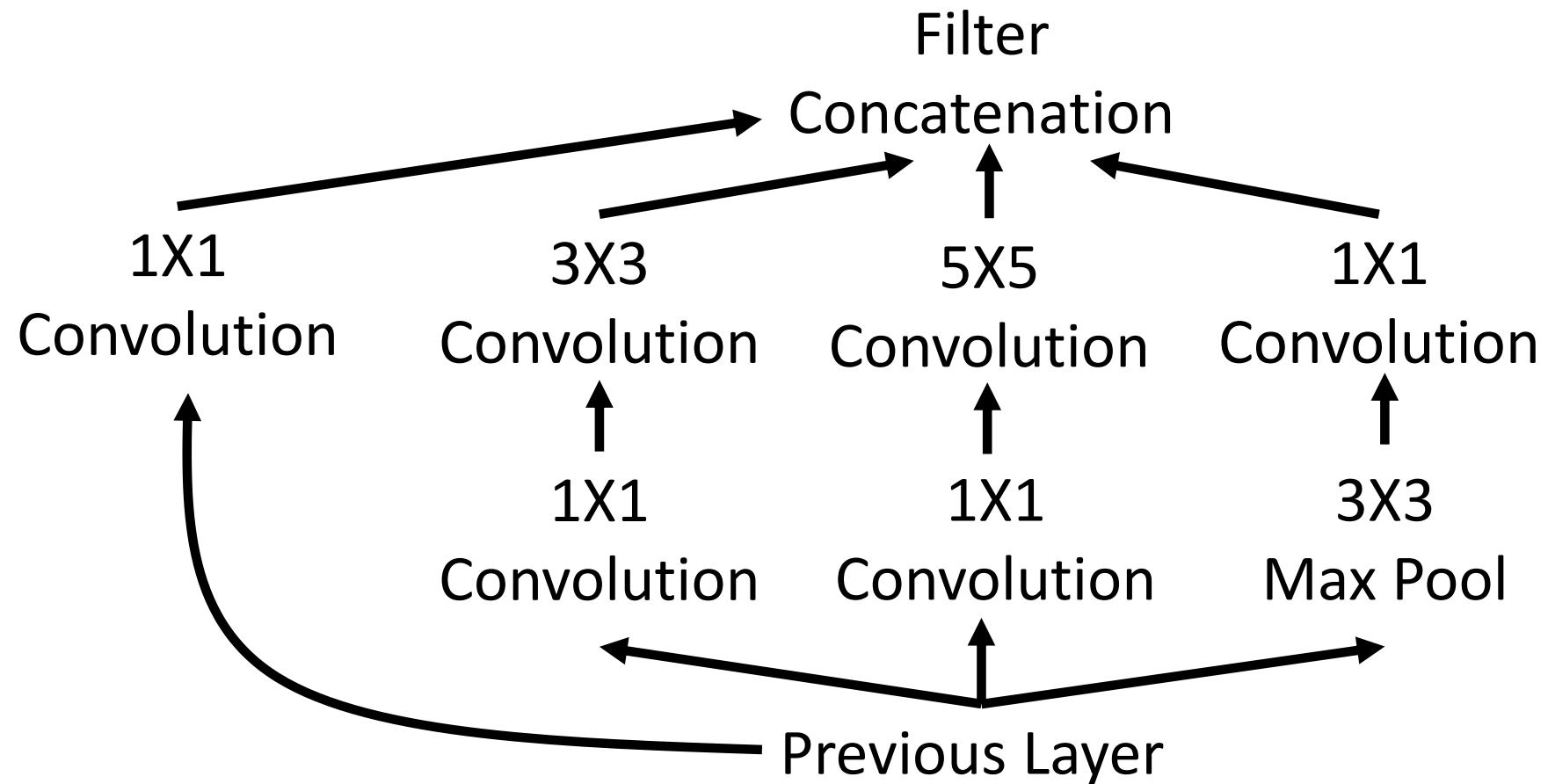
Problem?

- Heavy Compute: Too many feature maps because of feature concatenation

Solution?

- 1X1 Convolution: preserves spatial dimensions, reduces depth

GoogLeNet: Inception Module



ResNet (2015)

- What happens when we continue stacking deeper layers on a “plain” convolutional neural network?
- Hypothesis: the problem is an optimization problem, deeper models are harder to optimize
- The deeper model should be able to perform at least as well as the shallower model.
- A solution by construction is copying the learned layers from the shallower model and setting additional layers to identity mapping

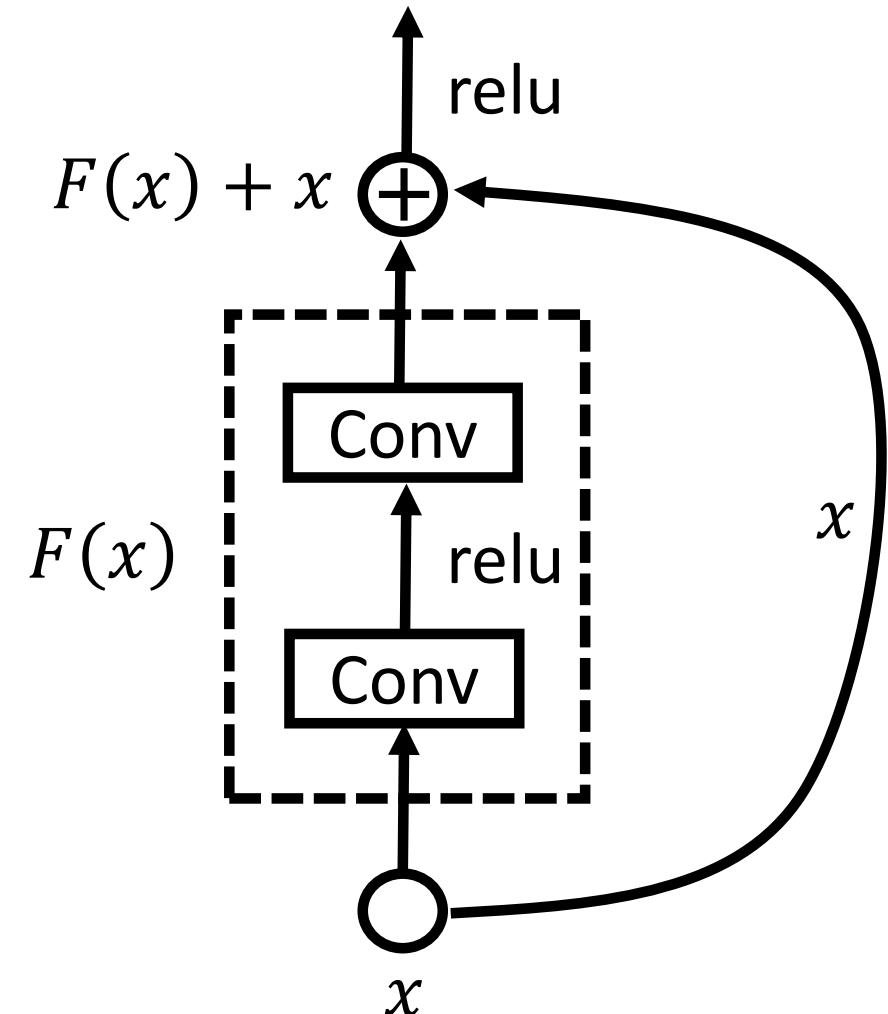


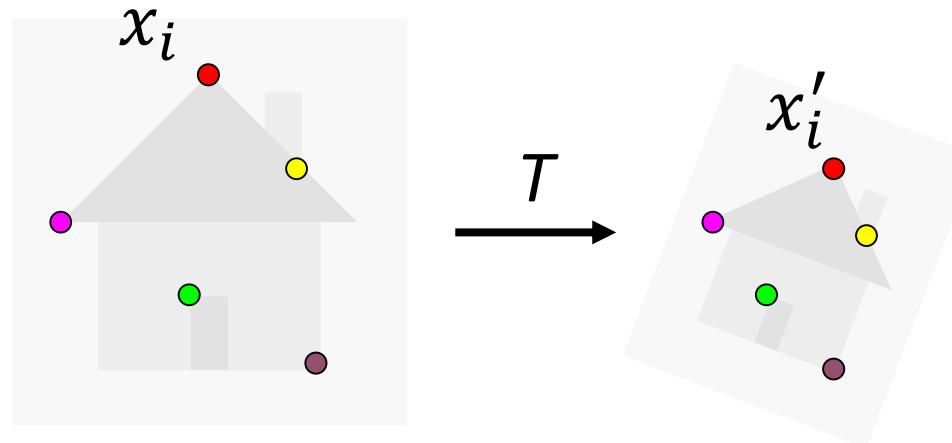
Image Classification Architectures: Summary

ILSVRC 2012	AlexNet		
ILSVRC 2013	ZFNet	Proposed by Zeiler-Fergus	improvement on AlexNet by tweaking the architecture hyperparameters, in particular by expanding the size of the middle convolutional layers.
ILSVRC 2014	GoogLe Net	Szegedy et al, Google	Development of an Inception Module that dramatically reduced the number of parameters in the network (4M, compared to AlexNet with 60M). Additionally, uses Average Pooling instead of Fully Connected layers at the top of the ConvNet, eliminating a large amount of parameters that do not seem to matter much.
ILSVRC 2014 Runner-up	VGGNet	Simonyan Zisserman	Showed that network depth is a critical component for good performance. Their final best network contains 16 CONV/FC layers and, appealingly, features an extremely homogeneous architecture that only performs 3x3 convolutions and 2x2 pooling from the beginning to the end. Outperformed GoogLeNet in multiple transfer learning tasks.
ILSVRC 2015	ResNets	He et al	Special skip connections and features heavy use of batch normalization. Also missing fully connected layers at the end of the network.

Object Detection

Object Detection Ideas: Geometrical Era

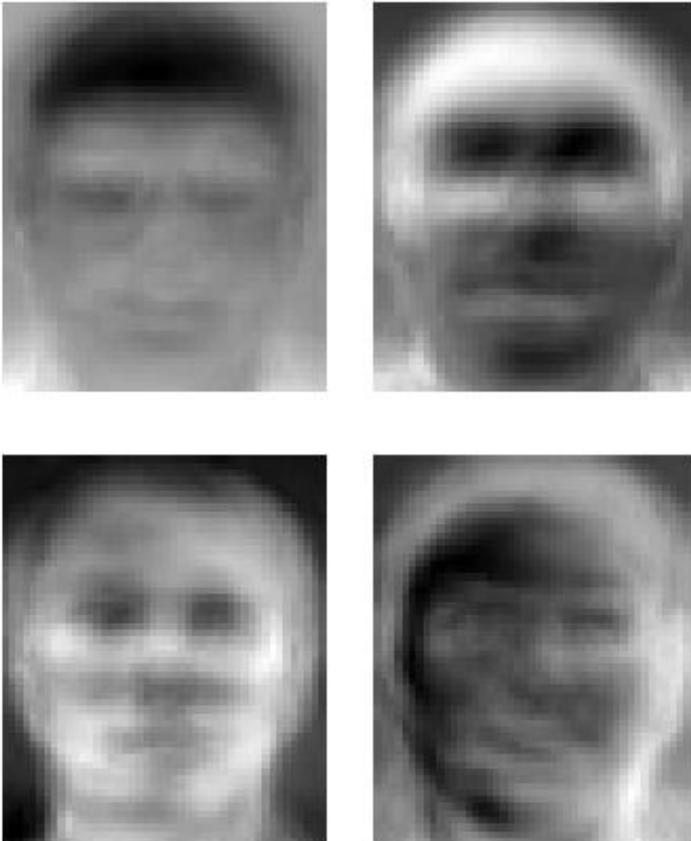
- Fit a model to a transformation between pairs of features (*matches*) in two images



Find transformation T
that minimizes:

$$\sum_i \text{residual}(T(x_i), x'_i)$$

Object Detection Ideas: Appearance Models



Eigenfaces
(Turk and
Pentland,
1991)

- Other global models based upon color histograms, appearance manifolds.
- Requires global registration of patterns. Not robust to clutter, occlusion, geometric transformations

Object Detection Ideas: Local Features

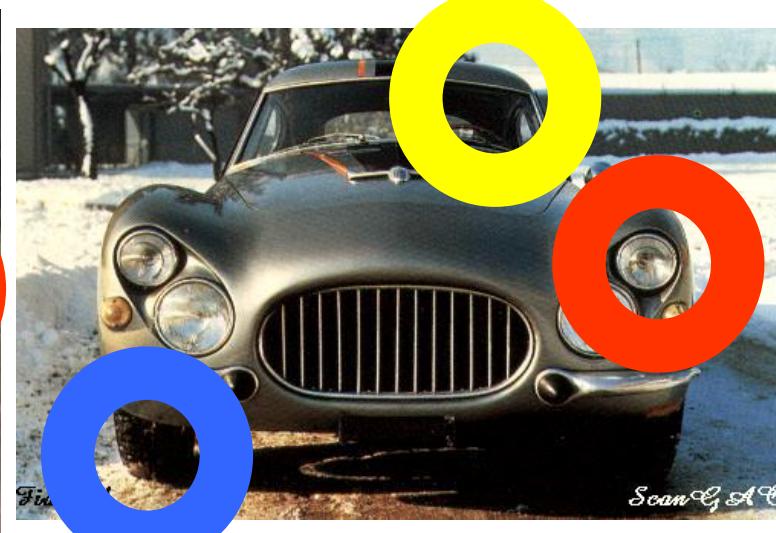


SIFT, D. Lowe (1999, 2004)



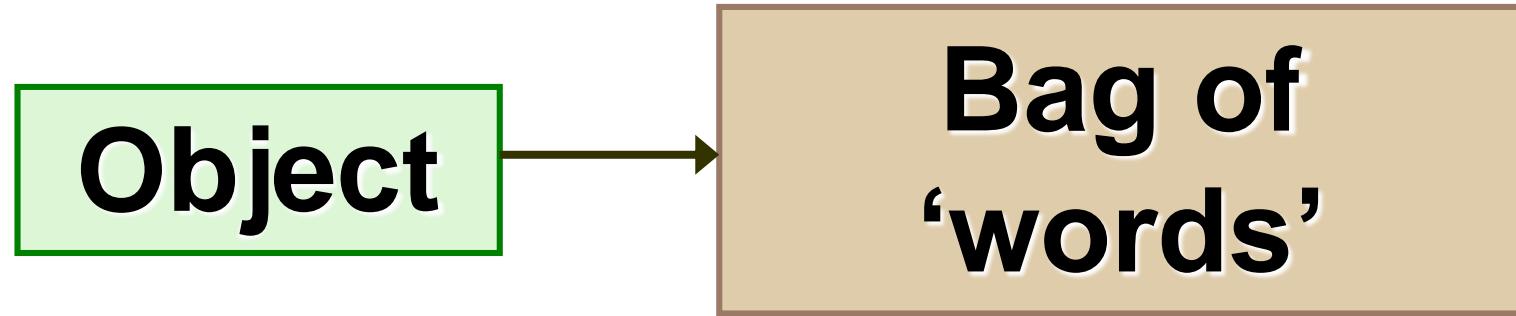
Large-scale image search, combining local features, indexing, and spatial constraints
Philbin et al. '07

Constellation Model



- Globally different objects, but have portions in common

Bag of Features Models



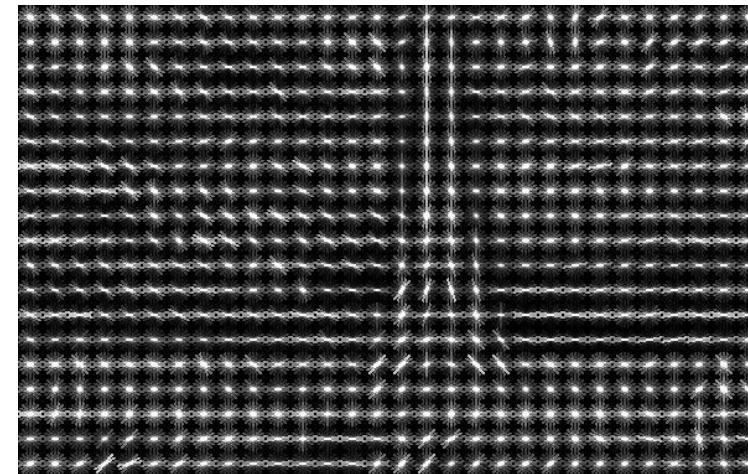
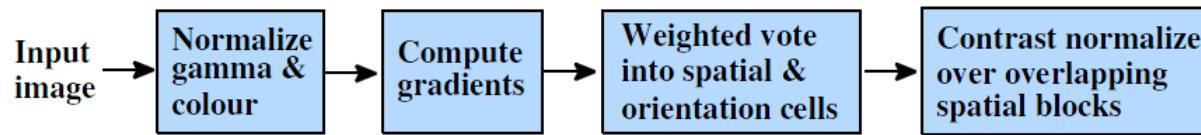
Object Detection Ideas: Sliding Window



- Turk and Pentland, 1991
- Belhumeur et al., 1997
- Schneiderman & Kanade 2004
- Viola and Jones, 2000
- Schneiderman & Kanade, 2004
- Agrawal and Roth, 2002

Histograms of oriented gradients (HOG)

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



N. Dalal and B. Triggs, [Histograms of Oriented Gradients for Human Detection](#), CVPR 2005

Image credit: N. Snavely

Pedestrian detection with HOG

- Train a pedestrian template using a linear support vector machine

positive training examples



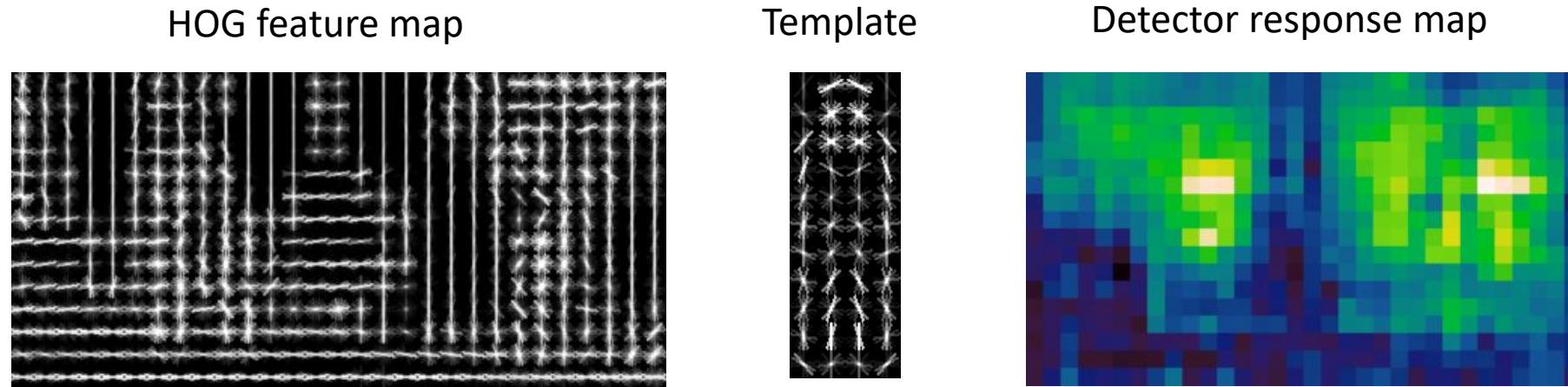
negative training examples



N. Dalal and B. Triggs, [Histograms of Oriented Gradients for Human Detection](#), CVPR 2005

Pedestrian detection with HOG

- Train a pedestrian template using a linear support vector machine
- 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*

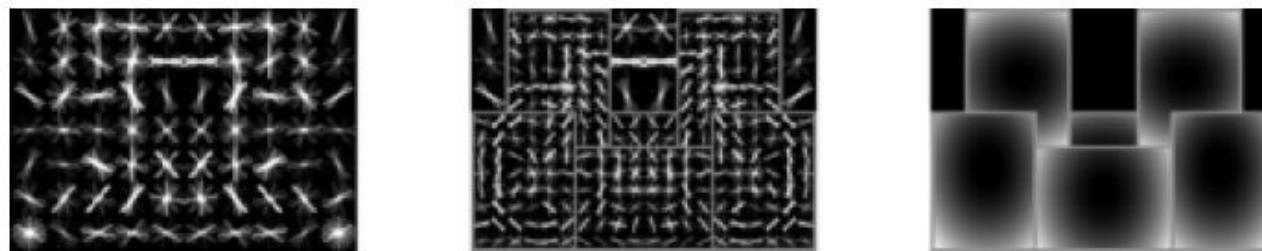
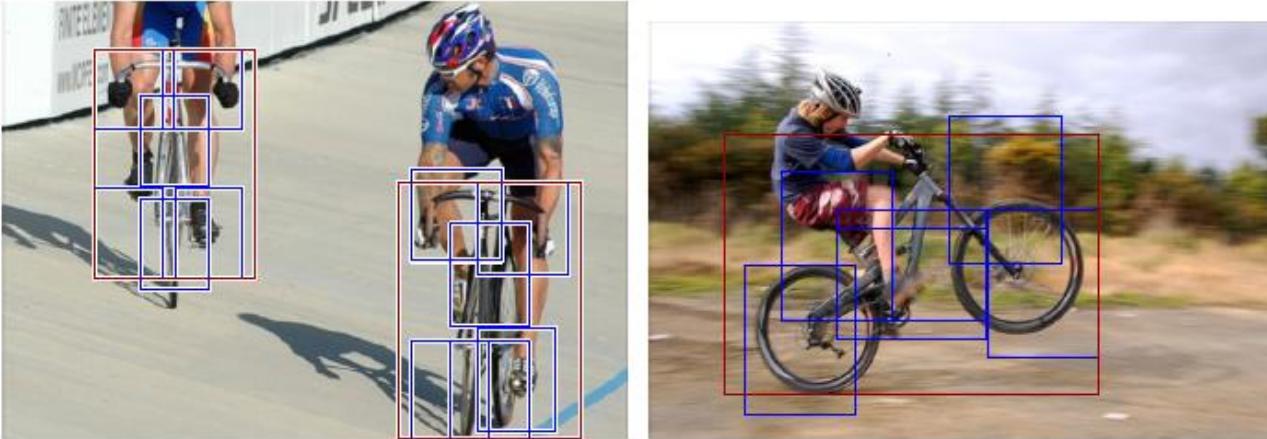


Example detections



[Dalal and Triggs, CVPR 2005]

Discriminatively trained part-based models



P. Felzenszwalb, R. Girshick, D. McAllester, D. Ramanan, PAMI 2009,
“Object Detection with
Discriminatively Trained Part-Based
Models”

Object Detection with Deep Neural Networks

HoG by Convolutional Layers

Steps of computing HOG:

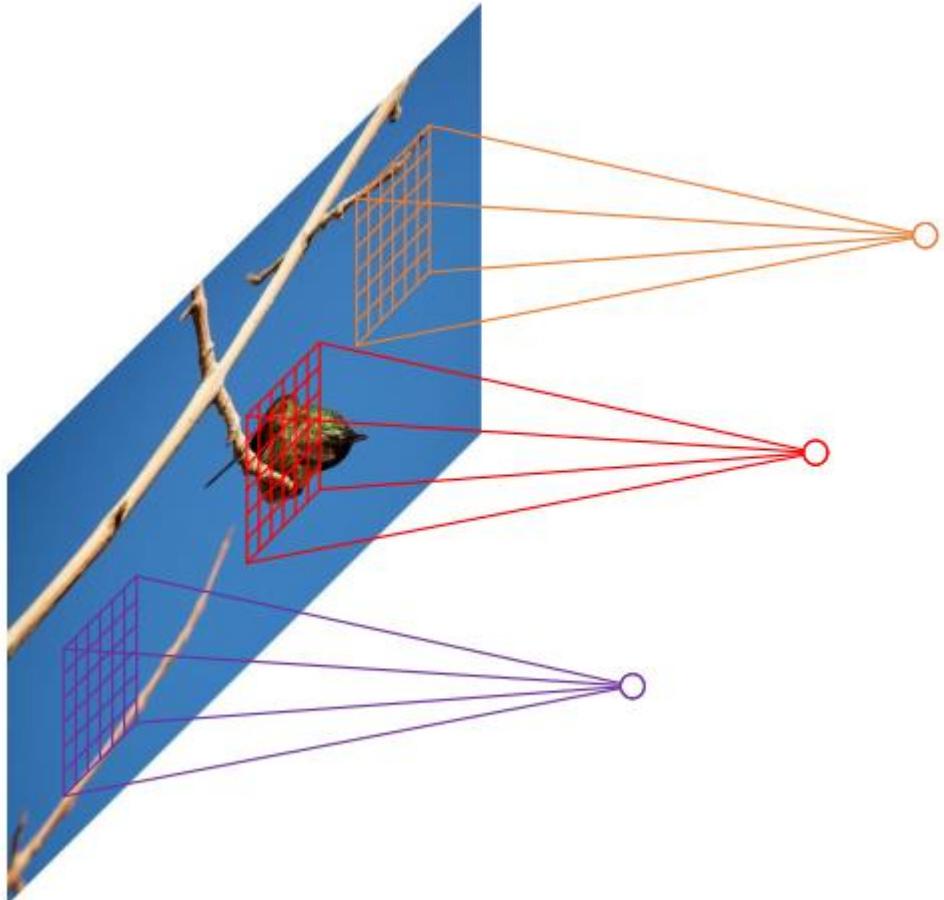
- Computing image gradients
- Binning gradients into 18 directions
- Computing cell histograms
- Normalizing cell histograms

Convolutional perspectives:

- Horizontal/vertical edge filters
- Directional filters + gating (non-linearity)
- Sum/average pooling
- Local response normalization (LRN)

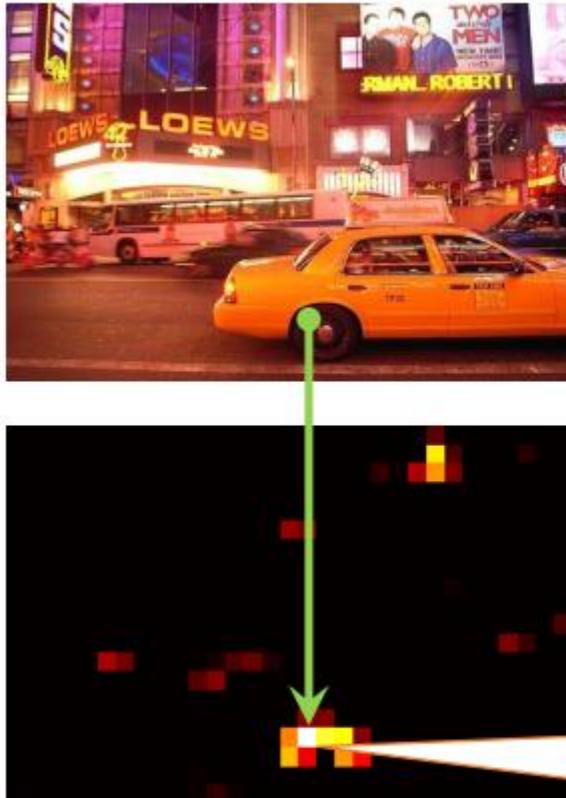
HOG, dense SIFT, and many other “hand-engineered” features are convolutional feature maps.

Convolutional Layers



- Convolutional layers are locally connected
- a filter/kernel/window slides on the image or the previous map
 - the position of the filter explicitly provides information for localizing
 - local spatial information w.r.t. the window is encoded in the channels

Feature Maps = features and their locations



ImageNet images with **strongest** responses of this channel



Intuition of *this* response:

There is a “circle-shaped” object (likely a tire) at this position.

What

Where

Feature Maps = features and their locations



one feature map of conv₅
(#66 in 256 channels of a model
trained on ImageNet)

ImageNet images with **strongest** responses of this channel



Intuition of *this* response:
There is a “λ-shaped” object (likely an underarm) **at this position**.

What

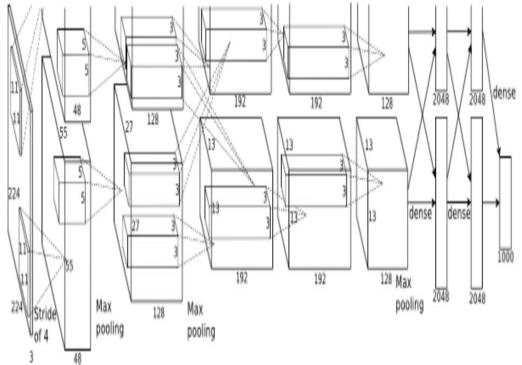
Where

Classification + Localization



Cat

Localization as Regression Problem



Output of last
fully connected
layer: 4096 dim
vector

4096 \rightarrow 1000

4096 \rightarrow 4

Class Scores

Cat: 0.9

Dog: 0.1

Car: 0.05

Correct Label

Softmax Loss

Total Loss

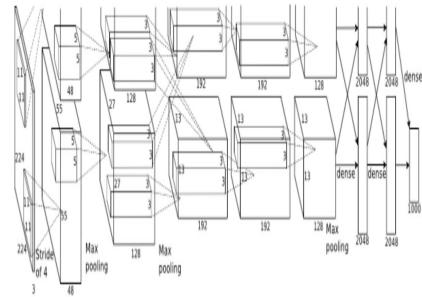
Box Coordinates

x, y, w, h

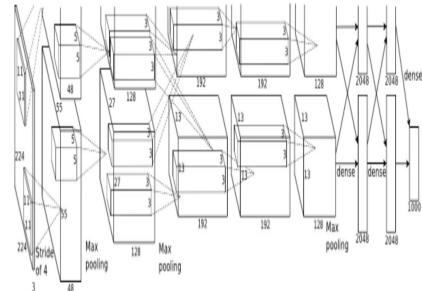
L2 Loss

Correct Box

Localization as Regression: Problem?



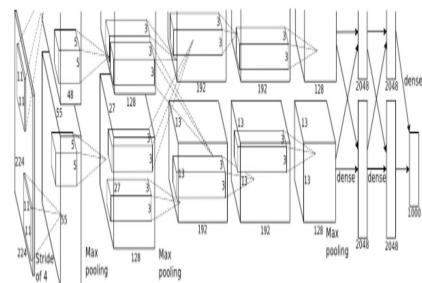
Cat: (x,y,w,h)



Dog: (x,y,w,h)

Cat: (x,y,w,h)

Each image needs a different number of outputs!



