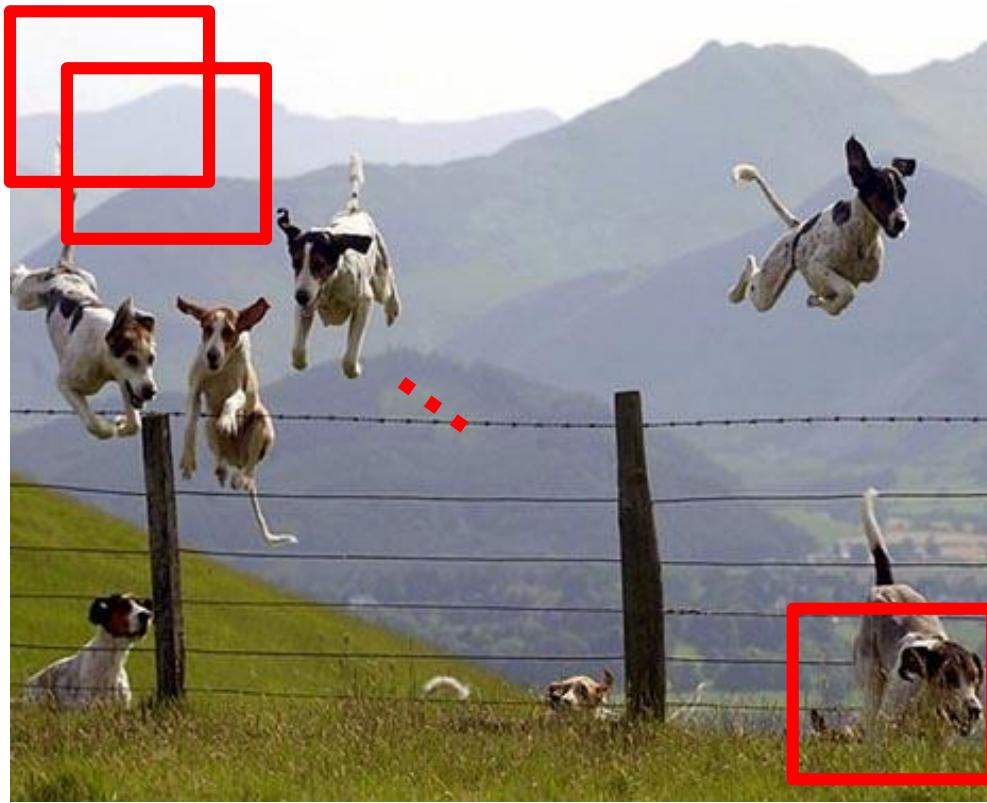


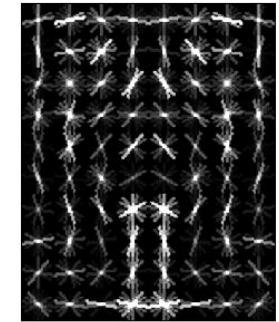
# Object Detection

# Object Category Detection

- Focus on object search: “Where is it?”
- Build templates that quickly differentiate object patch from background patch



Dog Model



**Object or  
Non-Object?**

# Challenges in modeling the object class



Illumination



Object pose



Clutter



Occlusions



Intra-class  
appearance



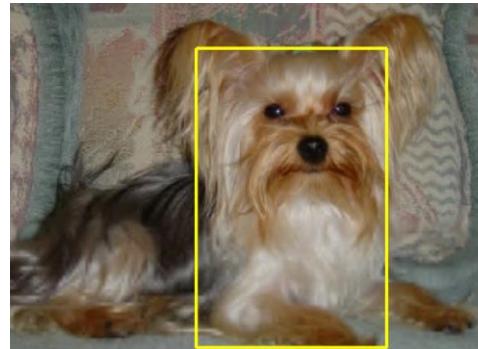
Viewpoint

# Challenges in modeling the non-object class

True  
Detections



Bad  
Localization



Confused with  
Similar Object



Misc. Background



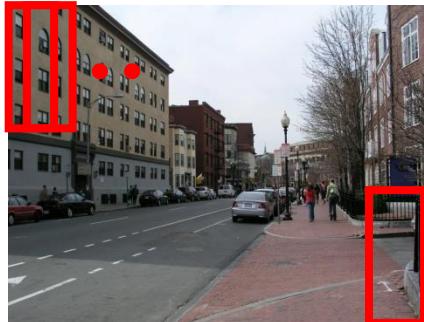
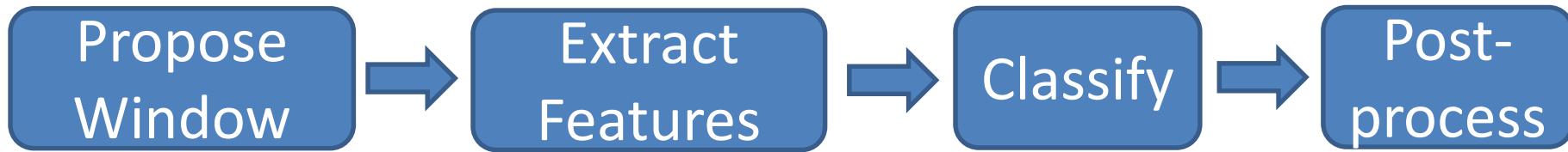
Confused with  
Dissimilar Objects



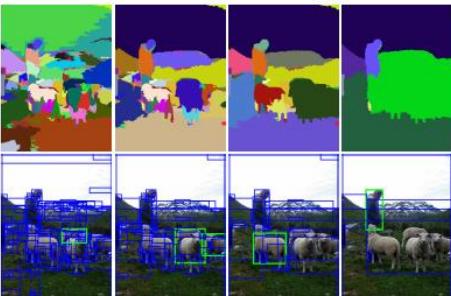
# Influential Works in Detection

- Sung-Poggio (1994, 1998) : ~2412 citations
  - Basic idea of statistical template detection (I think), bootstrapping to get “face-like” negative examples, multiple whole-face prototypes (in 1994)
- Rowley-Baluja-Kanade (1996-1998) : ~4953
  - “Parts” at fixed position, non-maxima suppression, simple cascade, rotation, pretty good accuracy, fast
- Schneiderman-Kanade (1998-2000,2004) : ~2600
  - Careful feature/classifier engineering, excellent results, cascade
- Viola-Jones (2001, 2004) : ~27,000
  - Haar-like features, Adaboost as feature selection, hyper-cascade, very fast, easy to implement
- Dalal-Triggs (2005) : ~18000
  - Careful feature engineering, excellent results, HOG feature, online code
- Felzenszwalb-Huttenlocher (2000): ~2100
  - Efficient way to solve part-based detectors
- Felzenszwalb-McAllester-Ramanan (2008,2010): ~7200
  - Excellent template/parts-based blend
- Girshick-Donahue-Darrell-Malik (2014) : ~4700
  - Region proposals + fine-tuned CNN features (marks significant advance in accuracy over hog-based methods)
- Redmon, Divvala, Girshick, Farhadi (2016): ~210
  - Refine and simplify RCNN++ approach to predict directly from last conv layer

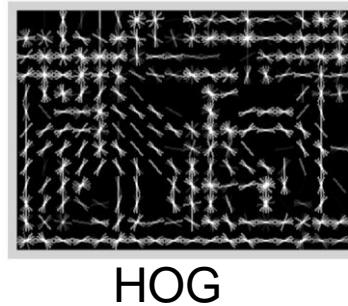
# Summary: statistical templates



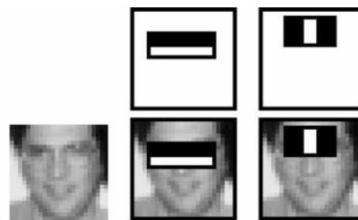
Sliding window: scan image pyramid



Region proposals:  
edge/region-based,  
resize to fixed window



HOG



Fast randomized features

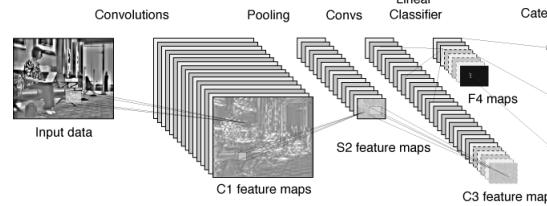
SVM

Boosted stubs

Neural network

Non-max suppression

Segment or refine localization



CNN features

# Generating hypotheses

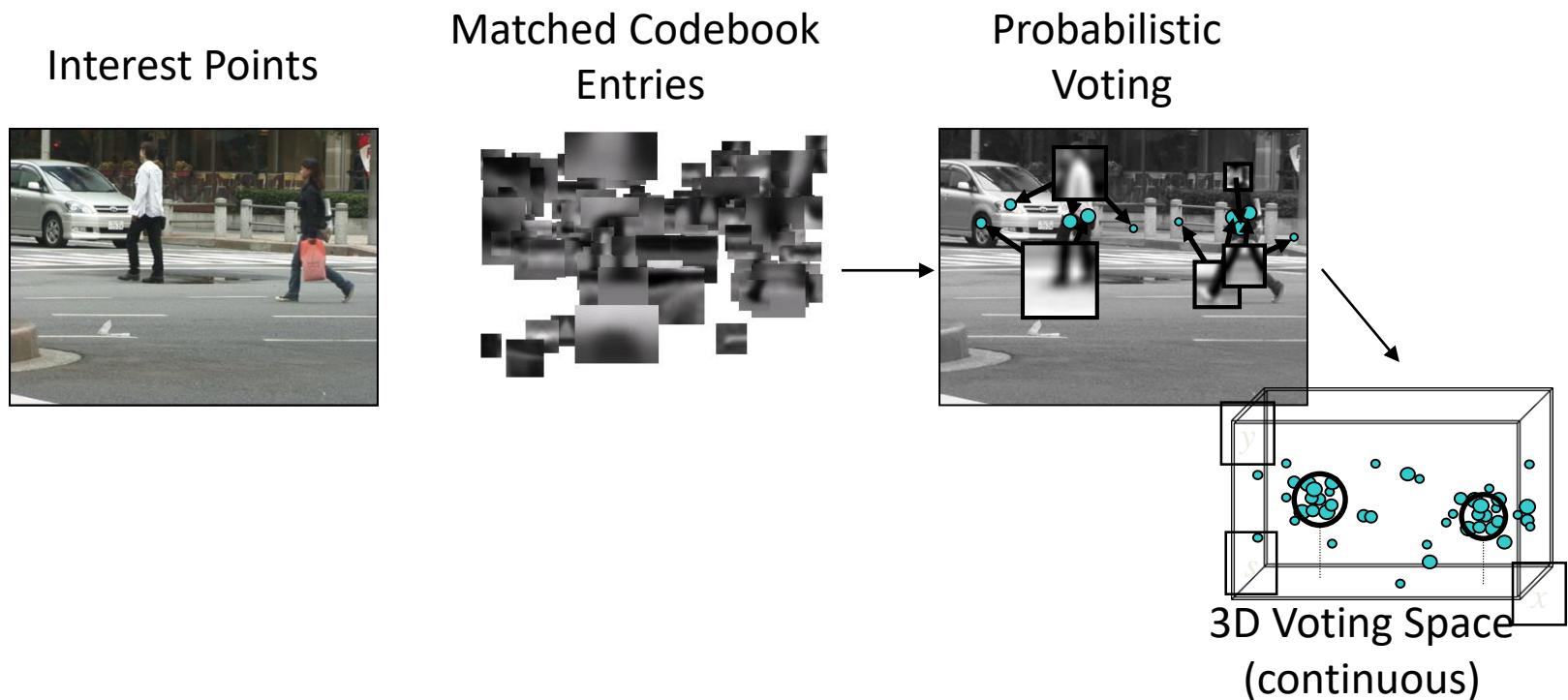
## 1. Sliding window

- Test patch at each location and scale



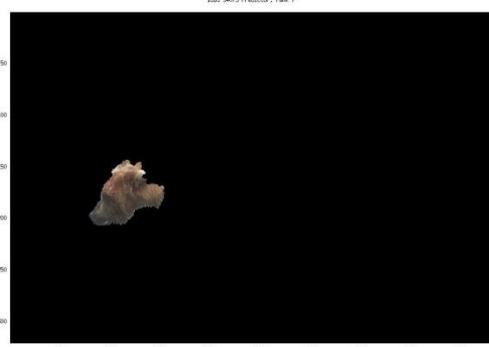
# Generating hypotheses

## 2. Voting from patches/keypoints



# Generating hypotheses

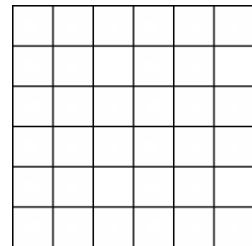
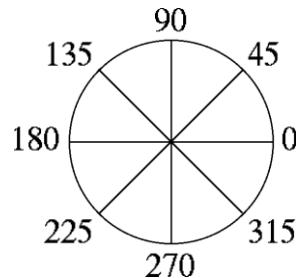
## 3. Region-based proposal



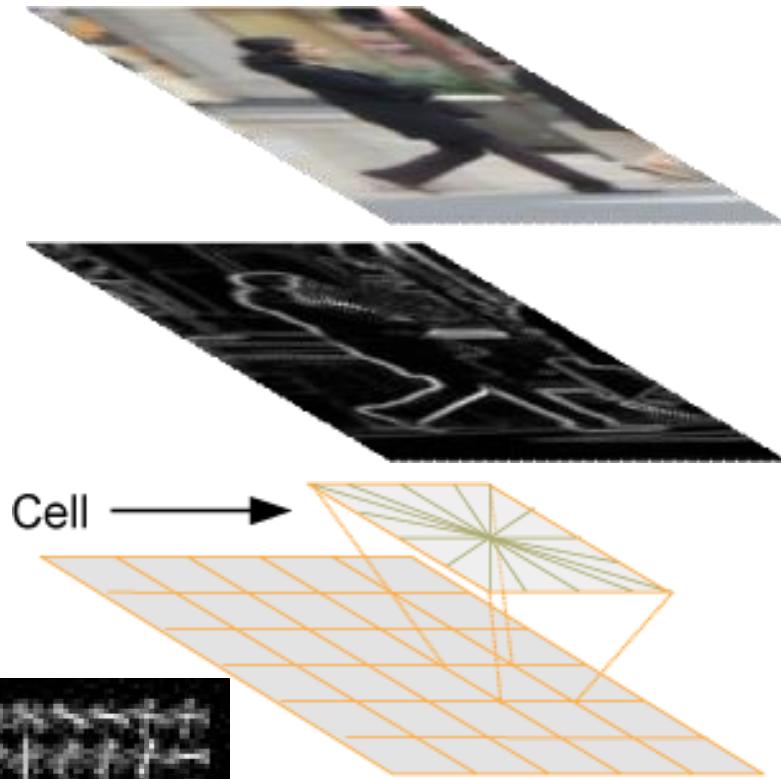
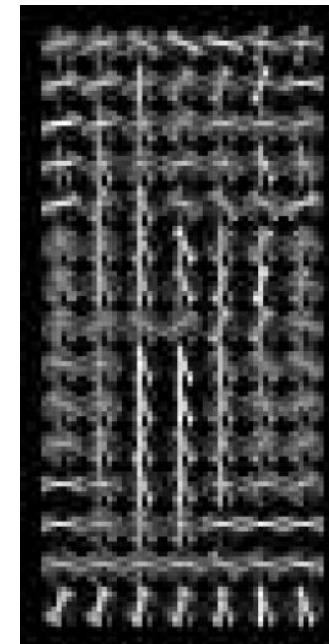
# Score Hypotheses

## Histograms of Oriented Gradients

1. Compute gradients on an image region of 64x128 pixels
  
2. Compute histograms on 'cells' of typically 8x8 pixels (i.e. 8x16 cells)
  - Histogram of gradient orientations
    - Orientation      -Position

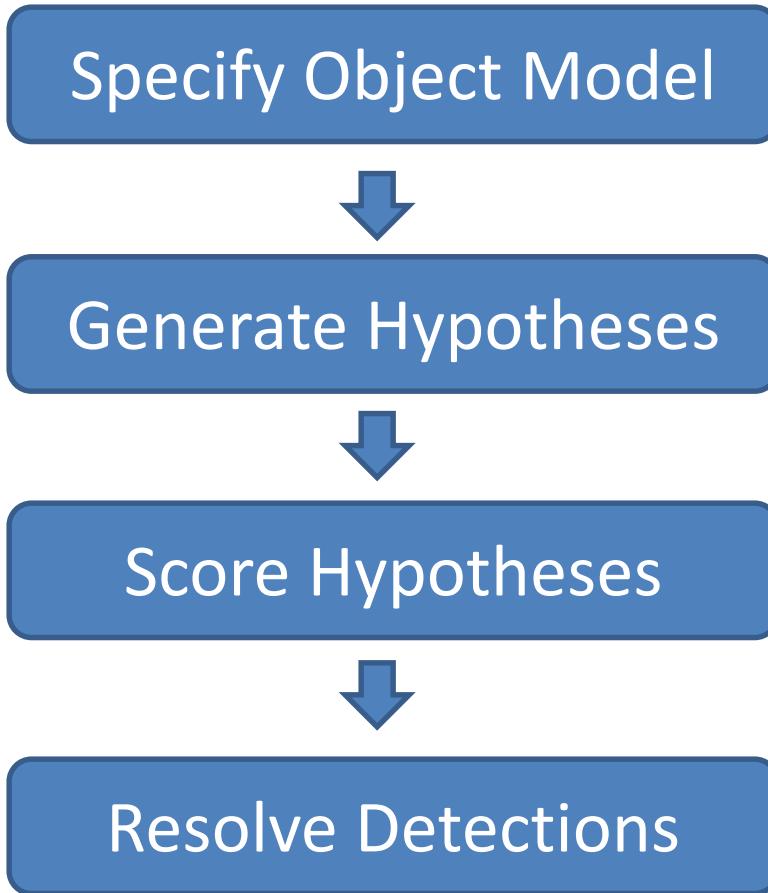


- Weighted by magnitude



7x15 ? Why ?

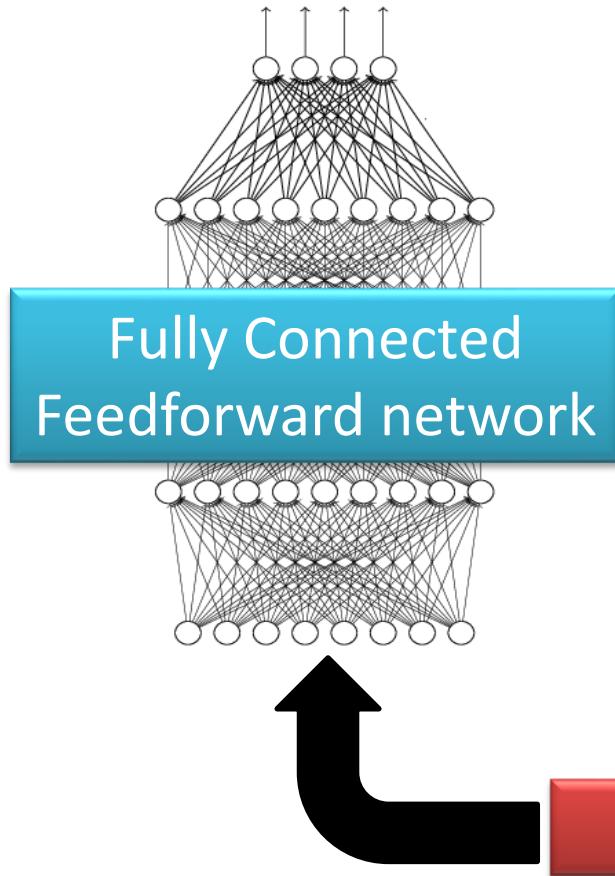
# General Process of Object Recognition



Currently CNN features and classifiers

# CNN

cat dog .....