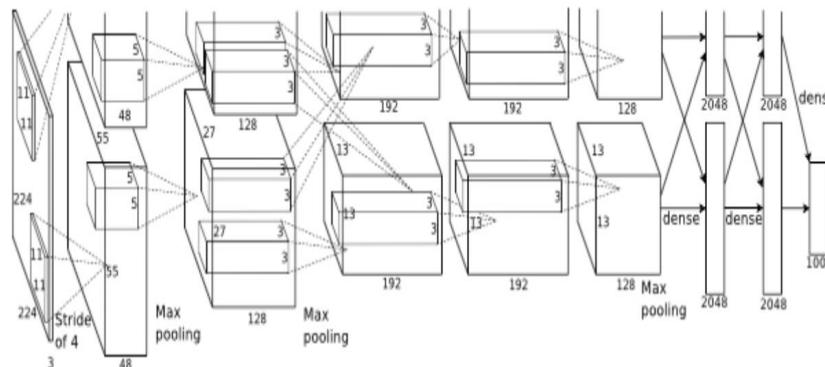
Duck: (x,y,w,h)

Duck: (x,y,w,h)

...

Object Detection as Classification

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



Dog: Yes

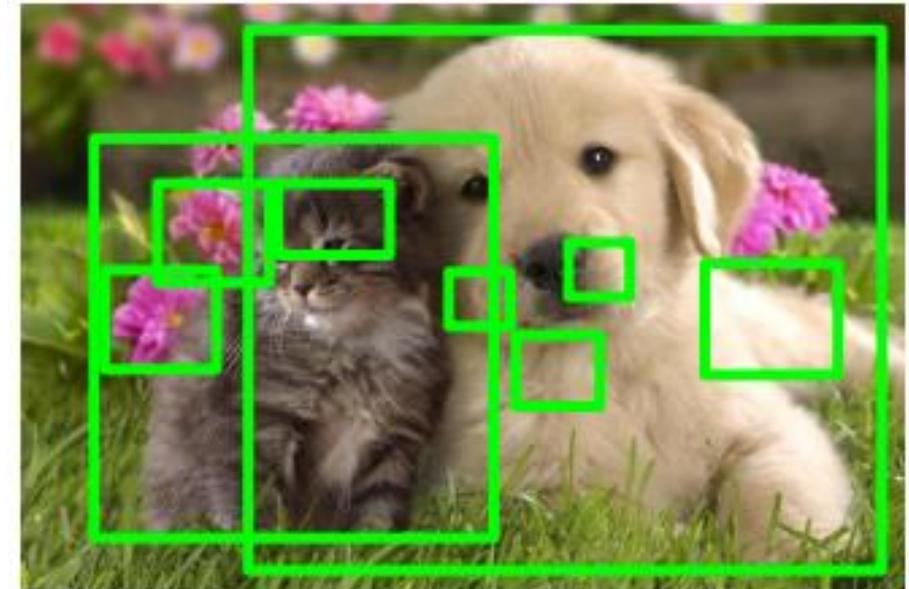
Cat: No

BG: Yes

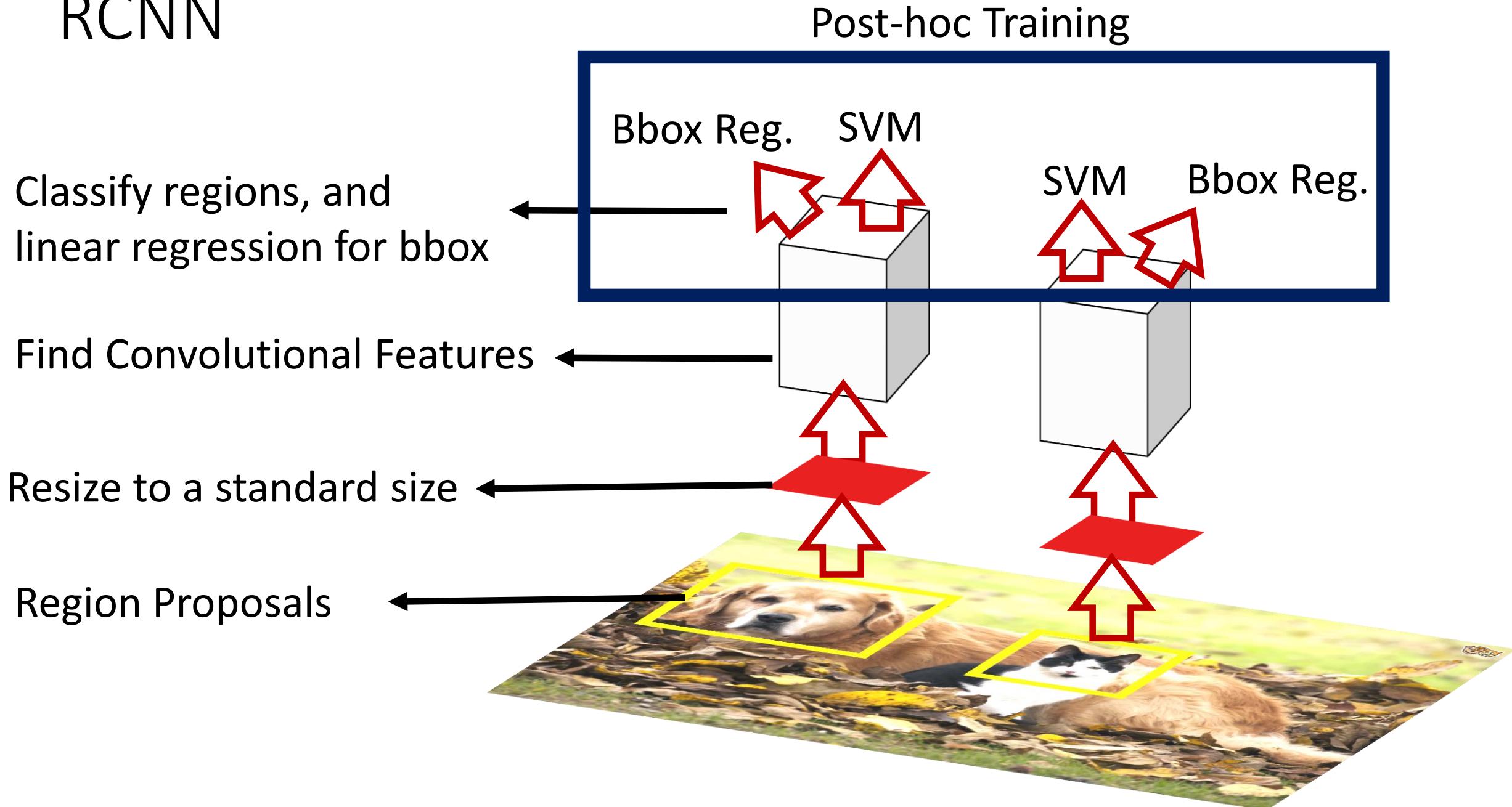
Problem: Need to apply CNN to huge number of locations and scales, very computationally expensive.

Region Proposals: “Objectness” Detection

- Find “blobby” image regions that are likely to contain object
- Relatively fast to run; e.g. Selective Search gives 1000 region proposals in a few seconds on CPU



RCNN



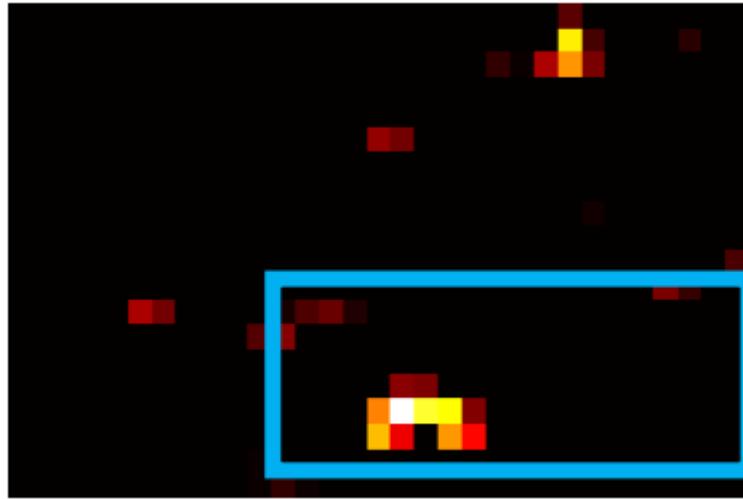
RCNN: Problems?

- 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
- Inference (detection) is slow
 - 47s / image with VGG16 [Simonyan & Zisserman. ICLR15]
 - ~2K proposals/passes through a Conv Net

Region-based CNN Features

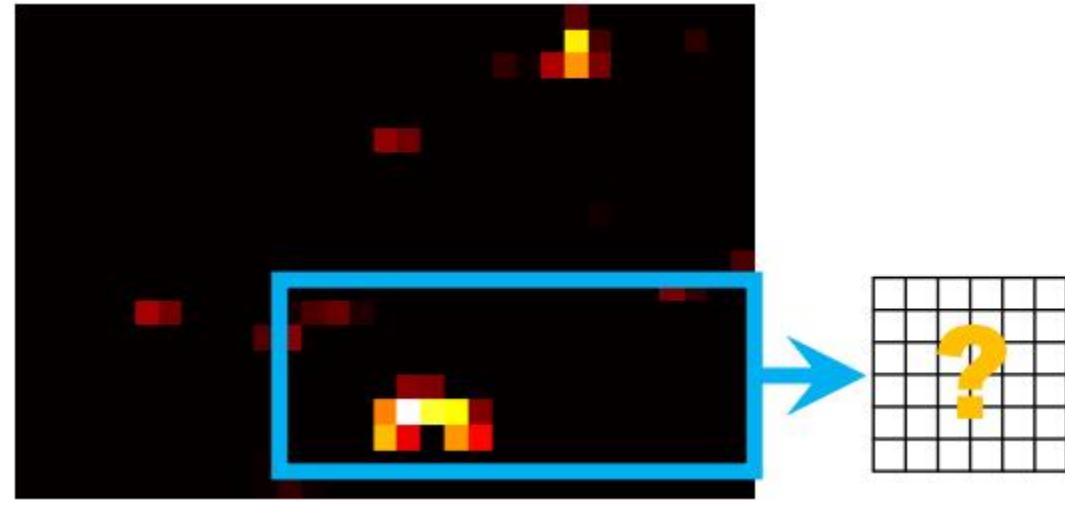
- Given proposal regions, what we need is a feature for each region
- R-CNN: cropping an image region + CNN on region, requires 2000 CNN computations
- What about cropping feature map regions?

Regions on Feature Maps



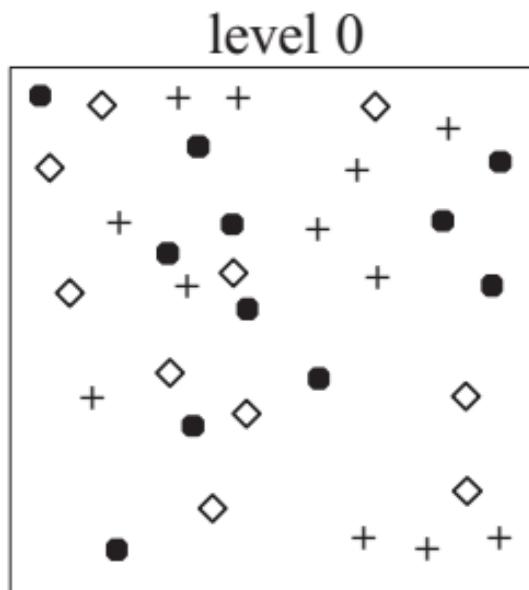
- Compute convolutional feature maps on the entire image only once.
- Project an image region to a feature map region (using correspondence of the receptive field center)
- Extract a region-based feature from the feature map region...

Regions on Feature Maps



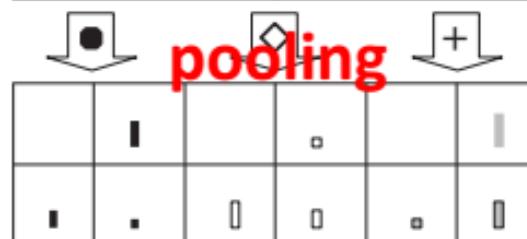
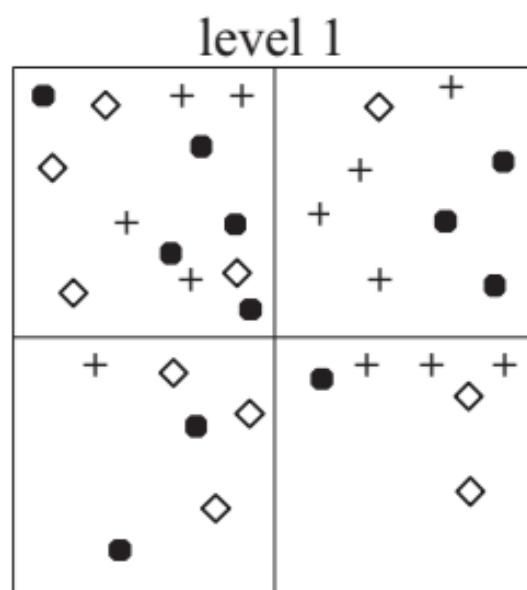
- Fixed-length features are required by fully-connected layers or SVM
- But how to produce a fixed-length feature from a feature map region?
- Solutions in traditional compute vision: Bag-of-words, SPM...

Bag-of-words & Spatial Pyramid Matching



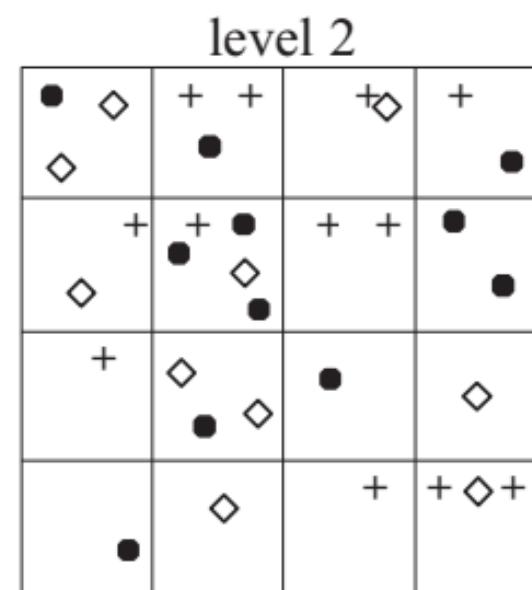
Bag-of-words

[J. Sivic & A. Zisserman, ICCV 2003]



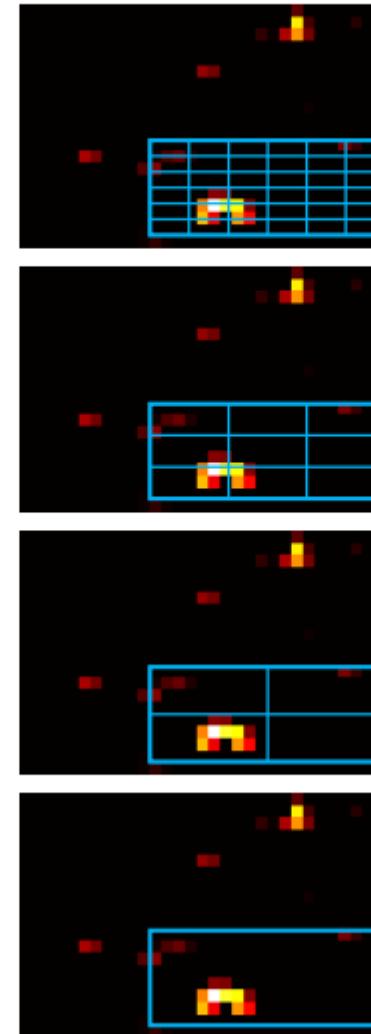
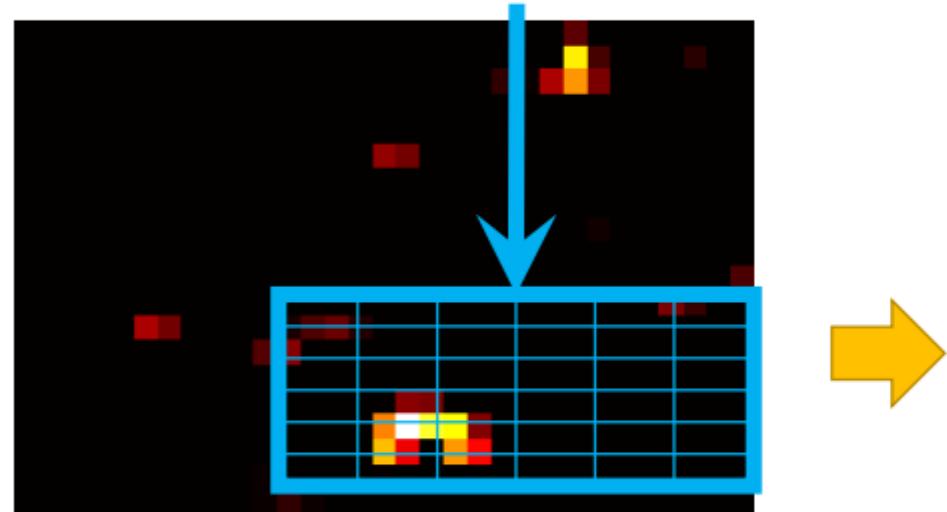
Spatial Pyramid Matching (SPM)

[K. Grauman & T. Darrell, ICCV 2005]
[S. Lazebnik et al, CVPR 2006]

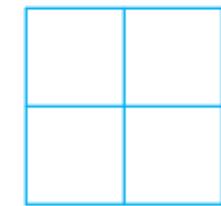
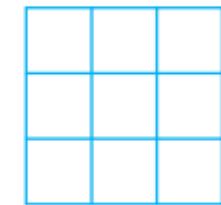
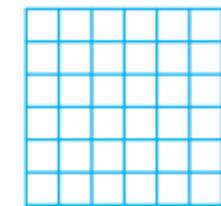


Spatial Pyramid Pooling (SPP) Layer

- fix the number of bins (instead of filter sizes)
- adaptively-sized bins



pooling



a finer level maintains explicit spatial information

concatenate, fc layers...

a coarser level removes explicit spatial information (bag-of-features)

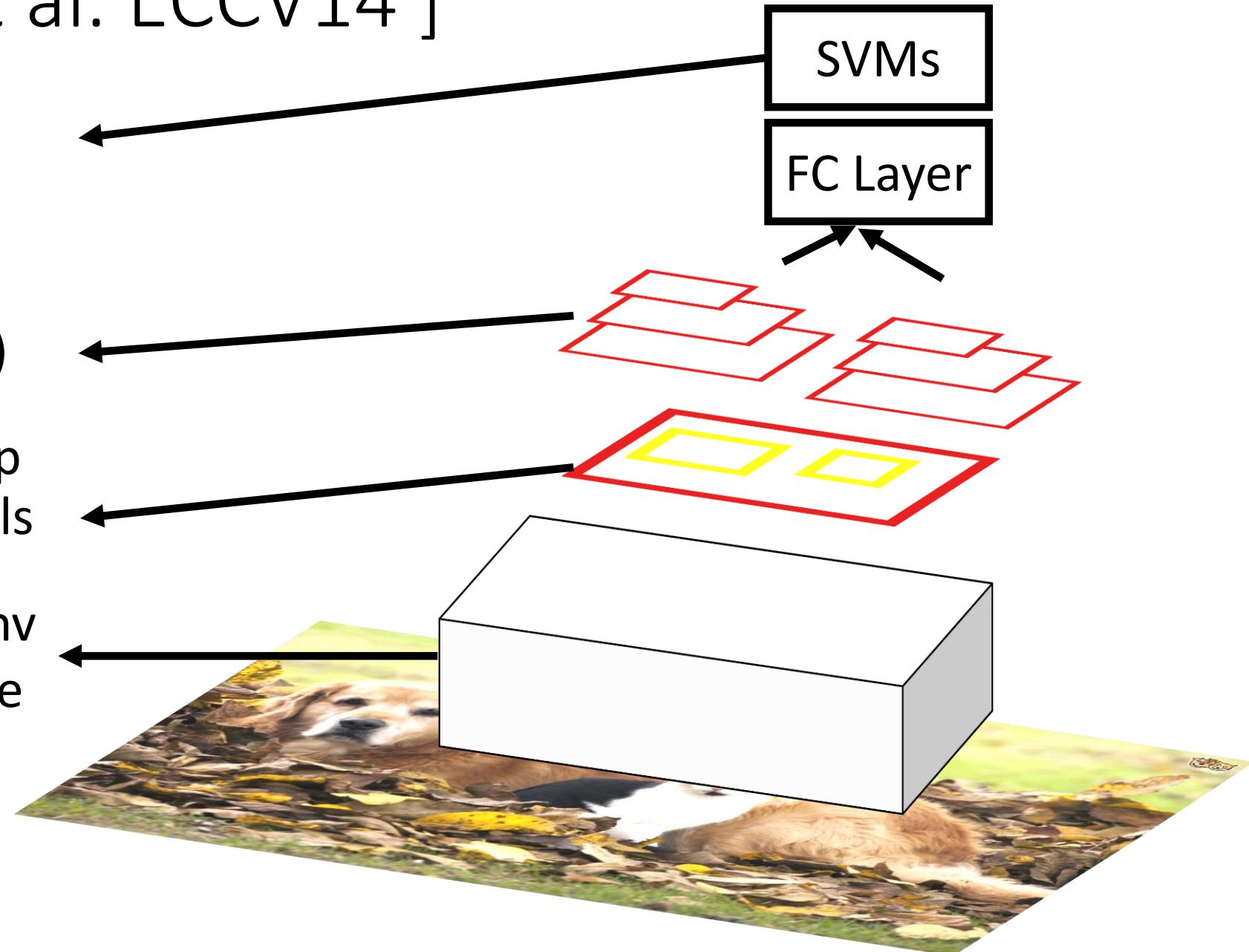
SPP Net [He et al. ECCV14]

Classify each region
with SVMs

Spatial Pyramid
Representation (SPP)

“Conv 5” Feature Map
with region proposals

Apply single pass conv
net over whole image



SPP Net: Good and Bad

- Good
 - Fixes the inference speed issue with RCNN
- Bad
 - Inherits adhoc training objectives from RCNN
 - Training still slow.
- New Problem
 - Parameters below SPP layers can not be updated during training

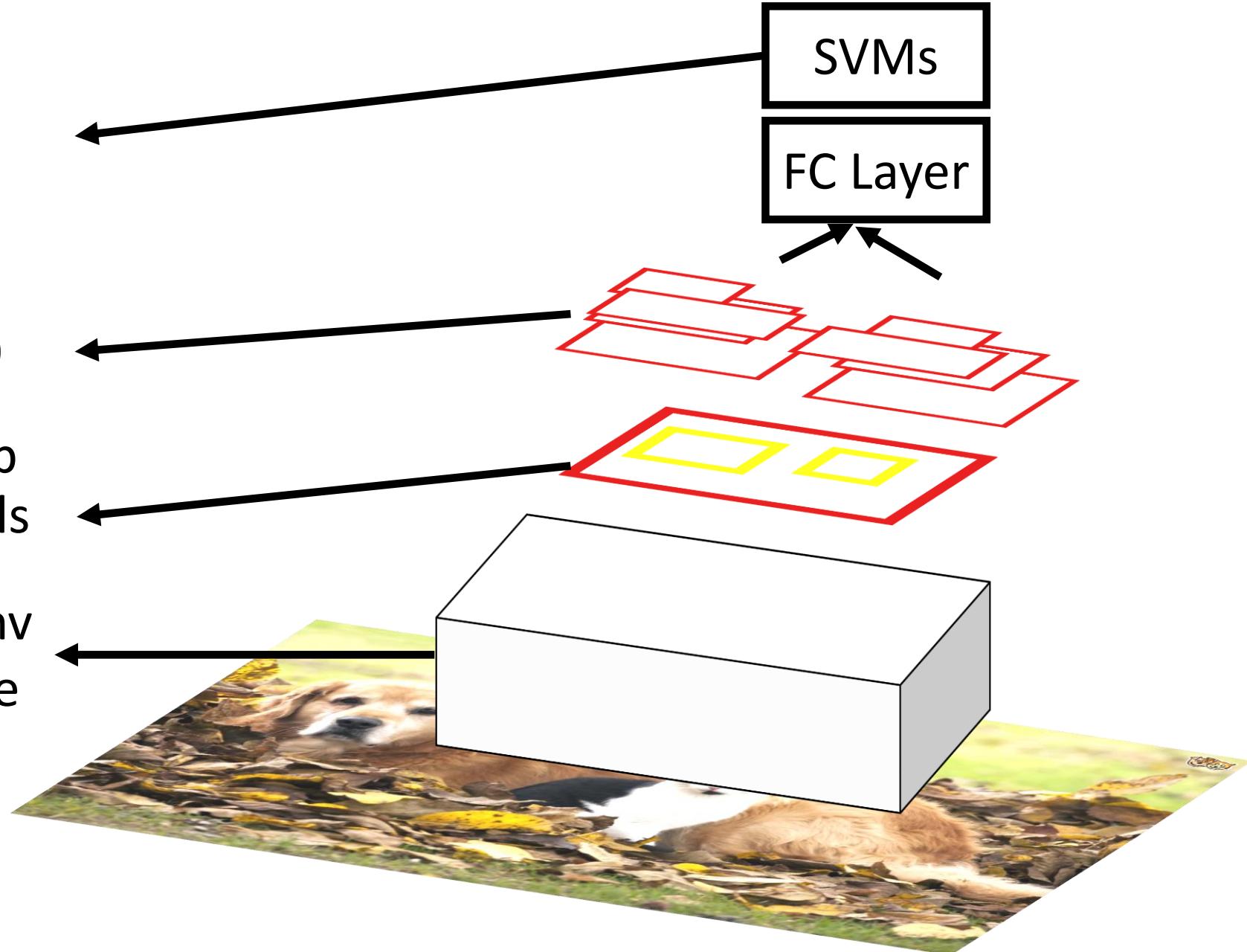
Fast RCNN

Classify each region
with SVMs

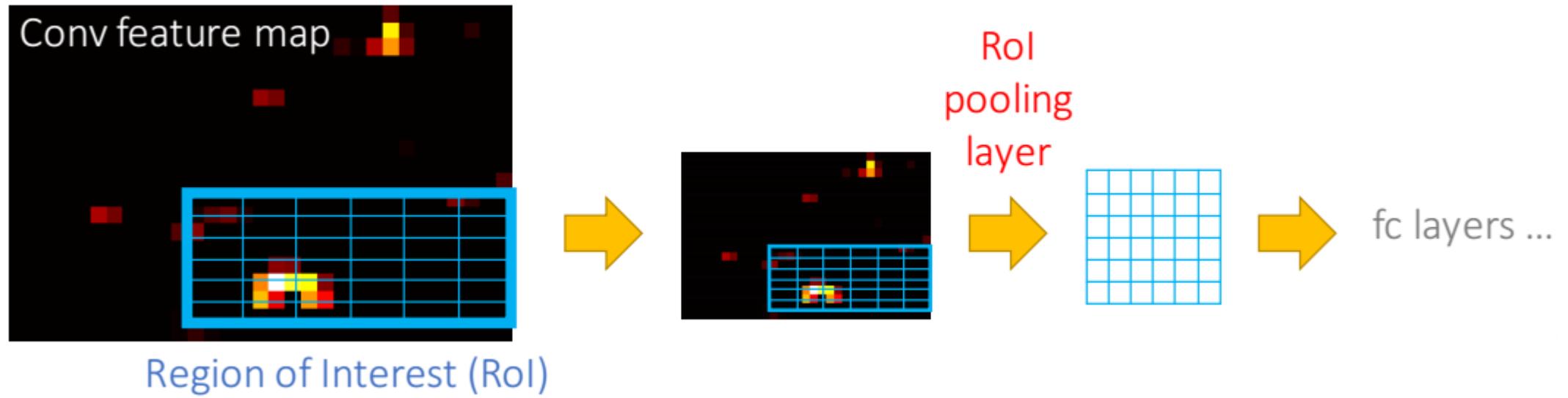
Spatial Pyramid
Representation (SPP)