Convolution

Max Pooling

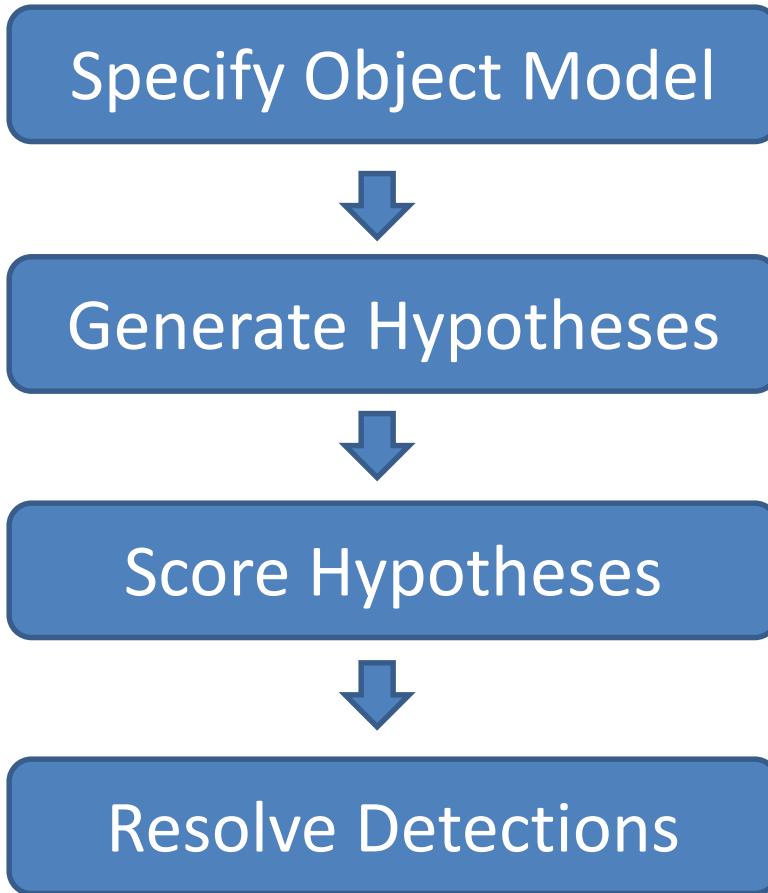
Convolution

Max Pooling

Can  
repeat  
many  
times

Flatten

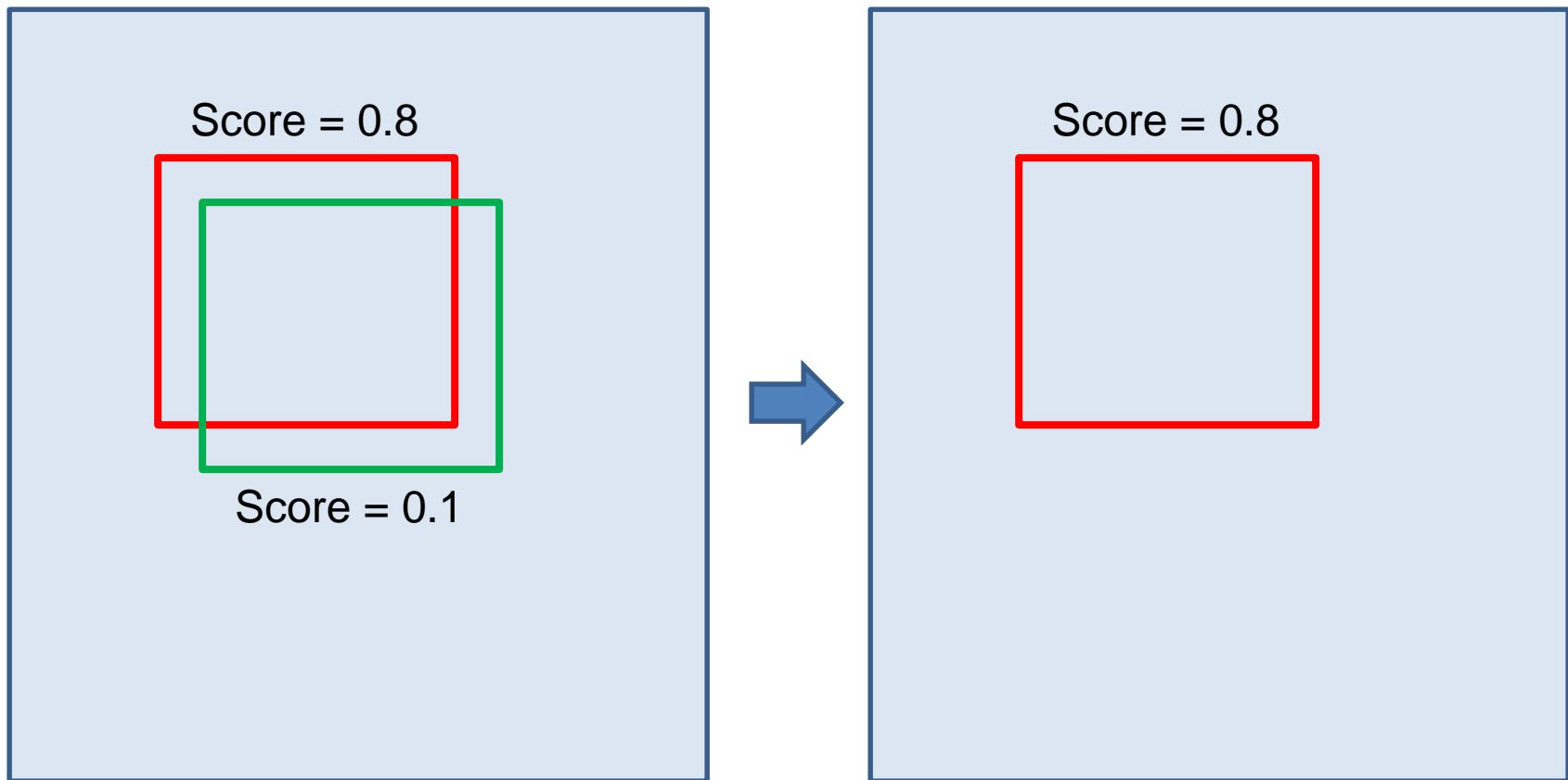
# General Process of Object Recognition



Optionally, rescore each proposed object based on whole set

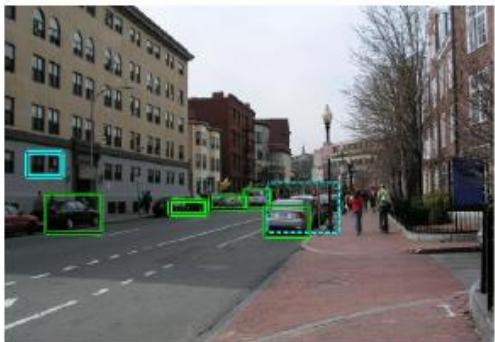
# Resolving detection scores

## 1. Non-max suppression



# Resolving detection scores

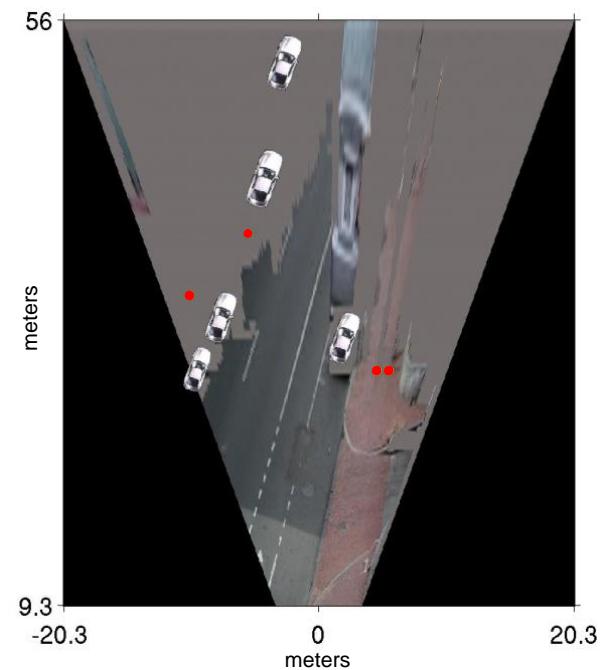
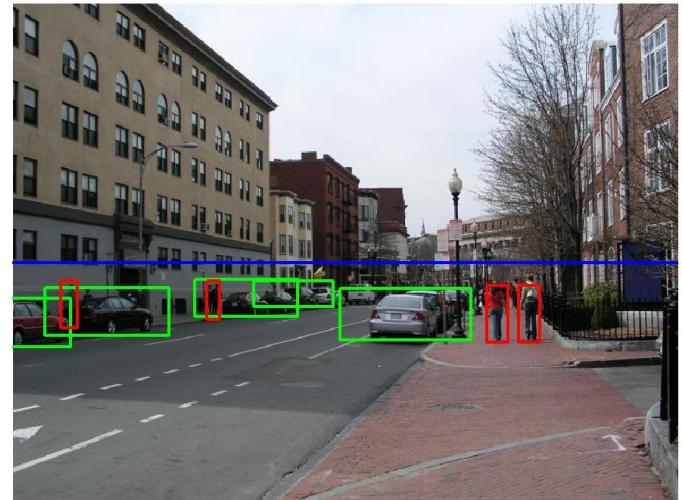
## 2. Context/reasoning



(g) Car Detections: Local



(h) Ped Detections: Local



# Object detection in computer vision

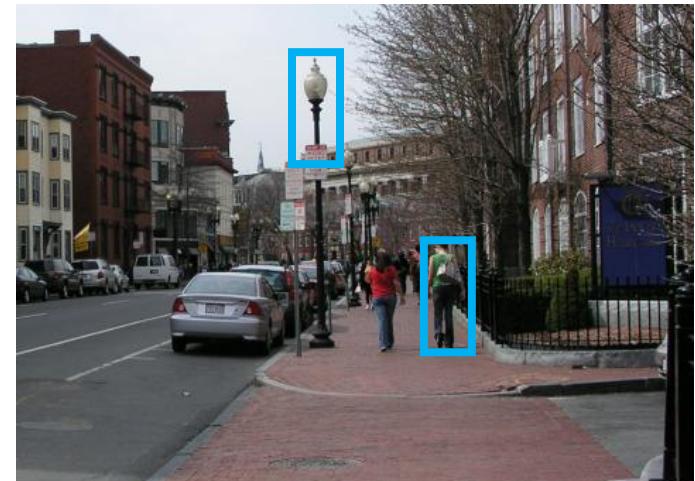
Goal: detect all pedestrians, cars, monkeys, etc in image



# Basic Steps of Category Detection

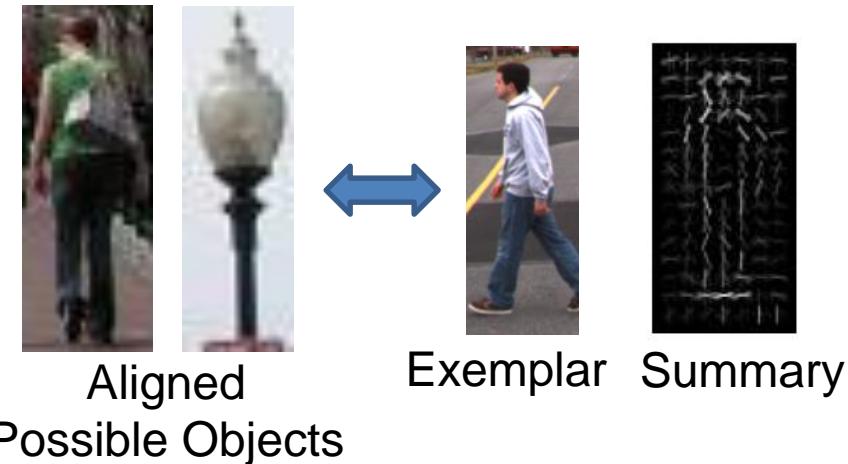
## 1. Align

- E.g., choose position, scale orientation
- How to make this tractable?

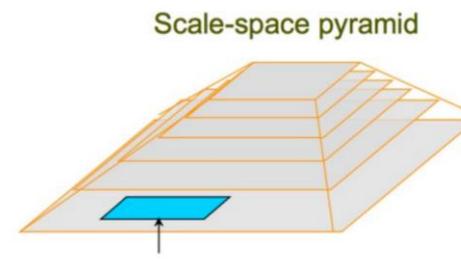


## 2. Compare

- Compute similarity to an example object or to a summary representation
- Which differences in appearance are important?



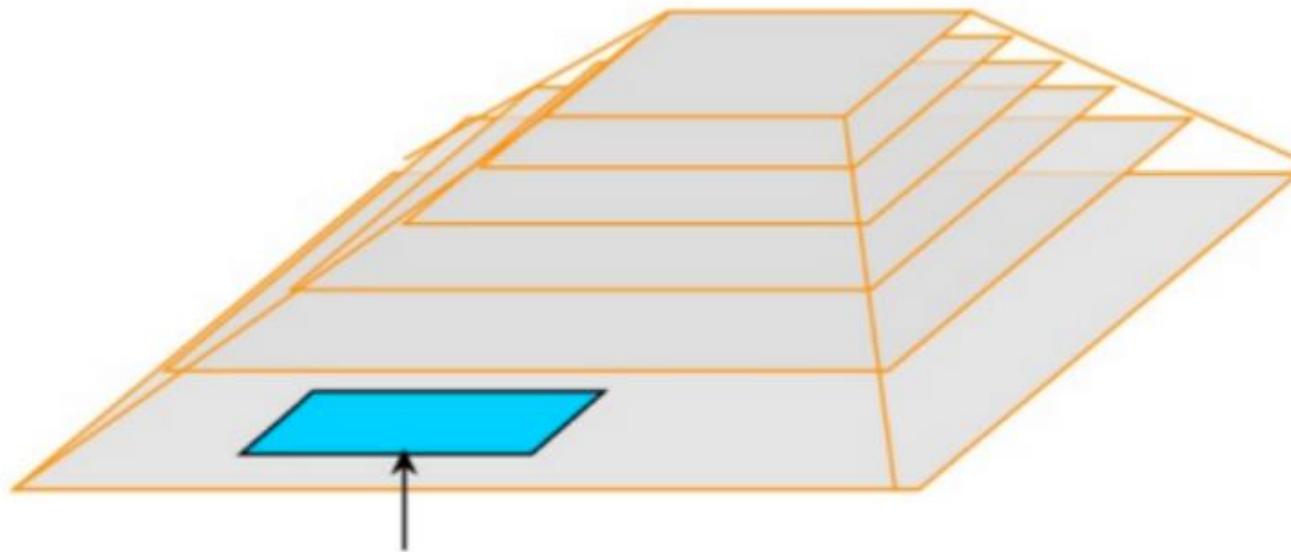
# Sliding window: a simple alignment solution



scales



# Scale-space Pyramid



Detection window

# Each window is separately classified



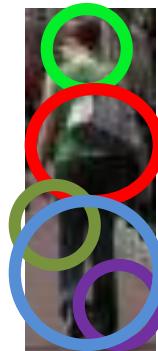
# Statistical Template

- Object model = sum of scores of features at fixed positions



$$+3 +2 -2 -1 -2.5 = -0.5 > 7.5 ?$$

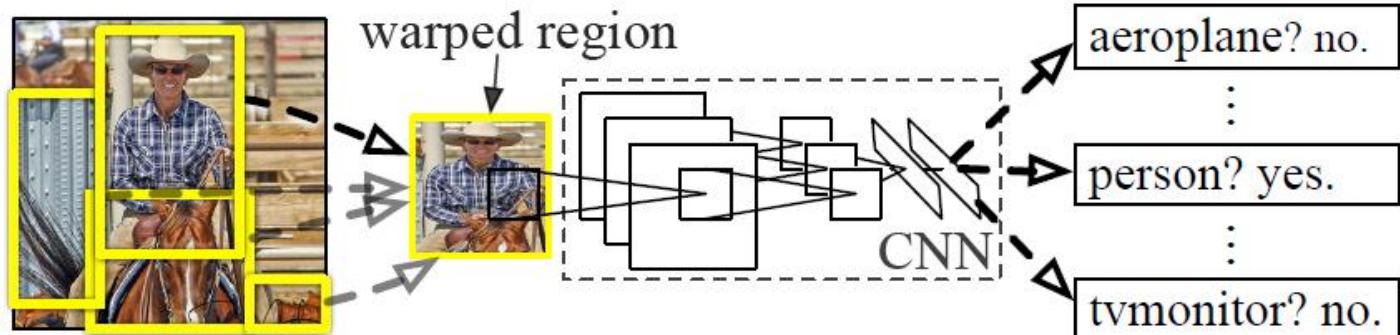
Non-object



$$+4 +1 +0.5 +3 +0.5 = 10.5 > 7.5 ?$$

Object

# R-CNN (Girshick et al. CVPR 2014)



1. Input image

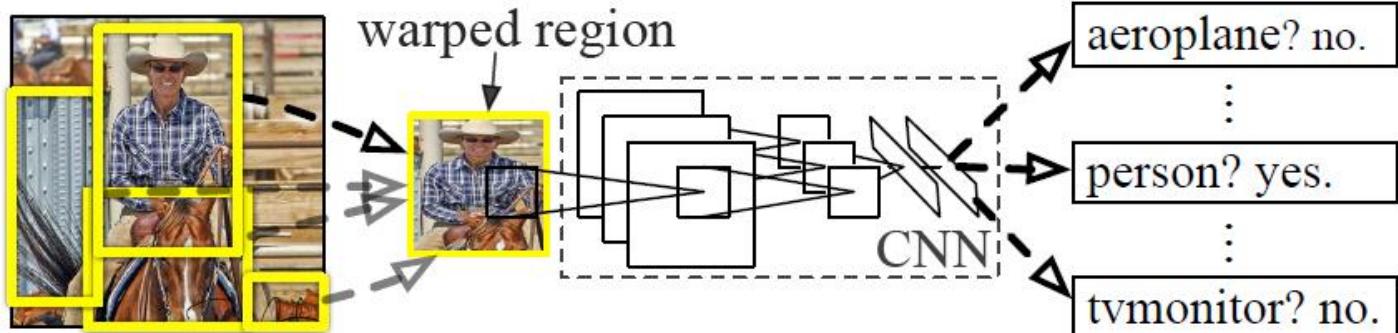
2. Extract region proposals (~2k)

3. Compute CNN features

4. Classify regions

- Replace sliding windows with “selective search” region proposals (Uijlings et al. IJCV 2013)
- Extract rectangles around regions and resize to 227x227
- Extract features with fine-tuned CNN (that was initialized with network trained on ImageNet before training)
- Classify last layer of network features with SVM

# R-CNN (Girshick et al. CVPR 2014)



1. Input image

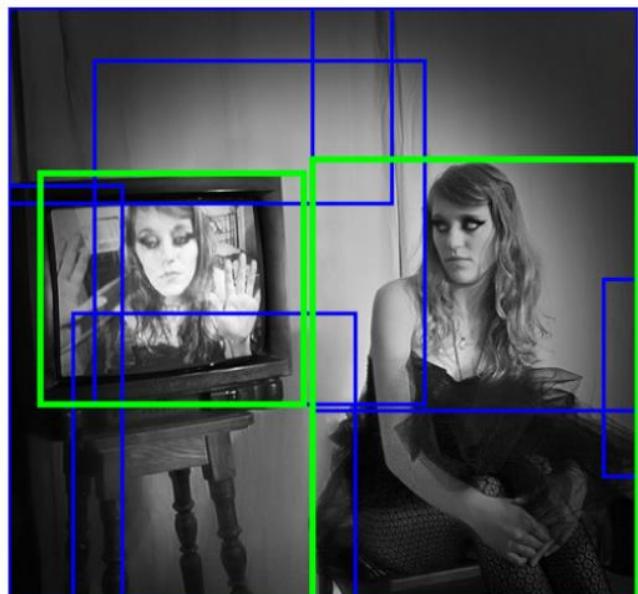
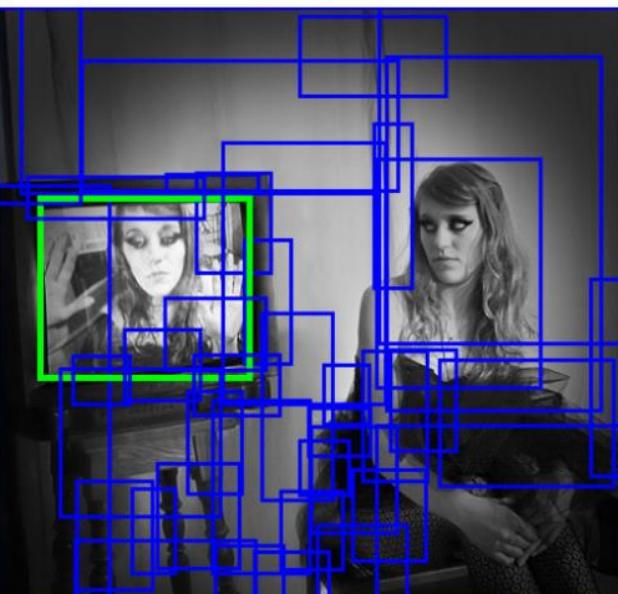
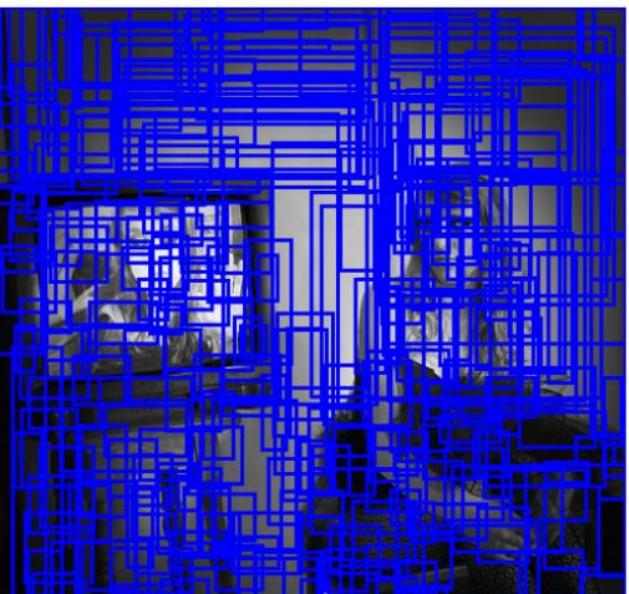
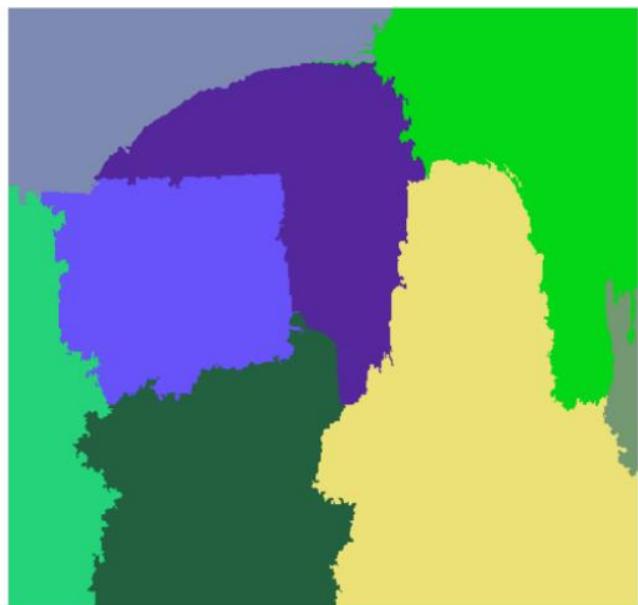
2. Extract region proposals (~2k)

3. Compute CNN features

4. Classify regions

- 產生一群約 2000 個可能的區域 (Region Proposals)
- 經由一個預先訓練好的 CNN 模型如 AlexNet 機取特徵，將結果儲存起來。
- 然後再以 SVM (Support Vector Machine) 分類器來區分是否為物體或者背景。
- 最後經由一個線性回歸模型來校正 bounding box 位置。

# Regions from selective search



# Fine-tuning example: ImageNet->VOC

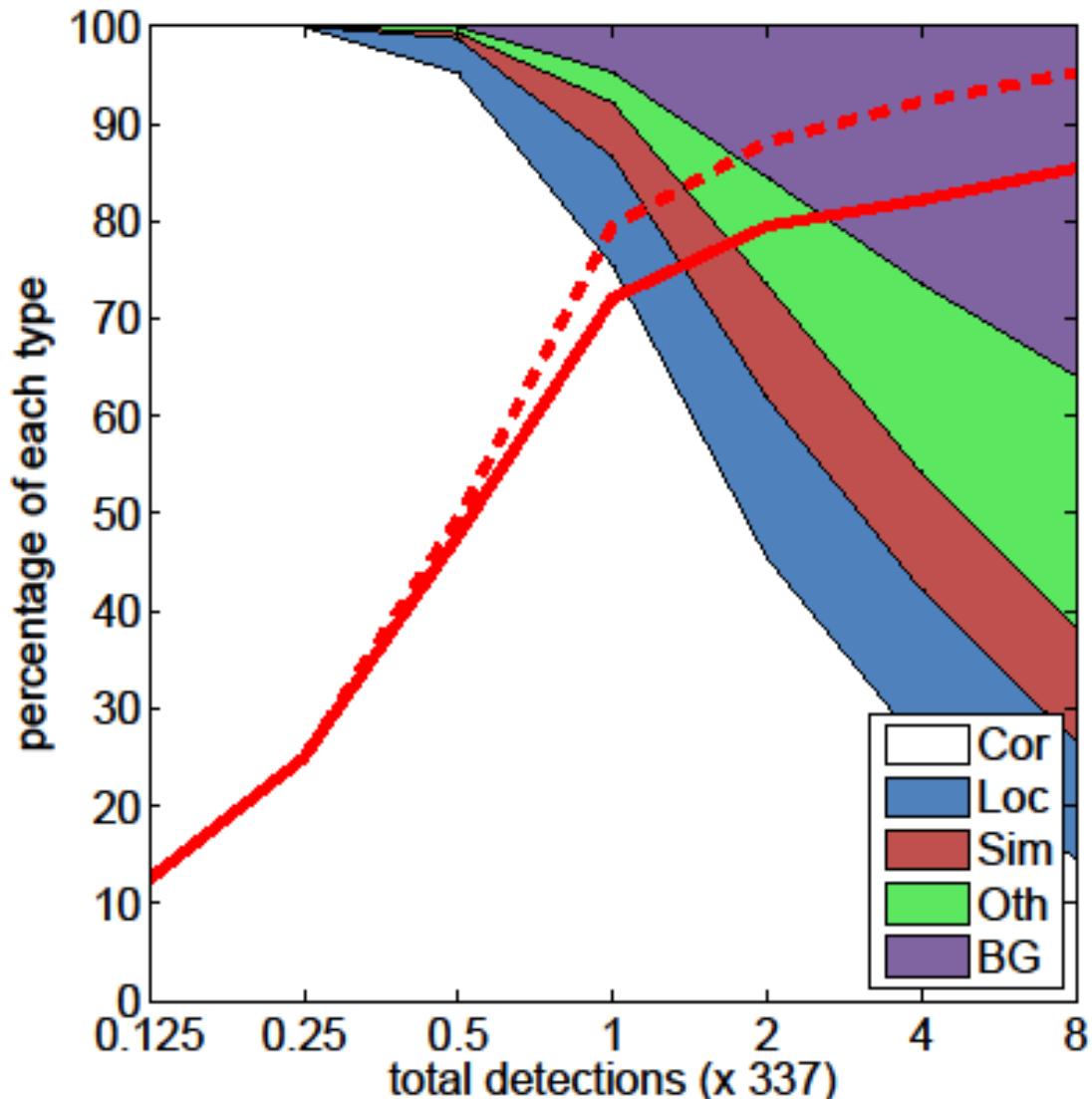
1. Train full network on ImageNet 1000-class classification
2. Replace classification layer with output layer for VOC (e.g. confidences for 20 classes)
3. Train on VOC pos/neg examples with low initial learning rate ( $1/10^{\text{th}}$  what is used for new network)

## Notes

- This usually works well if the “big data” task and target task are similar (object classification vs detection)
  - 0.45 AP without fine-tuning → 0.54 AP with fine tuning; training only on VOC does much worse
- Not necessary if target task is also very big

# Mistakes are often reasonable

Bicycle: AP = 0.73

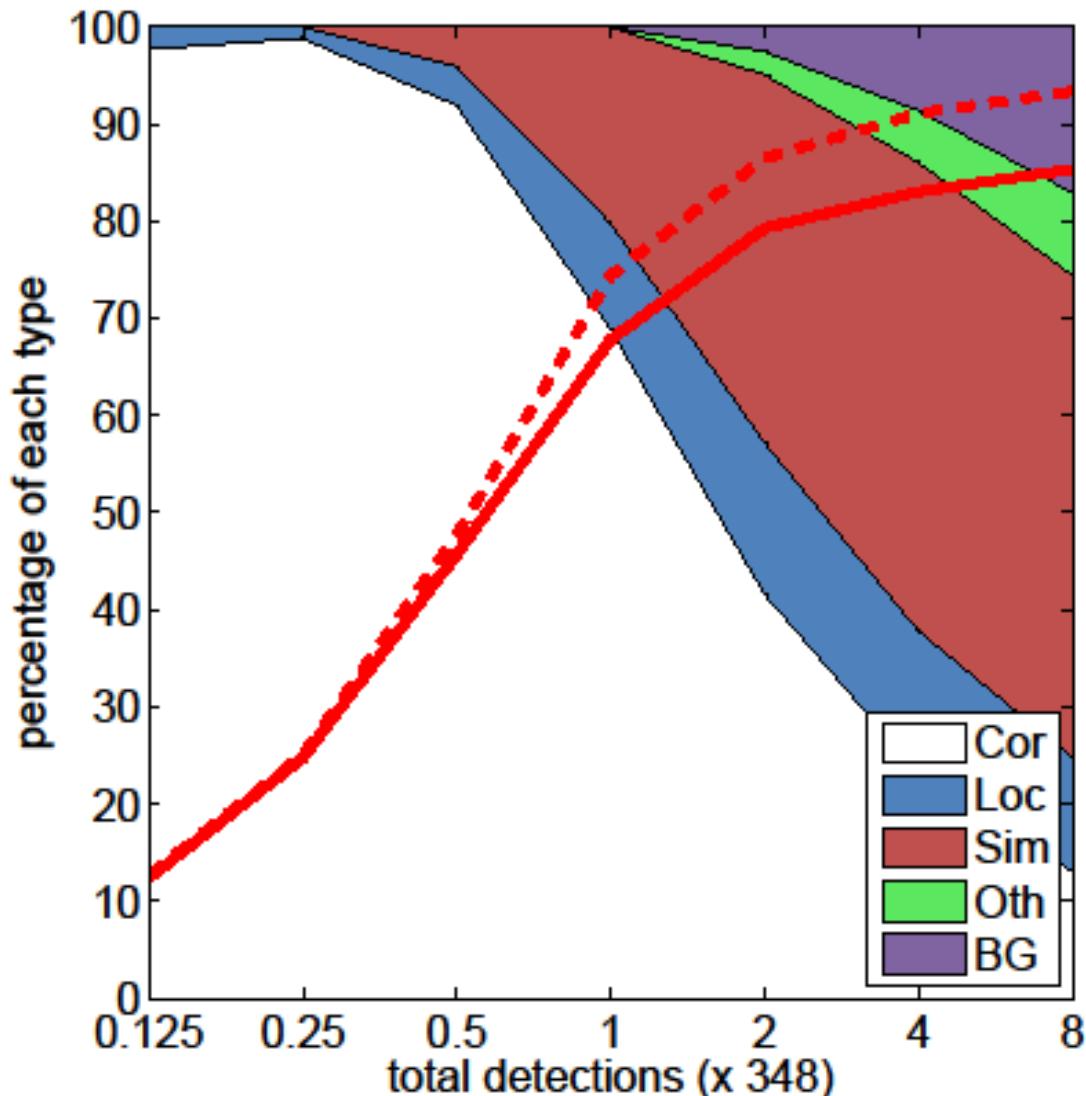


Confident Mistakes

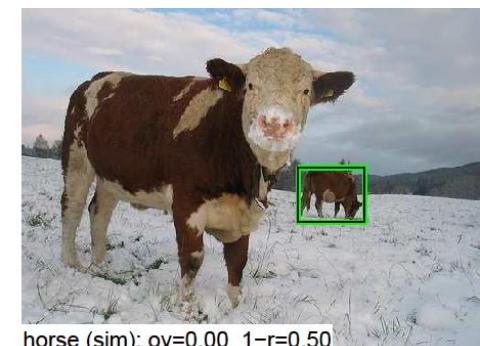


# Mistakes are often reasonable

Horse: AP = 0.69

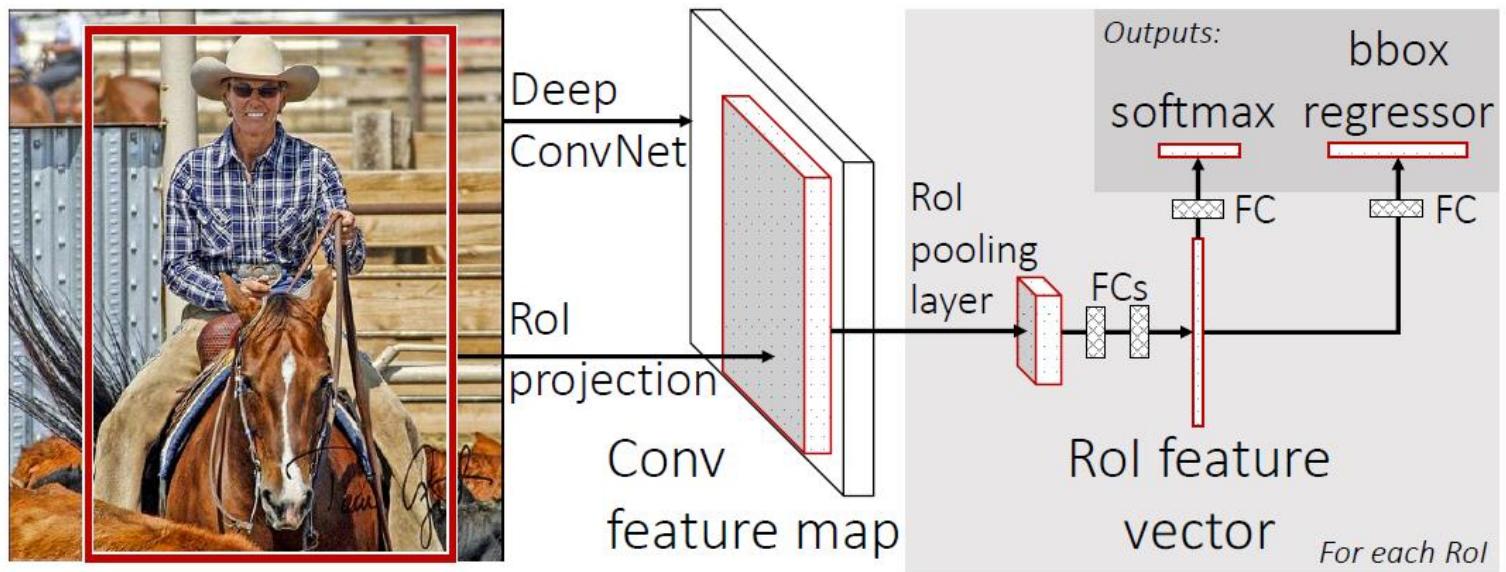


Confident Mistakes



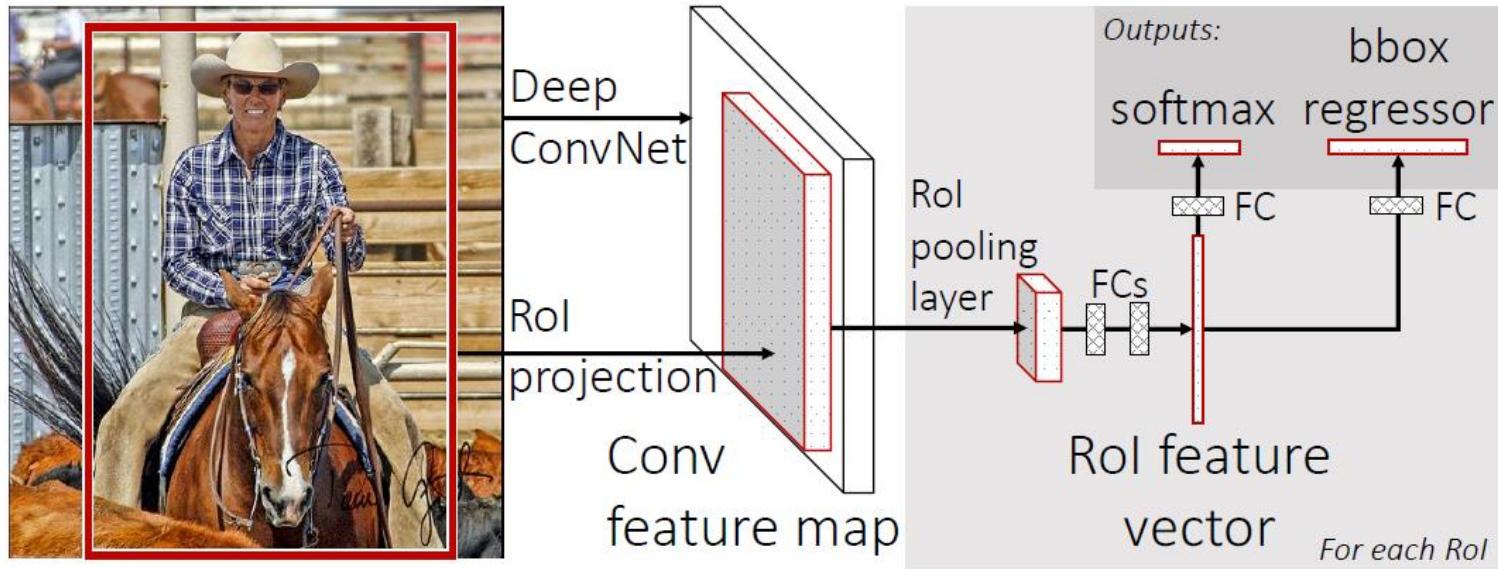
R-CNN results

# Fast R-CNN – Girshick 2015



- Compute CNN features for image once
- Pool into  $7 \times 7$  spatial bins for each region proposal, output class scores and regressed bboxes
- 100x speed up of R-CNN ( $0.02 - 0.1$  FPS  $\rightarrow$  0.5-20 FPS) with similar accuracy

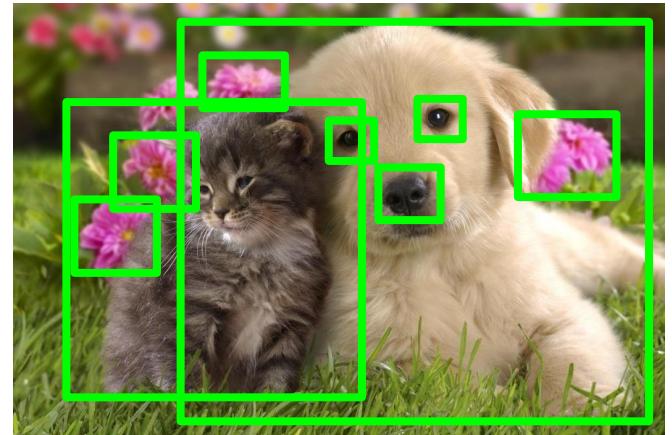
# Fast R-CNN – Girshick 2015



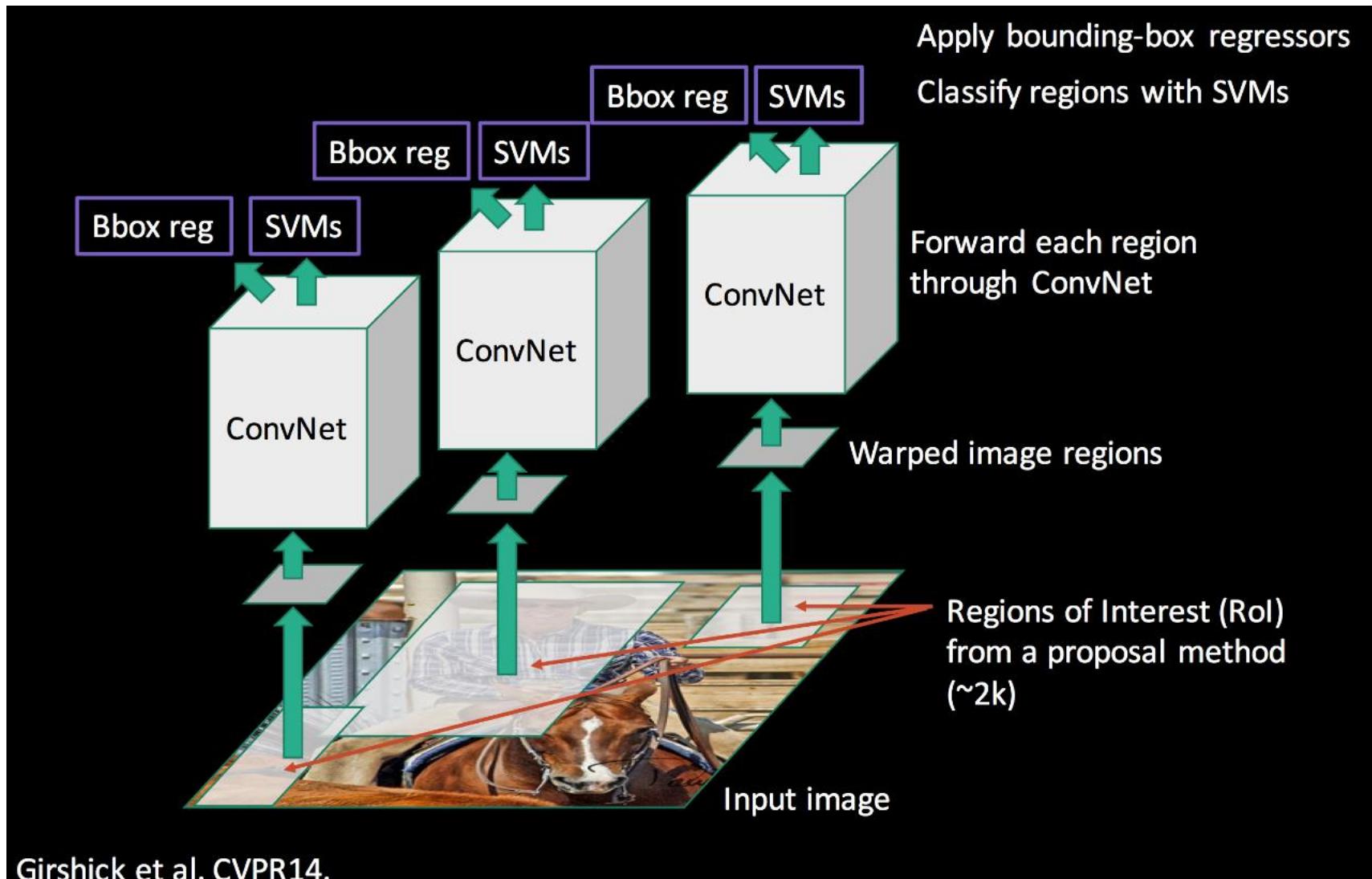
- **Fast R-CNN** : 在 R-CNN 中，2000 多個區域都要個別去運算 CNN，這些區域很多都是重疊的，也就是說這些重疊區域的 CNN 很多都是重複算的。所以 Fast R-CNN 的原則就是全部只算一次 CNN 就好，CNN 擷取出來的特徵可以讓這 2000 多個區域共用！
- Fast R-CNN: ROI Pooling (Region of Interest Pooling)

# Region Proposals

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



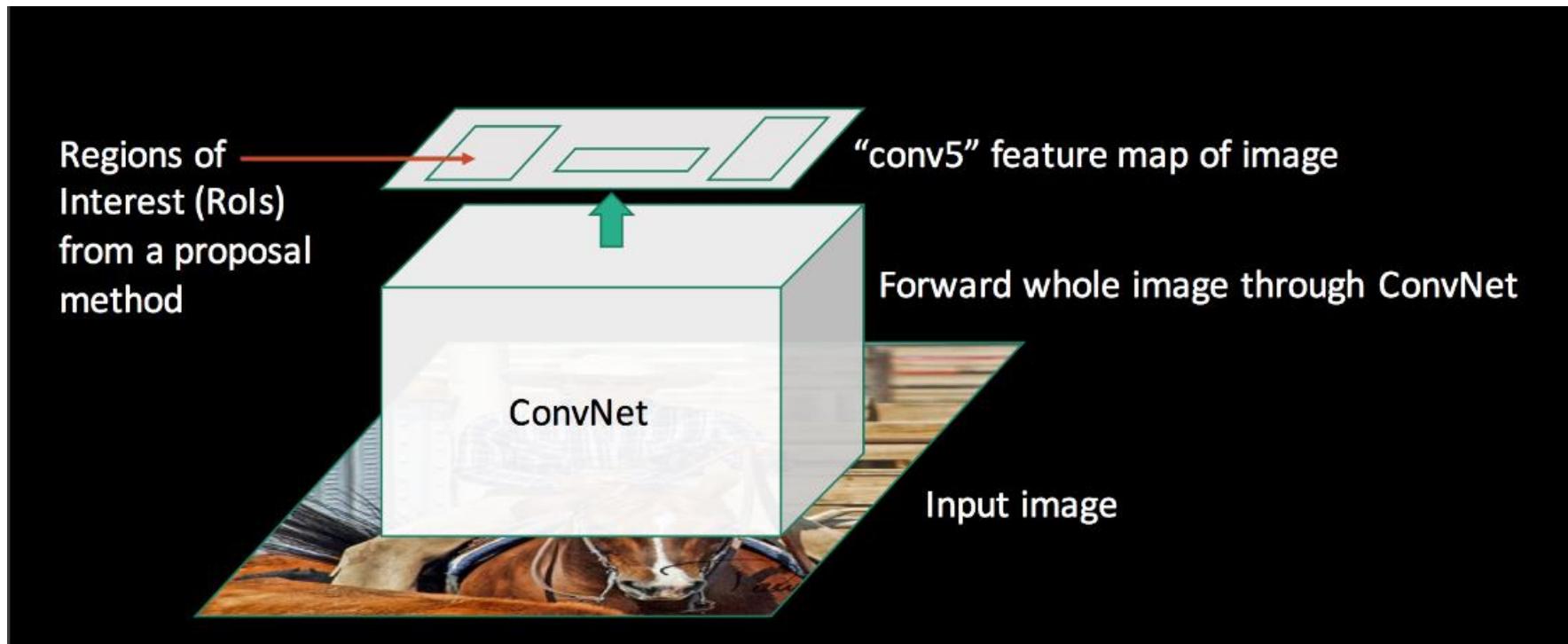
# R-CNN



Girshick et al. CVPR14.