“Conv 5” Feature Map
with region proposals

Apply single pass conv
net over whole image



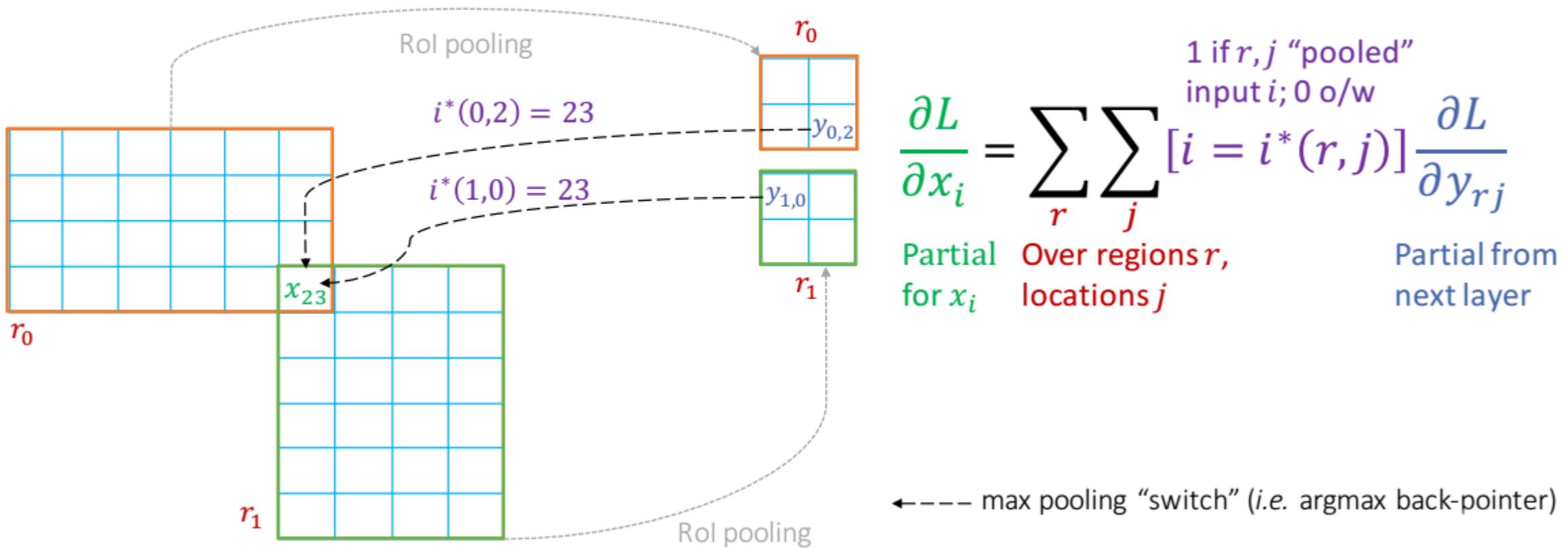
ROI Pooling Layer



Just a special case of the SPP layer with one pyramid level

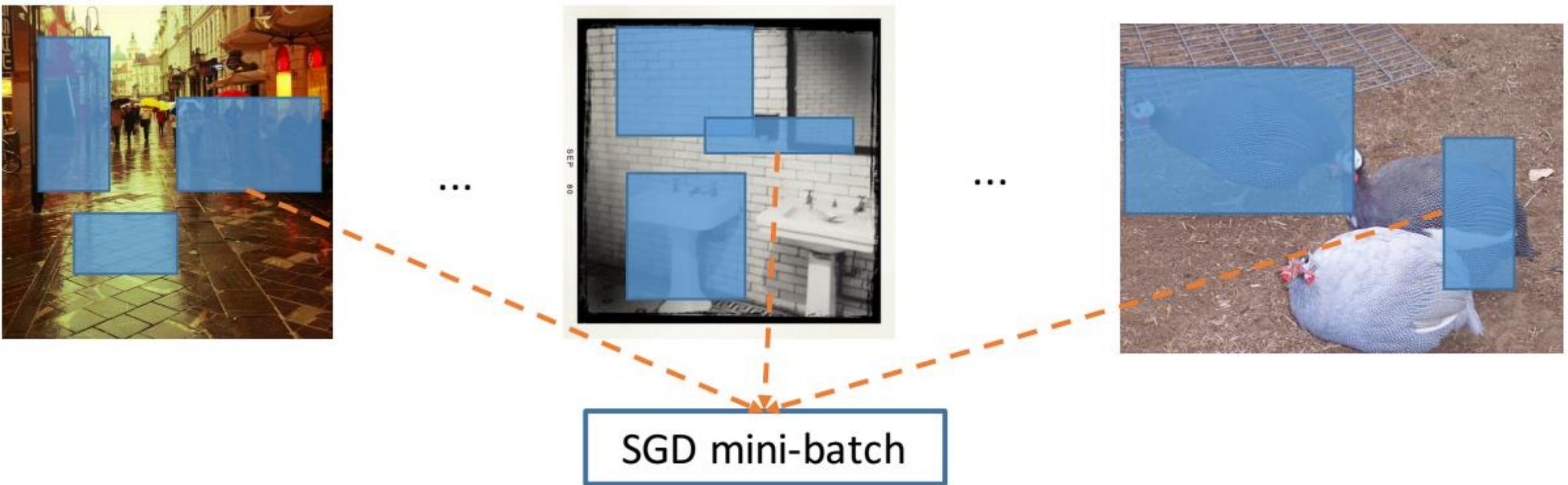
Differentiable ROI Pooling

ROI pooling / SPP is just like max pooling, except that pooling regions overlap



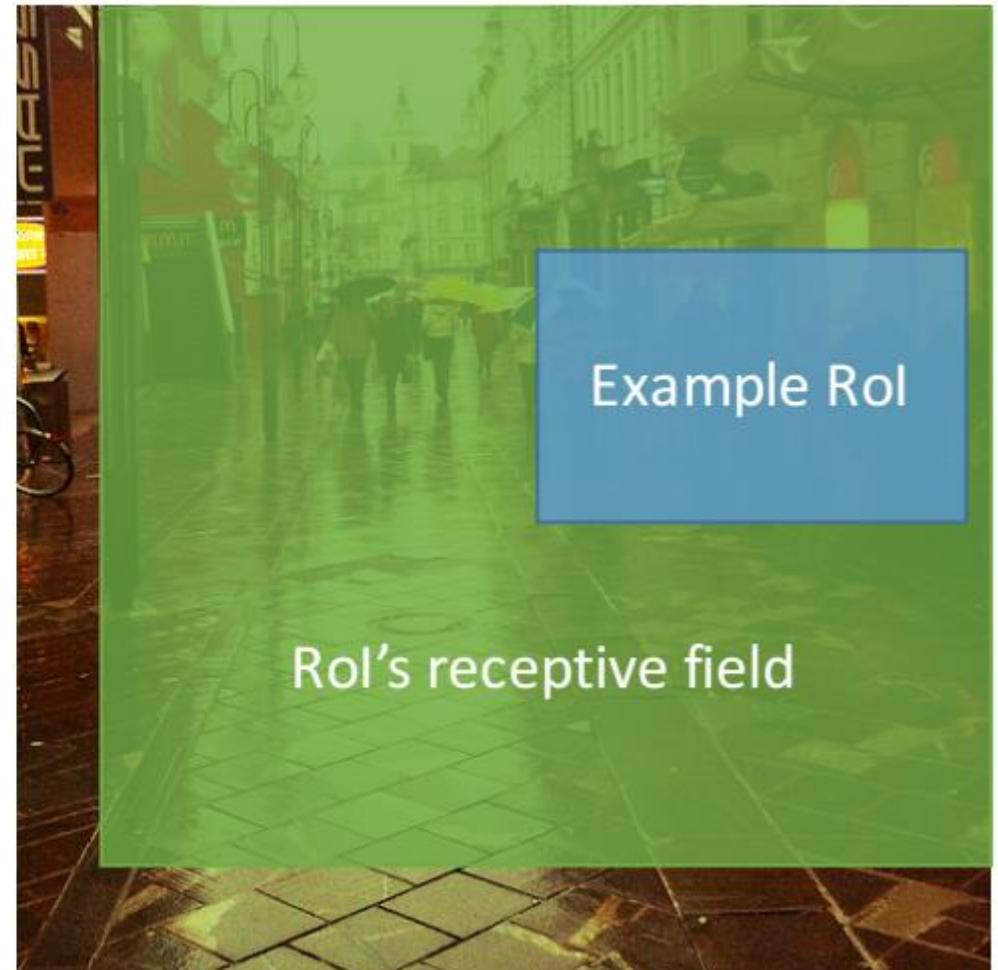
Efficient SGD

- R-CNN and SPP-net use region-wise sampling to make mini-batches
- Sample 128 example Rols uniformly at random
- Examples will come from different images with high probability



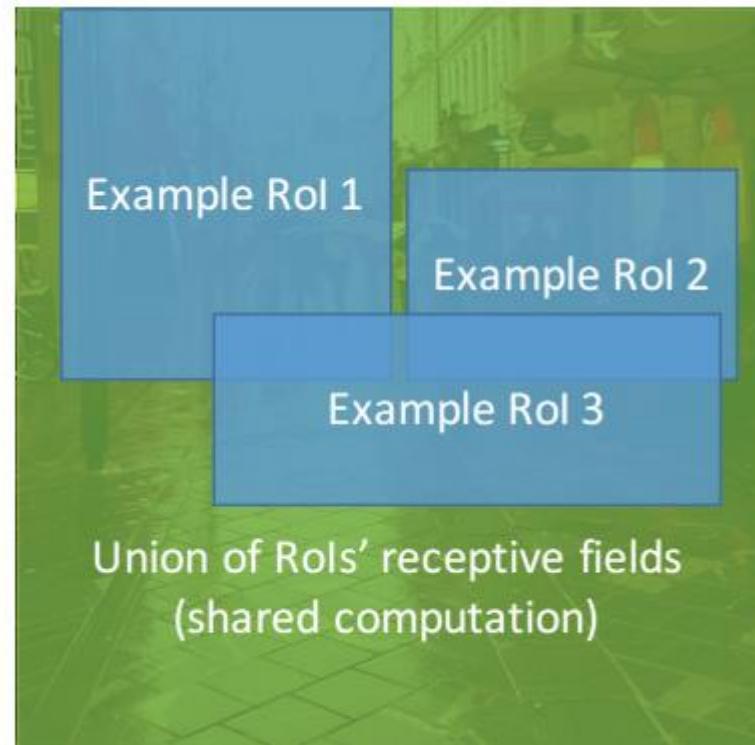
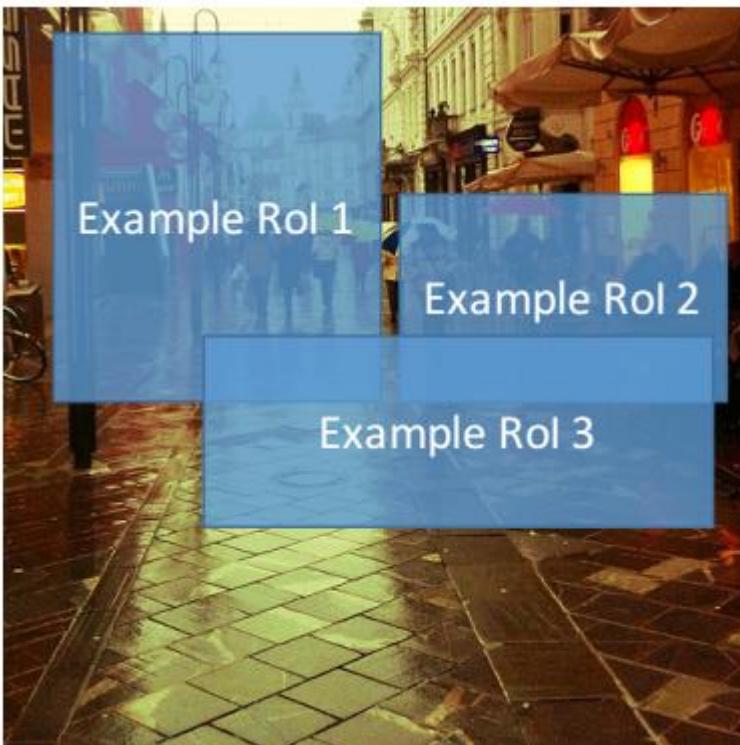
Efficient SGD

- The receptive field of one RoI could be very large (in worst case almost an whole image)

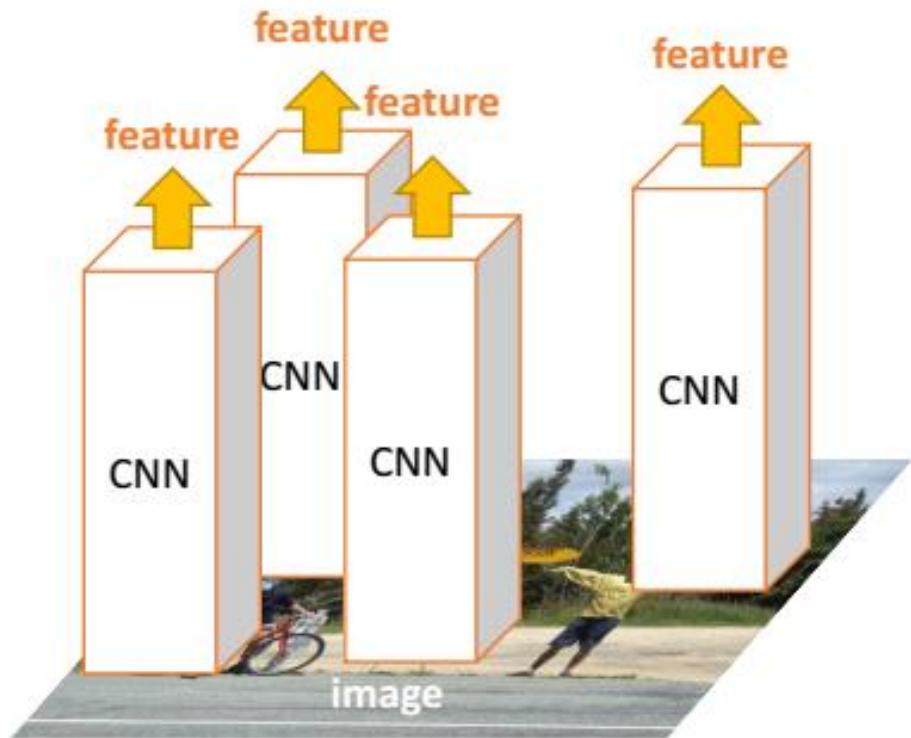


Efficient SGD

- Solution: use hierarchical sampling to build mini-batches
- Share computation between overlapping examples from the same image

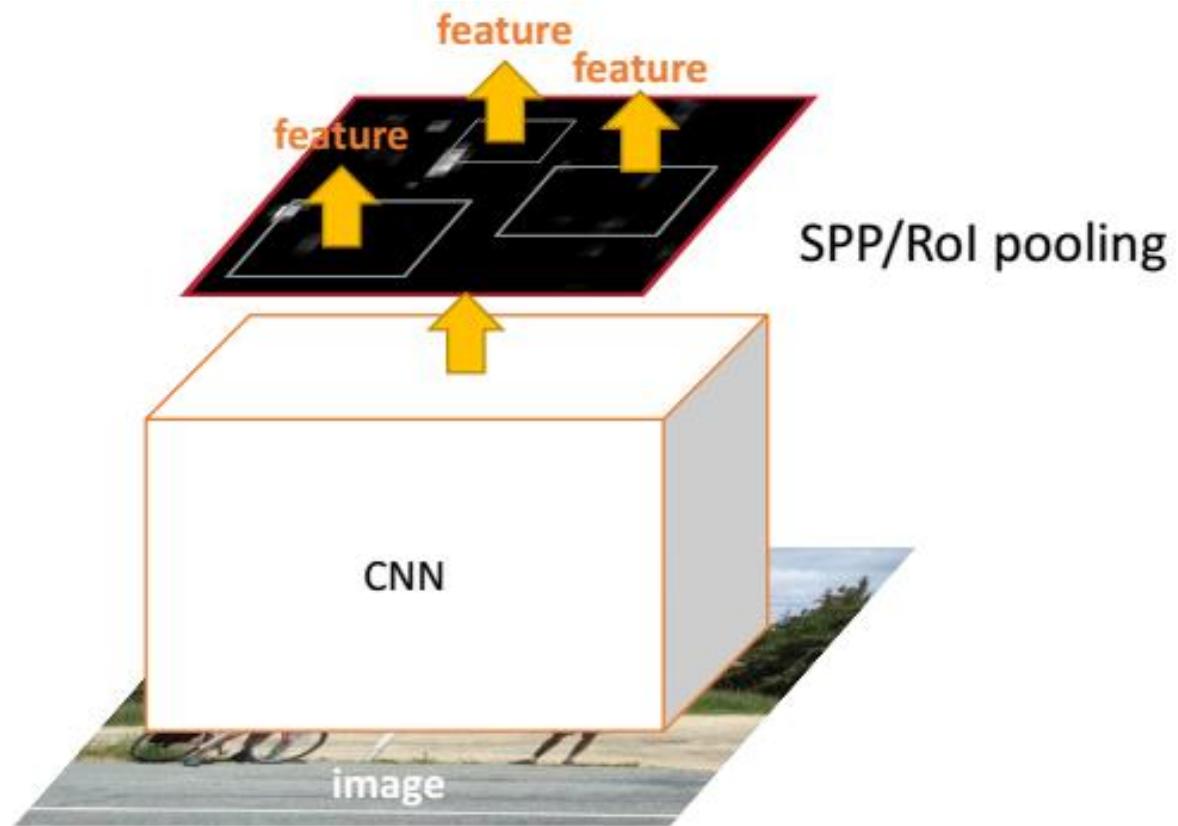


R-CNN vs. Fast R-CNN/SPP-Net



R-CNN

- Extract image regions
- 1 CNN per region (2000 CNNs)
- Classify region-based features



SPP-net & Fast R-CNN (the same forward pipeline)

- 1 CNN on the entire image
- Extract features from feature map regions
- Classify region-based features

Fast RCNN: Key Ideas

- Differentiable ROI Pooling
 - Improves accuracy
- Fast SGD by hierarchical sampling of ROIs
 - Improves speed
- What is still wrong
 - Out of network region proposals

Region Proposal from Feature Maps

- Object detection networks are fast (0.2s)...
- But what about region proposal?
 - Selective Search [Uijlings et al. ICCV 2011]: 2s per image
 - EdgeBoxes [Zitnick & Dollar. ECCV 2014]: 0.2s per image
- Can we do region proposal on the same set of feature maps?

Region Proposal from Feature Maps

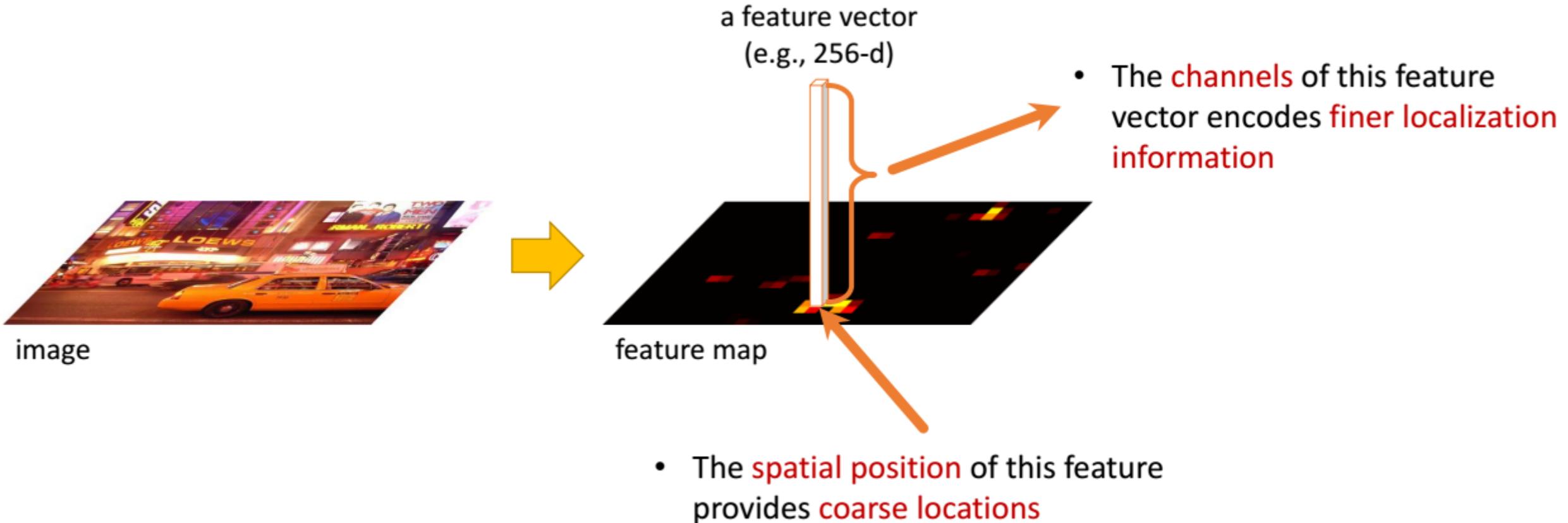
- By decoding **one response** at a single pixel, we can still roughly see the object outline*
- **Finer localization information** has been encoded in the channels of a convolutional feature response
- Extract this information for better localization...

* Zeiler & Fergus's method traces unpooling information so the visualization involves more than a single response. But other visualization methods reveal similar patterns.

Revisiting visualizations from Zeiler & Fergus

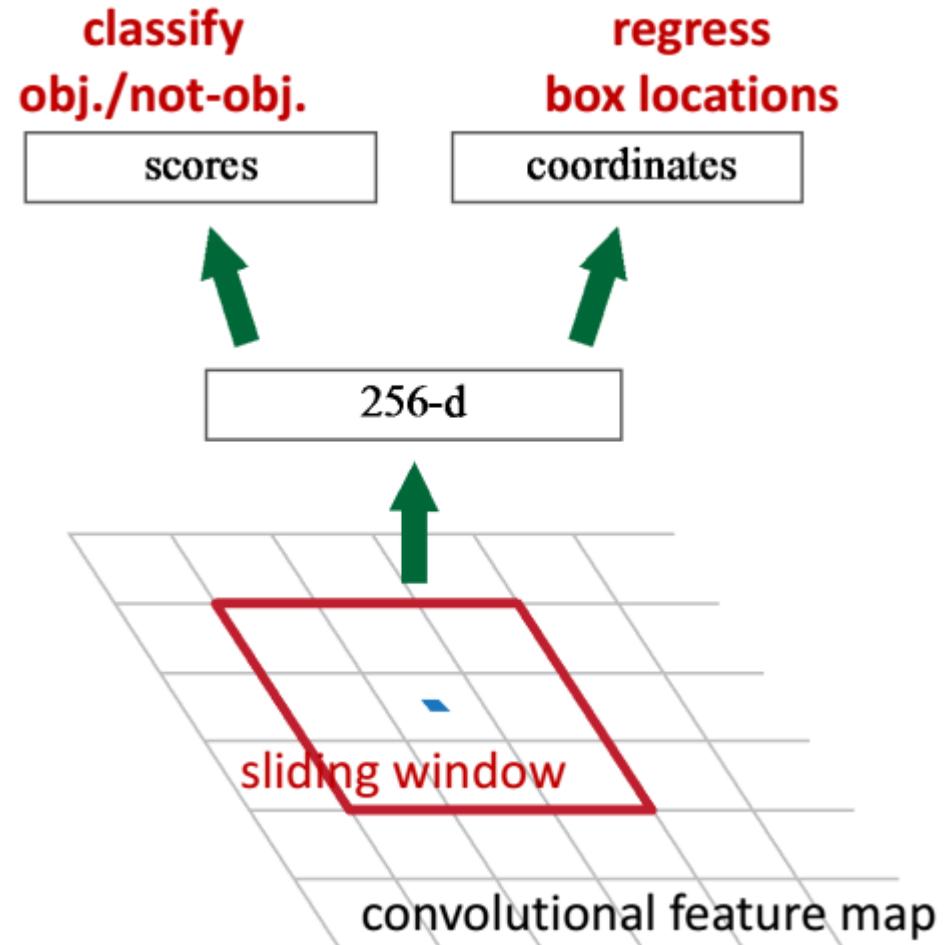


Region Proposal from Feature Maps



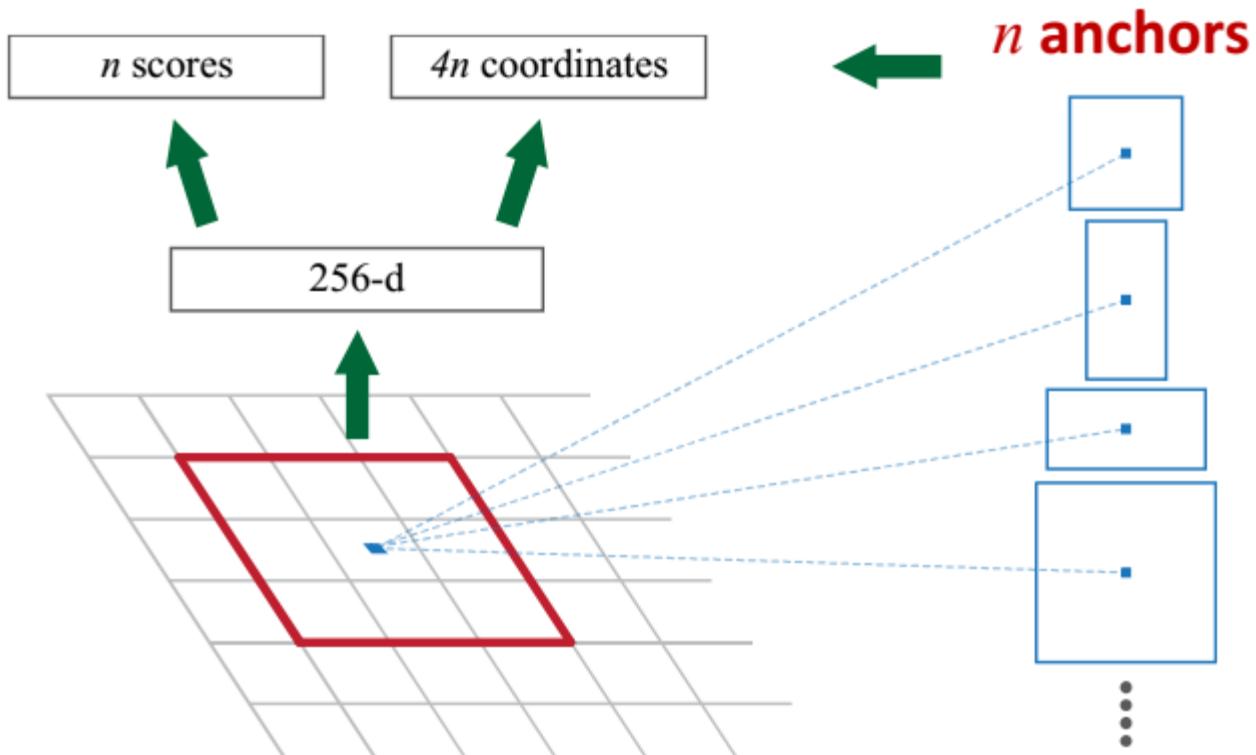
Region Proposal Network

- 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



Region Proposal Network

- **Anchors:** pre-defined reference boxes
- Box regression is with reference to anchors: regressing an anchor box to a ground-truth box
- Object probability is with reference to anchors, e.g.:
 - anchors as positive samples: if $\text{IoU} > 0.7$ or IoU is max
 - anchors as negative samples: if $\text{IoU} < 0.3$

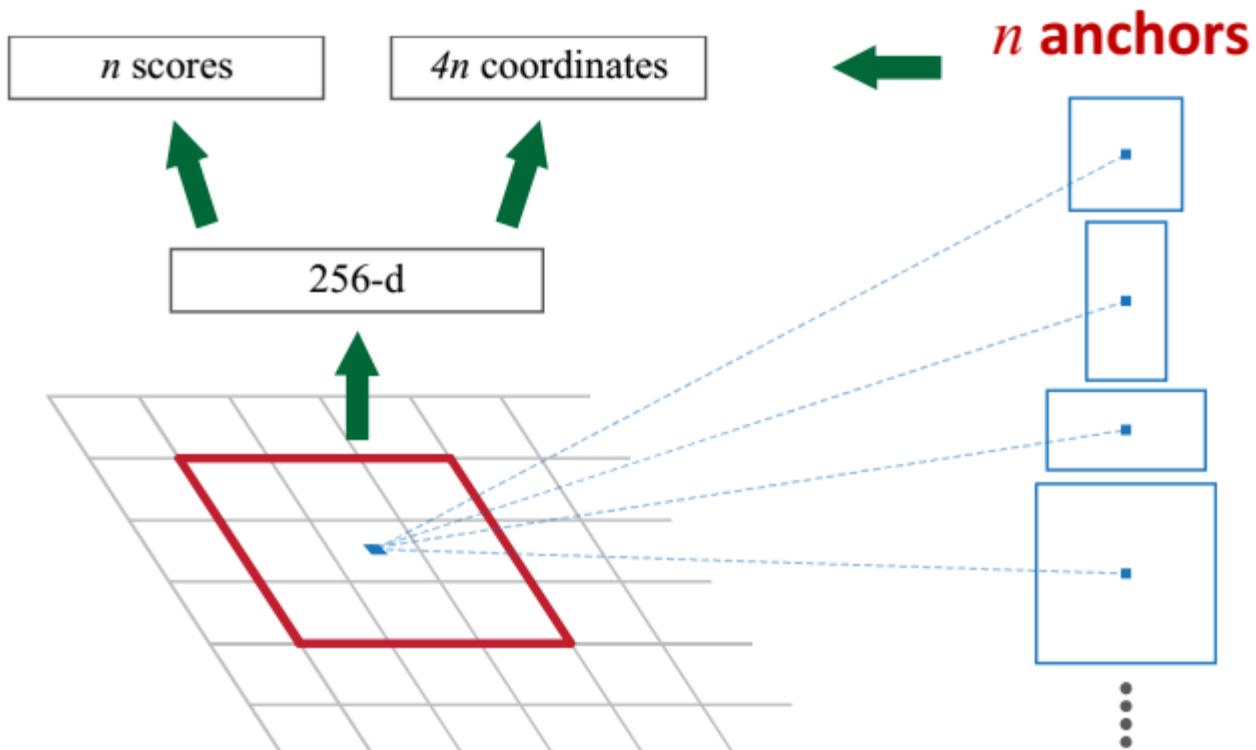


Region Proposal Network

- **Anchors:** pre-defined reference boxes

Translation-invariant anchors:

- the same set of anchors are used at each sliding position
- the same prediction functions (with reference to the sliding window) are used
- a translated object will have a translated prediction

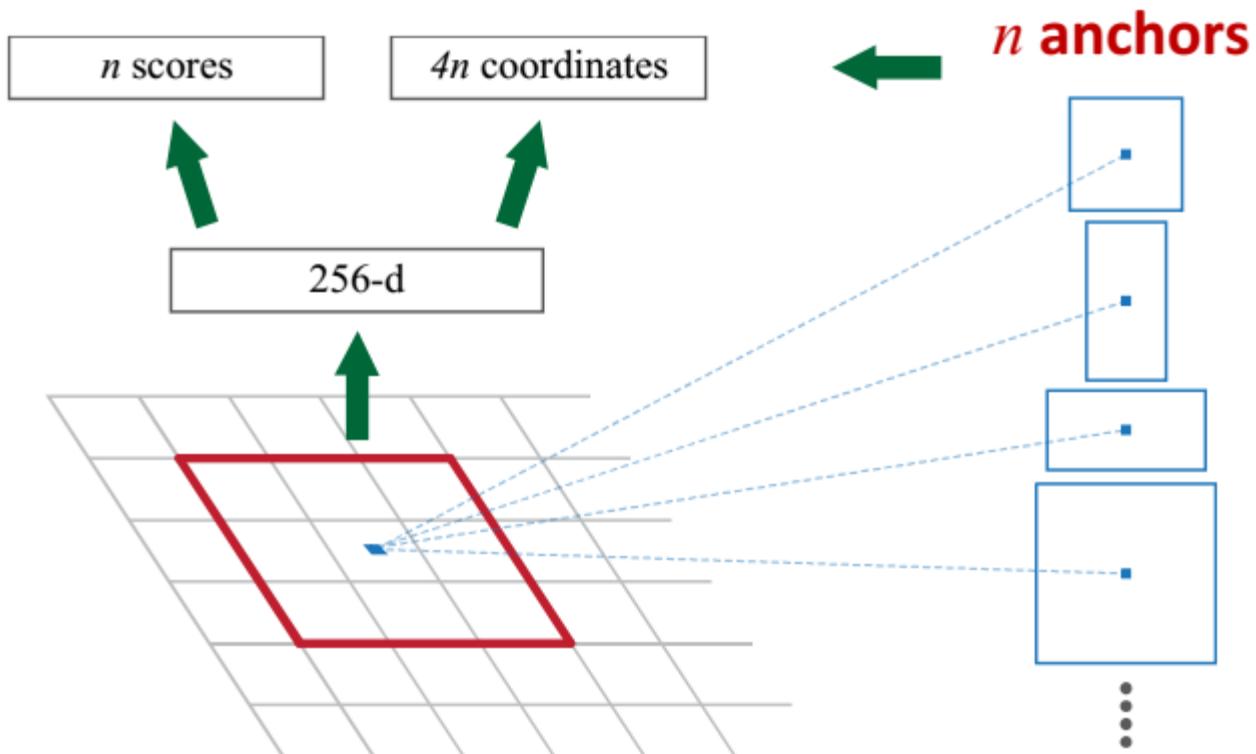


Region Proposal Network

- **Anchors:** pre-defined reference boxes

Multi-scale/size anchors:

- multiple anchors are used at each position: e.g., 3 scales ($128^2, 256^2, 512^2$) and 3 aspect ratios (2:1, 1:1, 1:2) yield 9 anchors
- each anchor has its own prediction function
- single-scale features, multi-scale predictions



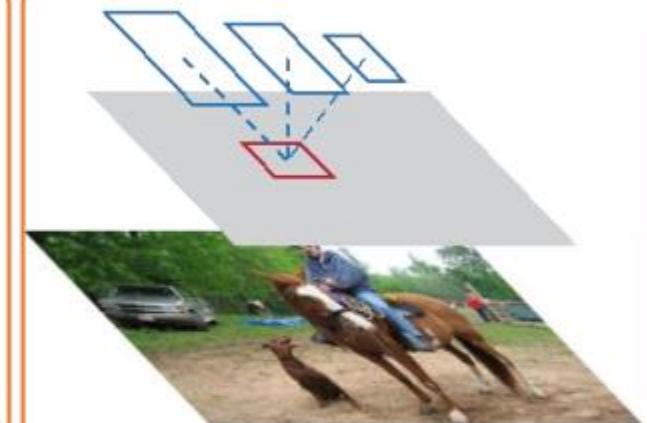
Comparisons of Multi Scale Strategies



Image/Feature Pyramid

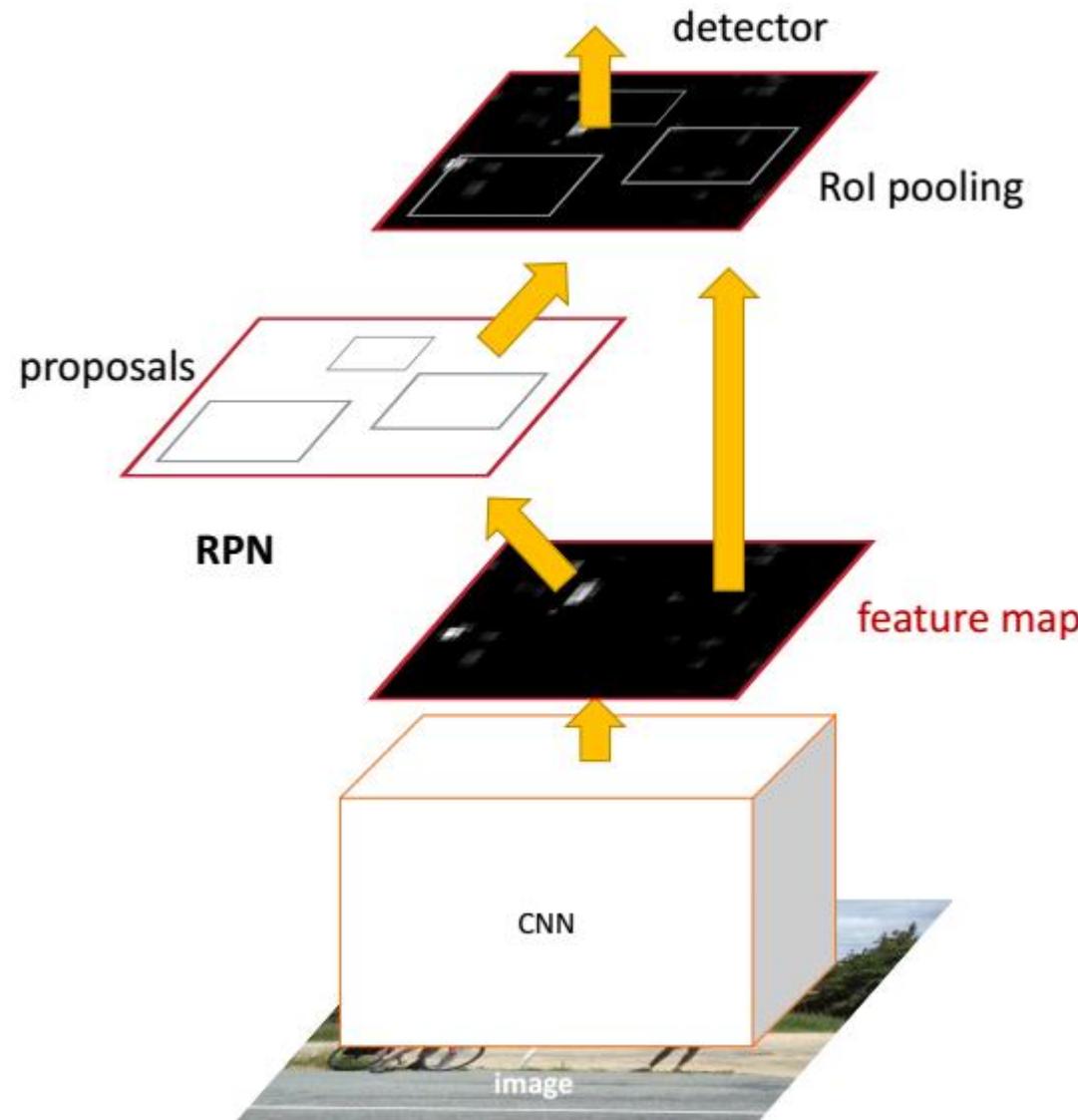


Filter Pyramid



Anchor Pyramid

Faster RCNN



YOLO



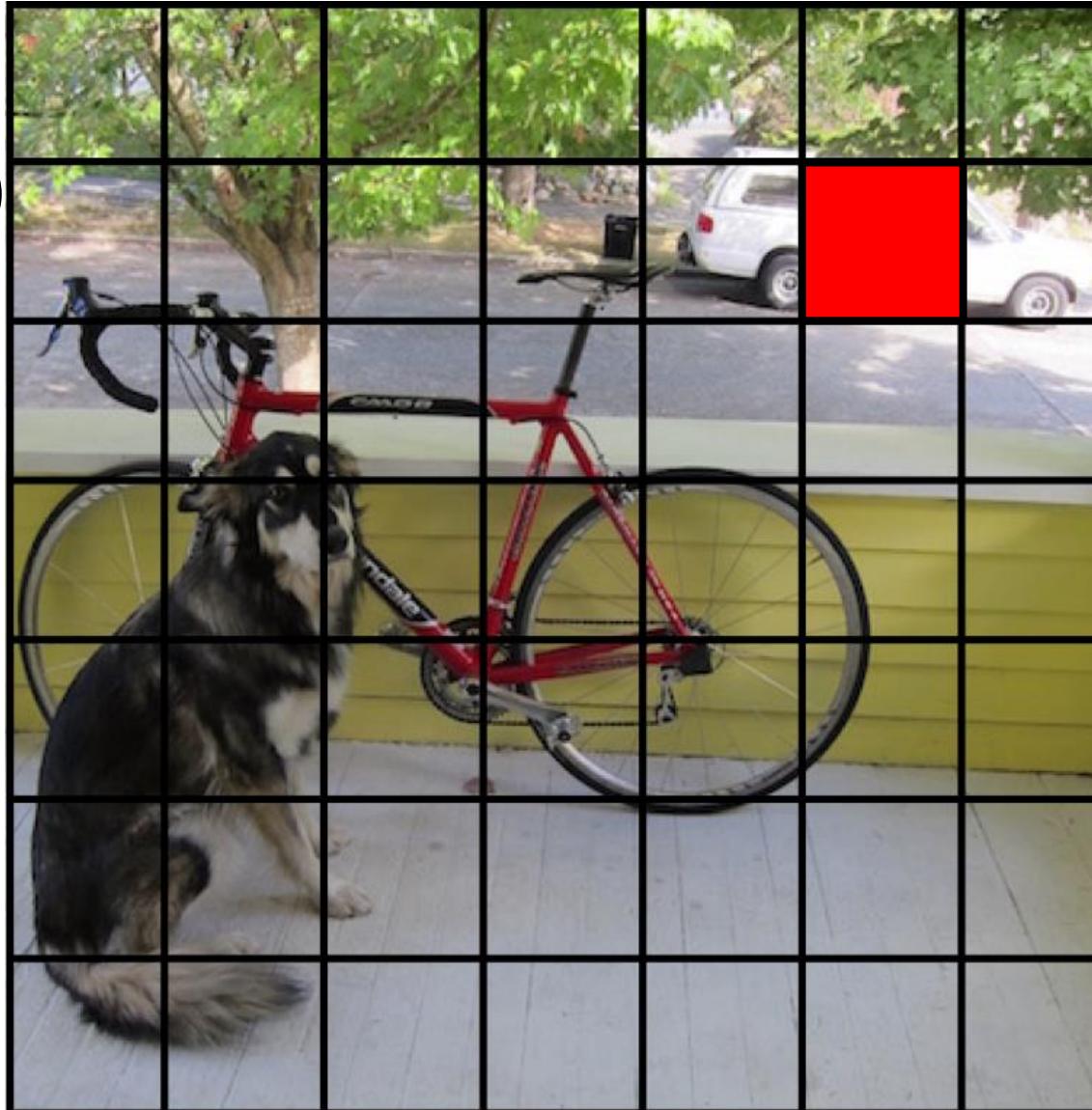
YOLO

Split the image
into a grid



YOLO

Each cell predicts
boxes and
confidences: $P(\text{Object})$



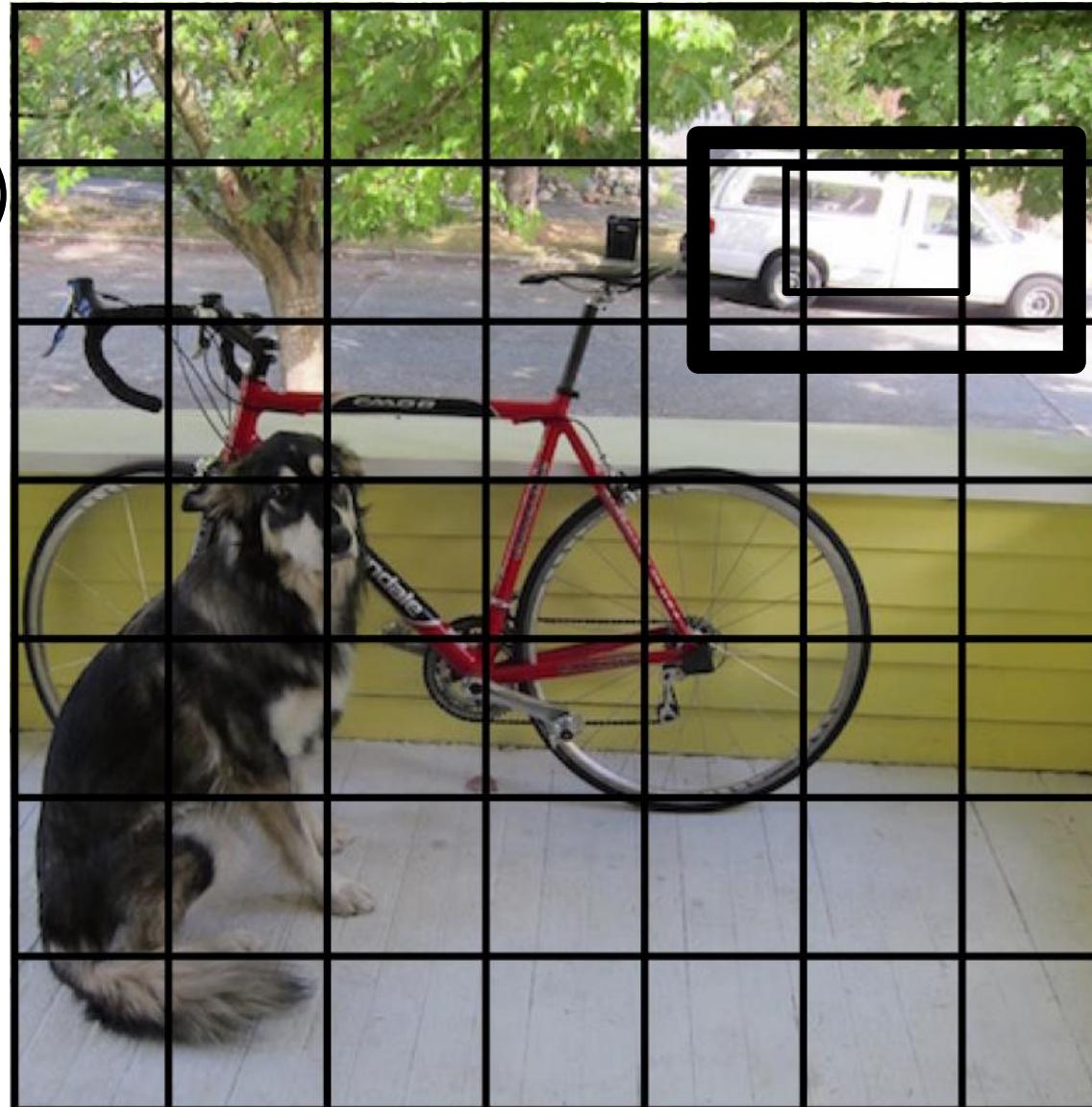
YOLO

Each cell predicts
boxes and
confidences: $P(\text{Object})$



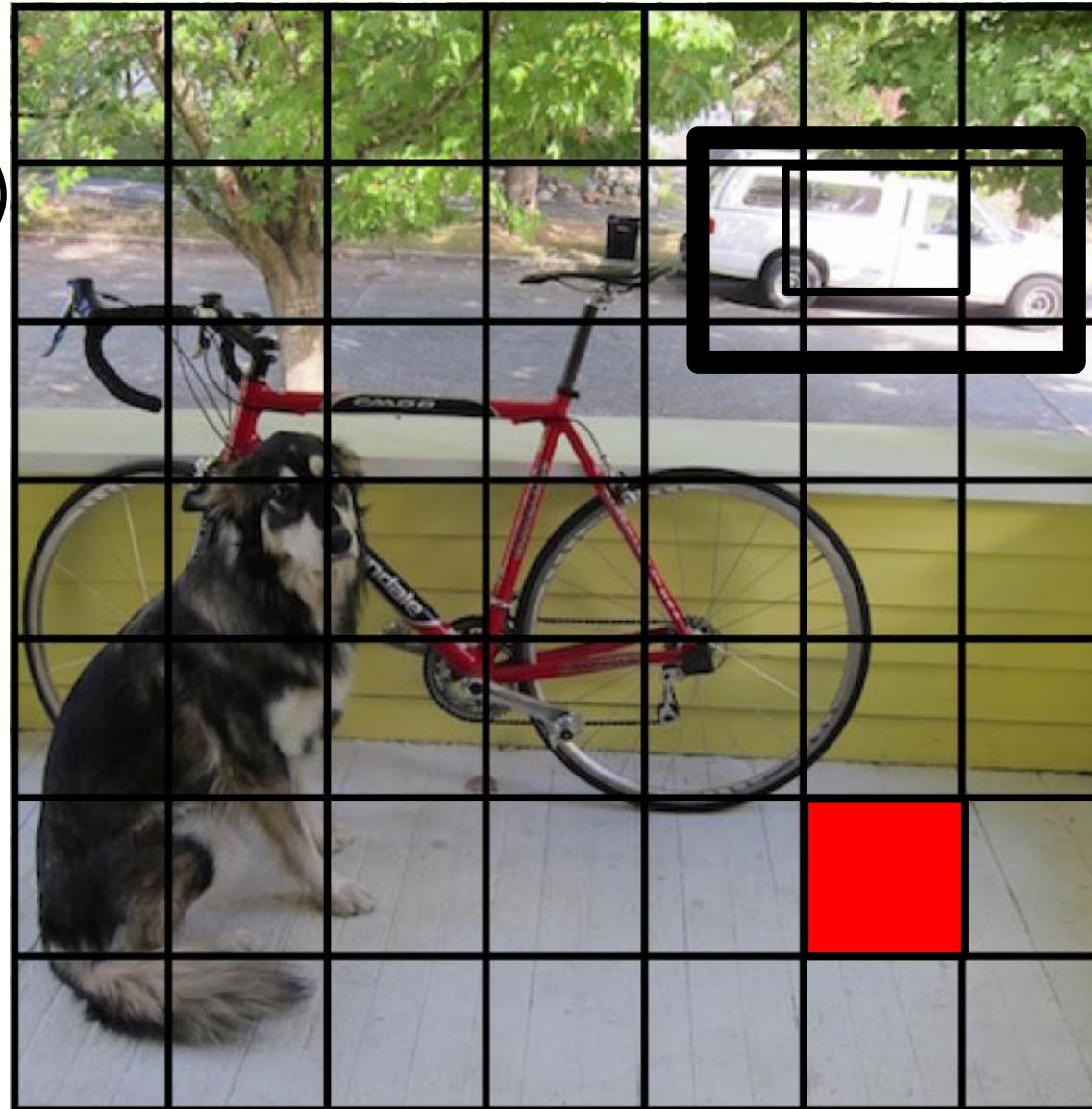
YOLO

Each cell predicts
boxes and
confidences: $P(\text{Object})$



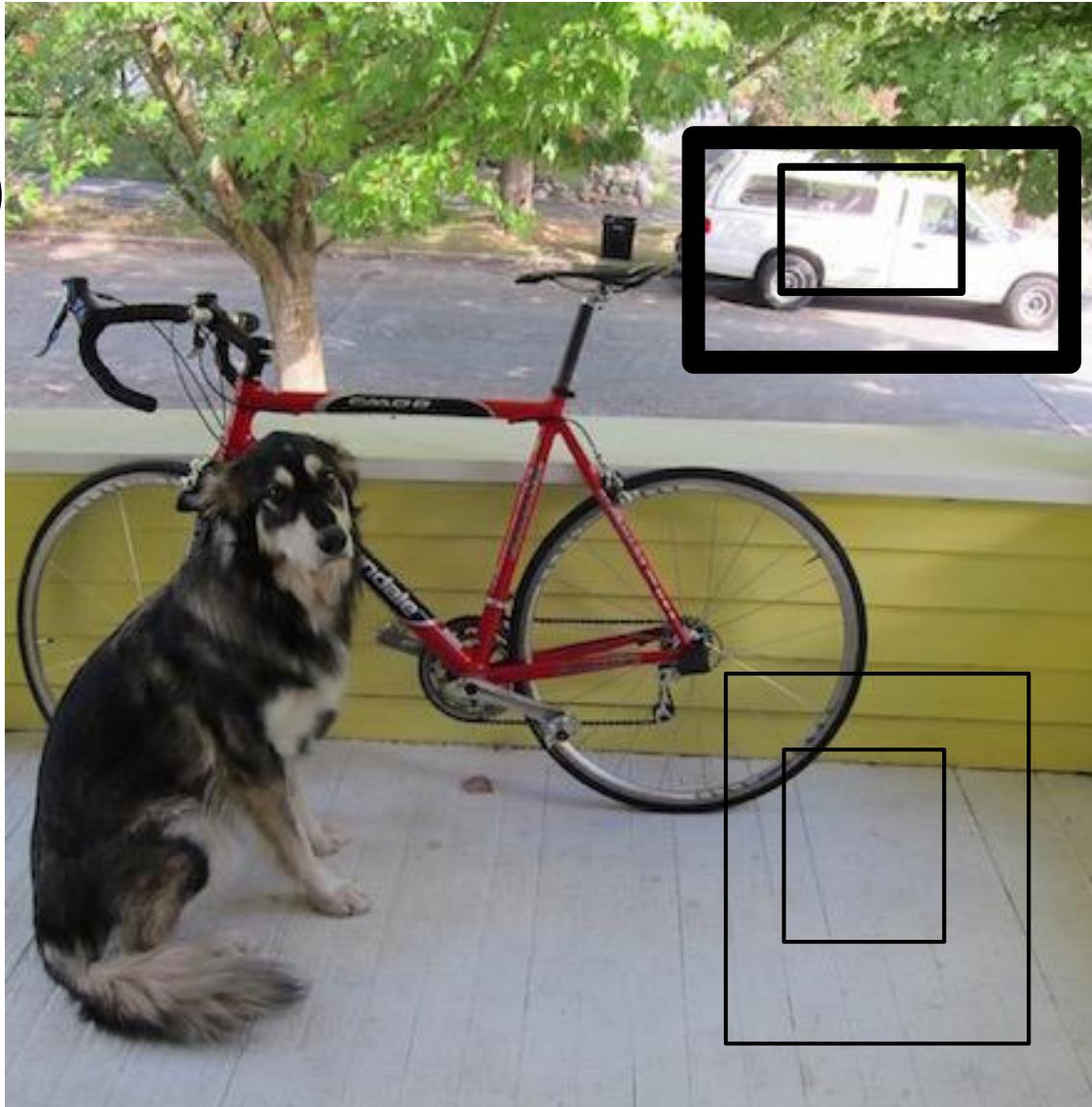
YOLO

Each cell predicts
boxes and
confidences: $P(\text{Object})$



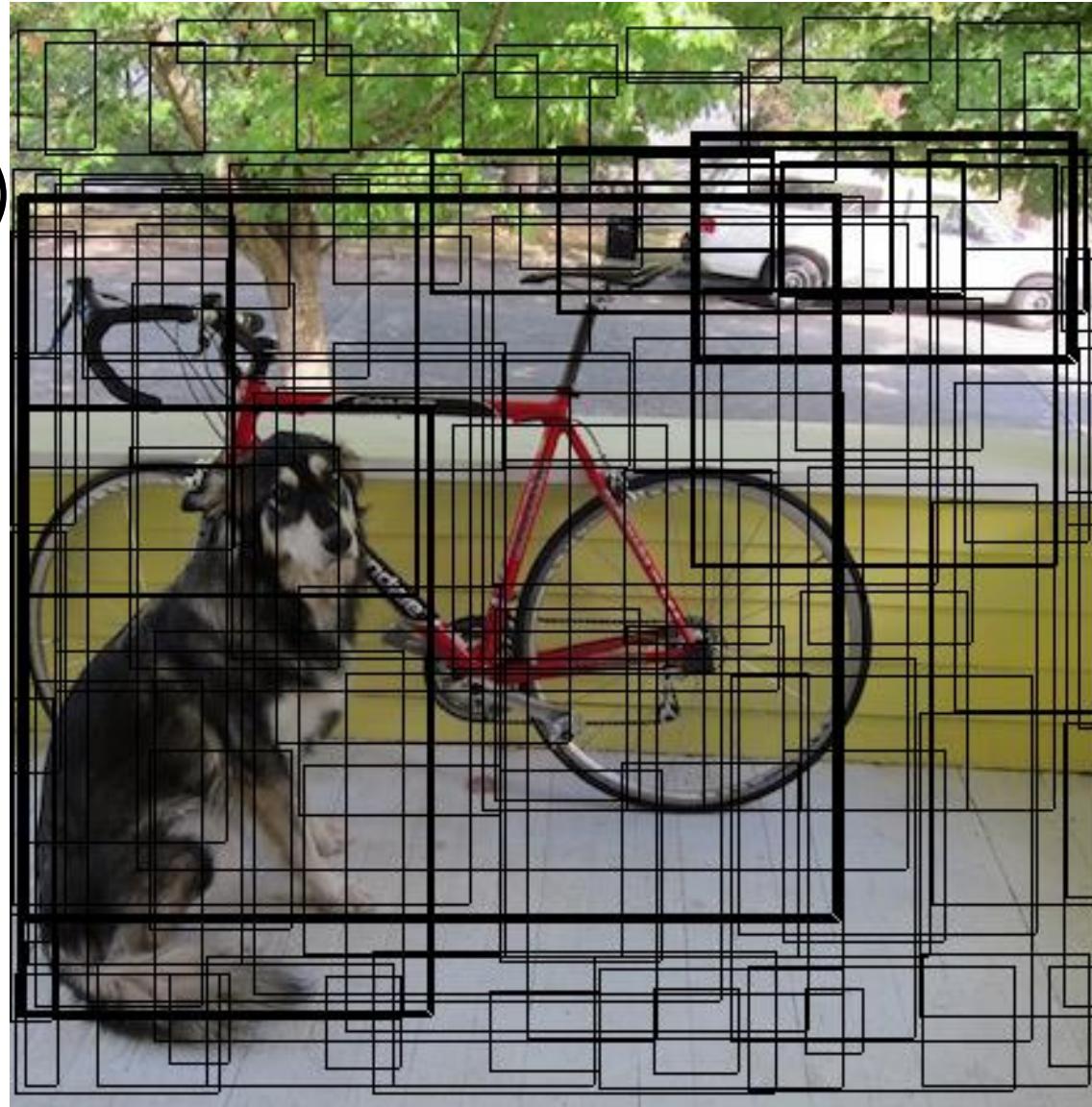
YOLO

Each cell predicts
boxes and
confidences: $P(\text{Object})$



YOLO

Each cell predicts
boxes and
confidences: $P(\text{Object})$



YOLO

Each cell also predicts a class probability.



YOLO

Each cell also predicts a class probability.
Conditioned on object: $P(\text{Car} | \text{Object})$

Bicycle



Car

Dog

Dining
Table

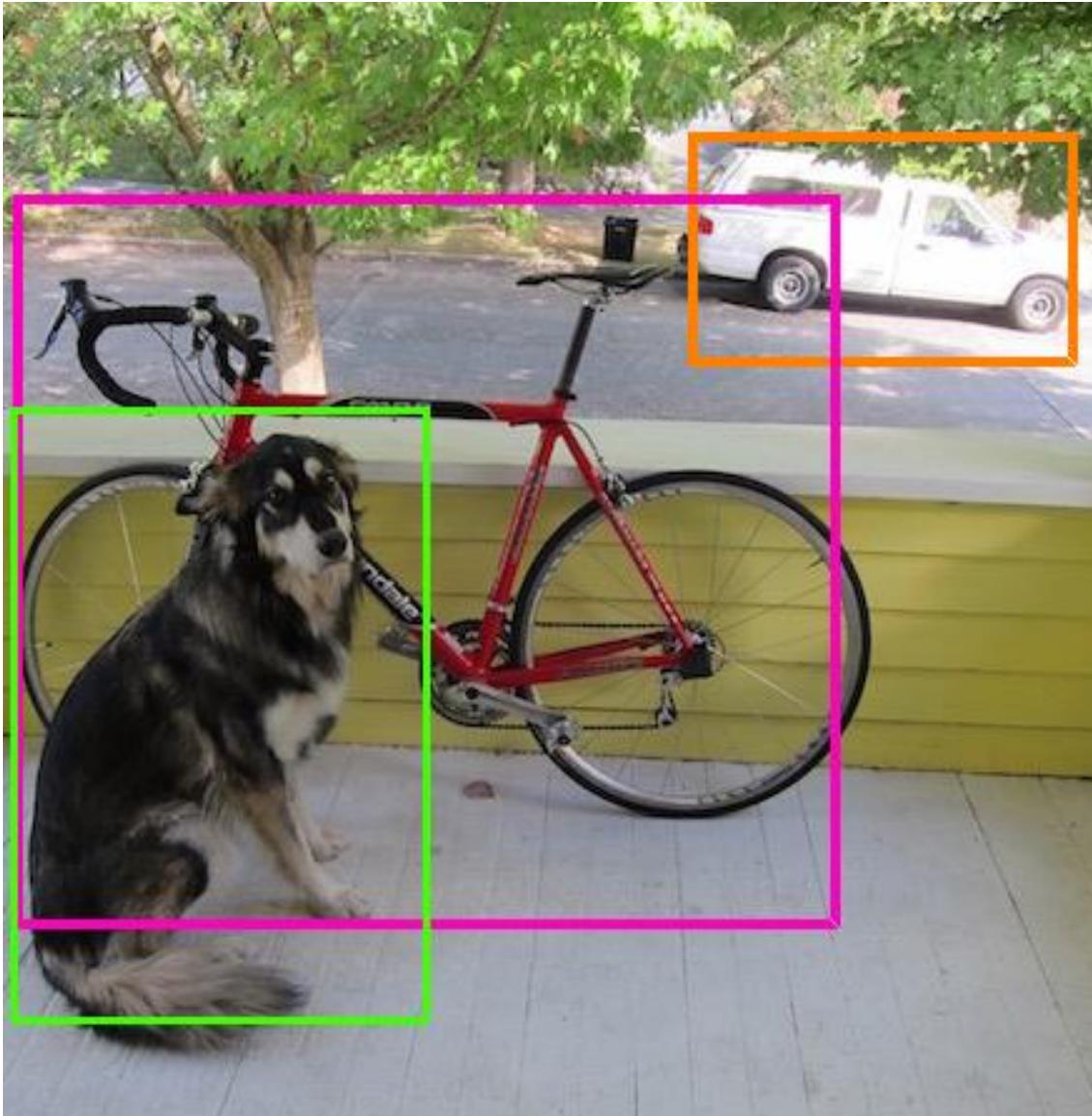
YOLO

Then combine the
box and class
predictions.