Number of CNN runs = number of proposals → high latency

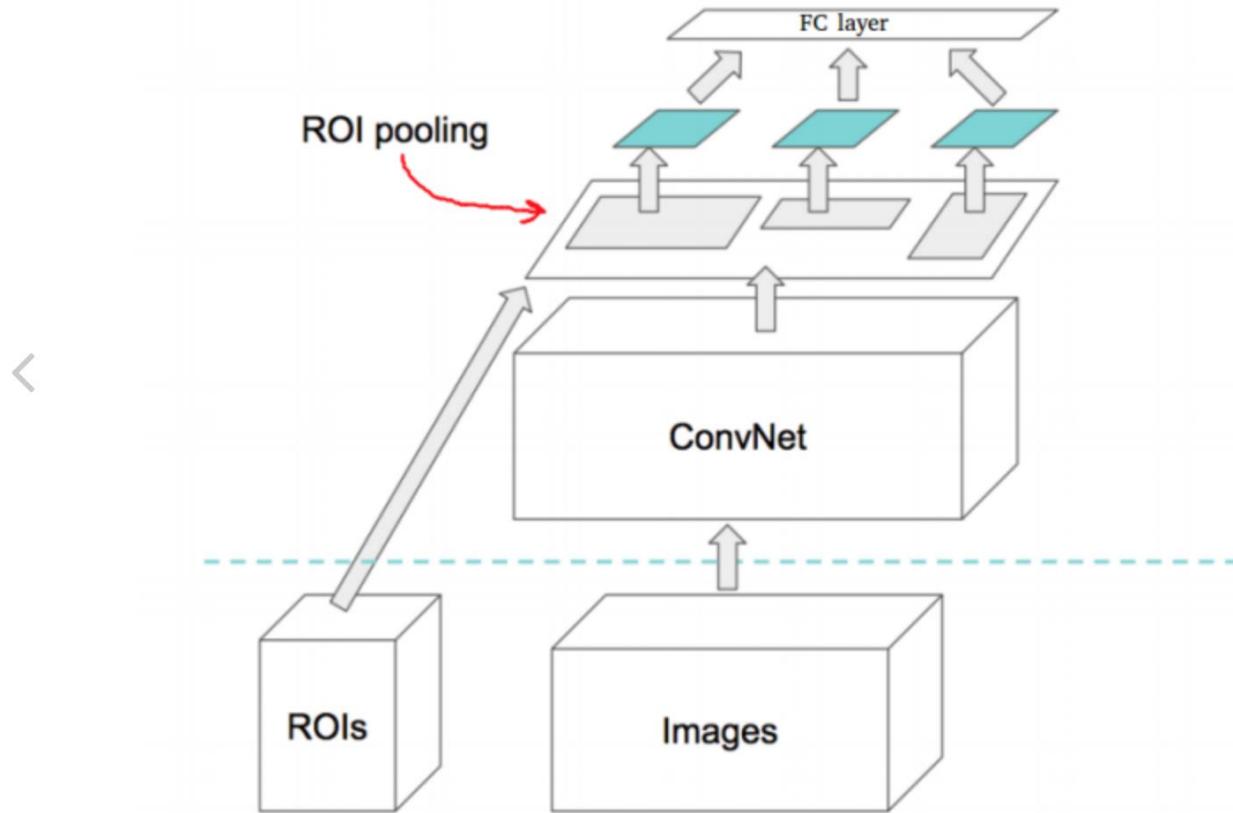
# Fast R-CNN (Girschick et. al. 2015)



1 CNN run for whole image

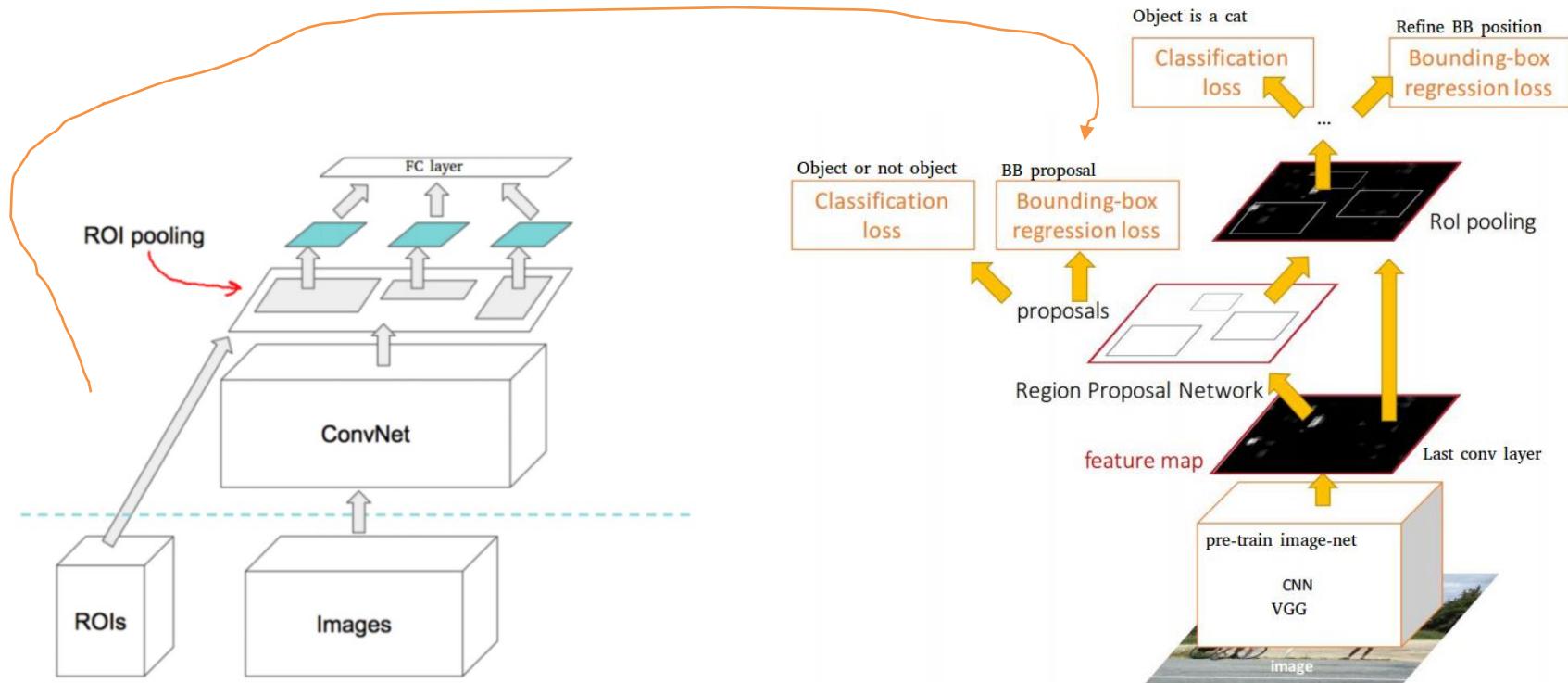
proposed regions extracted from CNN feature map, instead of image

# New operation: Region of Interest (ROI) pooling



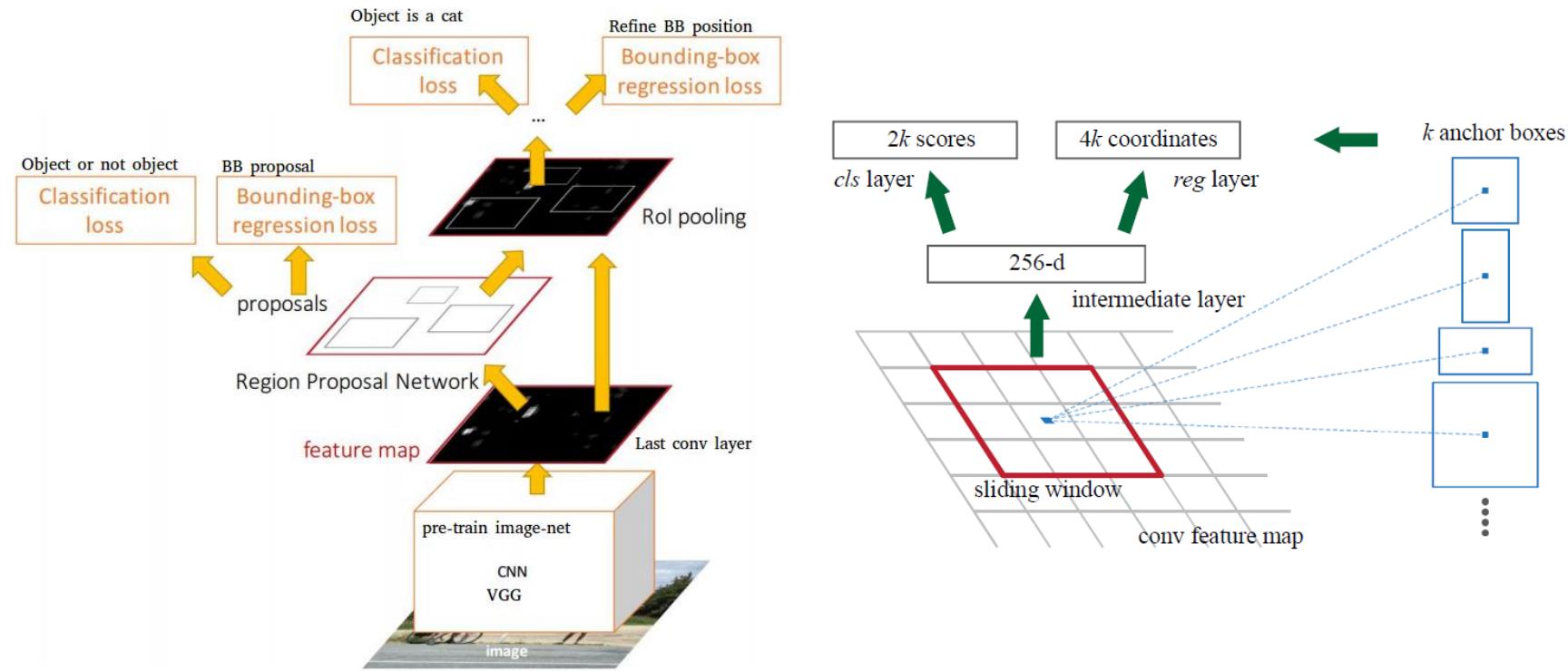
Necessary to convert variable sized ROIs to fixed size inputs  
expected by the following FC layers  
(similar to warping original image patches in R-CNN  
for fixed sized inputs to conv layers)

# Faster R-CNN – Ren et al. 2016



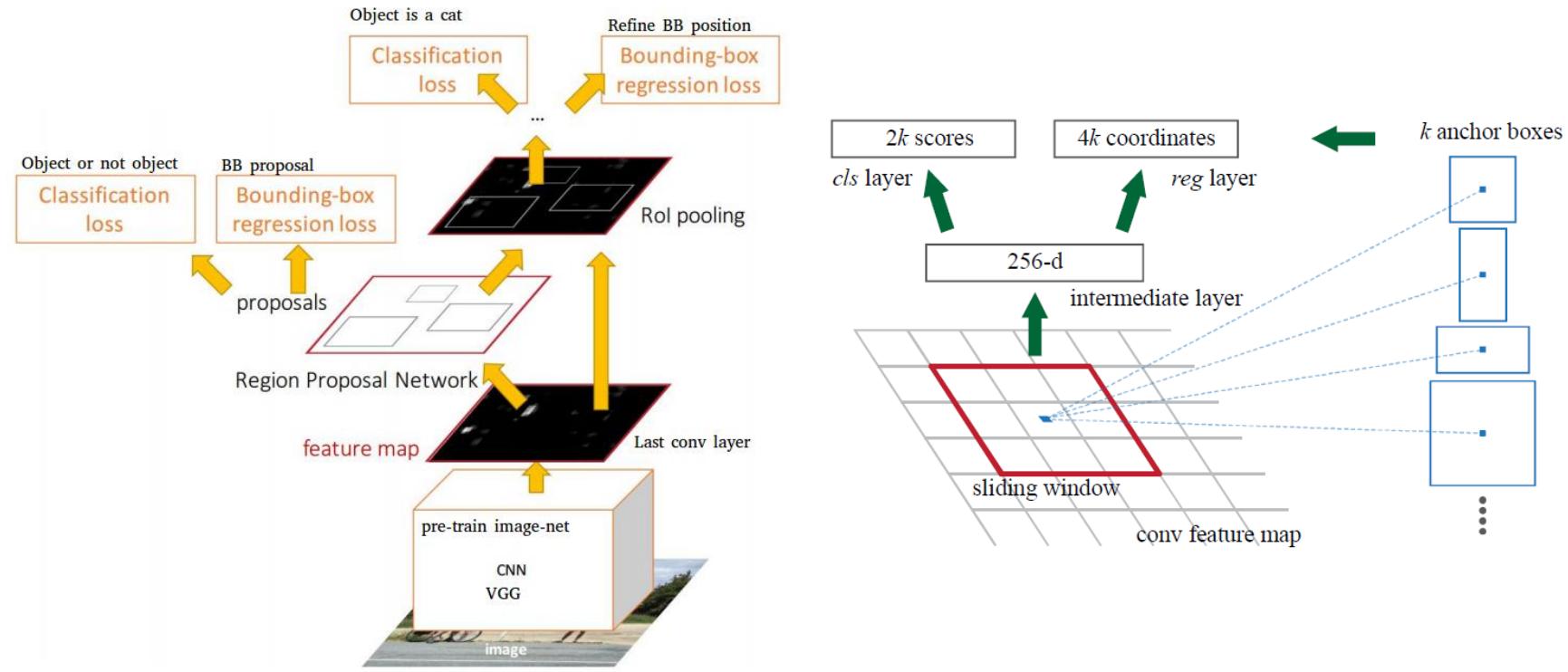
- Faster R-CNN 的直覺想法: 與其預先篩選 region proposals，到不如從 CNN 的 feature map 上選出 region proposals。

# Faster R-CNN – Ren et al. 2016



- Faster R-CNN 的直覺想法: 利用一個region proposal network 去產生region candidates

# Faster R-CNN – Ren et al. 2016



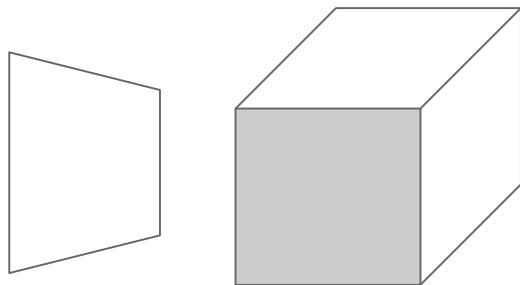
- Convolutional features used for generating proposals and scoring
  - Generate proposals with “objectness” scores and refined bboxes for each of  $k$  “anchors”
  - Score proposals in same way as Fast R-CNN
- Similar accuracy to Fast R-CNN with 10x speedup

# Fast R-CNN: Region of Interest Pooling

Convolution  
and Pooling

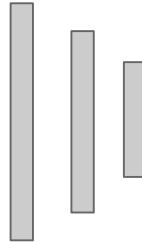


Hi-res input image:  
 $3 \times 800 \times 600$   
with region  
proposal



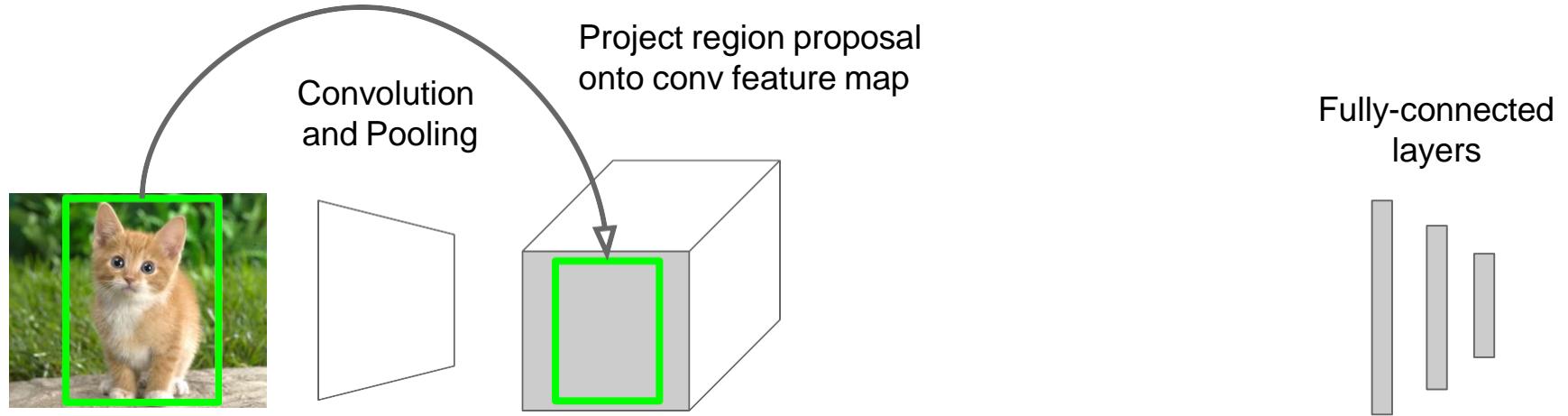
Hi-res conv features:  
 $C \times H \times W$   
with region proposal

Fully-connected  
layers



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

# Fast R-CNN: Region of Interest Pooling

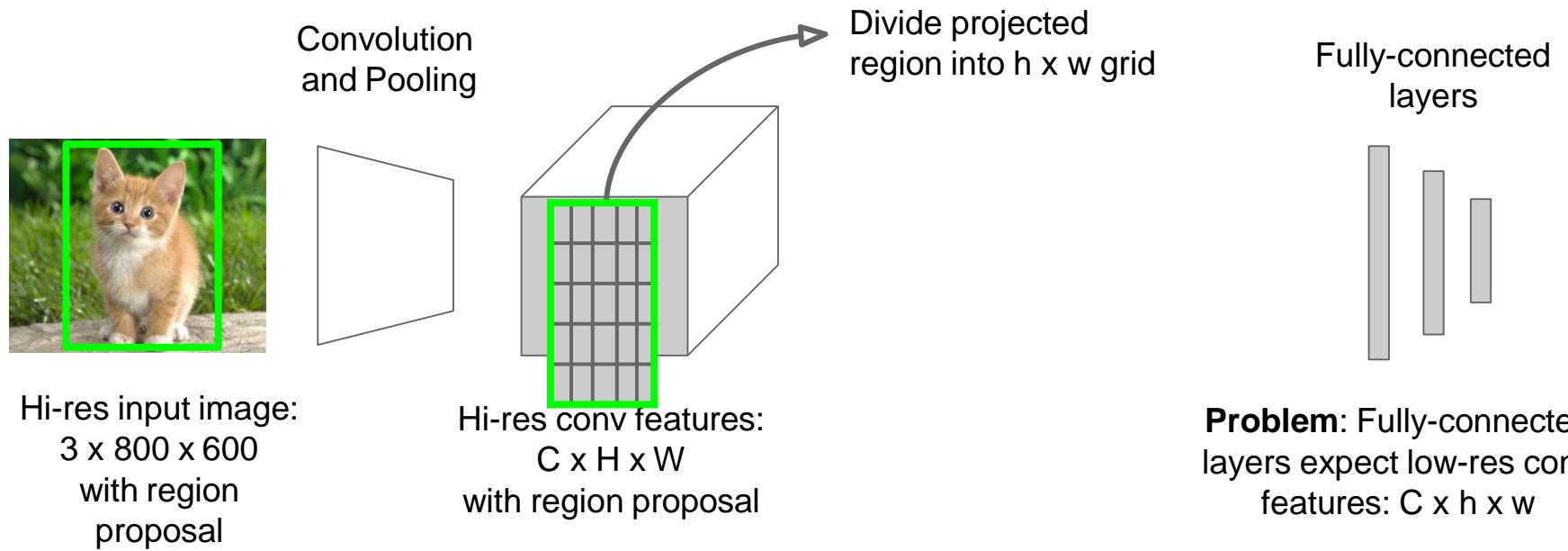


Hi-res input image:  
 $3 \times 800 \times 600$   
with region  
proposal

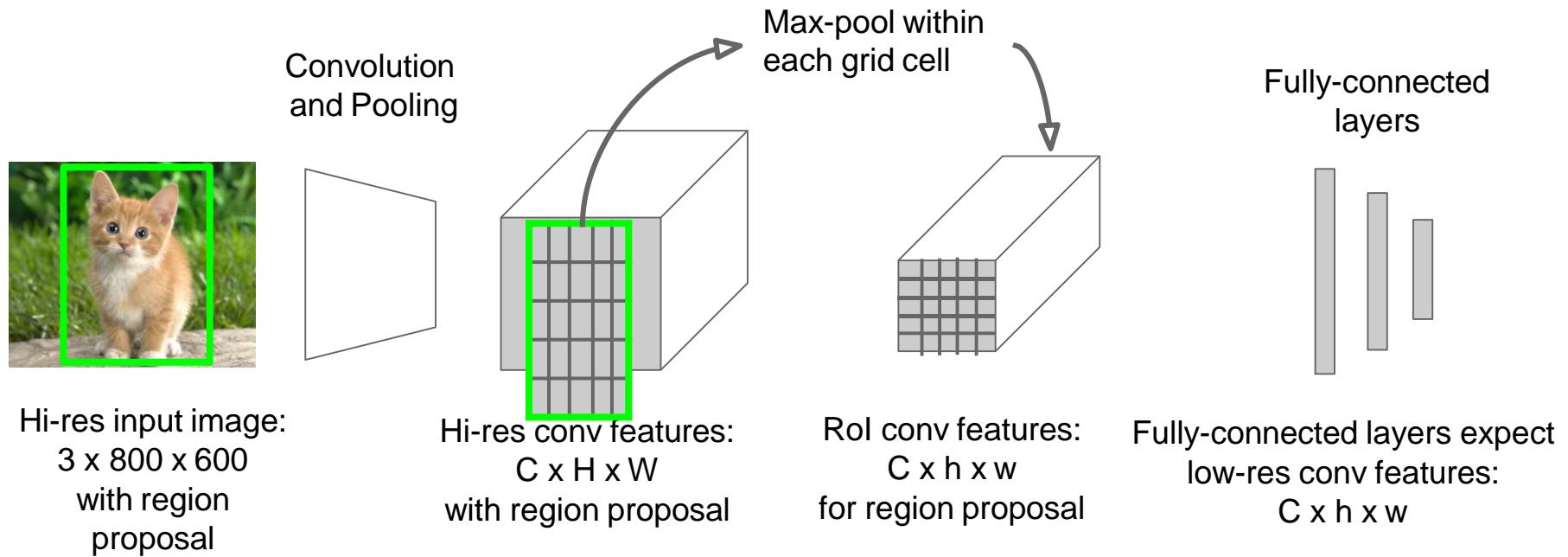
Hi-res conv features:  
 $C \times H \times W$   
with region proposal

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

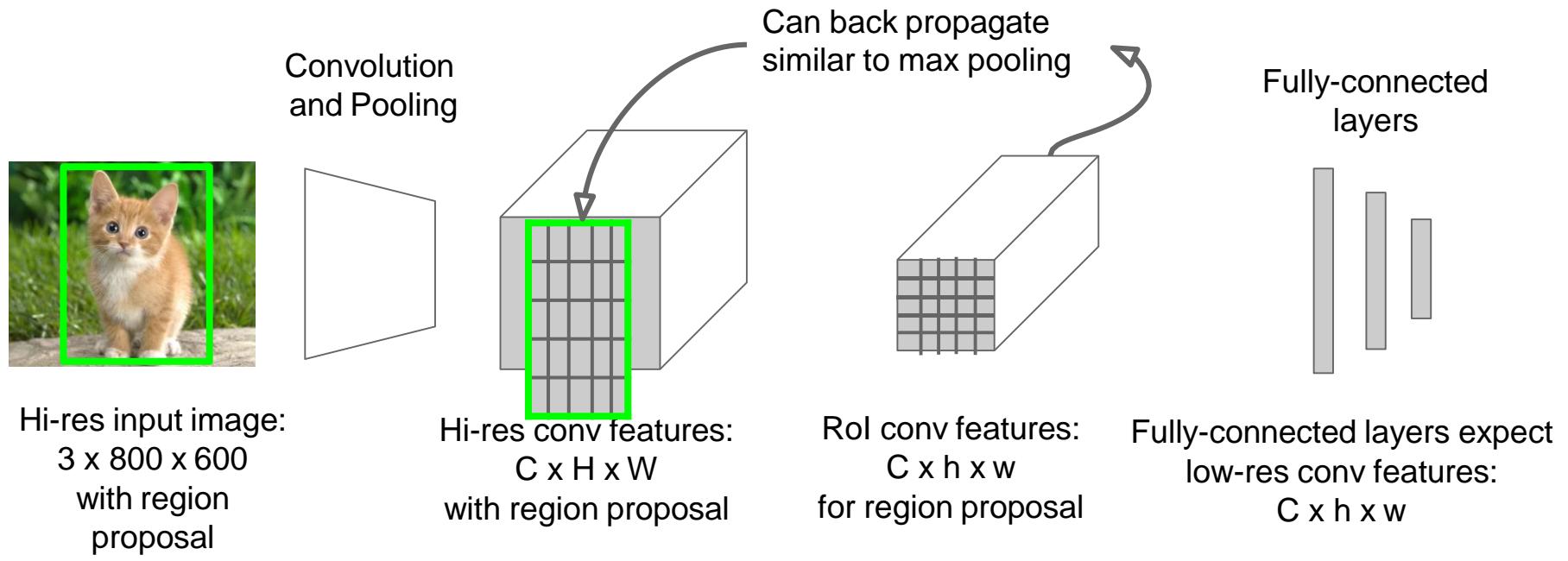
# Fast R-CNN: Region of Interest Pooling



# Fast R-CNN: Region of Interest Pooling



# Fast R-CNN: Region of Interest Pooling



# Fast R-CNN Results

Faster!

	R-CNN	Fast R-CNN
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

Faster!  
FASTER!

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

Faster!

FASTER!

Better!

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

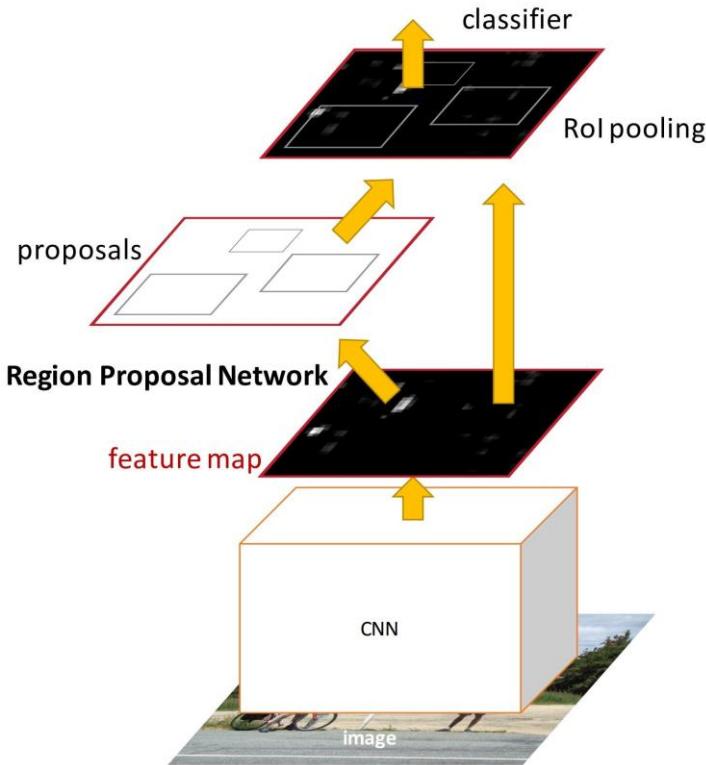
	R-CNN	Fast R-CNN
Test time per image	47 seconds	<b>0.32 seconds</b>
(Speedup)	1x	<b>146x</b>
Test time per image with Selective Search	50 seconds	<b>2 seconds</b>
(Speedup)	1x	<b>25x</b>

# 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

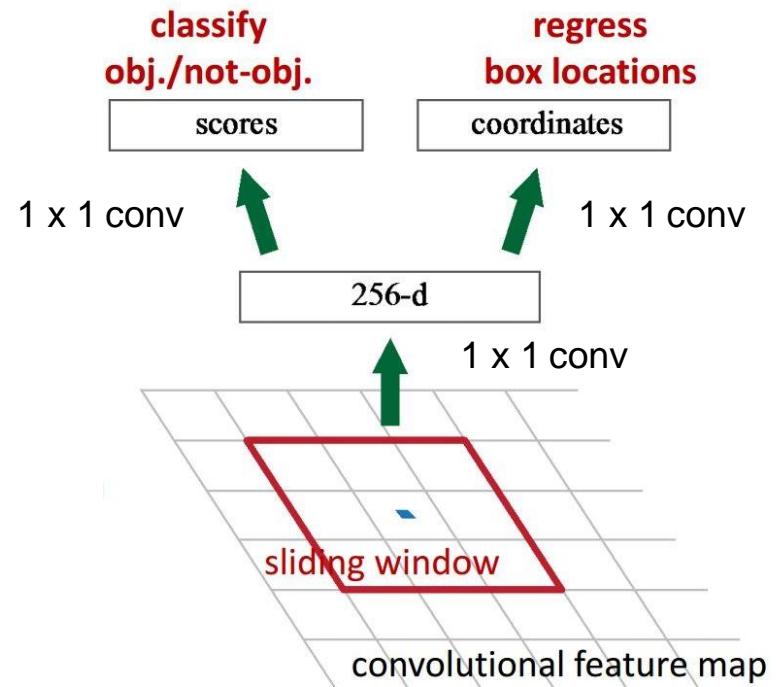
Slide a small window on the feature map

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

Fei-Fei Li & Andrej Karpathy & Justin Johnson

1 Feb 2016

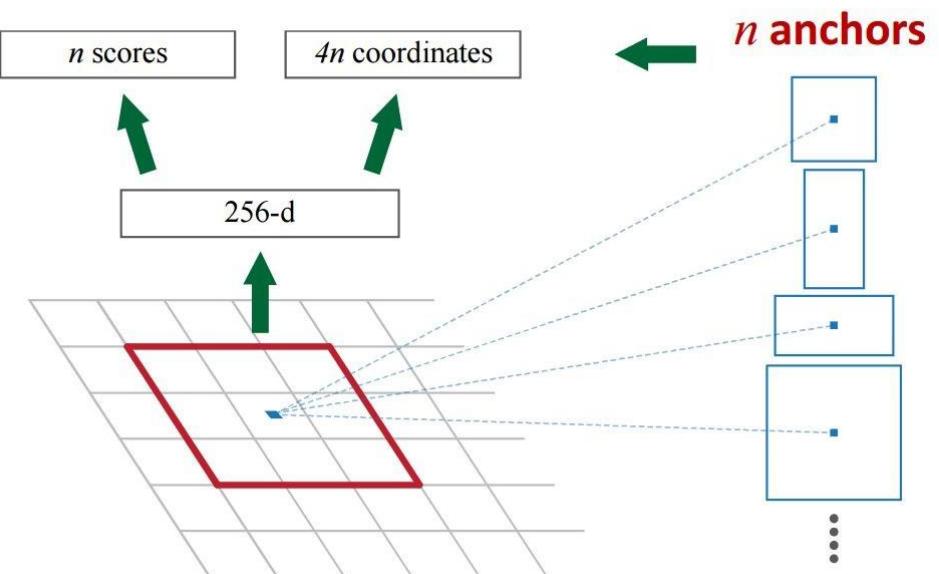
# Faster R-CNN: Region Proposal Network

Use **N anchor boxes** at each location

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

Regression gives offsets from anchor boxes

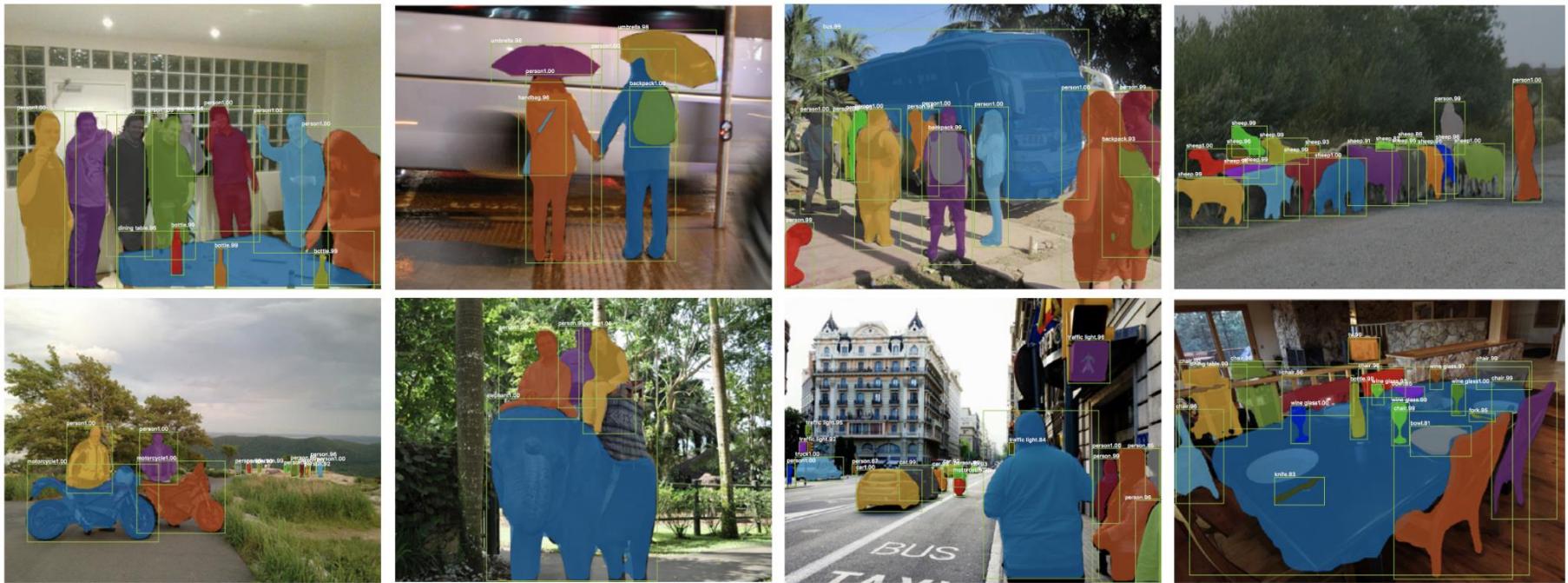
Classification gives the probability that each (regressed) anchor shows an object



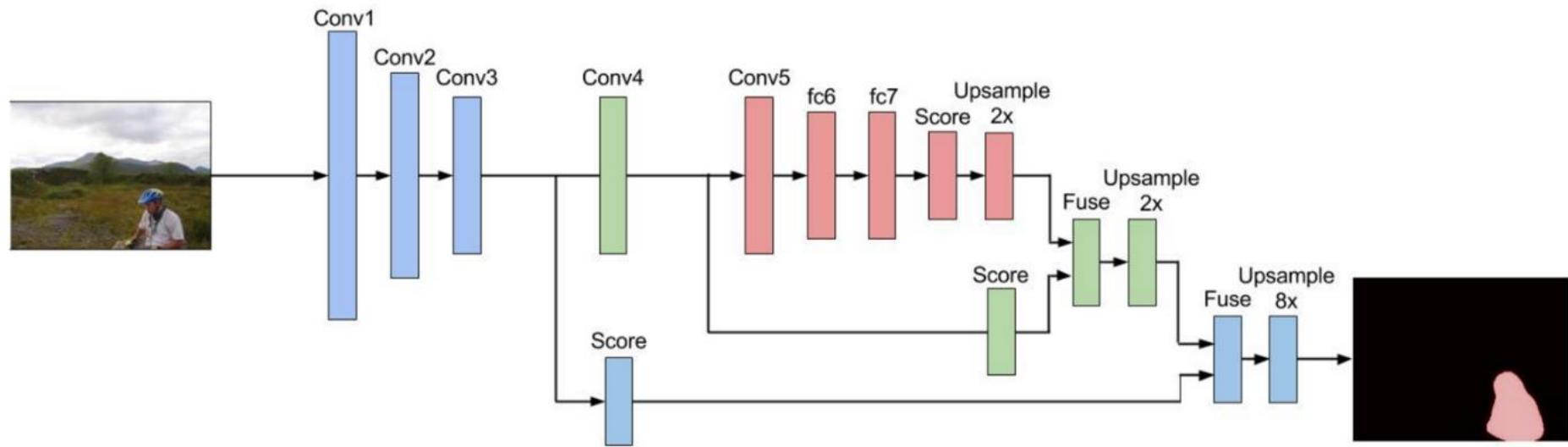
# Faster R-CNN: Results

	R-CNN	Fast R-CNN	Faster R-CNN
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>

# MASK R-CNN



# MASK R-CNN



# One stage method for object detection

JOSEPH

ROSS

SANTOSH

ALI

REDMON

GIRSHICK

DIVVALA

FARHADI

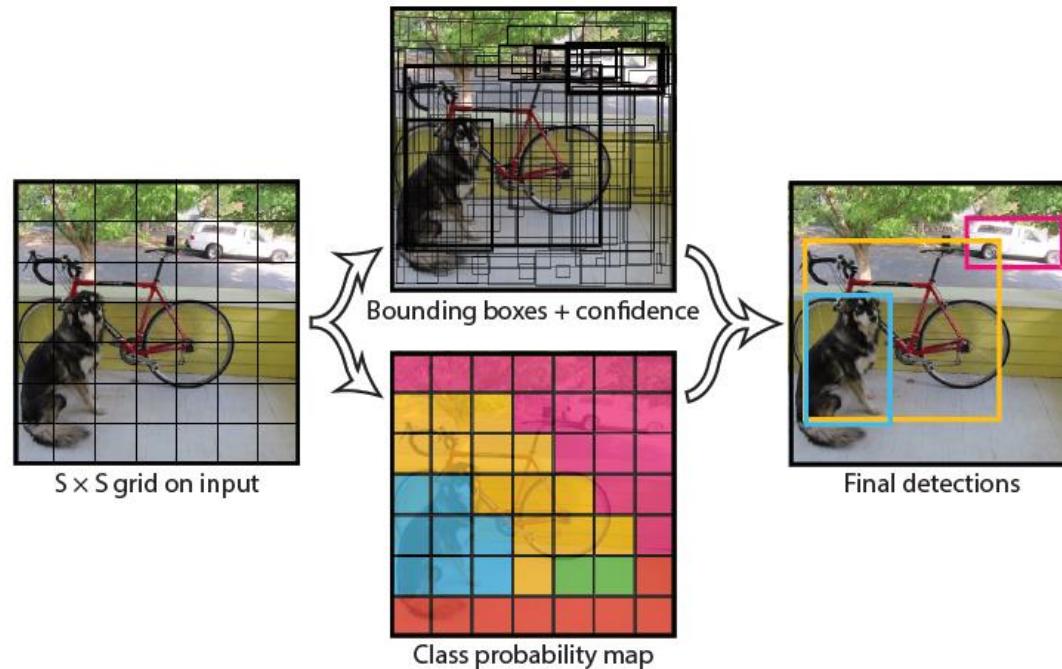
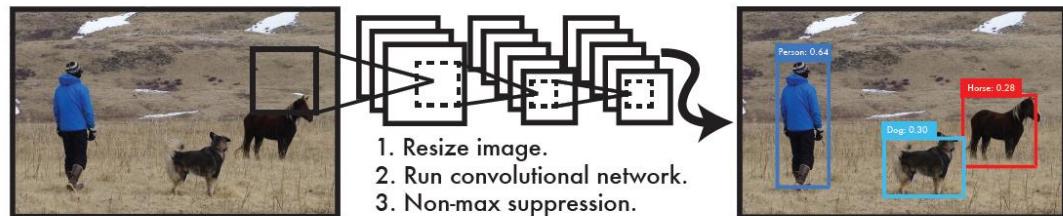
Dog



“YOU ONLY LOOK ONCE”  
REAL-TIME  
DETECTION

# YOLO – Redmon et al. 2016

1. CNN produces 4096 features for  $7 \times 7$  grid on image (fully conv.)
  2. Each cell produces a score for each object and 2 bboxes w/ conf
  3. Non-max suppression
- 
- 7x speedup over Faster RCNN (45-155 FPS vs. 7-18 FPS)
  - Some loss of accuracy due to lower recall, poor localization



# Accurate object detection is slow!

	Pascal 2007 mAP	Speed	
DPM v5	33.7	.07 FPS	14 s/img
R-CNN	66.0	.05 FPS	20 s/img
Fast R-CNN	70.0	.5 FPS	2 s/img



# Accurate object detection is slow!

	Pascal 2007 mAP	Speed	
DPM v5	33.7	.07 FPS	14 s/img
R-CNN	66.0	.05 FPS	20 s/img
Fast R-CNN	70.0	.5 FPS	2 s/img
Faster R-CNN	73.2	7 FPS	140 ms/img



# Accurate object detection is slow!

	Pascal 2007 mAP	Speed	
DPM v5	33.7	.07 FPS	14 s/img
R-CNN	66.0	.05 FPS	20 s/img
Fast R-CNN	70.0	.5 FPS	2 s/img
Faster R-CNN	73.2	7 FPS	140 ms/img
YOLO	63.4	45 FPS	22 ms/img



2 feet

# Accurate object detection is slow!

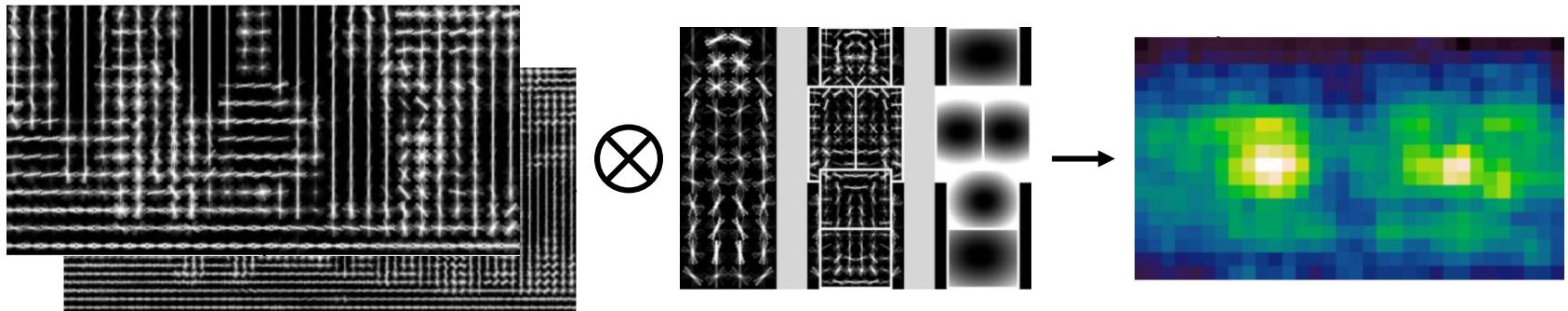
	Pascal 2007 mAP	Speed	
DPM v5	33.7	.07 FPS	14 s/img
R-CNN	66.0	.05 FPS	20 s/img
Fast R-CNN	70.0	.5 FPS	2 s/img
Faster R-CNN	73.2	7 FPS	140 ms/img
YOLO	<del>63.1</del> 69.0	45 FPS	22 ms/img



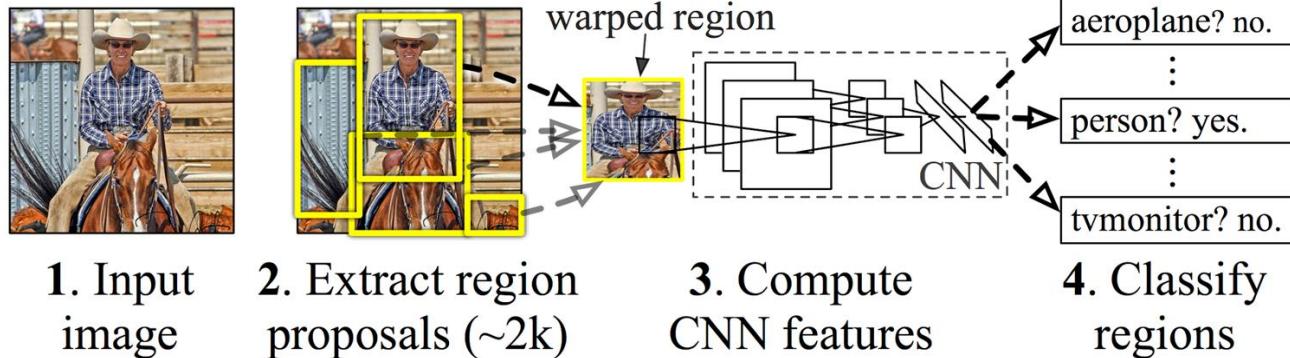
2 feet

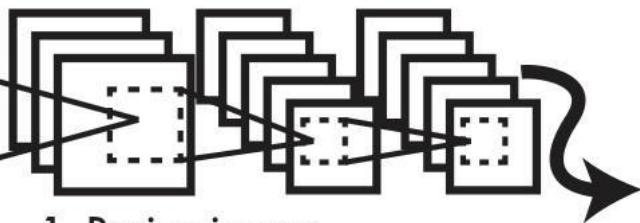
Sliding window, DPM, R-CNN all train region-based classifiers to perform detection