YOLO

Finally do NMS and threshold detections



YOLO

This parameterization fixes the output size

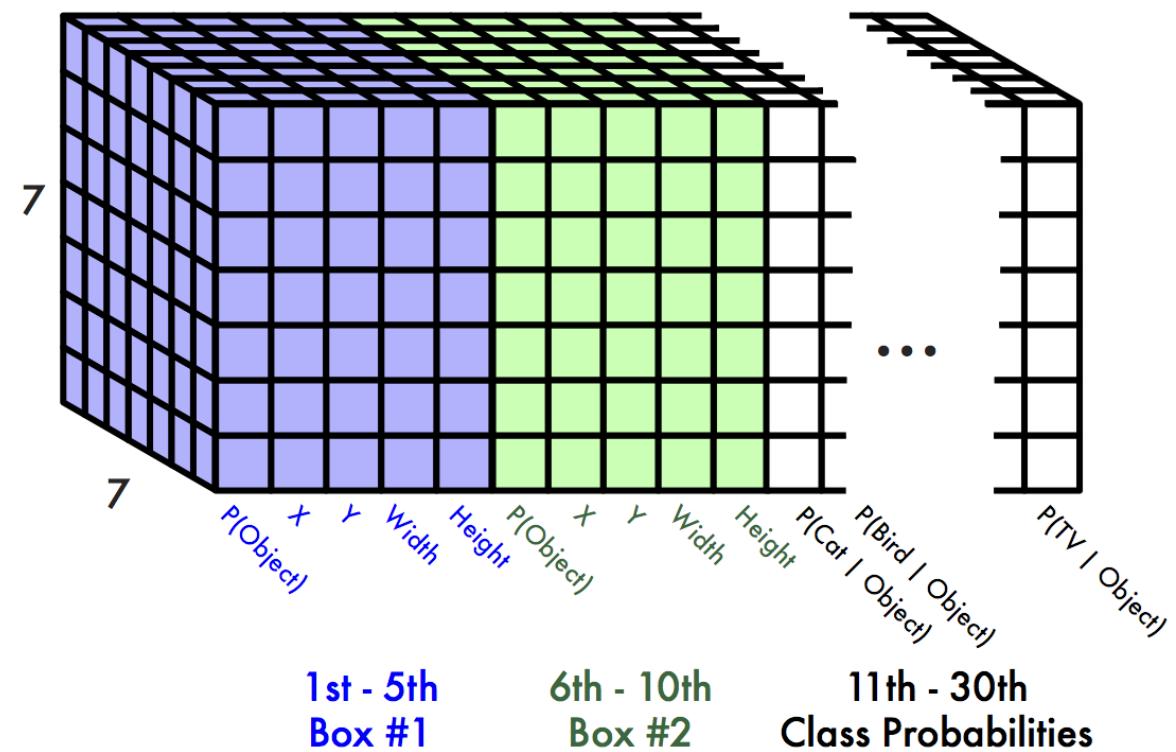
Each cell predicts:

- For each bounding box:
 - 4 coordinates (x, y, w, h)
 - 1 confidence value
- Some number of class probabilities

For Pascal VOC:

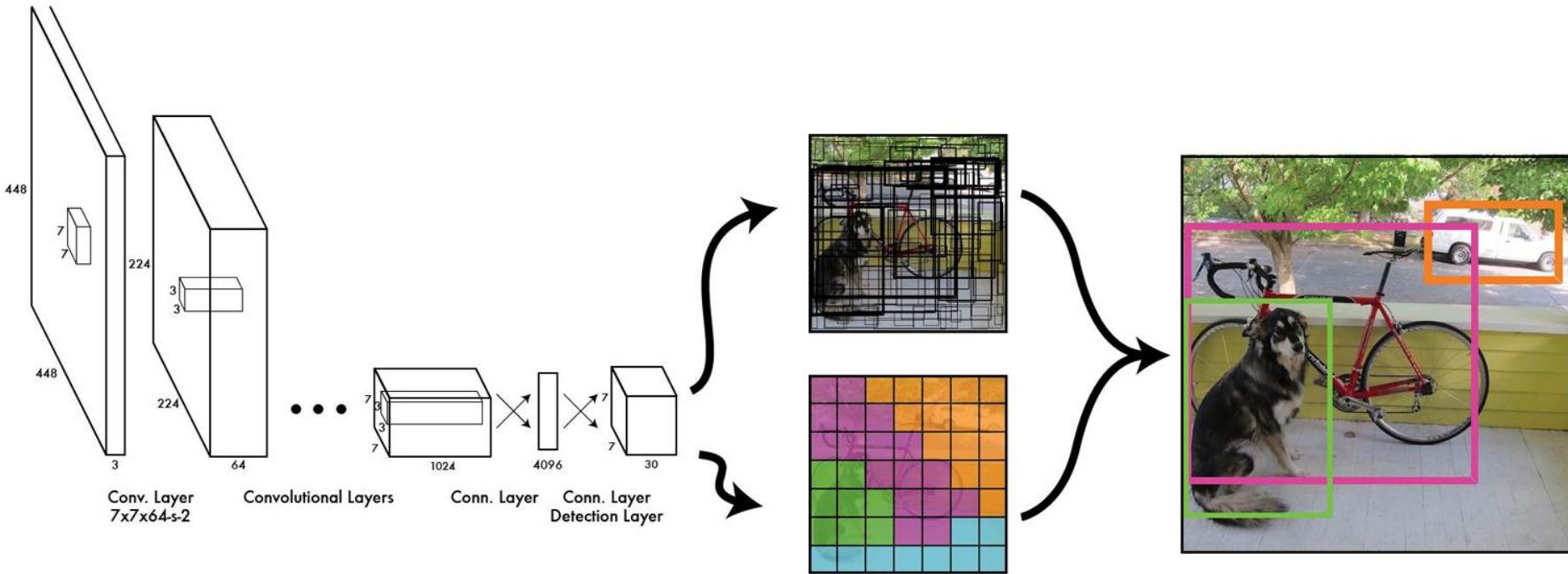
- 7x7 grid
- 2 bounding boxes / cell
- 20 classes

$$7 \times 7 \times (2 \times 5 + 20) = 7 \times 7 \times 30 \text{ tensor} = \mathbf{1470 \text{ outputs}}$$



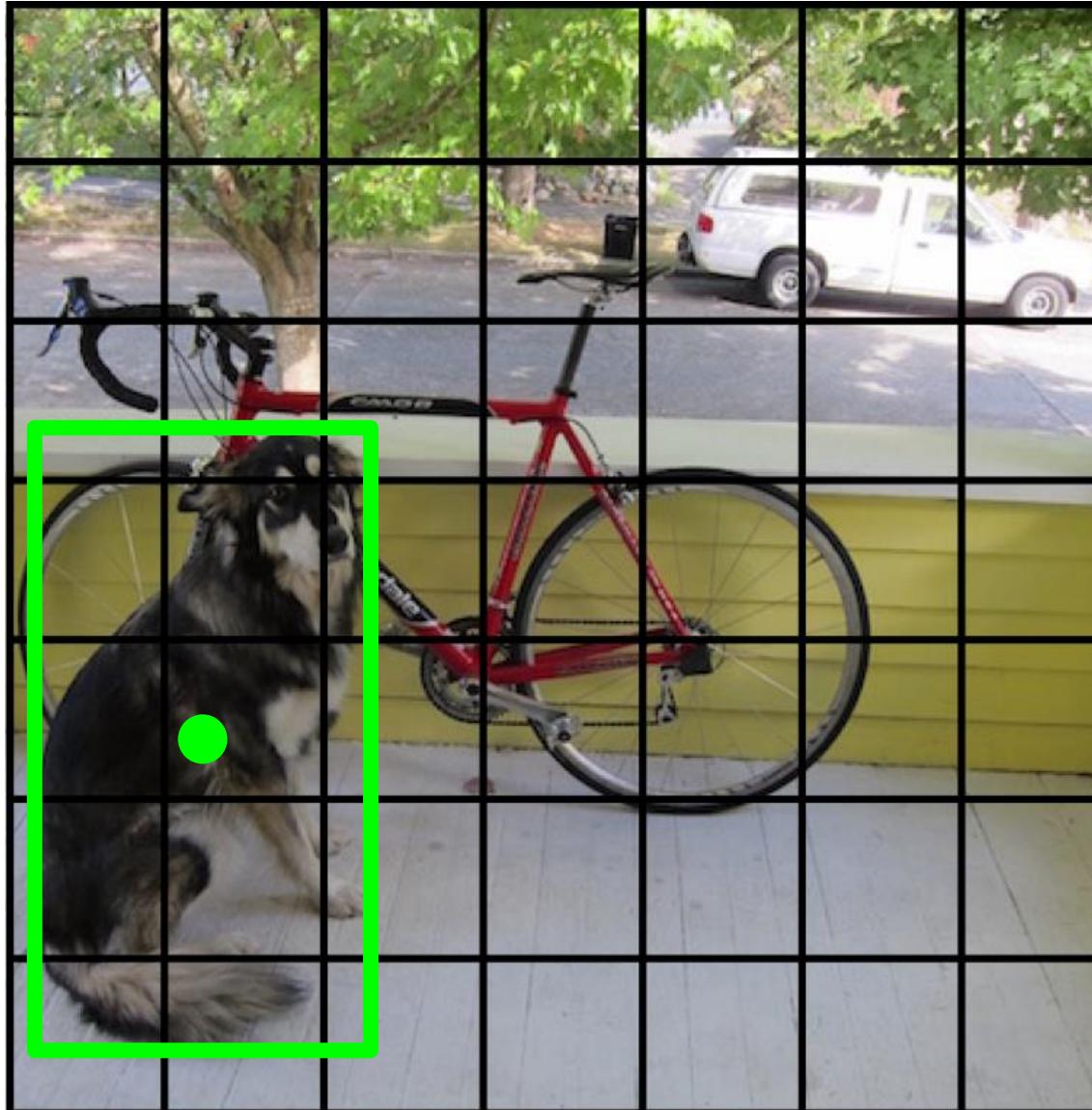
YOLO

Thus we can train one neural network to be a whole detection pipeline



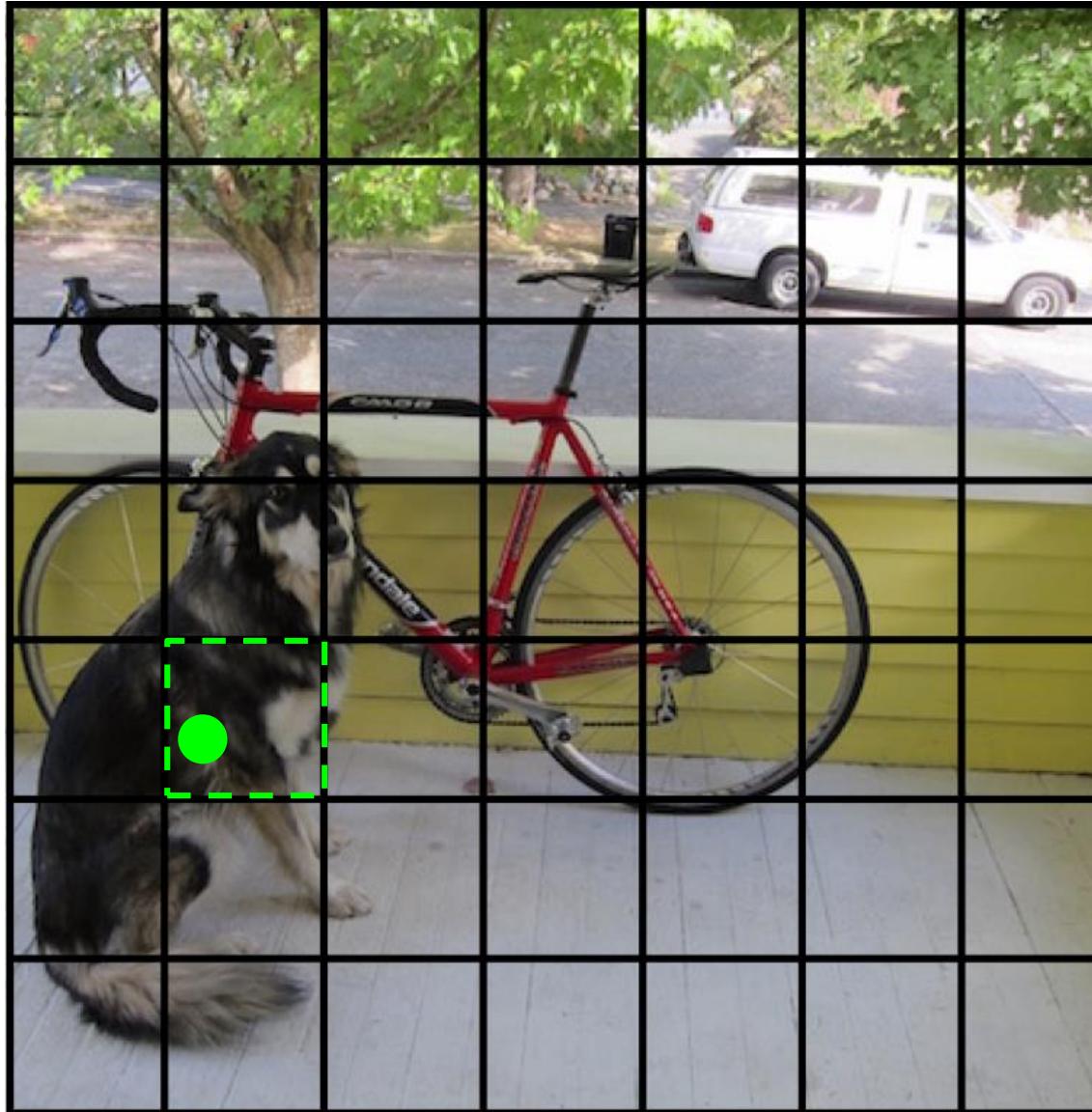
YOLO

During training, match example to the right cell



YOLO

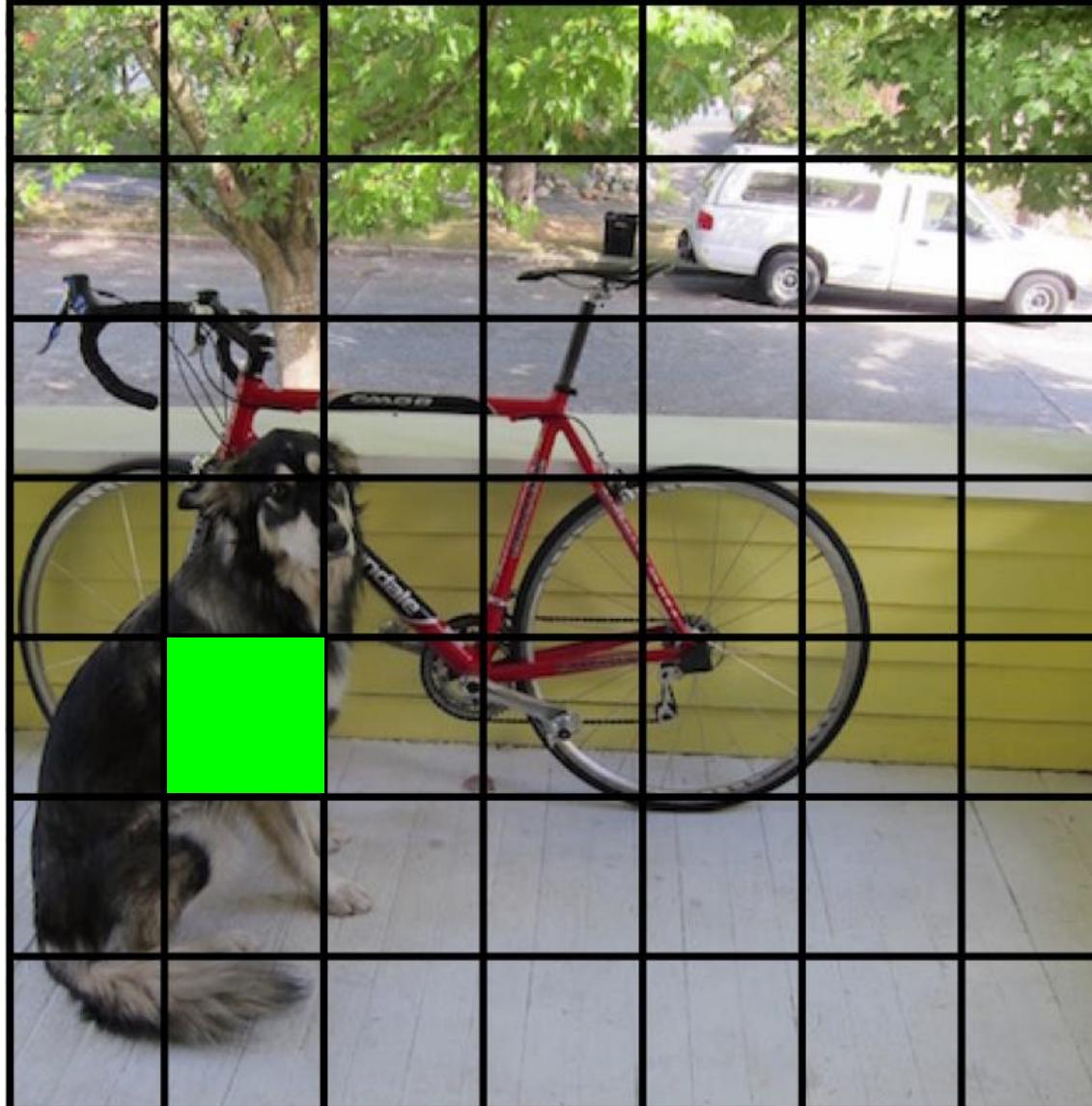
During training, match example to the right cell



YOLO

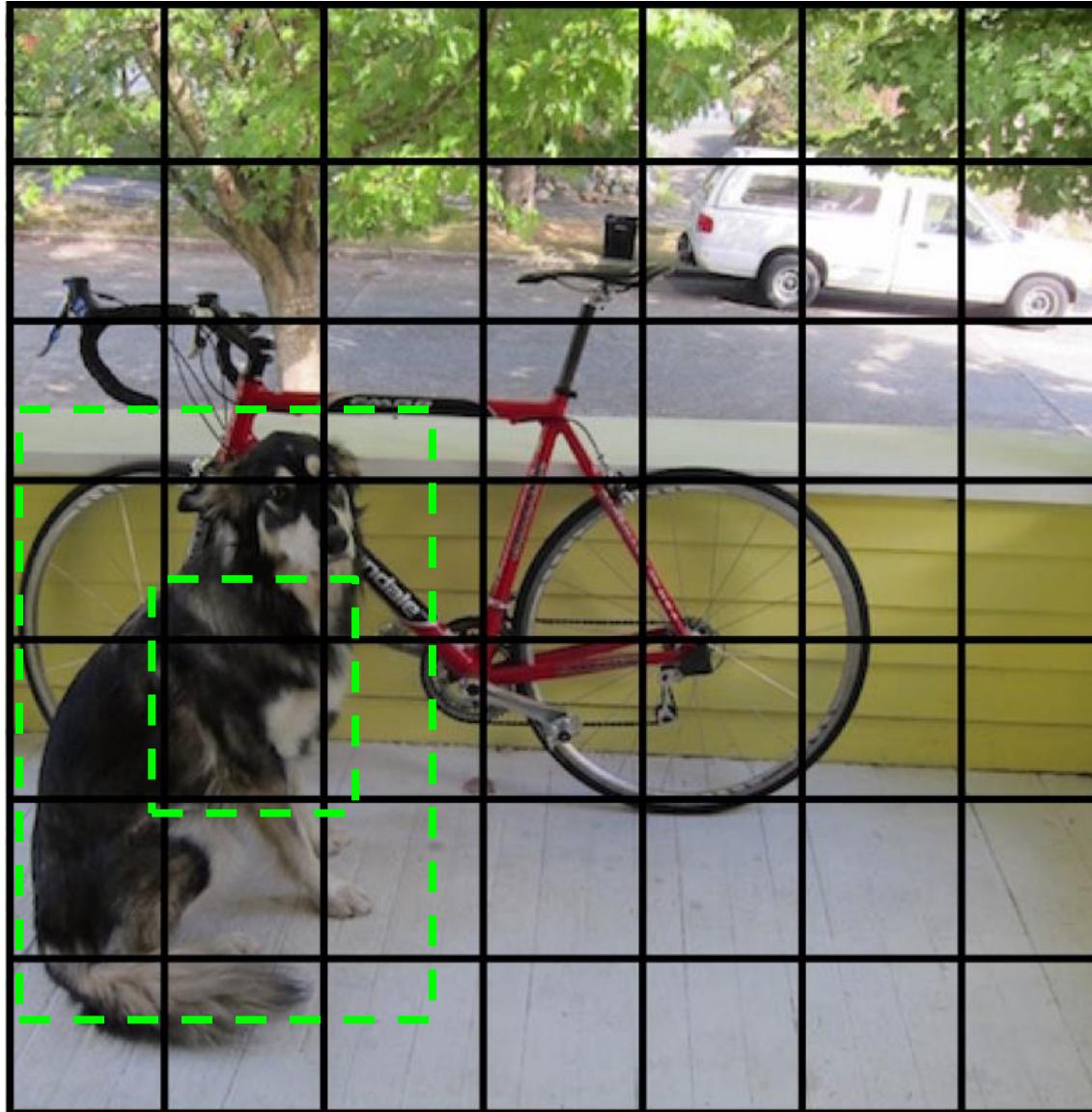
Adjust that cell's class prediction

Dog = 1
Cat = 0
Bike = 0
...



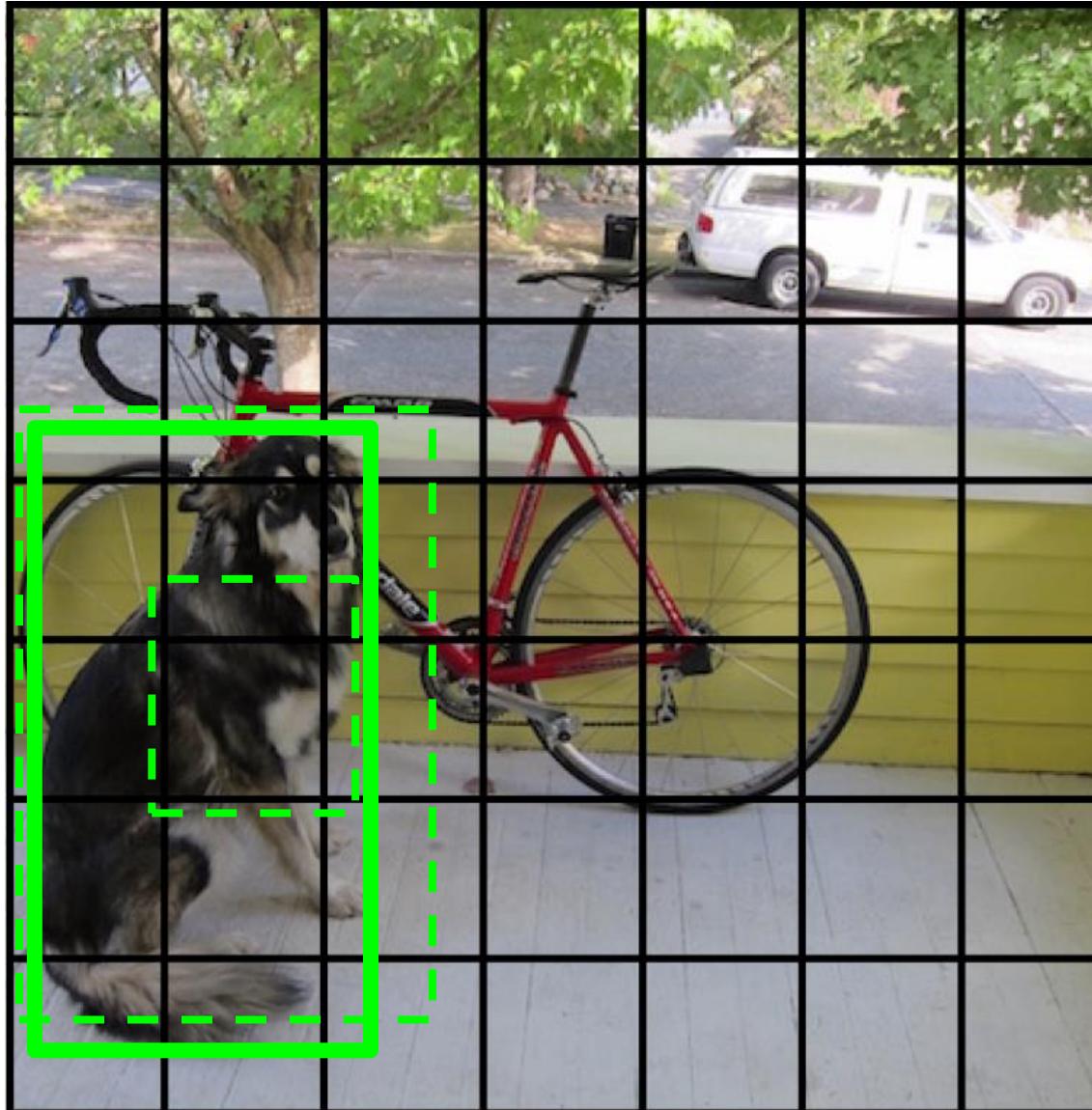
YOLO

Look at that cell's predicted boxes



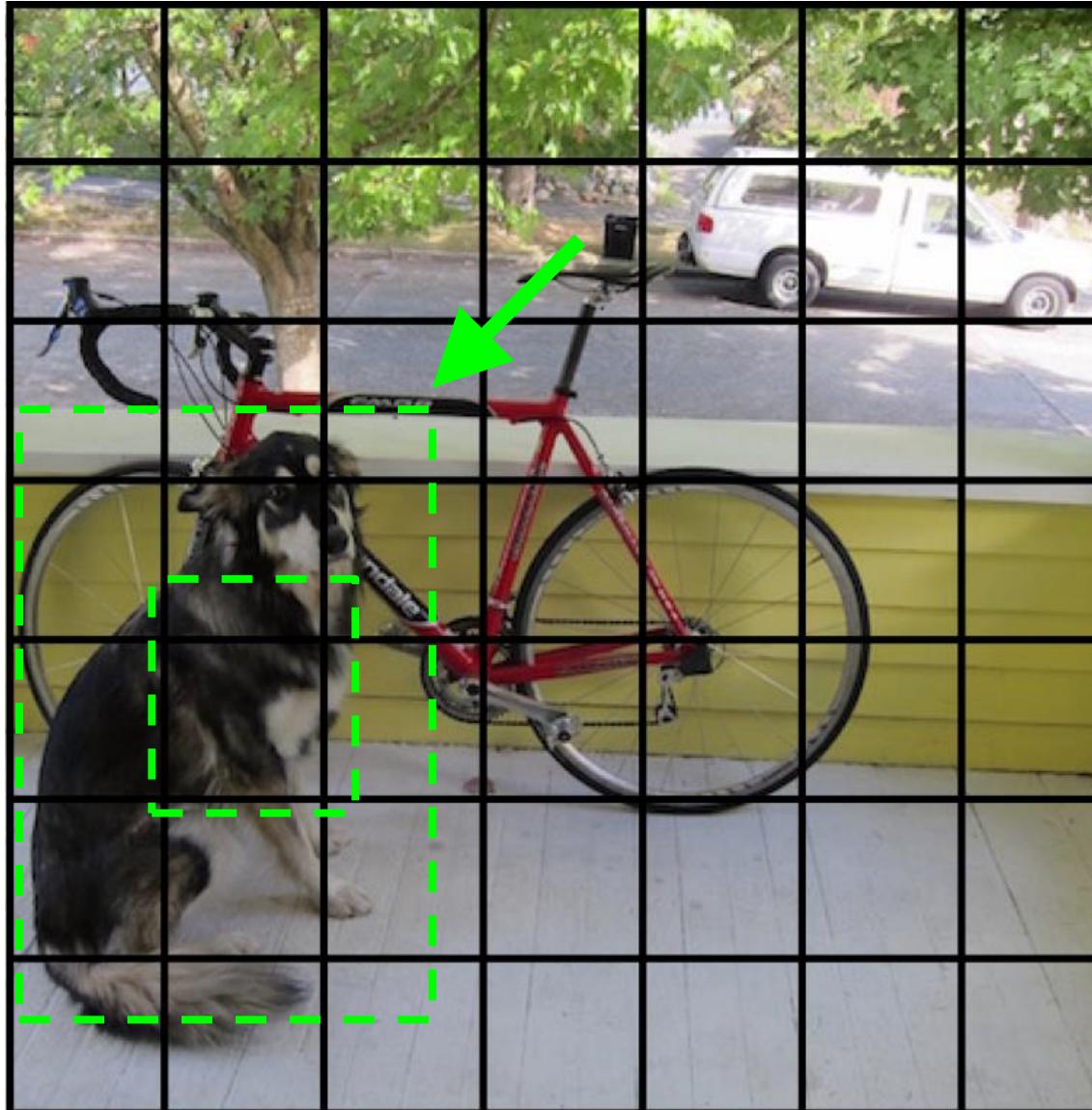
YOLO

Find the best one, adjust it, increase the confidence



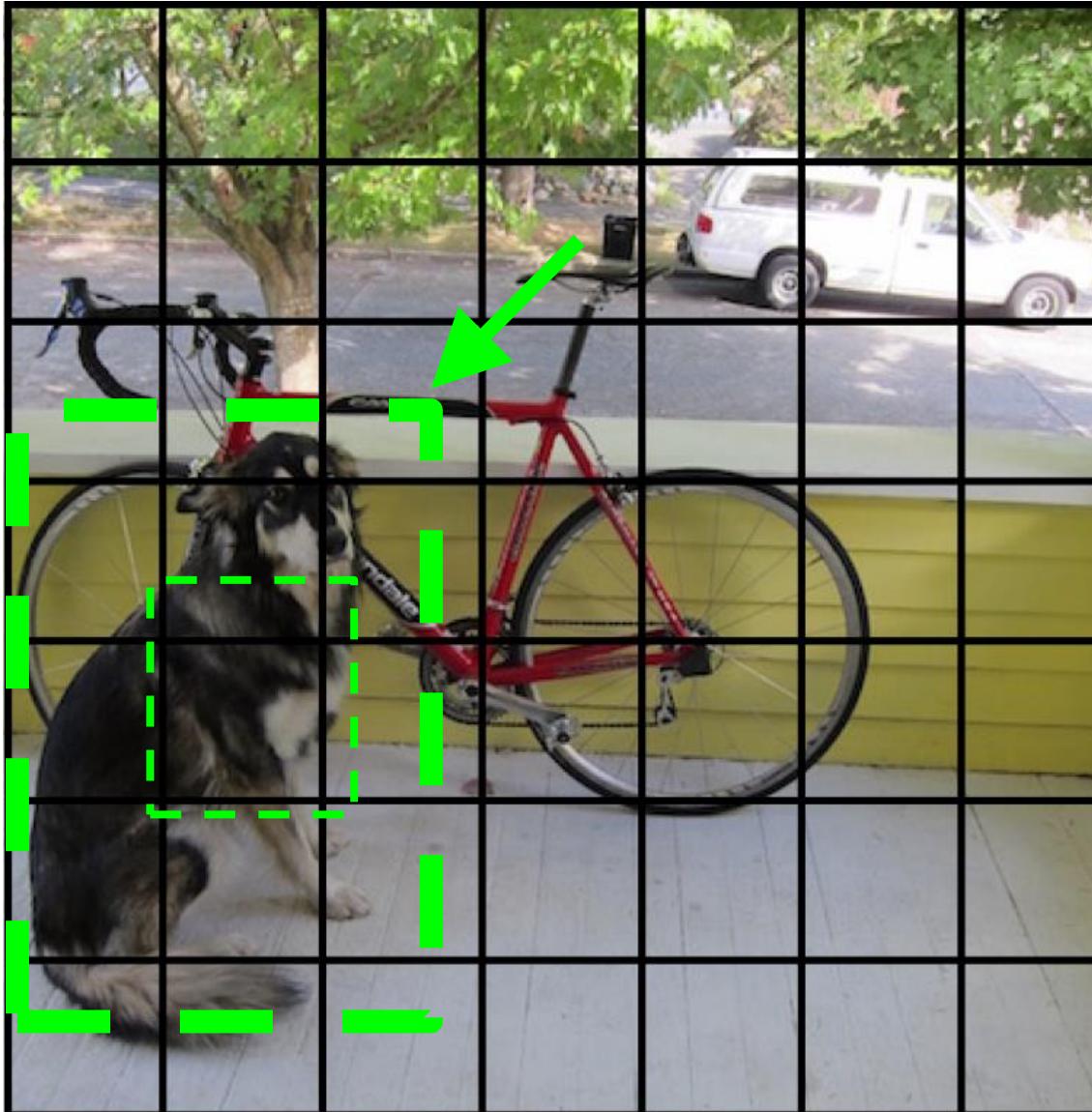
YOLO

Find the best one, adjust it, increase the confidence



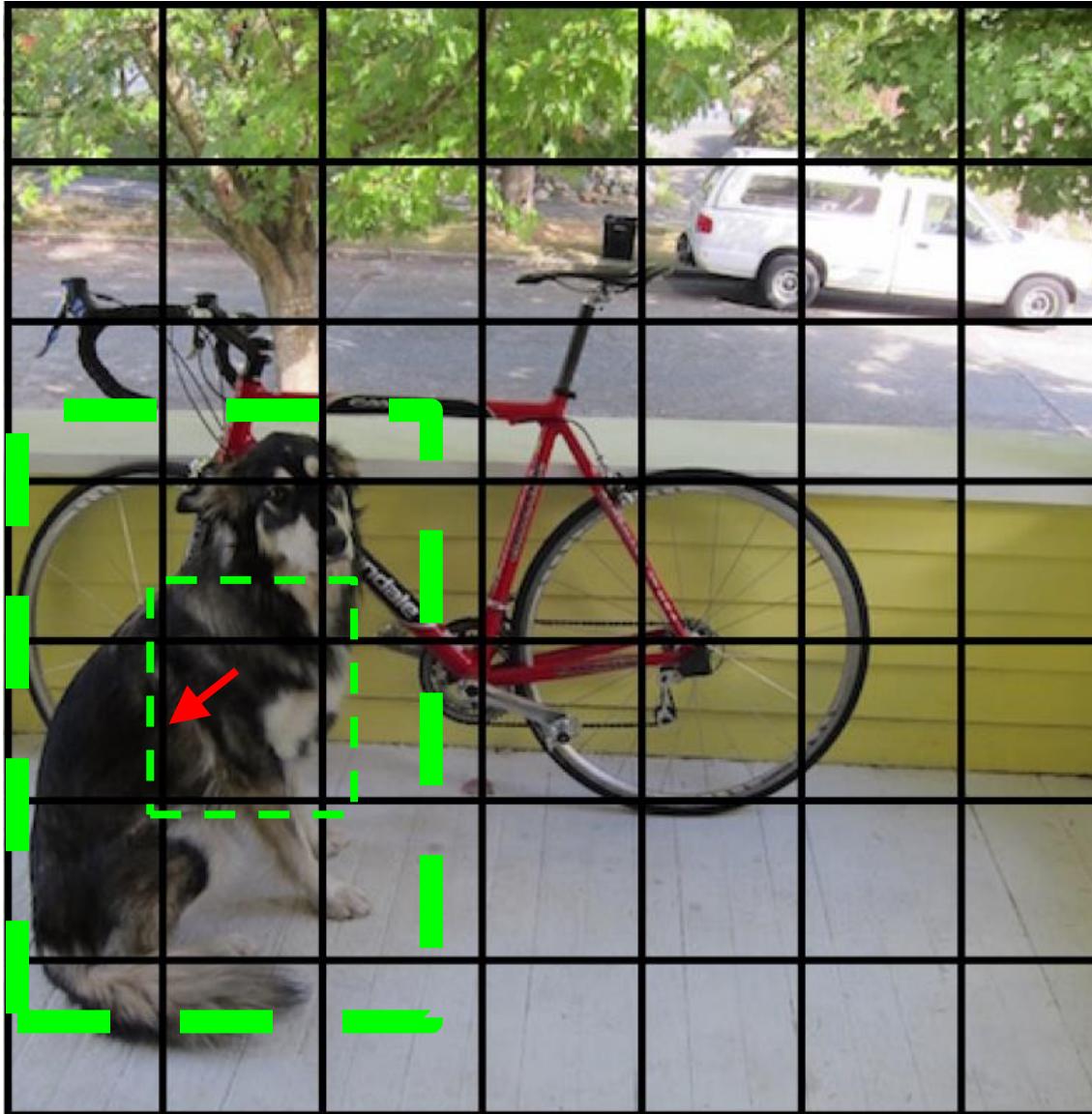
YOLO

Find the best one, adjust it, increase the confidence



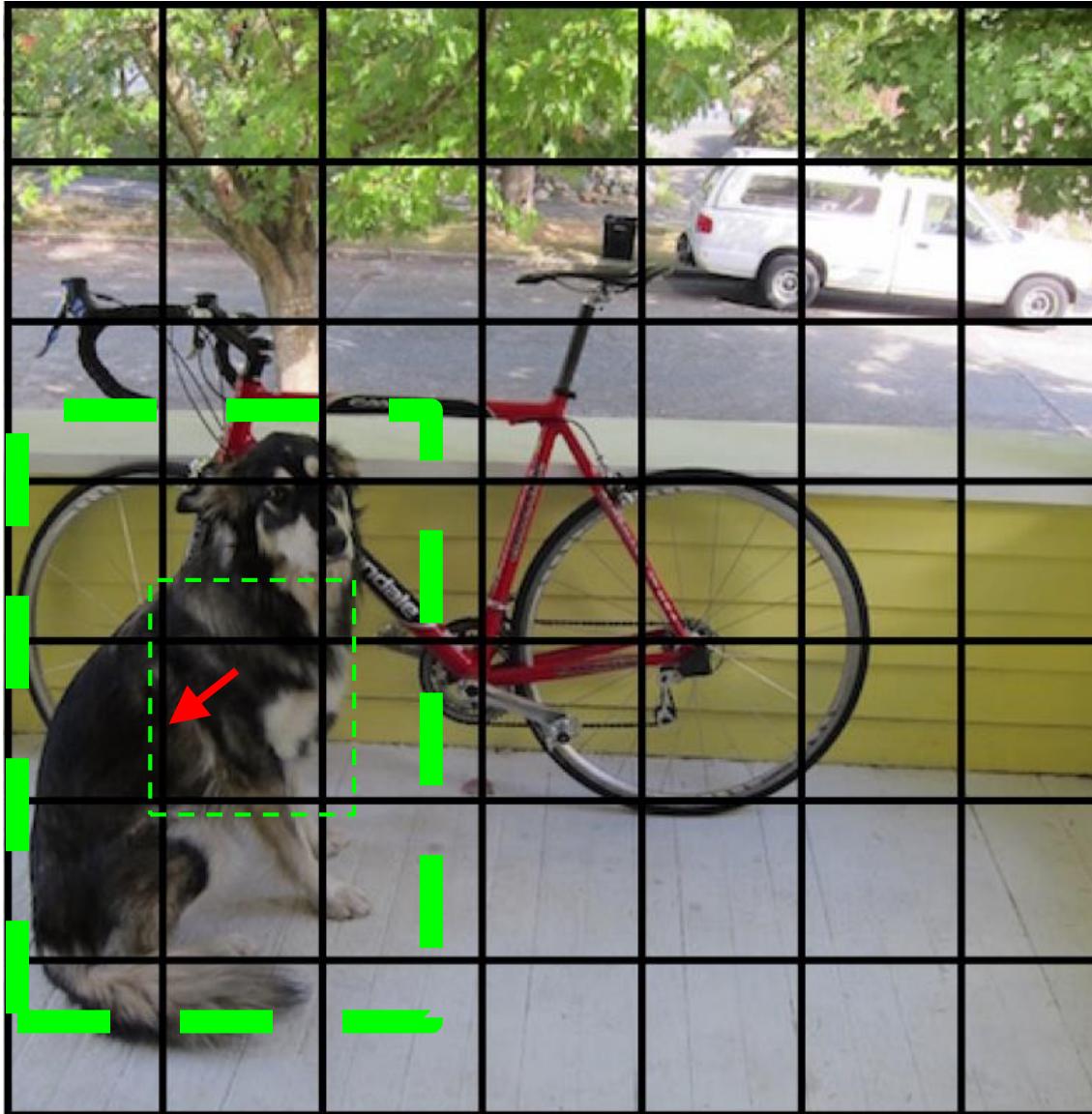
YOLO

Decrease the confidence of other boxes



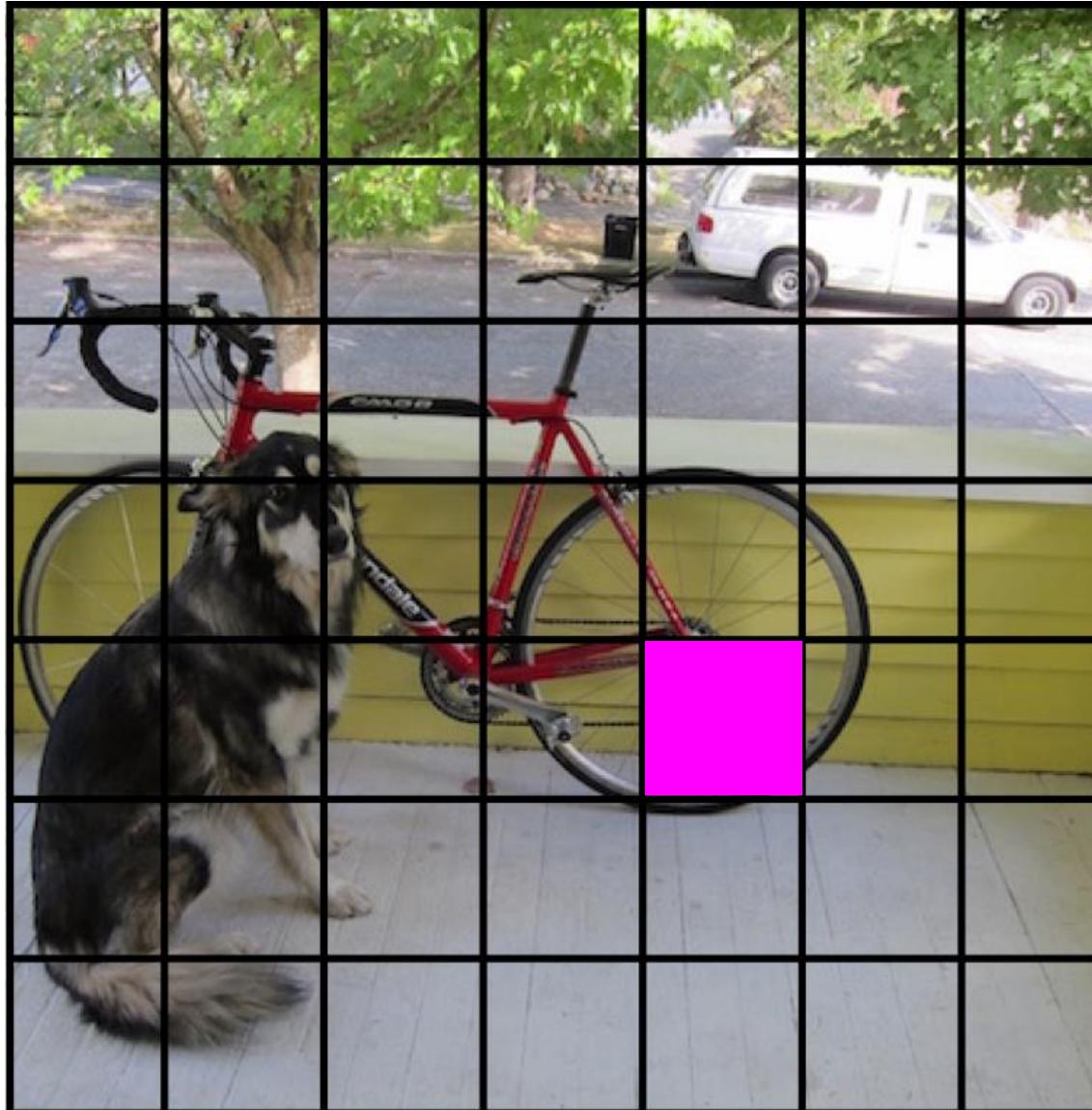
YOLO

Decrease the confidence of other boxes



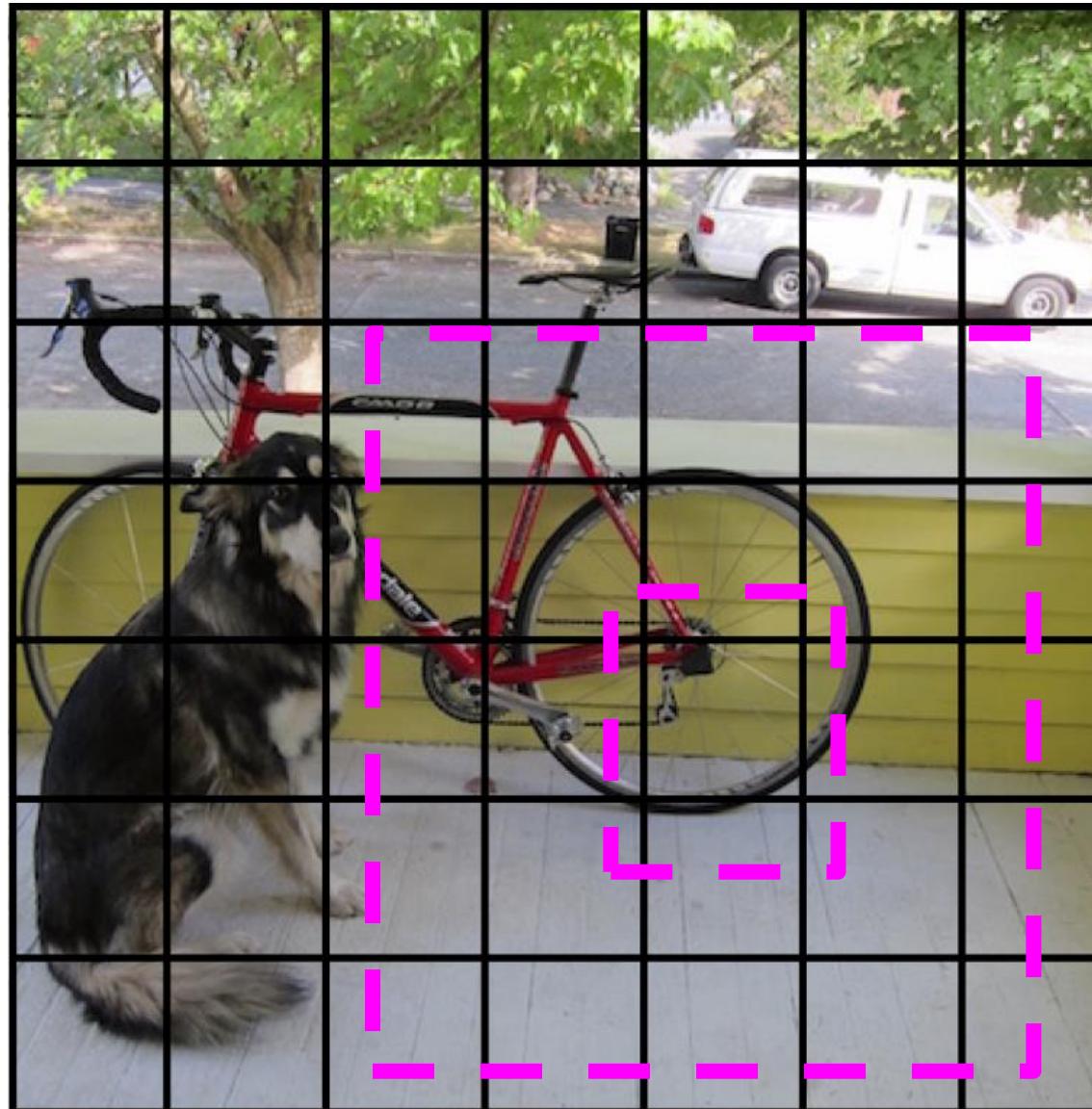
YOLO

Some cells don't have any ground truth detections!



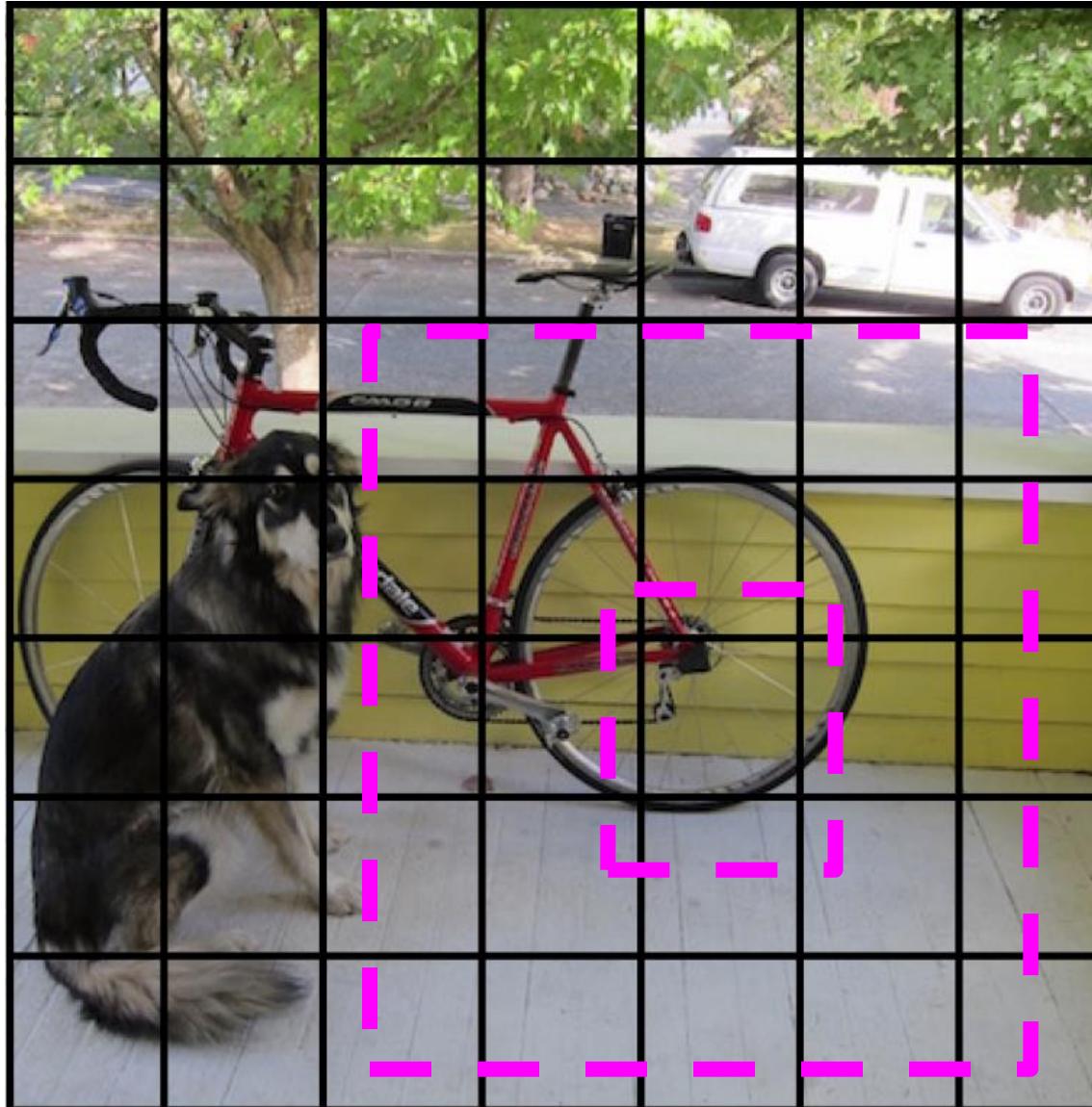
YOLO

Some cells don't have any ground truth detections!