### DPM: *Deformable Part Models*

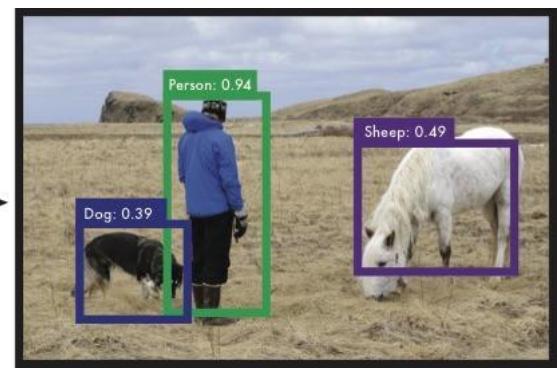


### R-CNN: *Regions with CNN features*



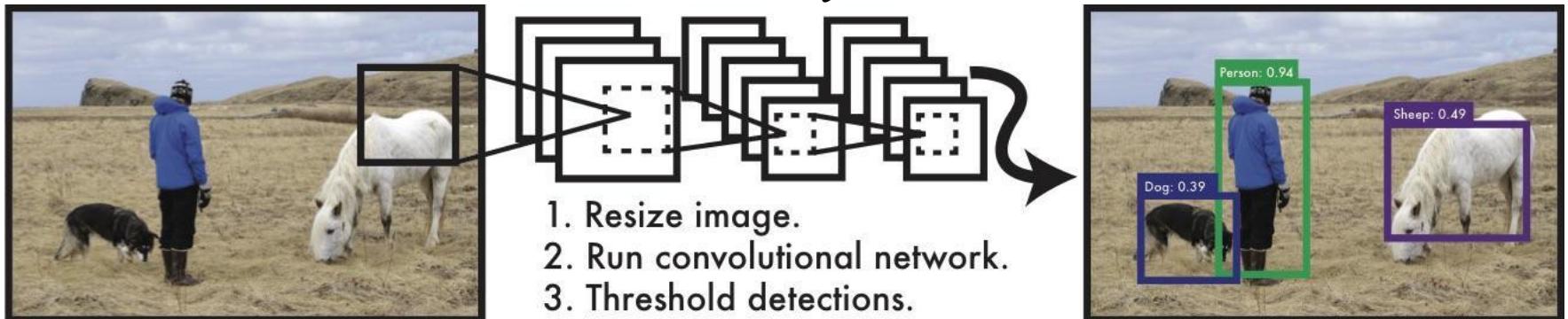


1. Resize image.
2. Run convolutional network.
3. Threshold detections.



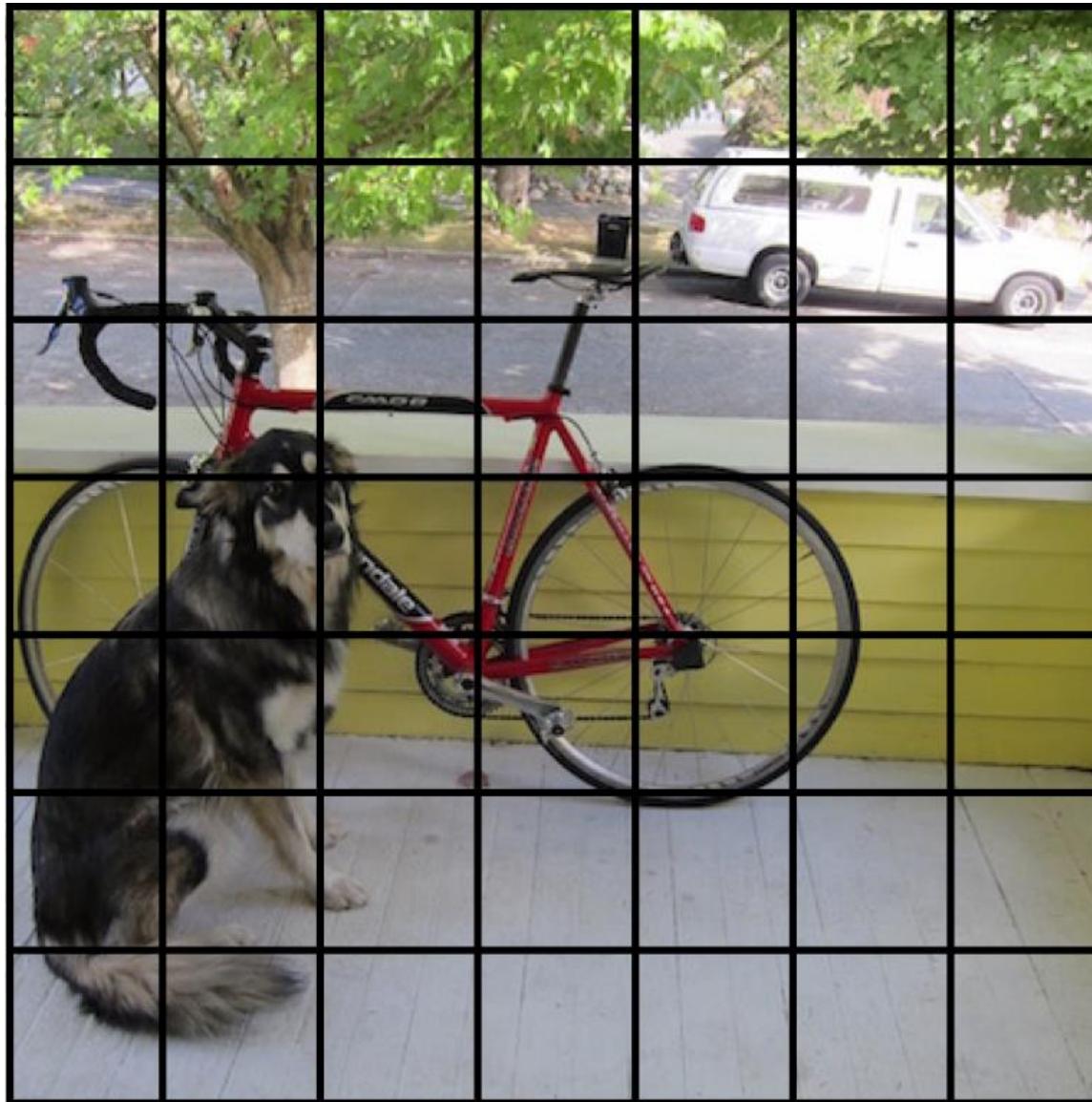
With YOLO, you only look once at an image to perform detection

### ***YOLO: You Only Look Once***

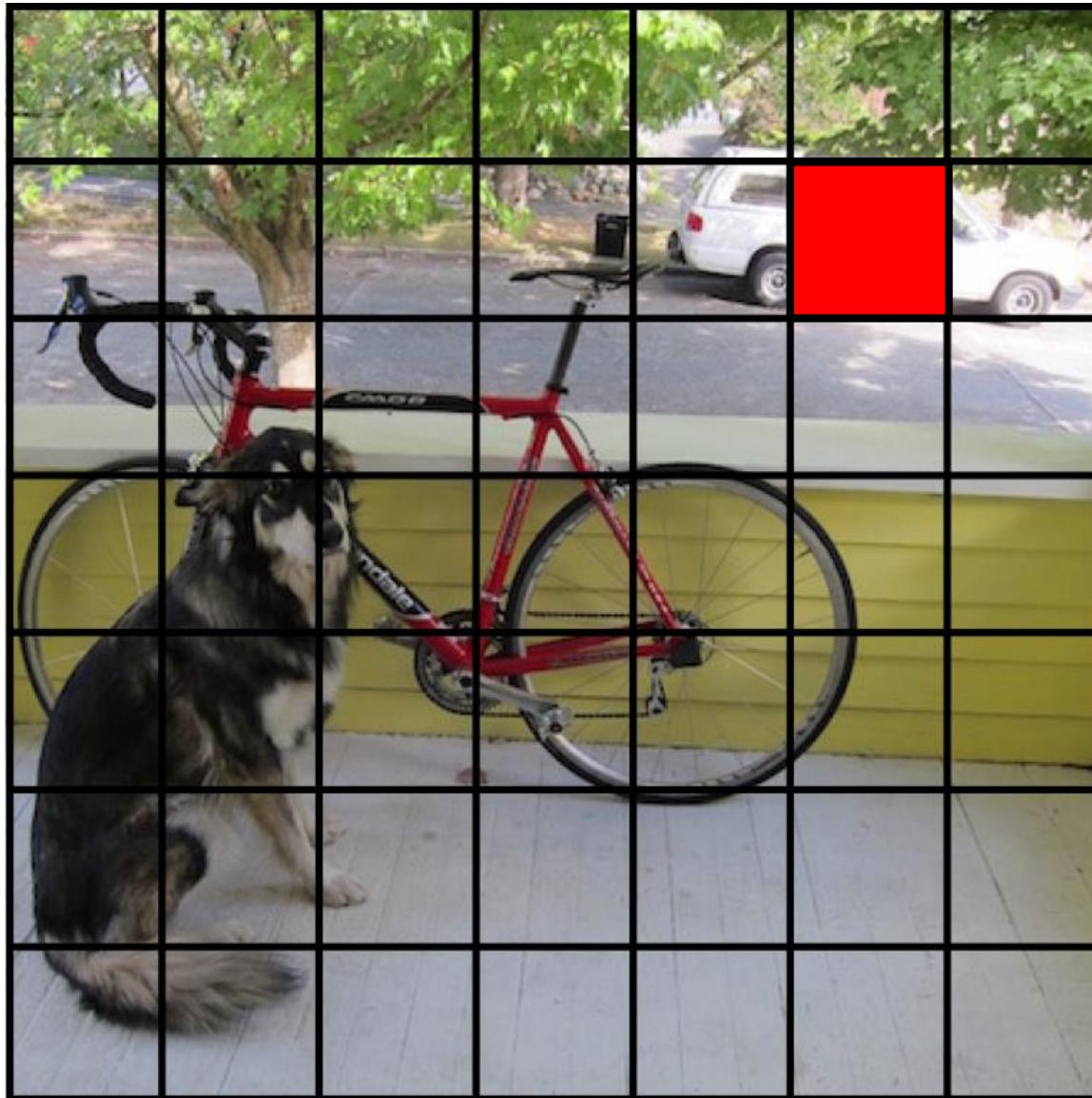




We split the image into a grid



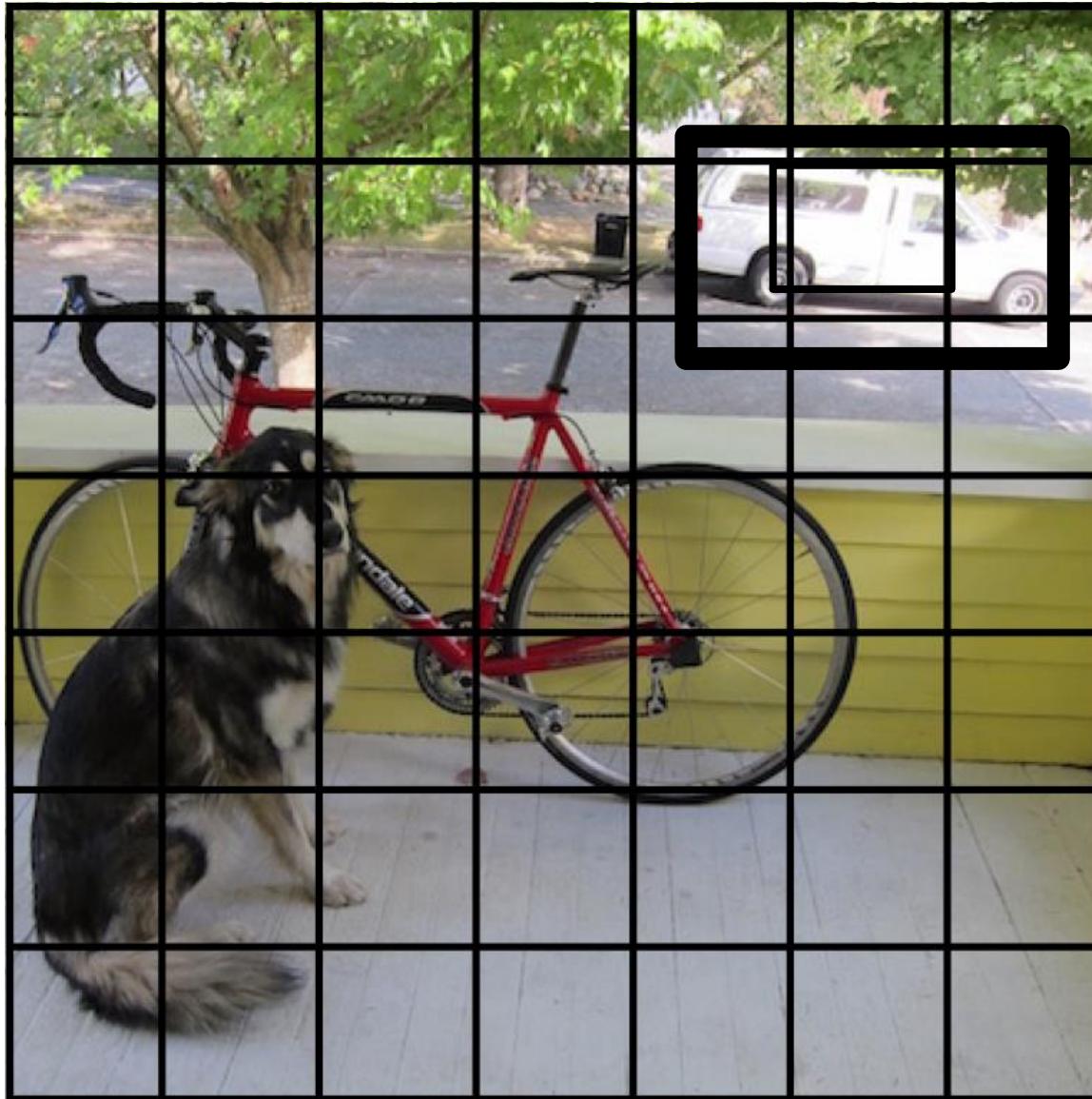
Each cell predicts boxes and confidences:  $P(\text{Object})$



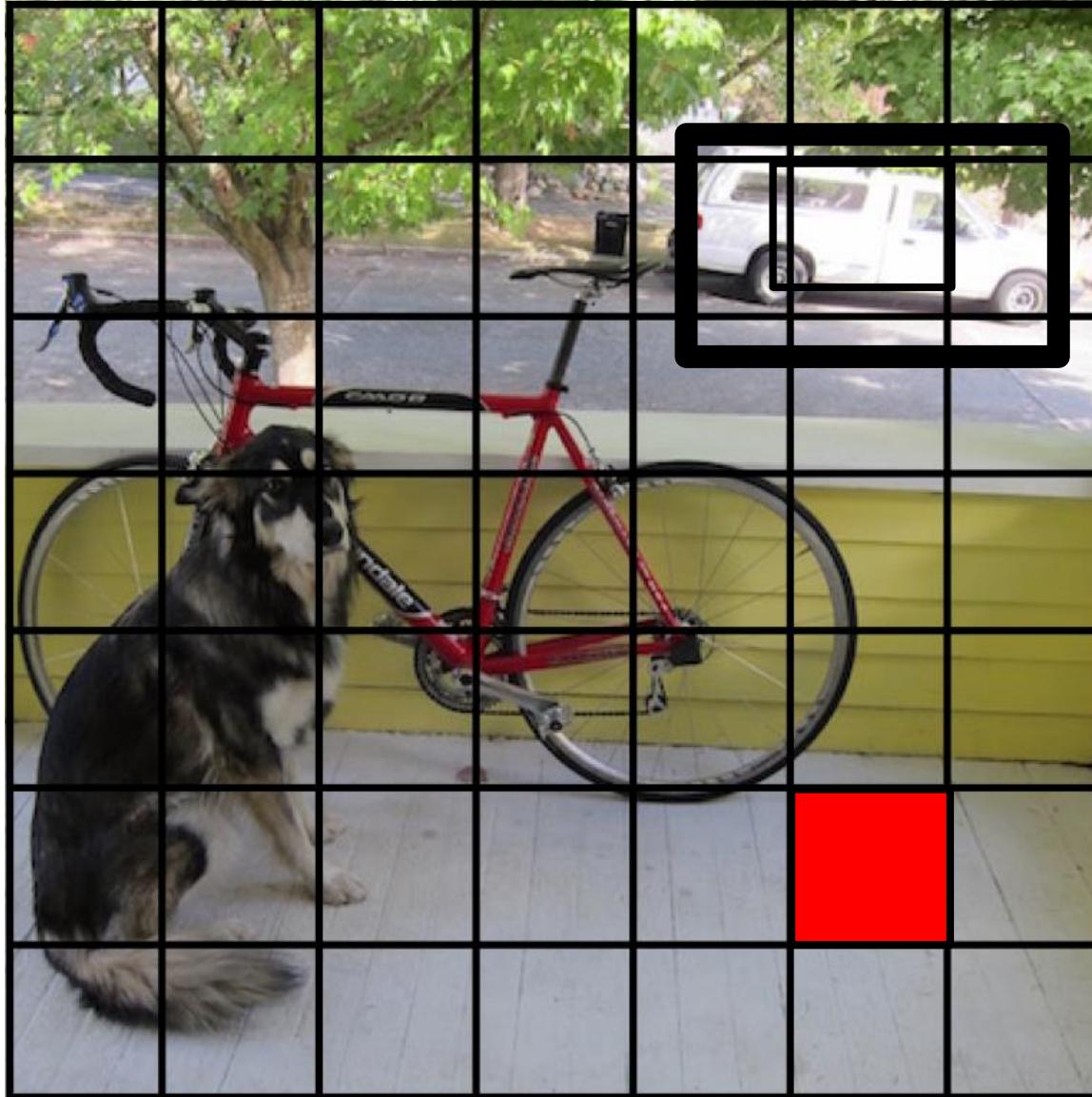
Each cell predicts boxes and confidences:  $P(\text{Object})$



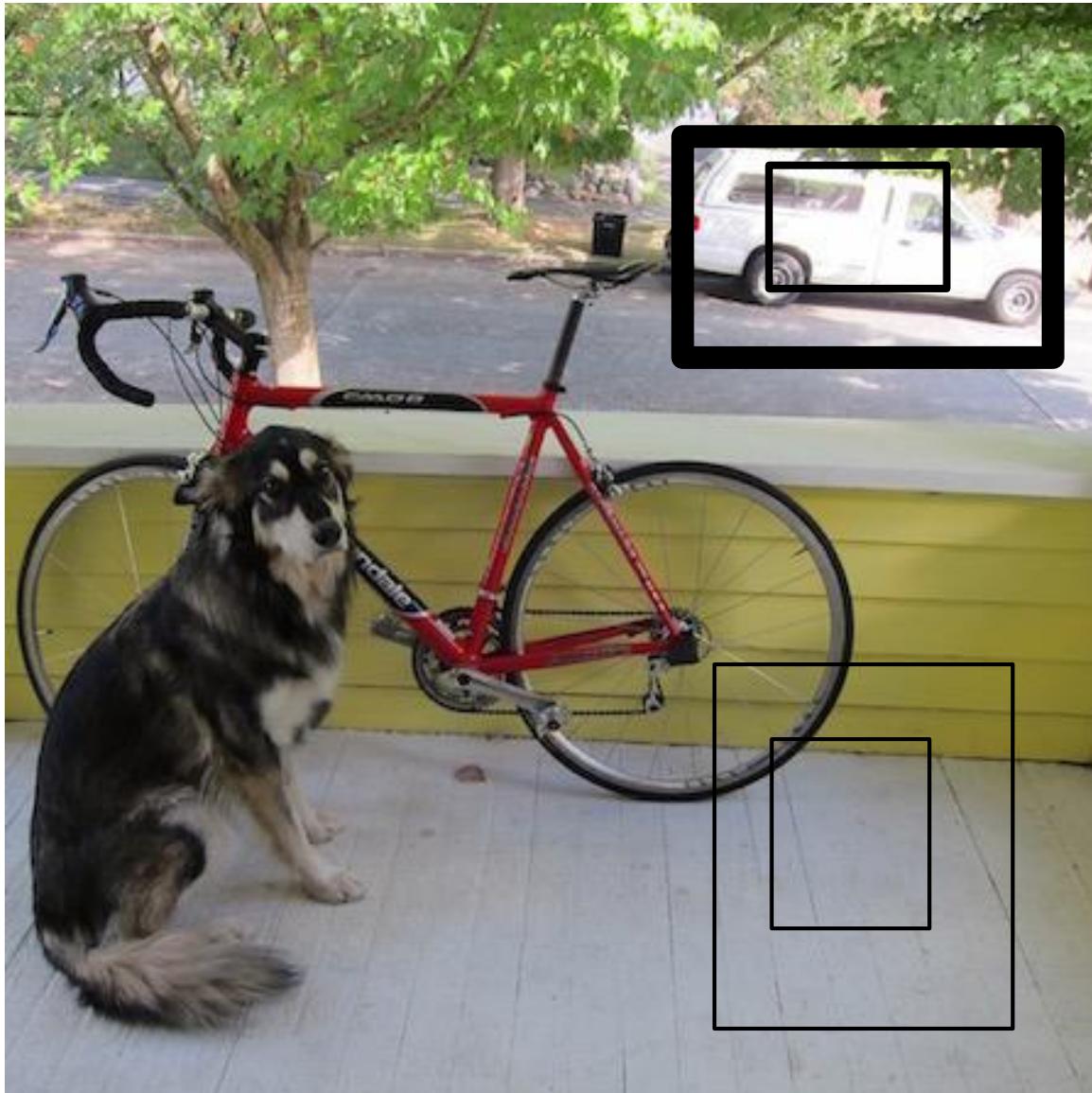
Each cell predicts boxes and confidences:  $P(\text{Object})$



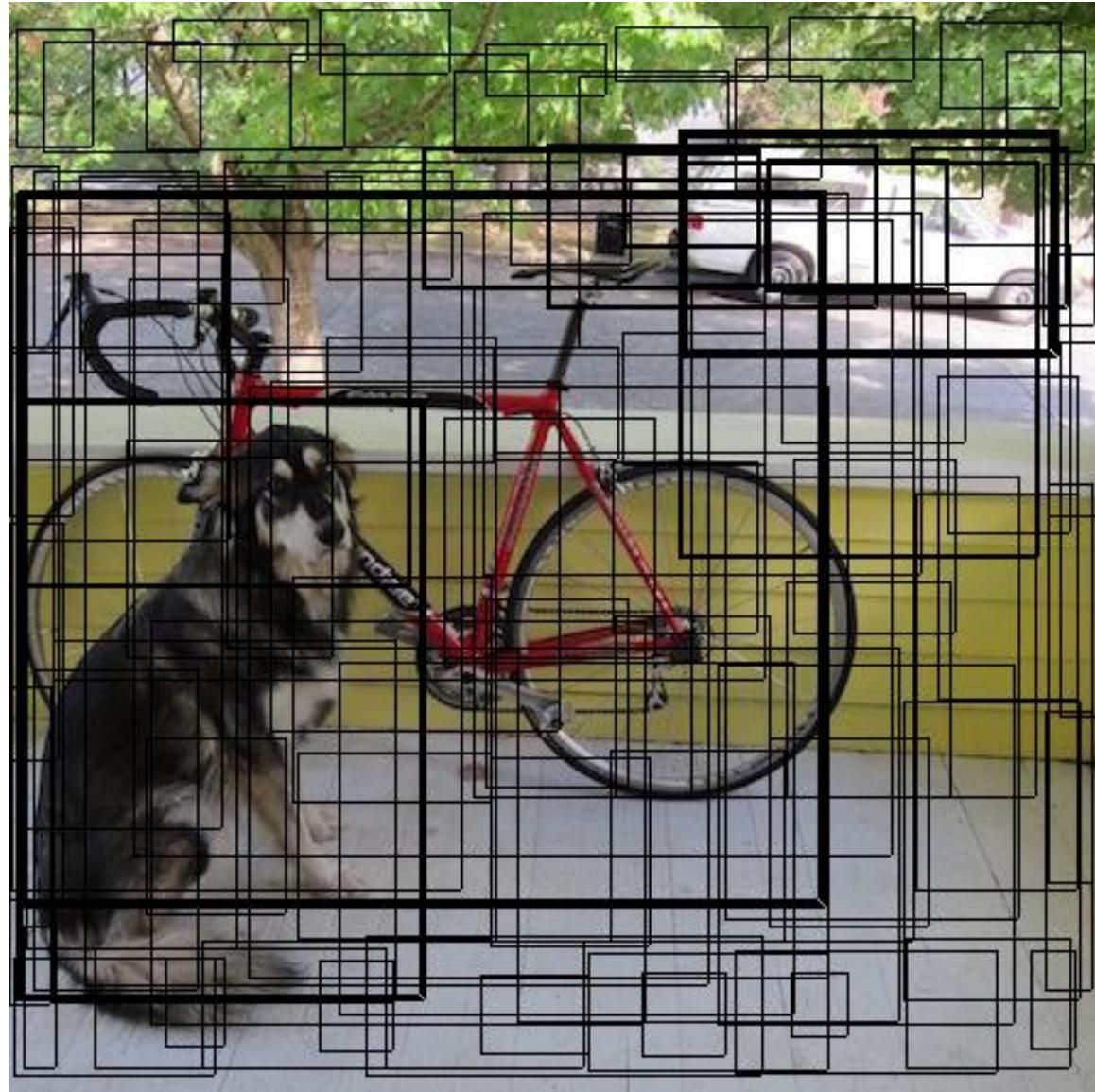
Each cell predicts boxes and confidences:  $P(\text{Object})$