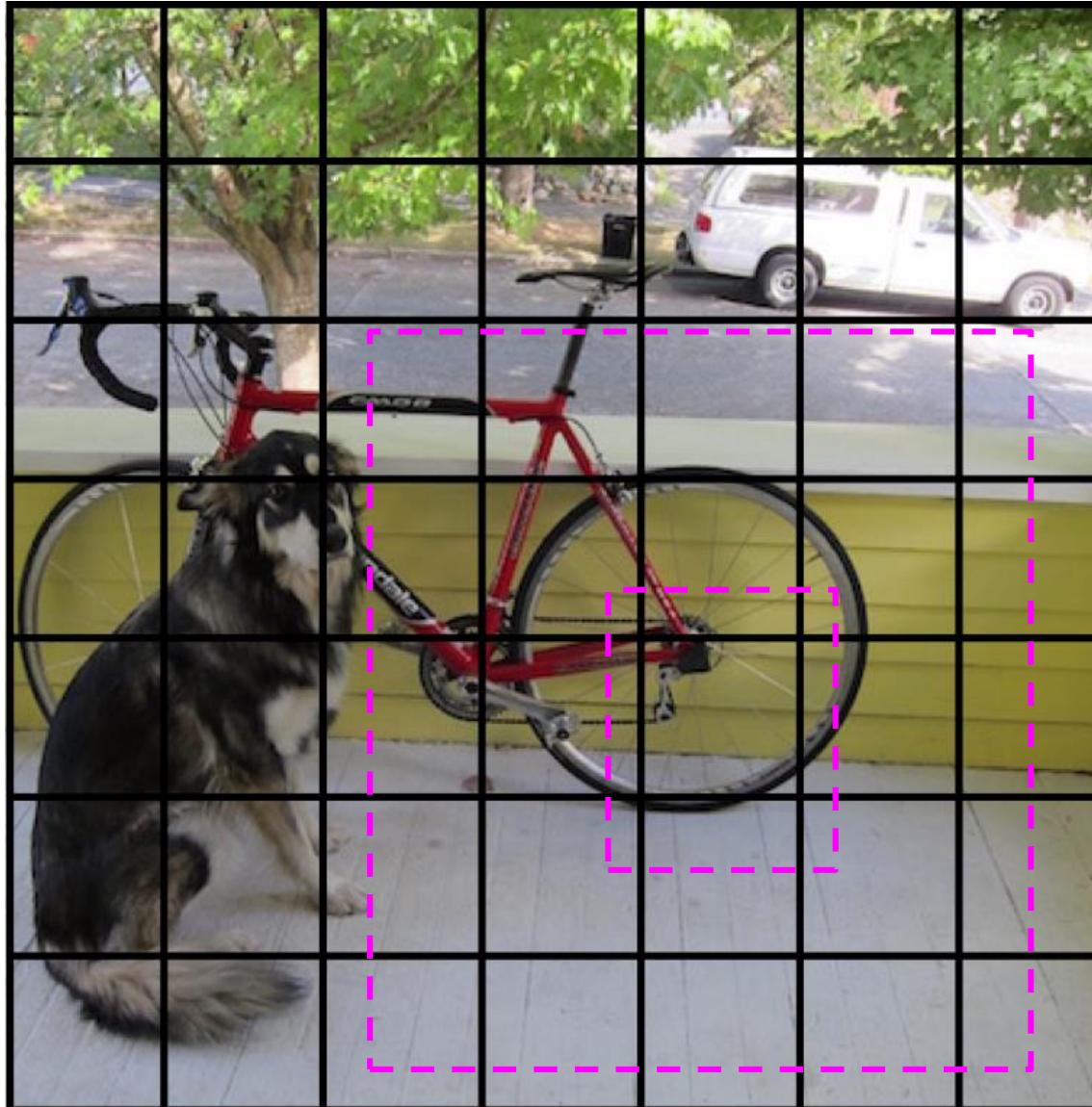
YOLO

Decrease the confidence of these boxes



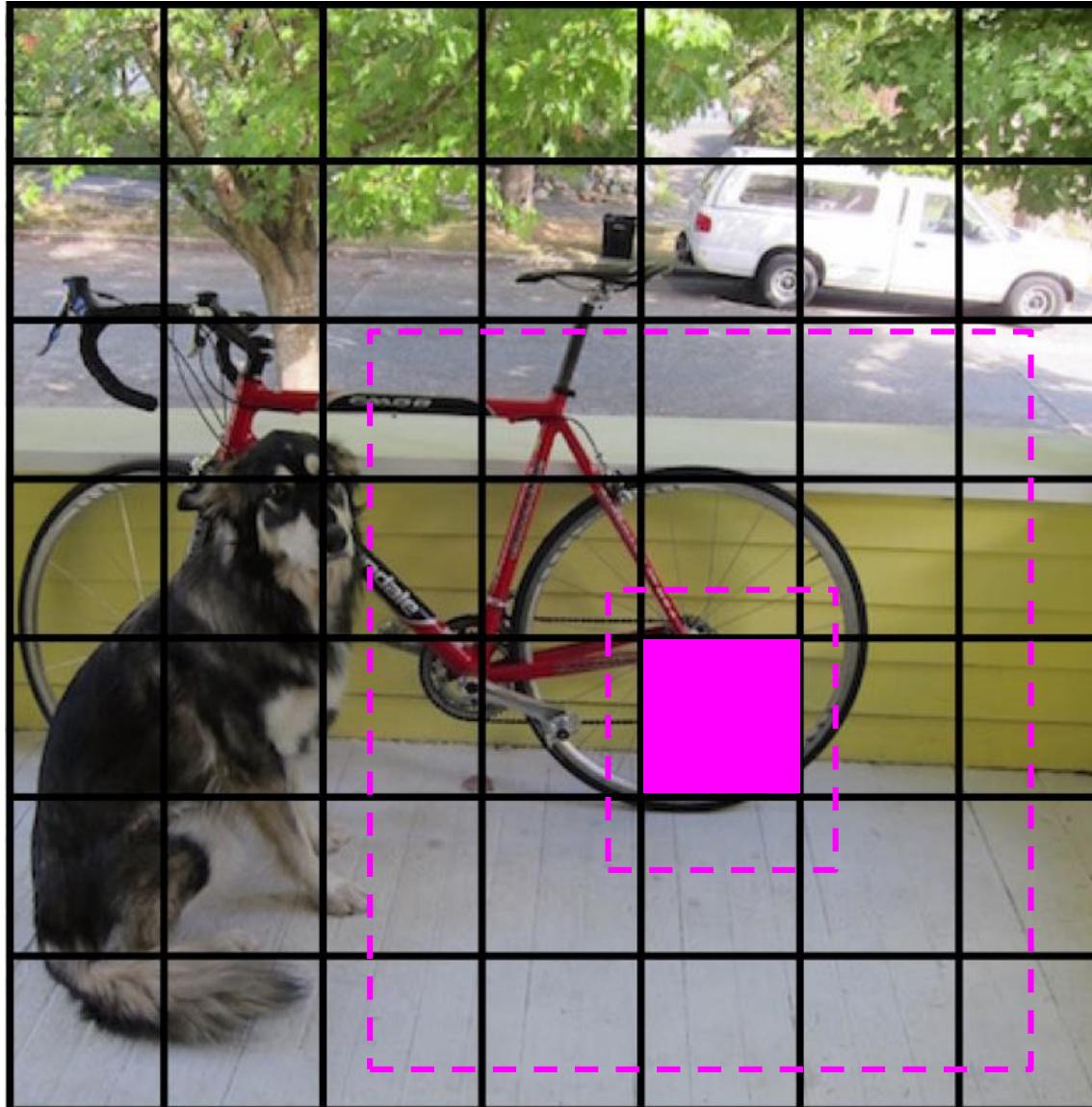
YOLO

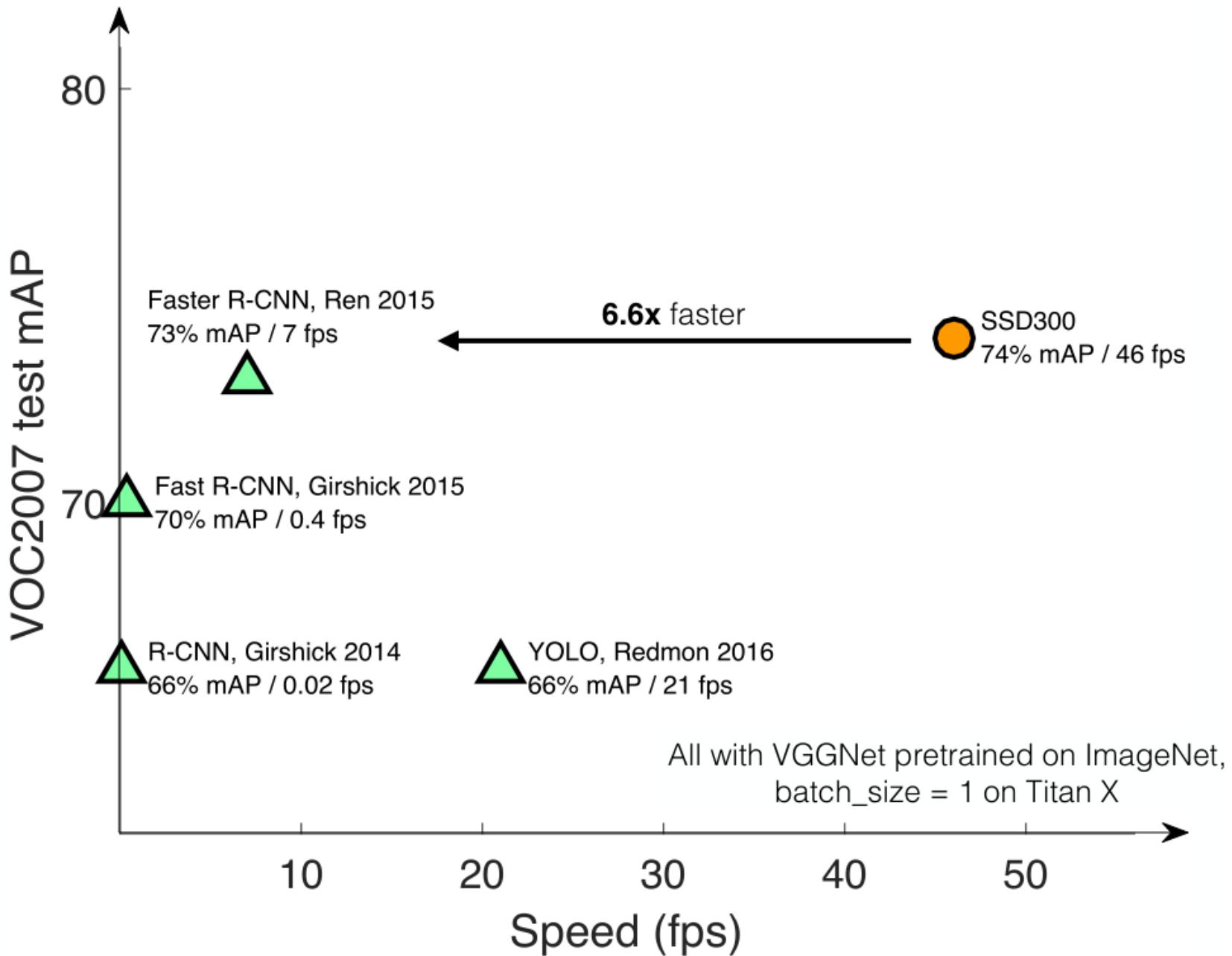
Decrease the confidence of these boxes



YOLO

Don't adjust the class probabilities or coordinates





SSD: Single Shot MultiBox Detector

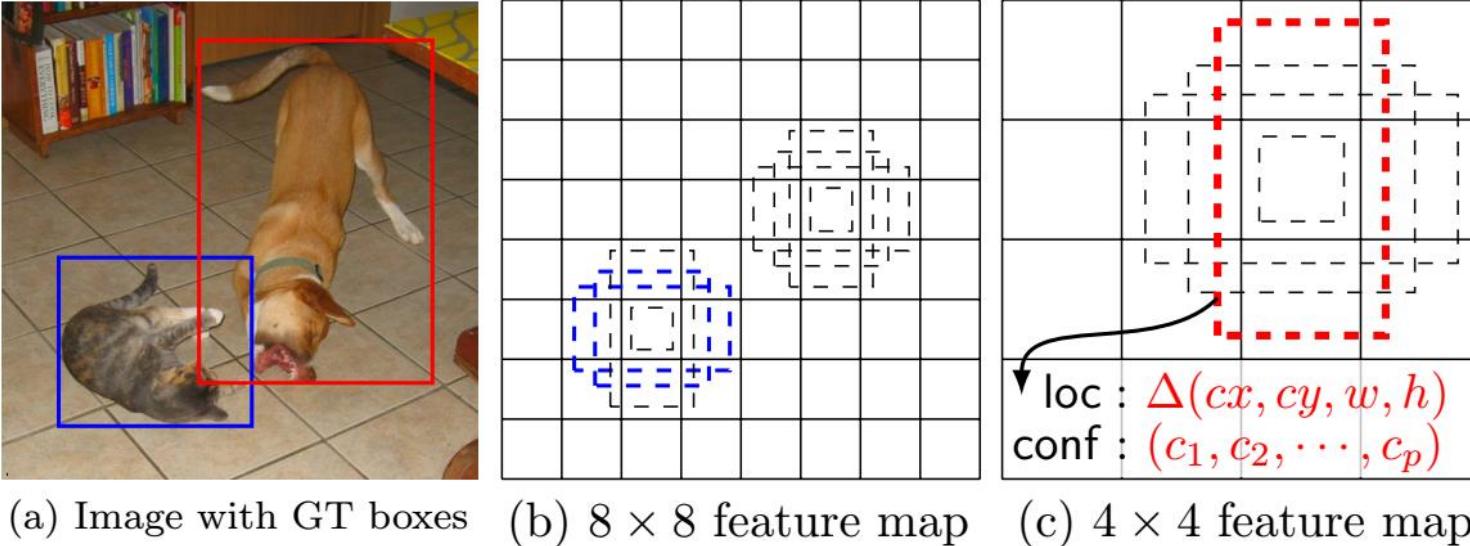


Fig. 1: **SSD framework.** (a) SSD only needs an input image and ground truth boxes for each object during training. In a convolutional fashion, we evaluate a small set (e.g. 4) of default boxes of different aspect ratios at each location in several feature maps with different scales (e.g. 8×8 and 4×4 in (b) and (c)). For each default box, we predict both the shape offsets and the confidences for all object categories $((c_1, c_2, \dots, c_p))$. At training time, we first match these default boxes to the ground truth boxes. For example, we have matched two default boxes with the cat and one with the dog, which are treated as positives and the rest as negatives. The model loss is a weighted sum between localization loss (e.g. Smooth L1 [6]) and confidence loss (e.g. Softmax).

SSD: Single Shot MultiBox Detector

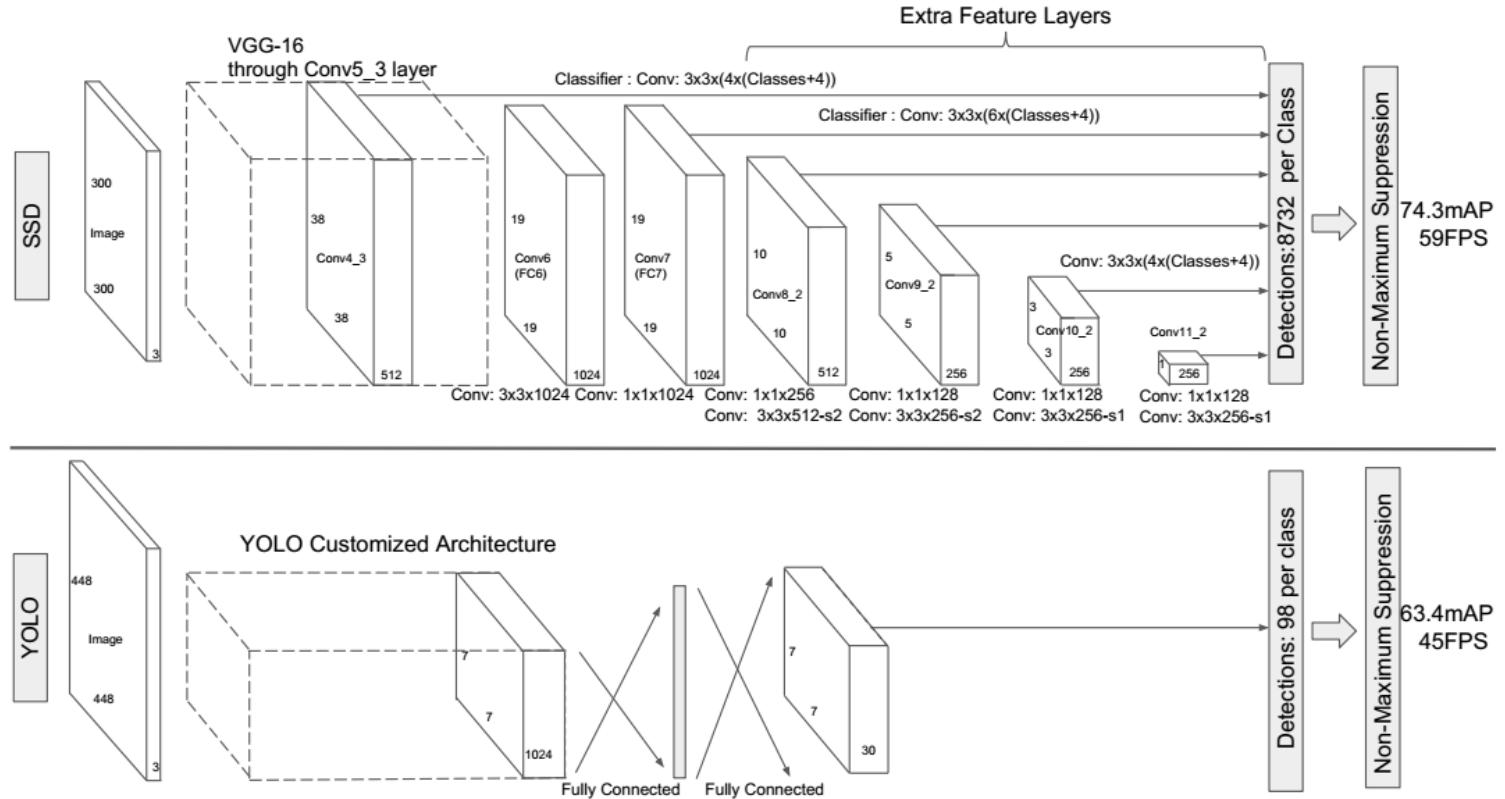
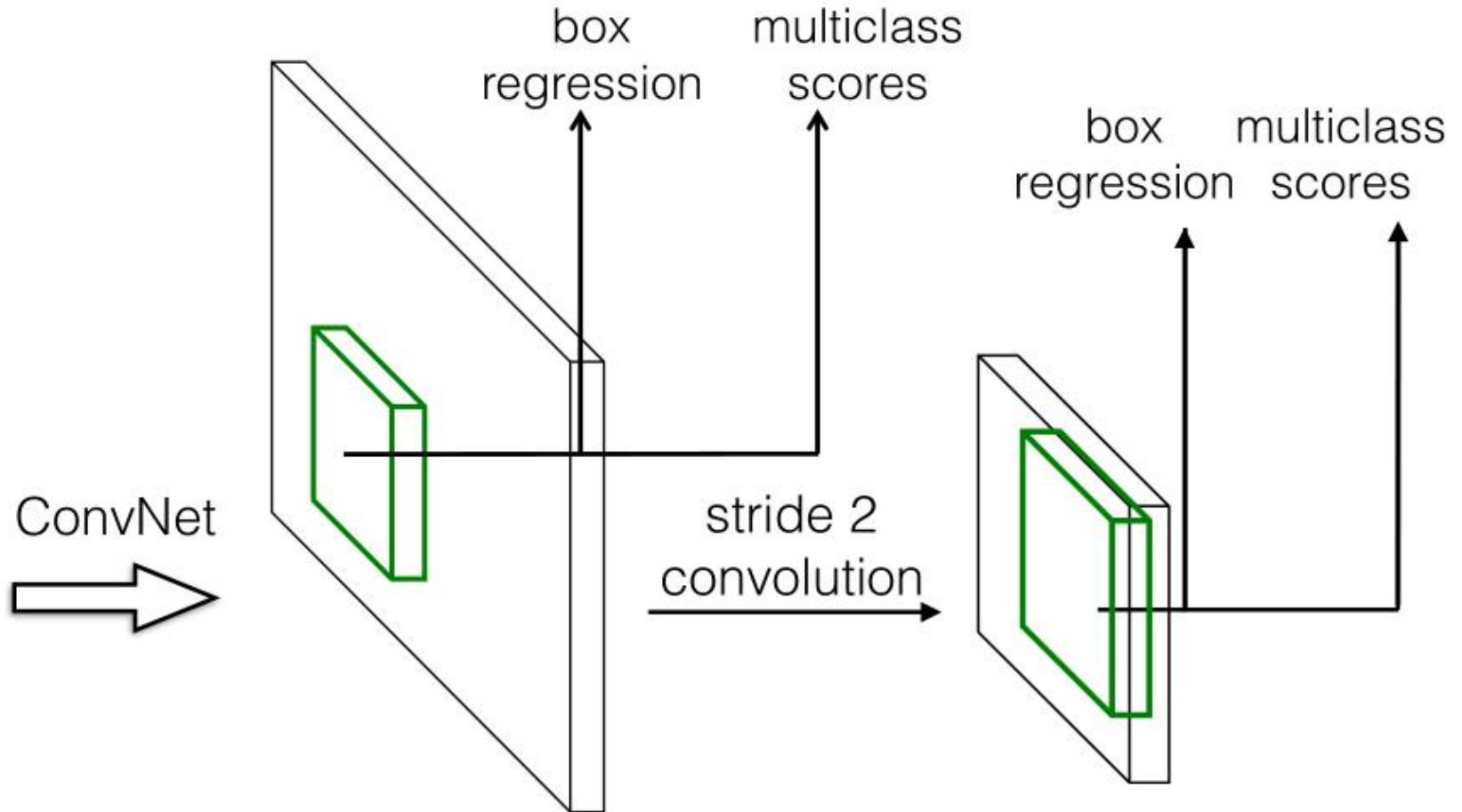


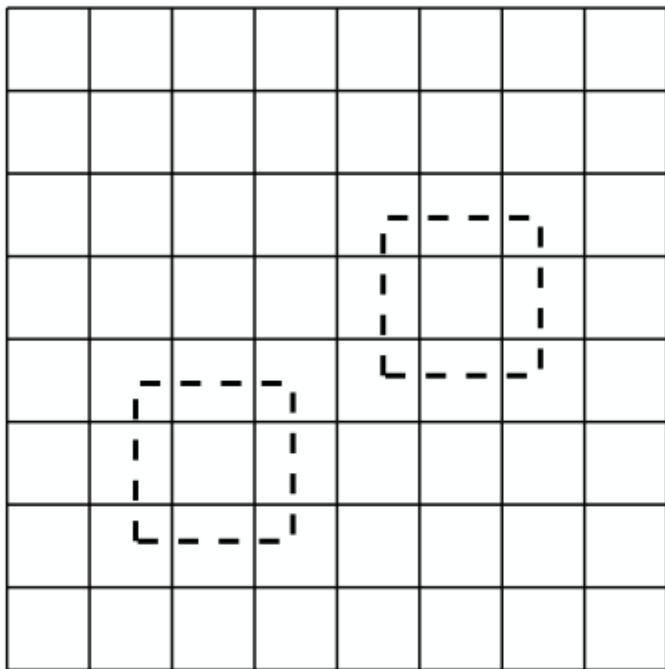
Fig. 2: A comparison between two single shot detection models: SSD and YOLO [5]. Our SSD model adds several feature layers to the end of a base network, which predict the offsets to default boxes of different scales and aspect ratios and their associated confidences. SSD with a 300×300 input size significantly outperforms its 448×448 YOLO counterpart in accuracy on VOC2007 test while also improving the speed.

Key Ideas: Multi Scale Feature Maps

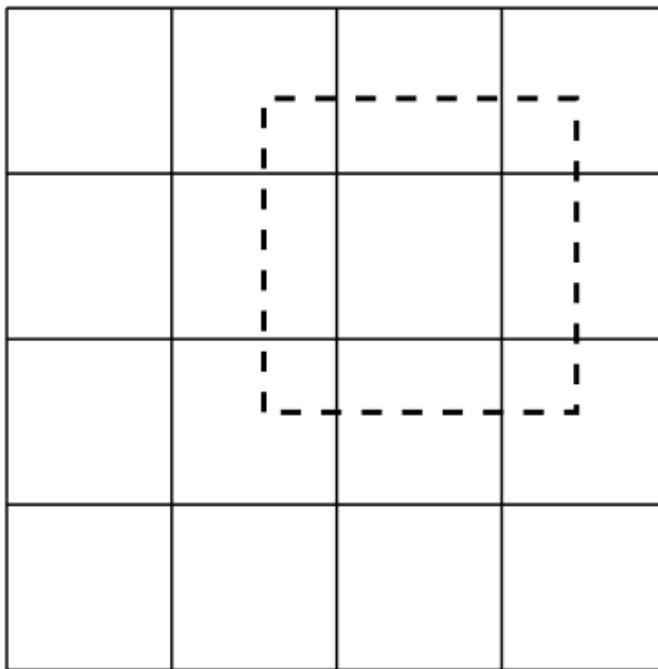


Key Ideas: Multi Scale Feature Maps

SSD



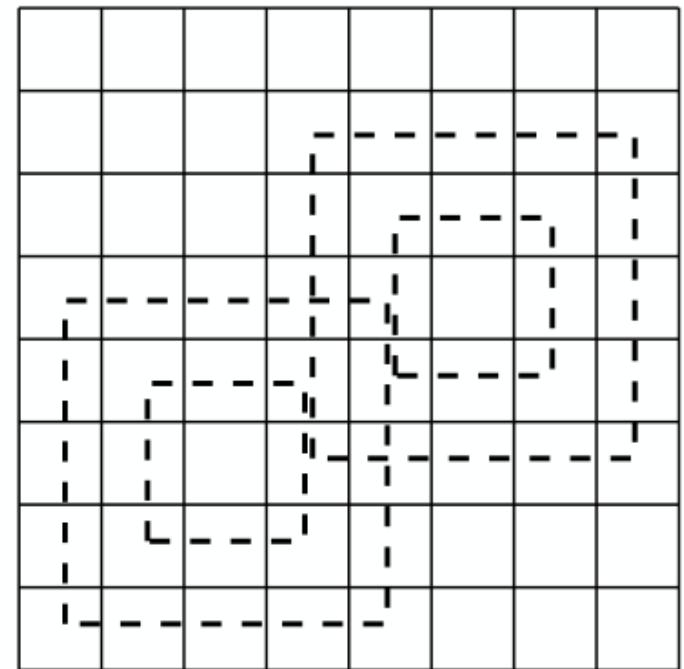
8×8 feature map



4×4 feature map

Faster R-CNN Objectness
Proposal, Ren 2015

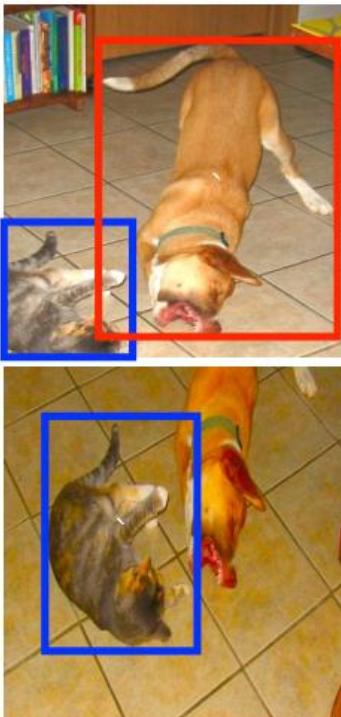
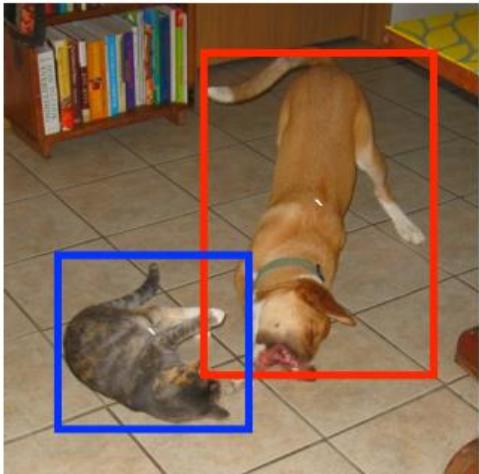
vs.



8×8 feature map

Key Ideas: Devil is in the Details!

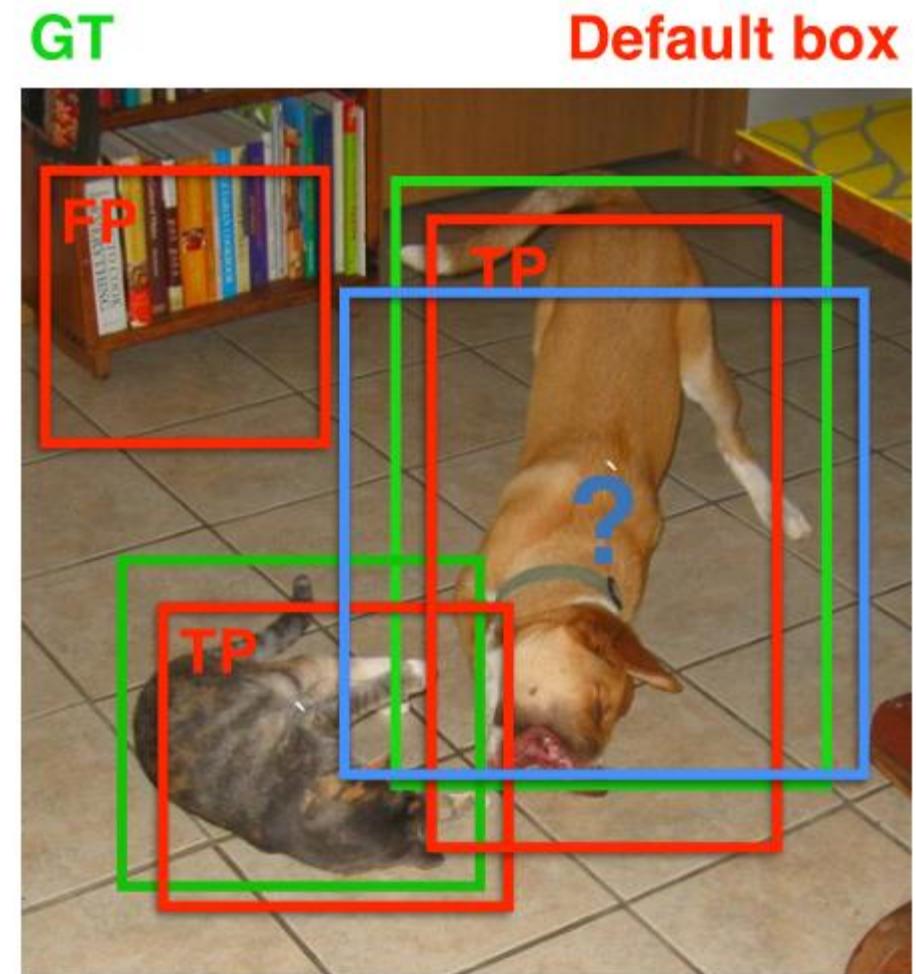
- Data Augmentation



data augmentation	SSD300	
horizontal flip	✓	✓
random crop & color distortion		✓
VOC2007 test mAP	65.5	74.3

Multiple Default Boxes to Ground Truth

- Matching ground truth and default boxes
 - Match each GT box to closest default box
 - Also match each GT box to all unassigned default boxes with $\text{IoU} > 0.5$

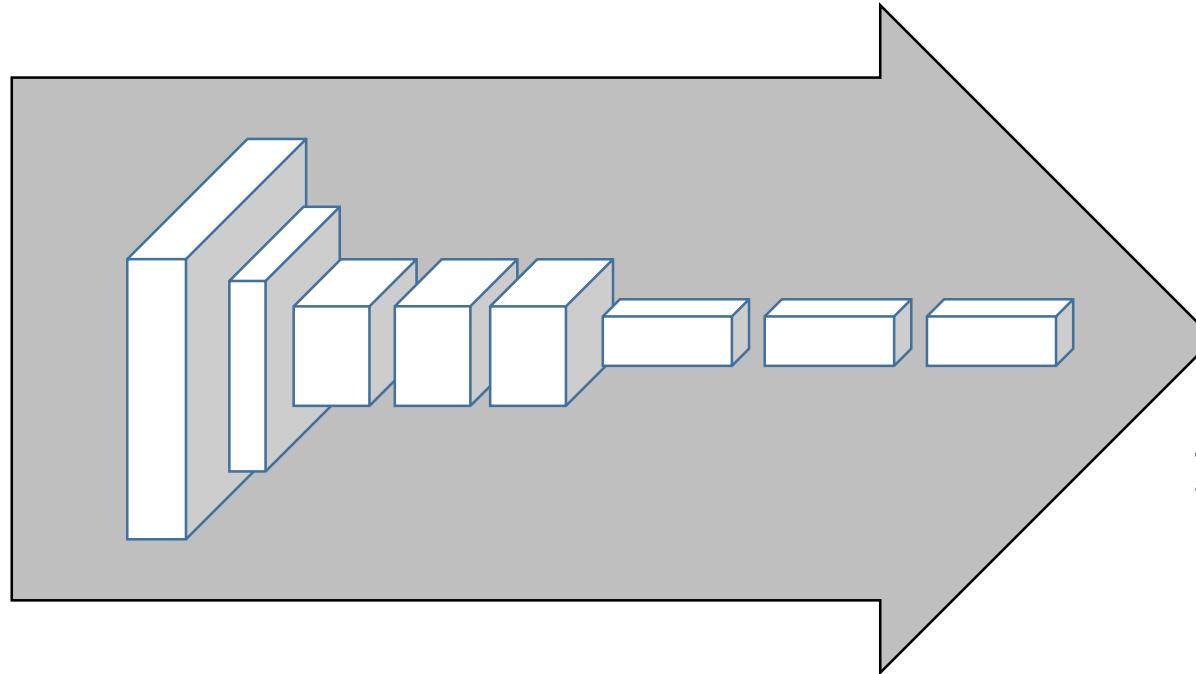


Keys to efficient CNN-based object detection

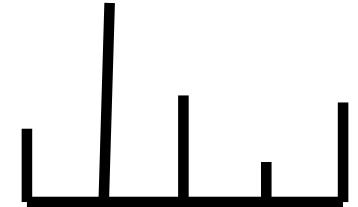
- Feature sharing
 - R-CNN => SPP-net & Fast R-CNN: sharing features among proposal regions
 - Fast R-CNN => Faster R-CNN: sharing features between proposal and detection
 - All are done by shared convolutional feature maps
- Efficient multi-scale solutions
 - Single-scale convolutional feature maps are good trade-offs
 - Multi-scale anchors are fast and flexible.
 - Multi-scale grid cells are even better.
- End to End learning helps
 - RCNN => Fast RCNN => Faster RCNN => YOLO

Semantic Segmentation

Regular CNNs



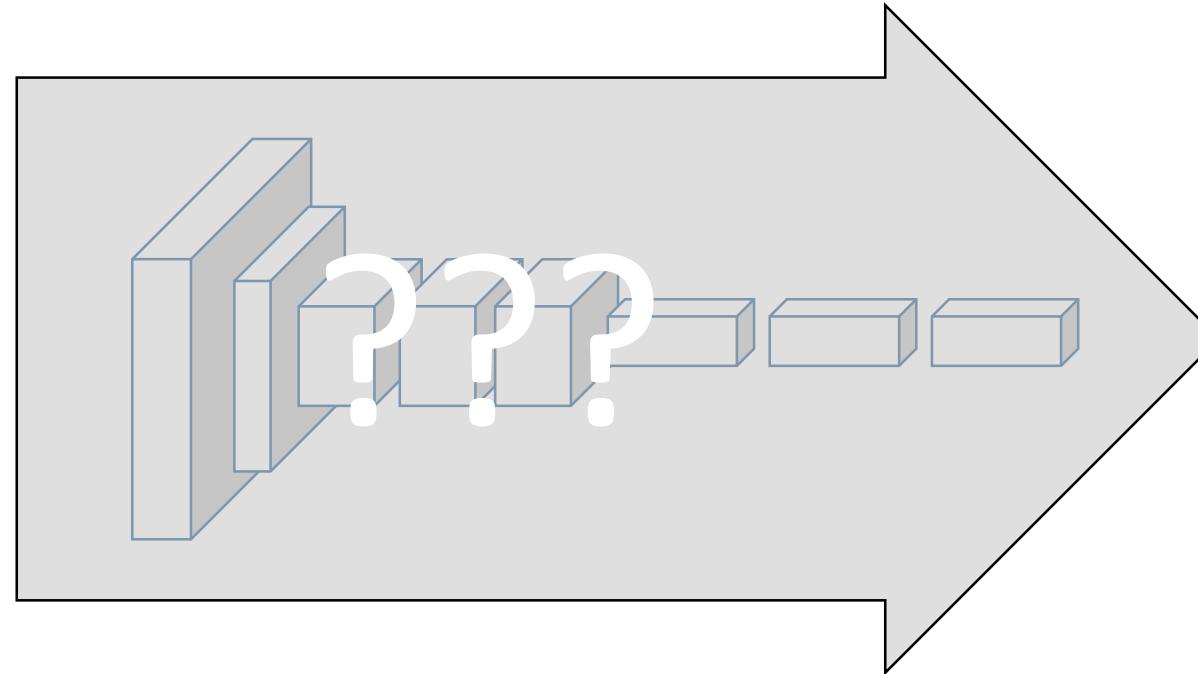
“Tabby Cat”



1000 dimensional
vector

End to end learning

How to do semantic segmentation?



End to end learning

Fully Convolutional Networks for Semantic Segmentation

- Using Convolutional Net for Segmentation (Long, Shelhamer, Darrell, 2015)
 - Determine what objects are where
 - “what” depends on global information
 - “where” depends on local information

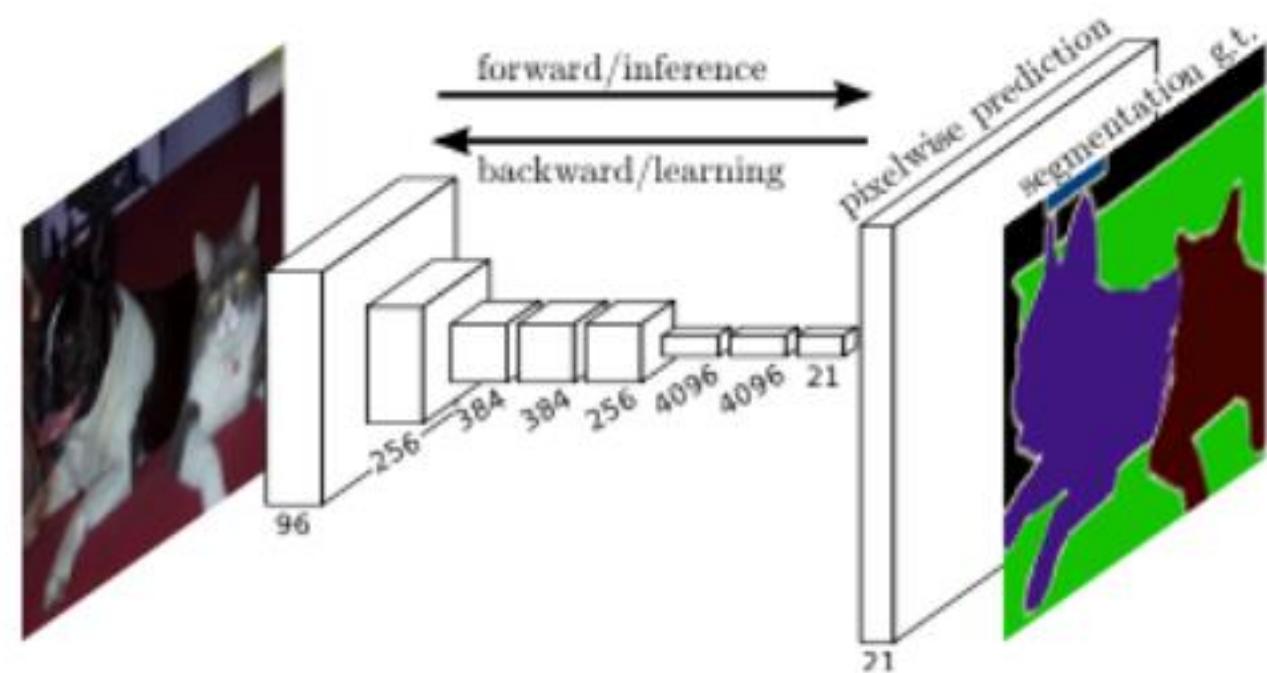


Figure 1. Fully convolutional networks can efficiently learn to make dense predictions for per-pixel tasks like semantic segmentation.

Fully Convolutional Networks for Semantic Segmentation

- Higher layers of typical convnets are nonspatial
- If every layer remains spatial, then output can specify where as well as what (albeit coarsely)

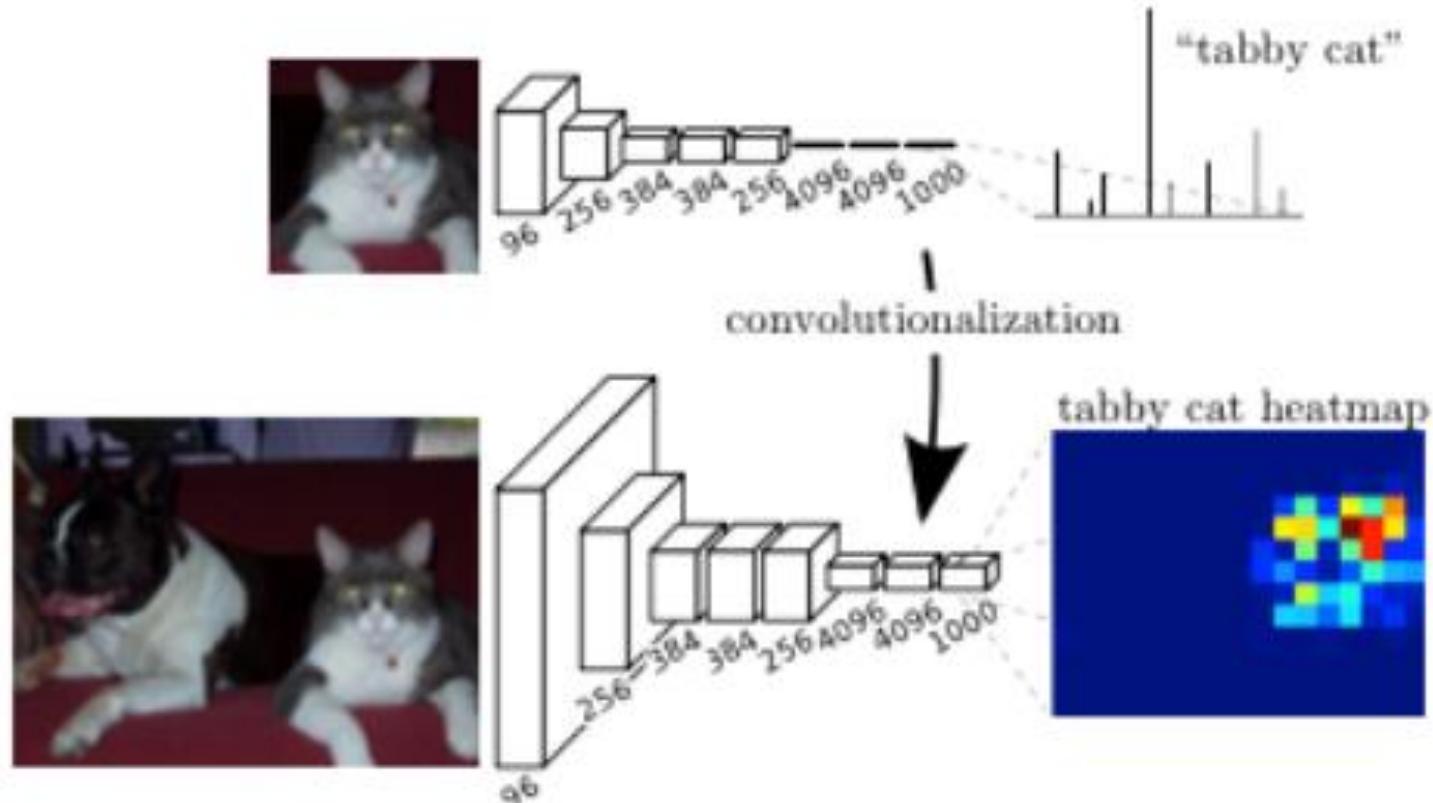
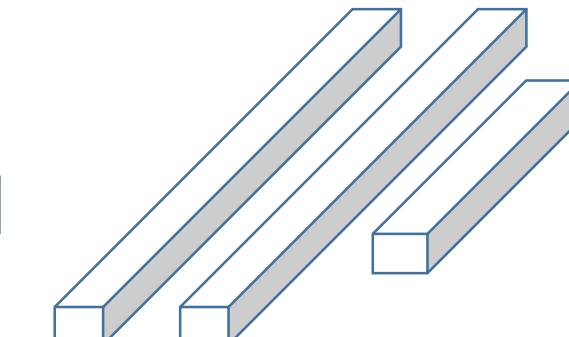
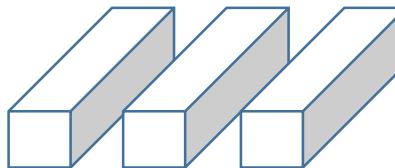
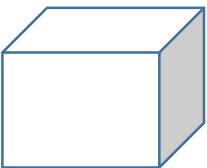


Figure 2. Transforming fully connected layers into convolution layers enables a classification net to output a heatmap. Adding layers and a spatial loss (as in Figure 1) produces an efficient machine for end-to-end dense learning.

Making Fully Convolutional + Upsampling Network



$H \times W$

$$\frac{H}{4} \times \frac{W}{4}$$

$$\frac{H}{8} \times \frac{W}{8}$$

$$\frac{H}{16} \times \frac{W}{16}$$

$$\frac{H}{32} \times \frac{W}{32}$$

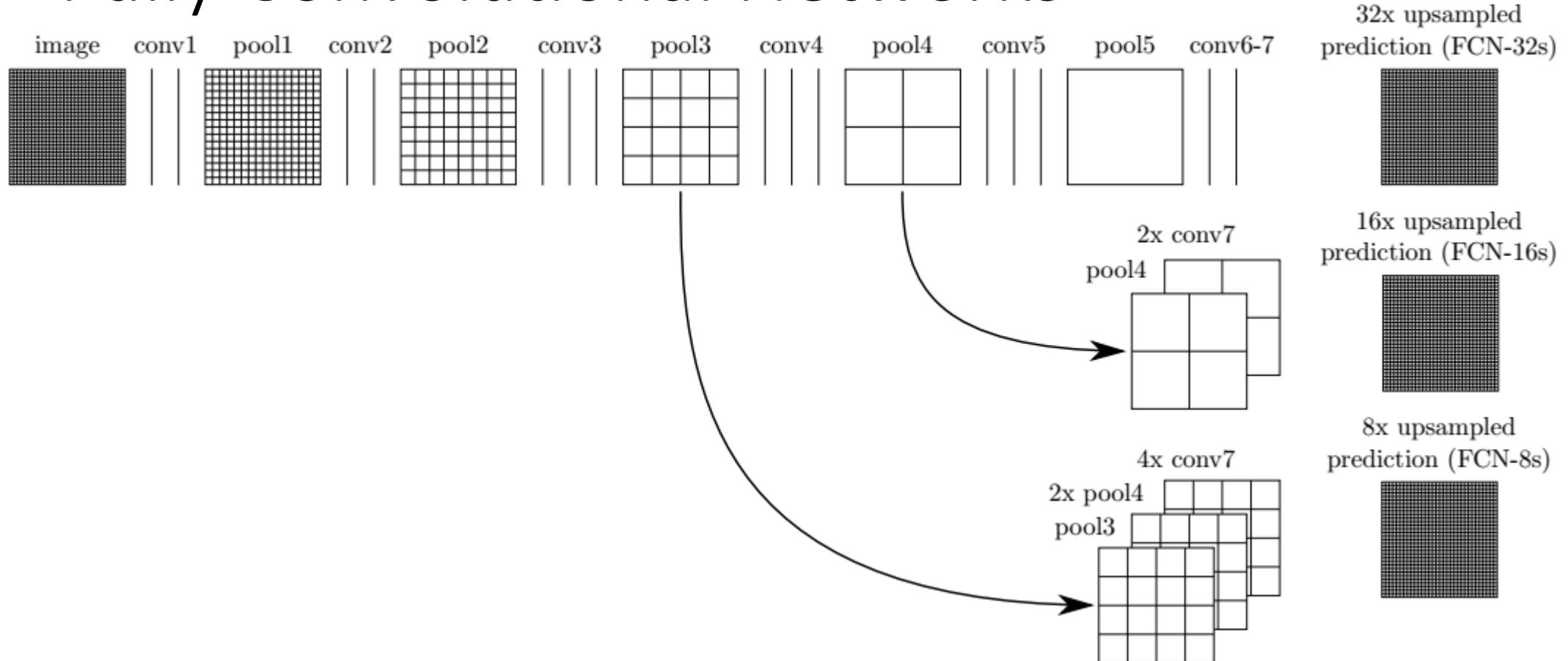


Conv, Pool,

$H \times W$
Pixelwise
Output + Loss

Upsampling

Fully Convolutional Networks



Our DAG nets learn to combine coarse, high layer information with fine, low layer information. Pooling and prediction layers are shown as grids that reveal relative spatial coarseness, while intermediate layers are shown as vertical lines. First row (FCN-32s): Our singlestream net, described in Section 4.1, upsamples stride 32 predictions back to pixels in a single step. Second row (FCN-16s): Combining predictions from both the final layer and the pool4 layer, at stride 16, lets our net predict finer details, while retaining high-level semantic information. Third row (FCN-8s): Additional predictions from pool3, at stride 8, provide further precision.

Upsampling/Unpooling

Nearest Neighbor

1	2
3	4



1	1	2	2
1	1	2	2
3	3	4	4
3	3	4	4

Input: 2 x 2

Output: 4 x 4

“Bed of Nails”

1	2
3	4



1	0	2	0
0	0	0	0
3	0	4	0
0	0	0	0

Input: 2 x 2

Output: 4 x 4

Upsampling: Max Unpooling

Max Pooling

Remember which element was max!

1	2	6	3
3	5	2	1
1	2	2	1
7	3	4	8

Input: 4 x 4

5	6
7	8

Output: 2 x 2

Max Unpooling

Use positions from pooling layer

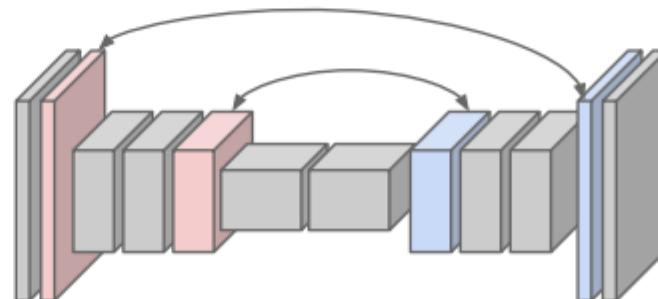
1	2
3	4

Input: 2 x 2

0	0	2	0
0	1	0	0
0	0	0	0
3	0	0	4

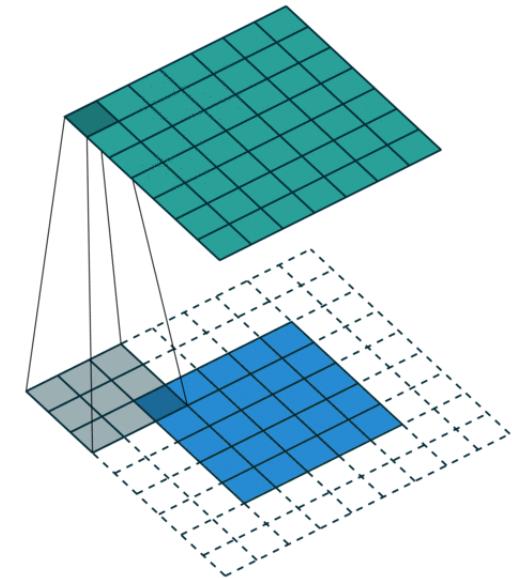
Output: 4 x 4

Corresponding pairs of
downsampling and
upsampling layers



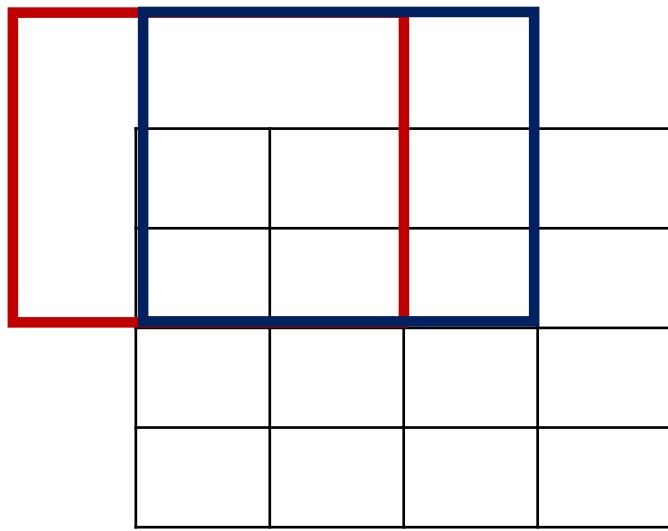
Key Idea: Learnable Upsampling

- Simple bilinear interpolation computes each output y_{ij} from the nearest four inputs by a linear map that depends only on the relative positions of the input and output cells.
- In normal convolution you take each pixel of your input image (or feature map) and calculate the dot product with your let's say 3×3 kernel. The number you get goes into the output pixel, and then you shift to the next pixel



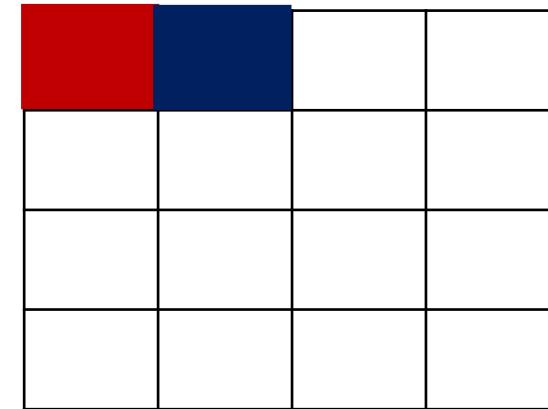
Deconvolution or Transposed Convolution

Typical 3X3 convolution: Stride 1 Pad 1



Input: 4X4

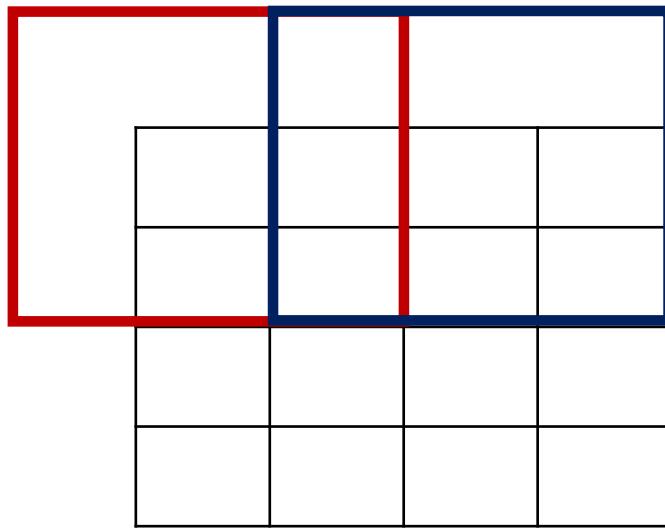
Dot product
between filter
and input



Output: 4X4

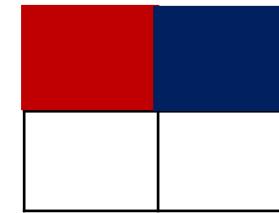
Deconvolution or Transposed Convolution

Typical 3X3 convolution: **Stride 2** Pad 1



Input: 4X4

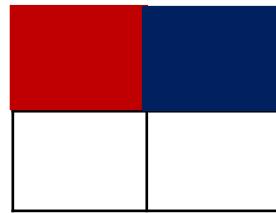
Dot product
between filter
and input



Output: 2X2

Deconvolution or Transposed Convolution

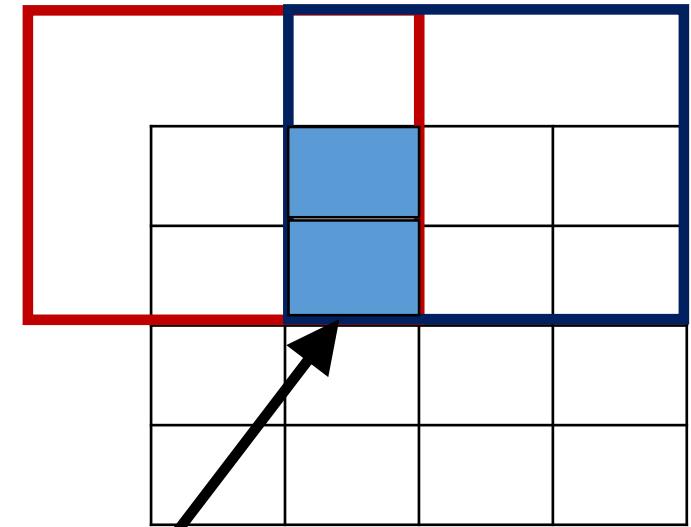
3X3 **deconvolution**: Stride 2 Pad 1



Input: 2X2



Input gives
weight for filter

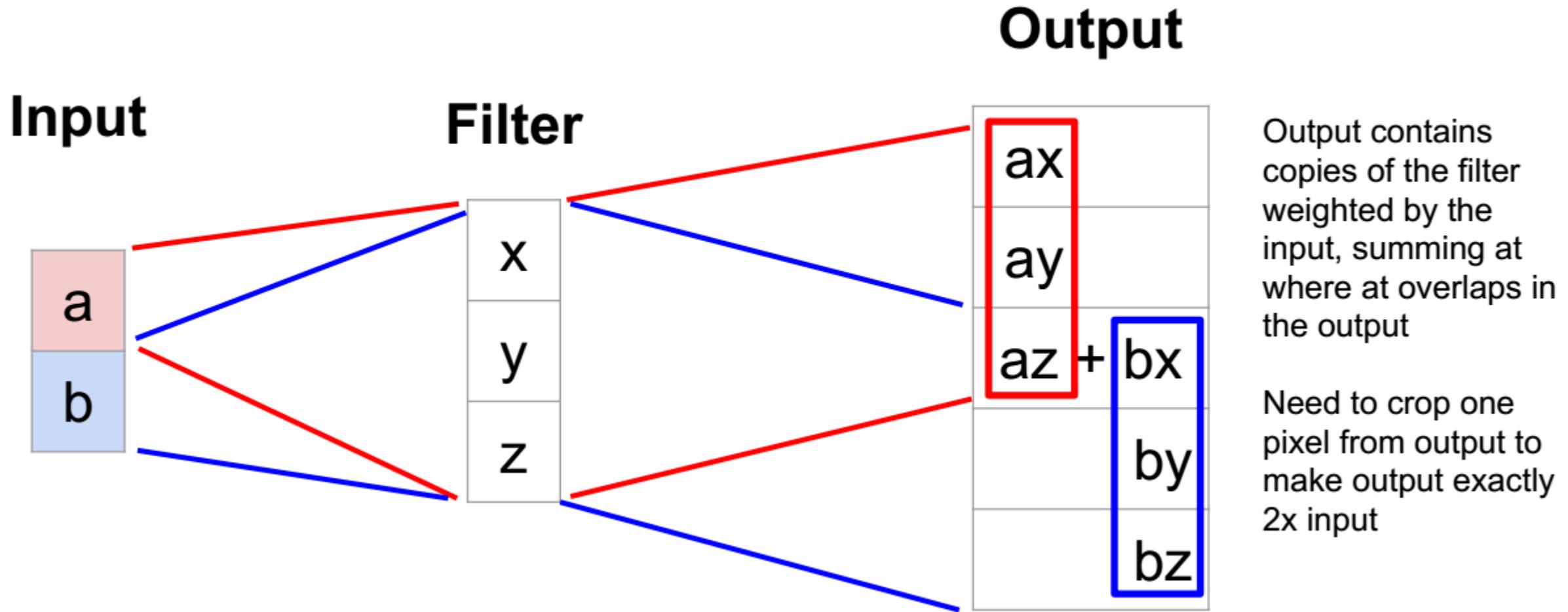


Output: 4X4

Same as backward pass for normal convolution!

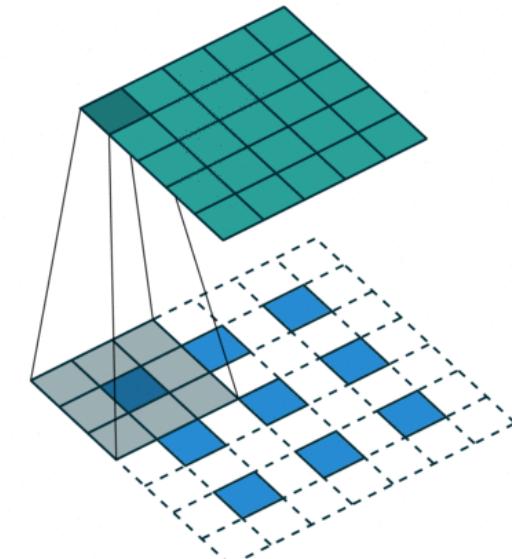
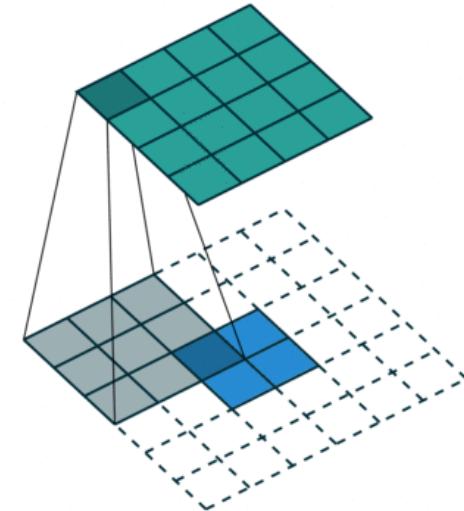
Sum where
output overlaps

Deconvolution: 1D Example

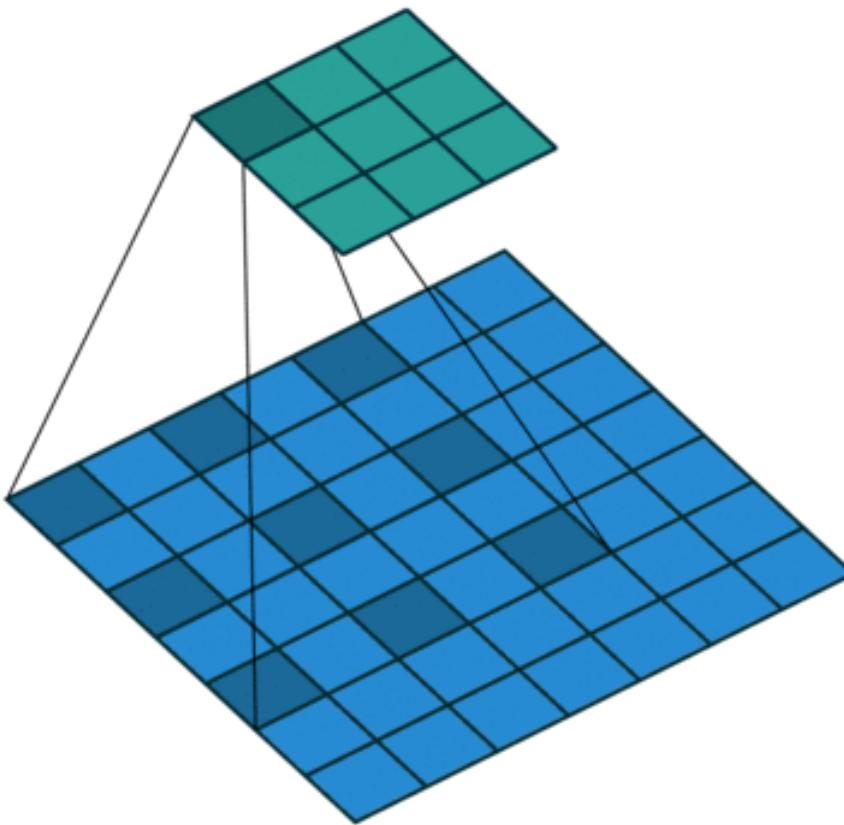


Transposed Convolution: Fractional Stride

- Normal Stride 1: You take each pixel of your input image, you multiply each pixel of your say 3×3 kernel with the input pixel to get a 3×3 weighted kernel, then you insert this in your output image. Where the outputs overlap you sum them.
- You can also use a stride larger than one to increase the upsampling effect (or some would call this a convolution with a fractional stride $1/f$)



Dilated Convolution



Deconvolution or Transposed Convolution

- “Deconvolution” is a bad name, already defined as “inverse of convolutions
- Better names:
 - Convolution Transpose
 - Backward Strided Convolution
 - $\frac{1}{2}$ Strided Convolution
 - UpConvolution

Deconvolution or Transposed Convolution

- It is more proper to say “convolutional transpose operation” rather than “deconvolutional” operation. Hence, we will be using “convolutional transpose” from now
 - Im et al., Generating images with recurrent adversarial networks. arXiv 2016.
- A series of four fractionally-strided convolutions (in some recent papers, these are wrongly called deconvolutions)
 - Radford et al., Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks. ICLR 2016

Key Idea: Skip Connection

- The 32 pixel stride at the final prediction layer limits the scale of detail in the upsampled output
- This is resolved by adding skips that combine the final prediction layer with lower layers with finer strides. This turns a line topology into a DAG, with edges that skip ahead from lower layers to higher ones
- Combining fine layers and coarse layers lets the model make local predictions that respect global structure

Results

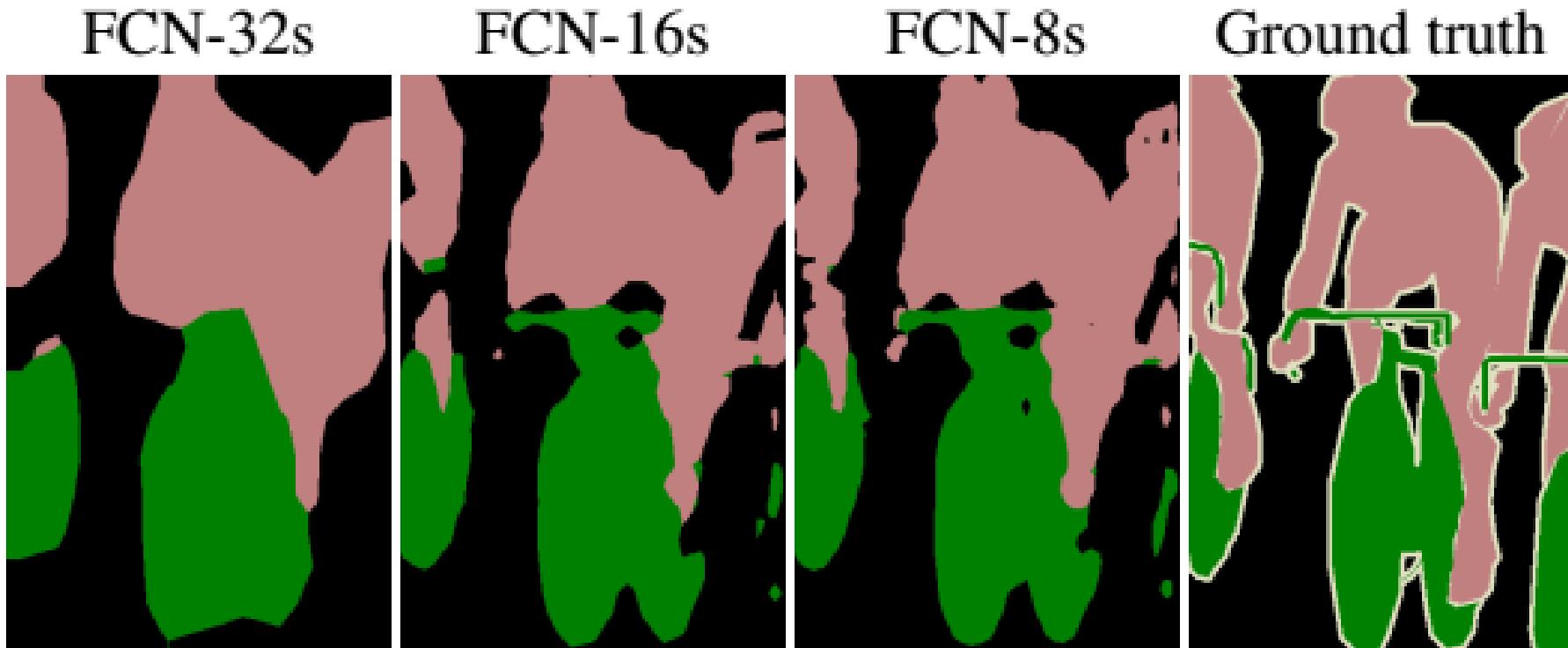


Figure 4. Refining fully convolutional nets by fusing information from layers with different strides improves segmentation detail. The first three images show the output from our 32, 16, and 8 pixel stride nets (see Figure 3).

Results on Pascal VOC 2011

Frequency
weighted IU

	pixel acc.	mean acc.	mean IU	f.w. IU
FCN-32s-fixed	83.0	59.7	45.4	72.0
FCN-32s	89.1	73.3	59.4	81.4
FCN-16s	90.0	75.7	62.4	83.0
FCN-8s	90.3	75.9	62.7	83.2

Upsampling: Other Works