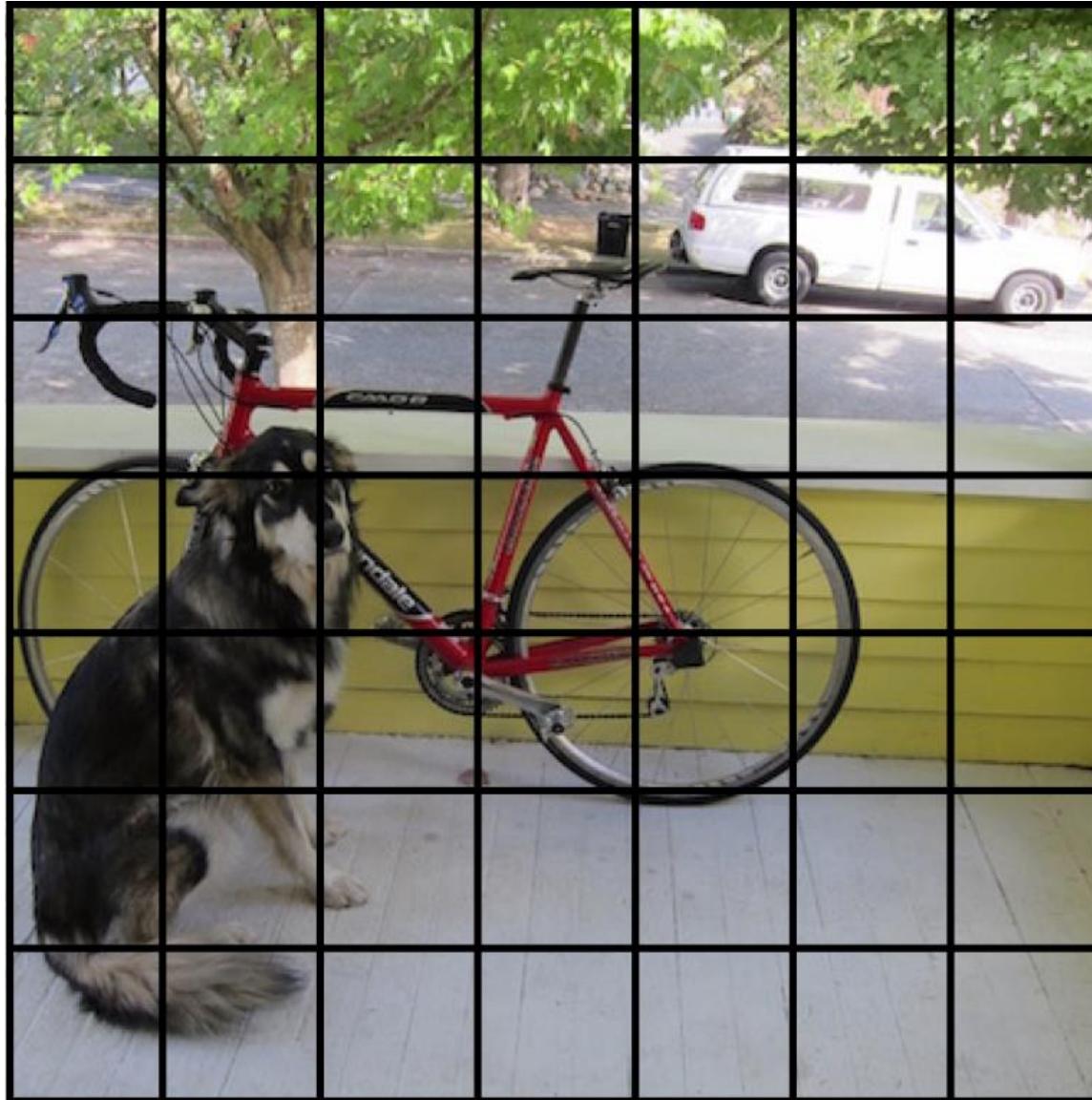
Each cell predicts boxes and confidences:  $P(\text{Object})$



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



Each cell also predicts a class probability.



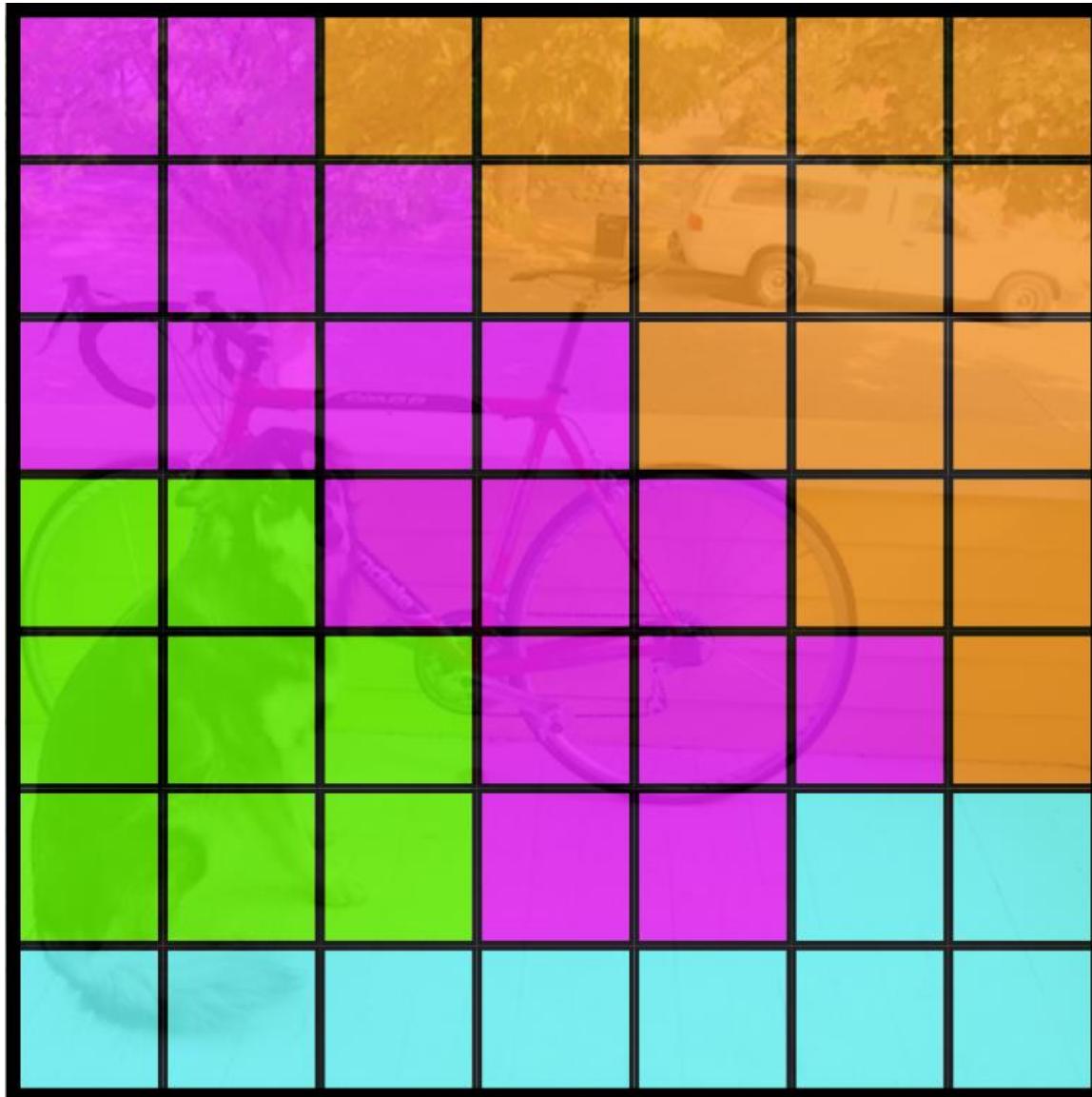
Each cell also predicts a class probability.

Bicycle

Car

Dog

Dining  
Table



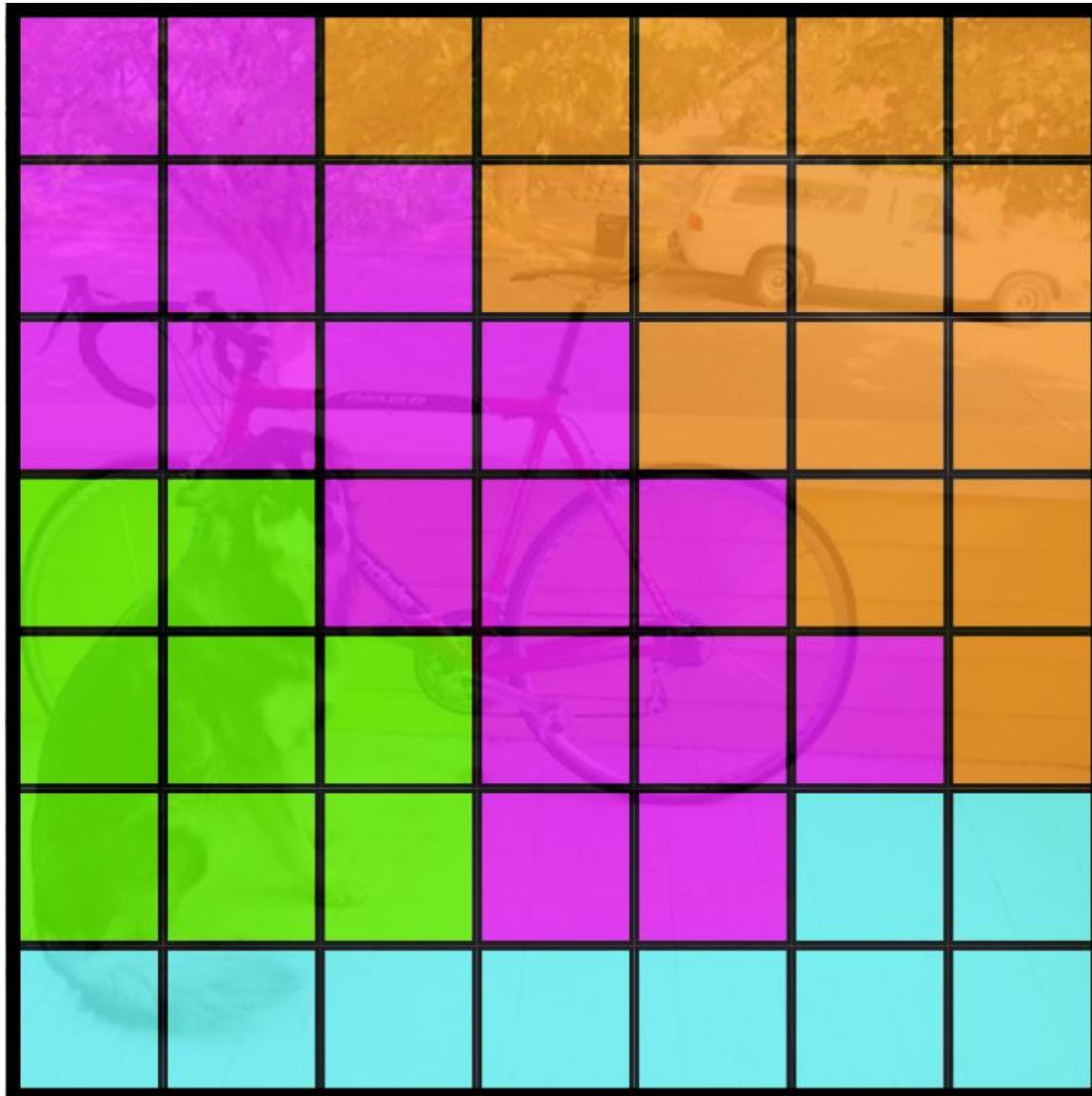
Conditioned on object:  $P(\text{Car} \mid \text{Object})$

Bicycle

Car

Dog

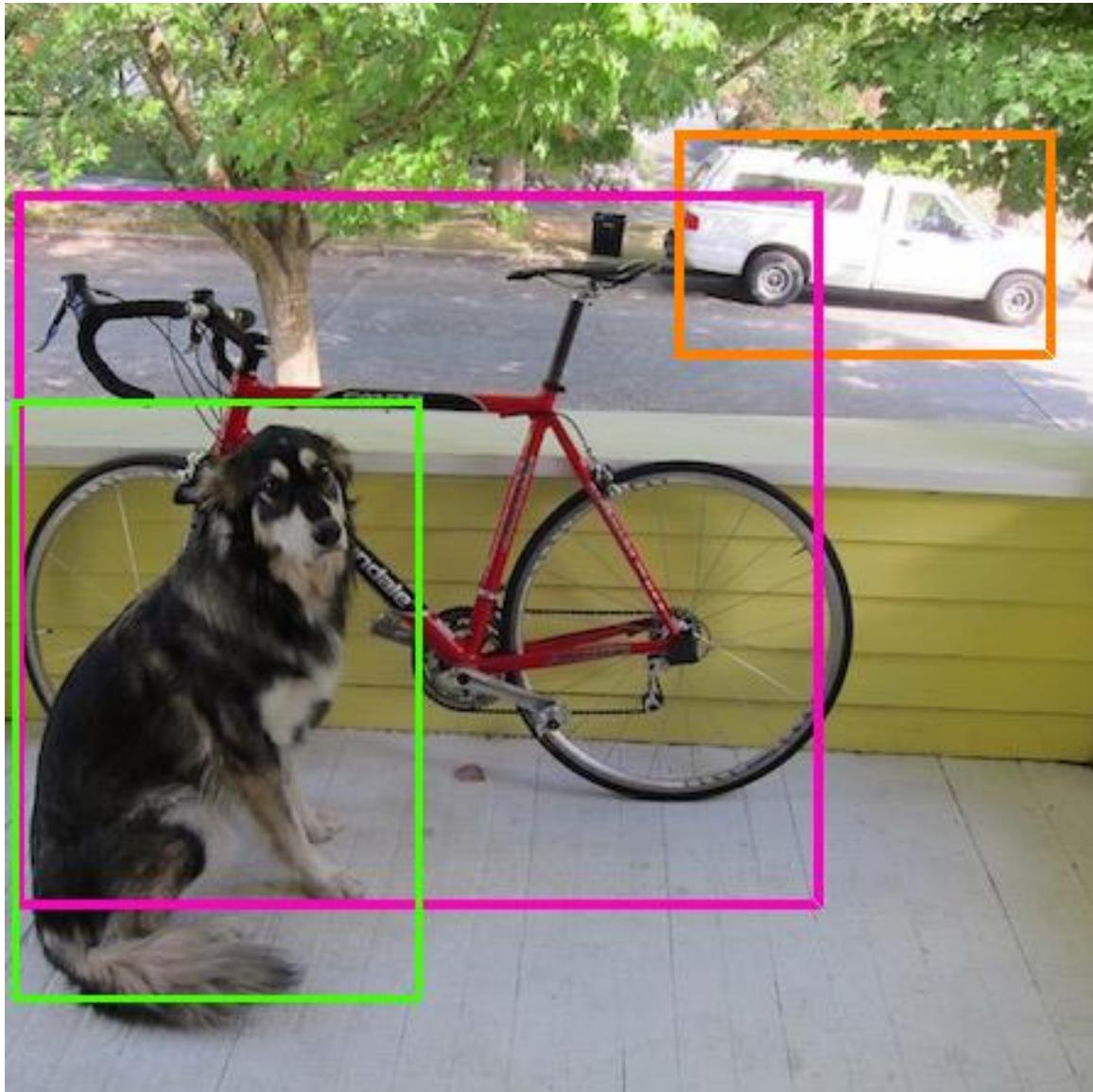
Dining  
Table



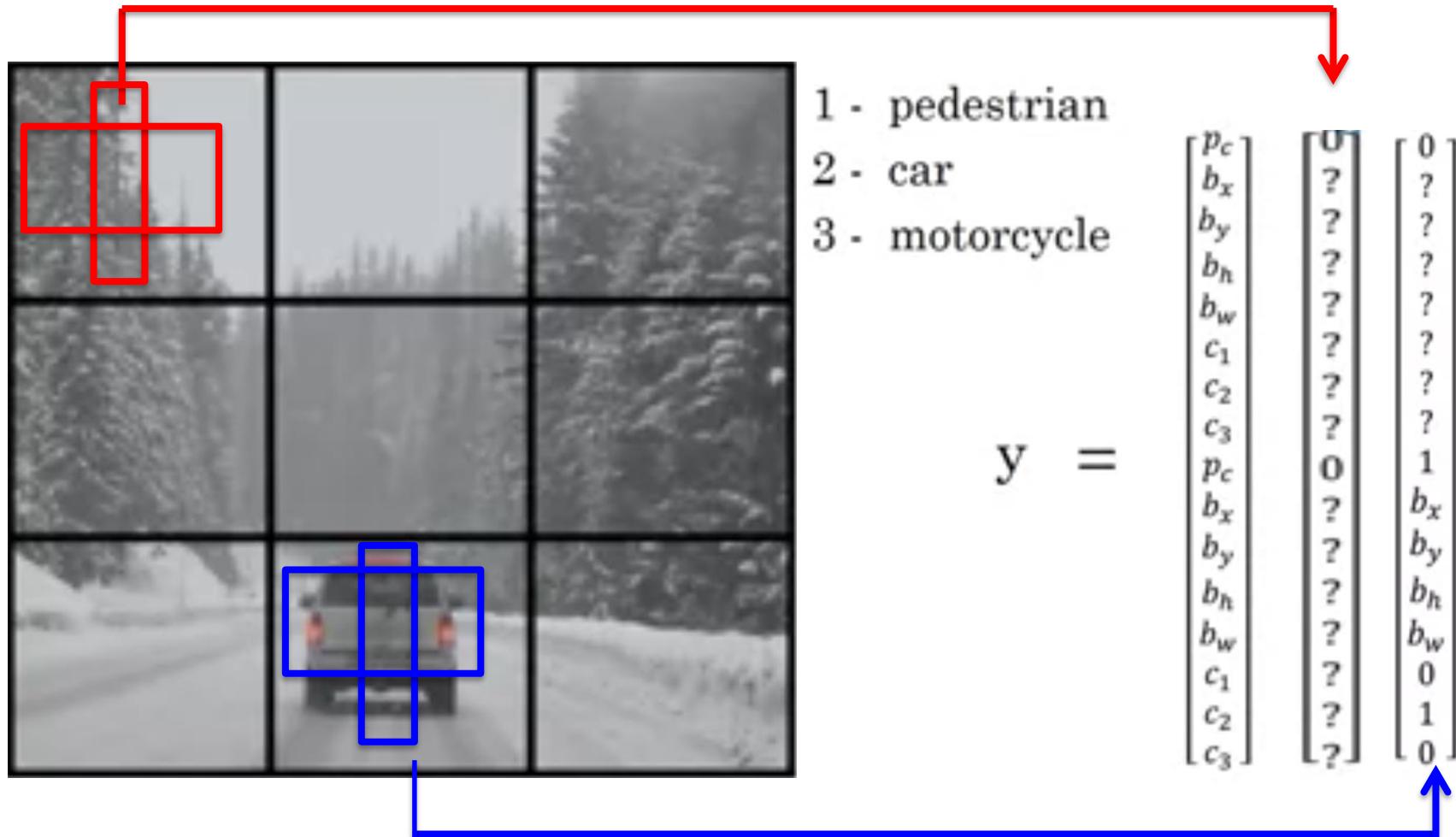
Then we combine the box and class predictions.



Finally we do NMS and threshold detections



# Yolo: CNN output for 2 bounding boxes (similar to anchor boxes in RPN)

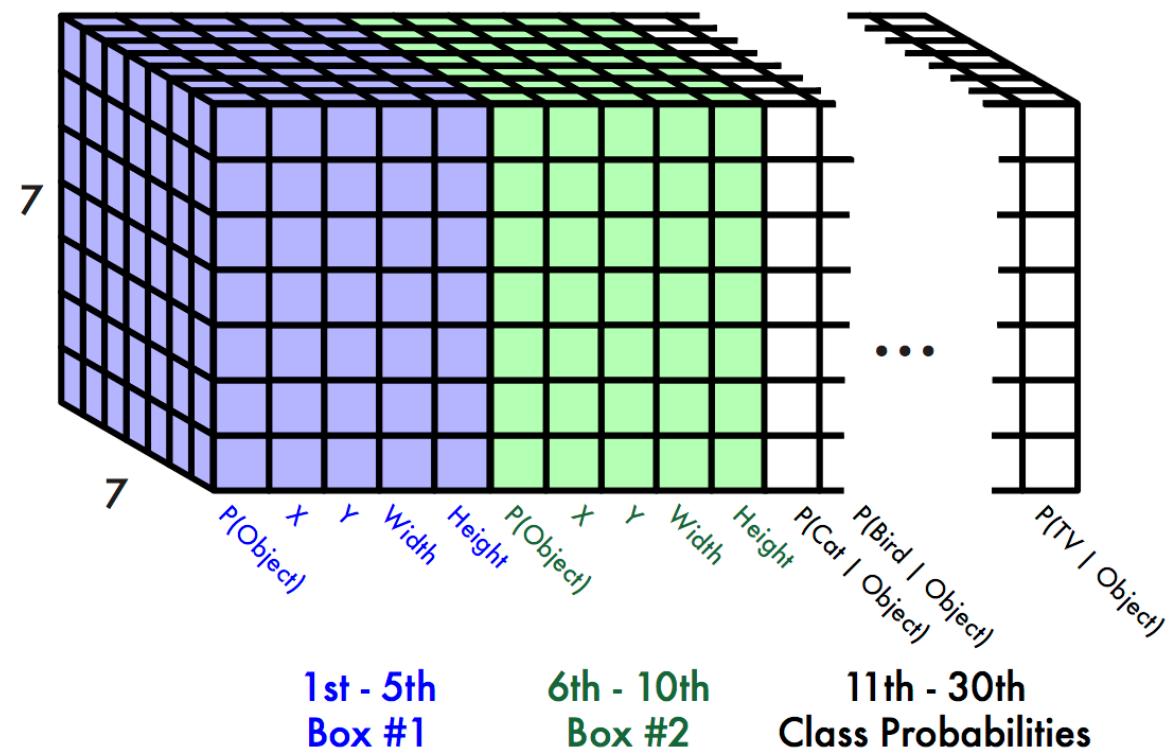


y size: 3 (#grid rows) x 3 (#grid columns) x 2 (#bounding boxes) x 8 (confidence 1, coord 4, class 3)

# This parameterization fixes the output size

Each cell predicts:

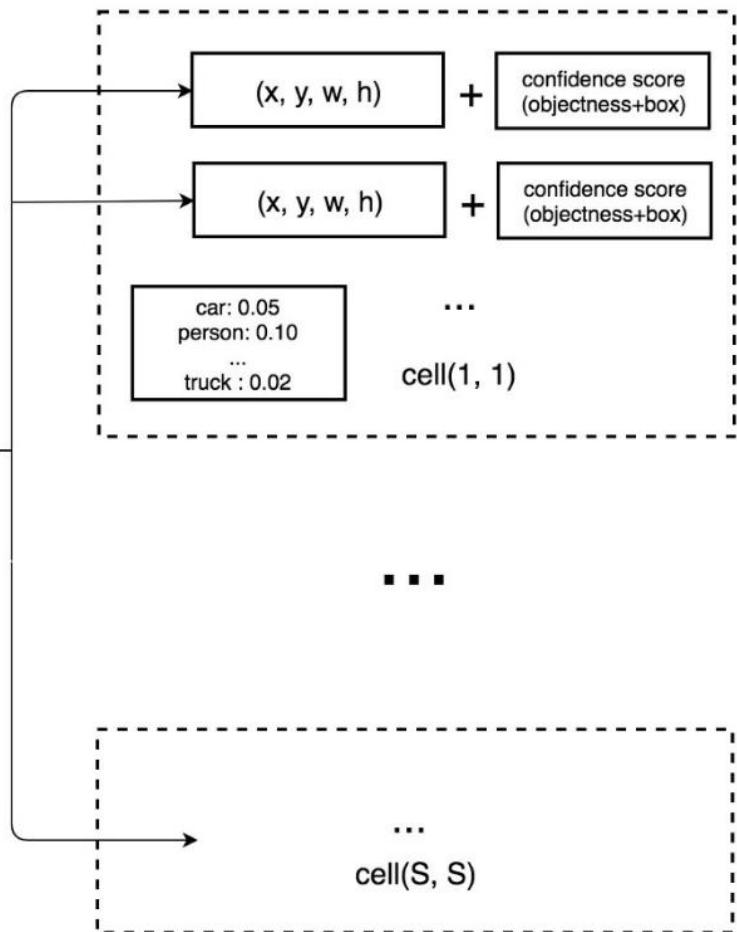
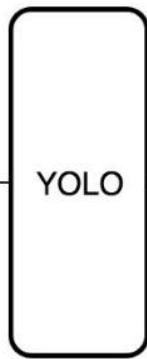
- For each bounding box:
  - 4 coordinates (x, y, w, h)
  - 1 confidence value
- Some number of class probabilities



For Pascal VOC:

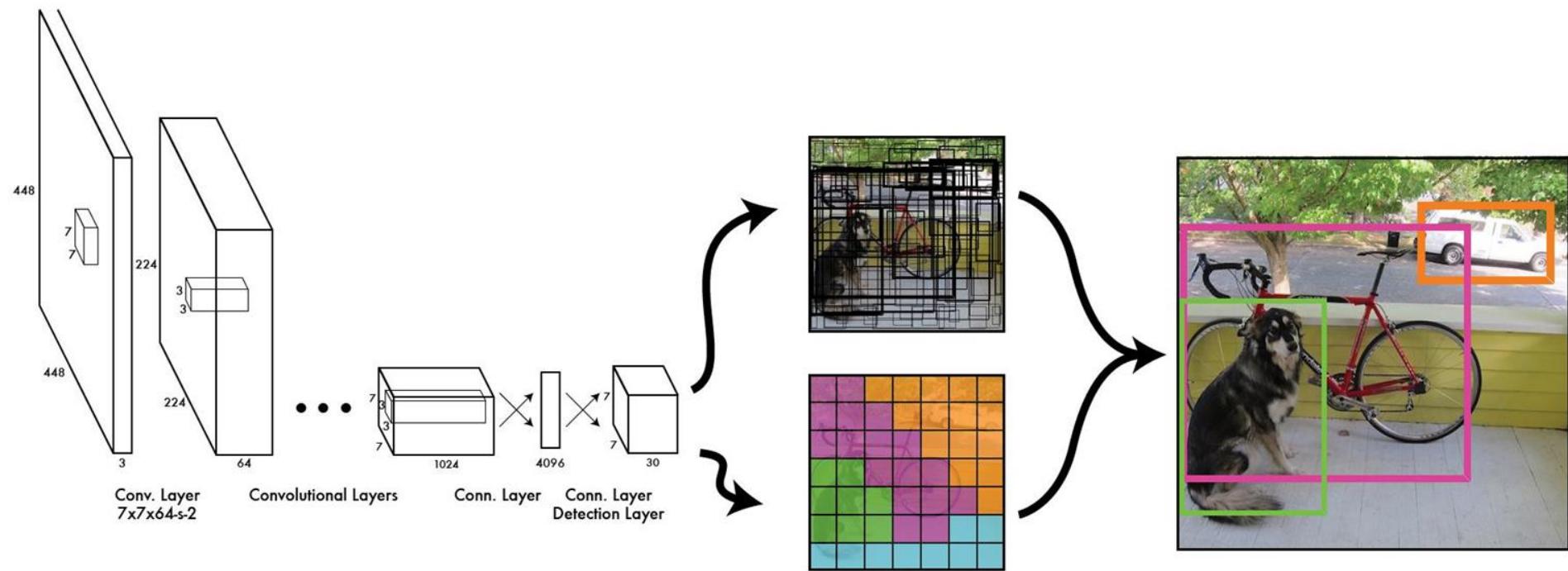
- 7x7 grid
- 2 bounding boxes / cell
- 20 classes

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

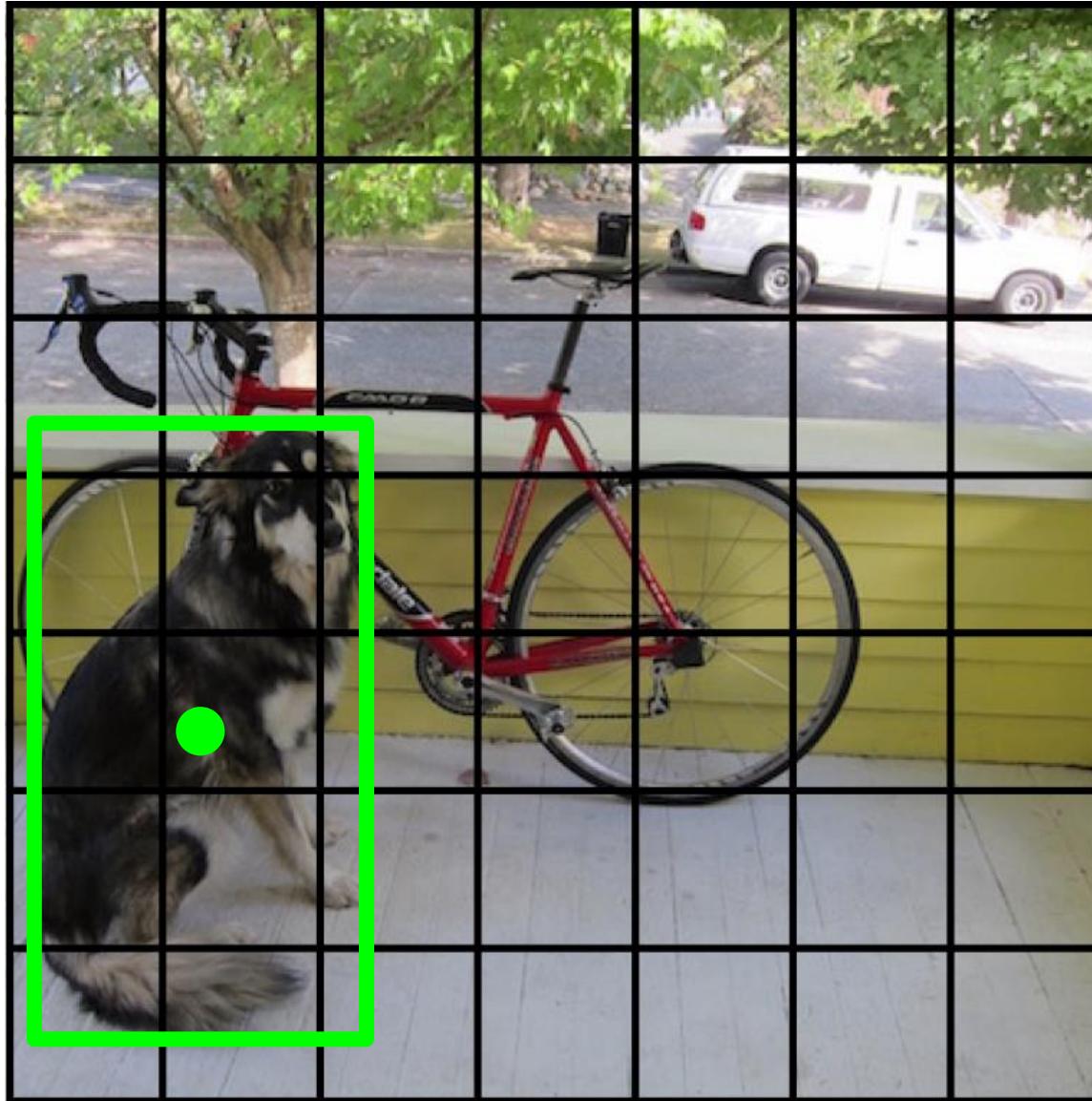


$$\begin{aligned}
& \lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} \left[ (x_i - \hat{x}_i)^2 + (y_i - \hat{y}_i)^2 \right] \\
& + \lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} \left[ \left( \sqrt{w_i} - \sqrt{\hat{w}_i} \right)^2 + \left( \sqrt{h_i} - \sqrt{\hat{h}_i} \right)^2 \right] \\
& + \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} \left( C_i - \hat{C}_i \right)^2 \\
& + \lambda_{\text{noobj}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{noobj}} \left( C_i - \hat{C}_i \right)^2 \\
& + \sum_{i=0}^{S^2} \mathbb{1}_i^{\text{obj}} \sum_{c \in \text{classes}} (p_i(c) - \hat{p}_i(c))^2
\end{aligned}$$

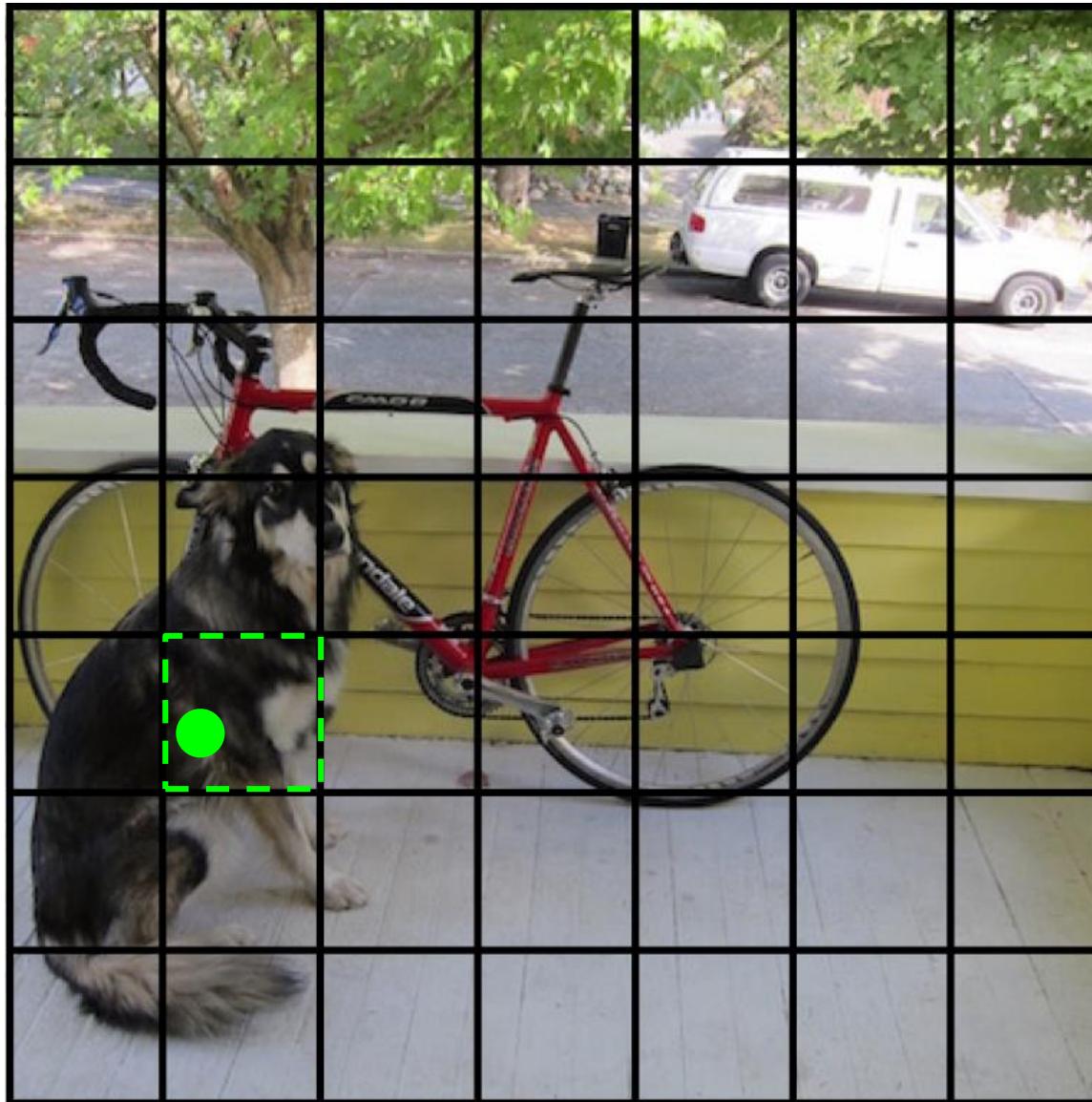
Thus we can train one neural network to be a whole detection pipeline



During training, match example to the right cell



During training, match example to the right cell



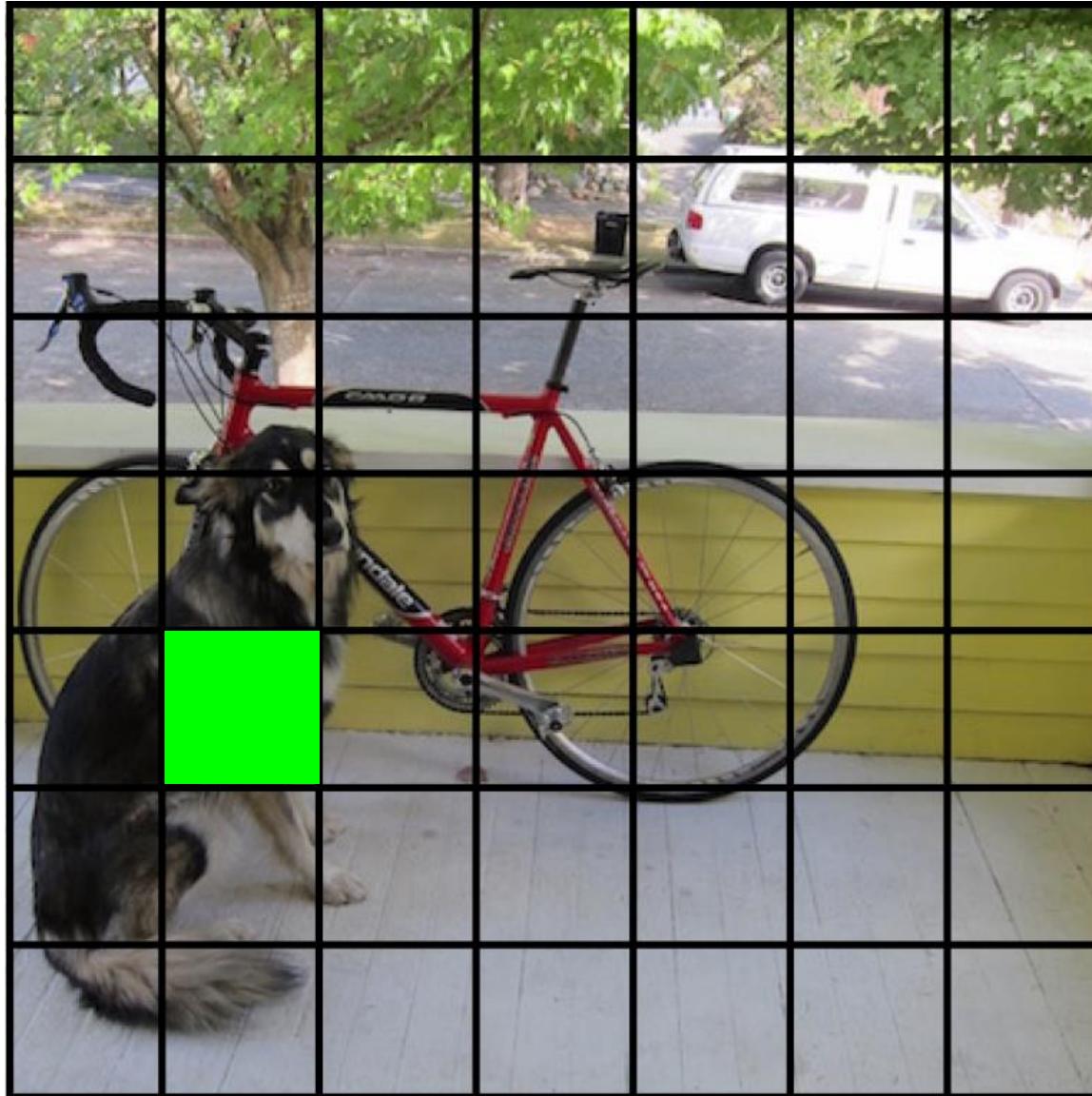
Adjust that cell's class prediction

Dog = 1

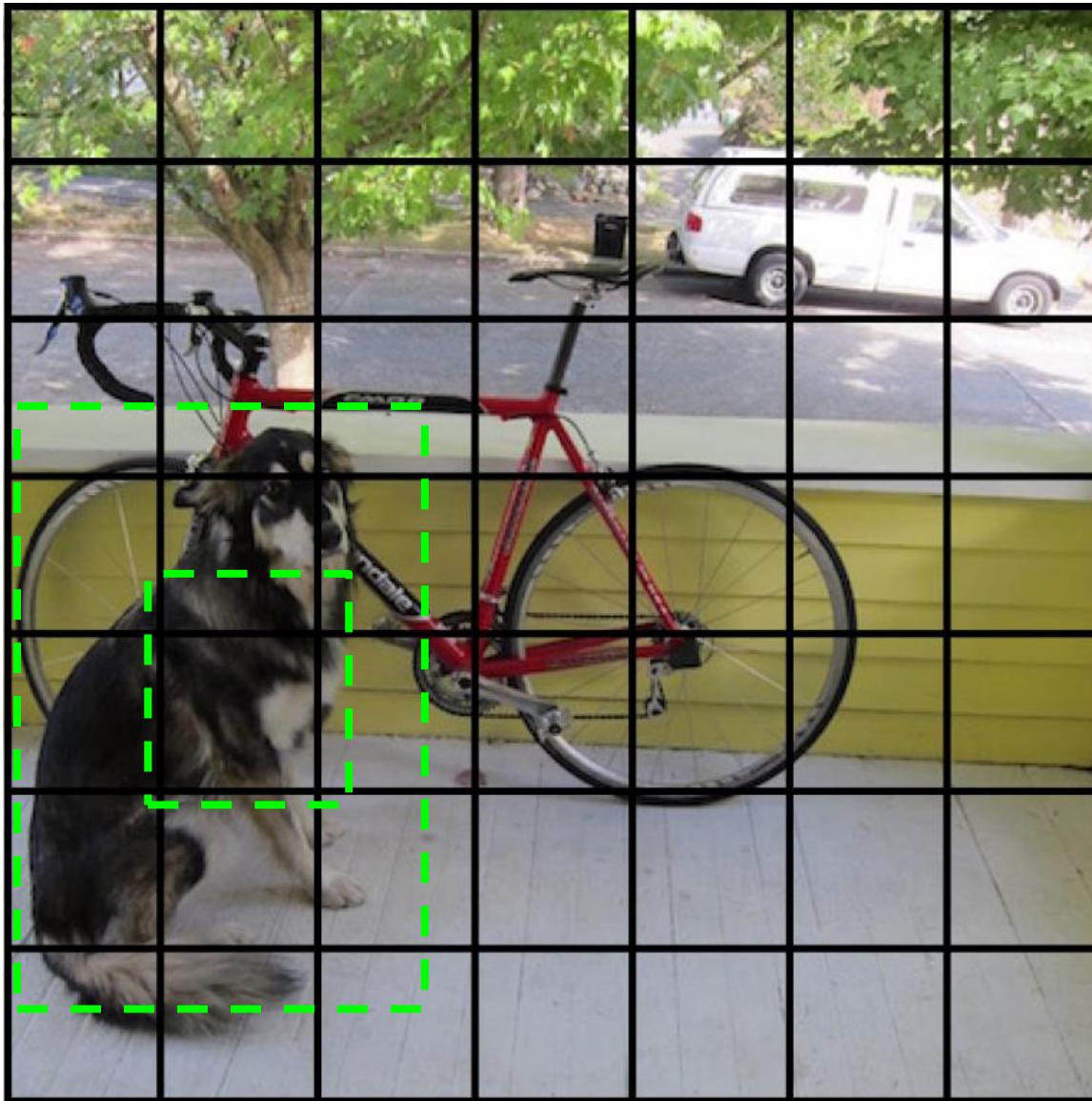
Cat = 0

Bike = 0

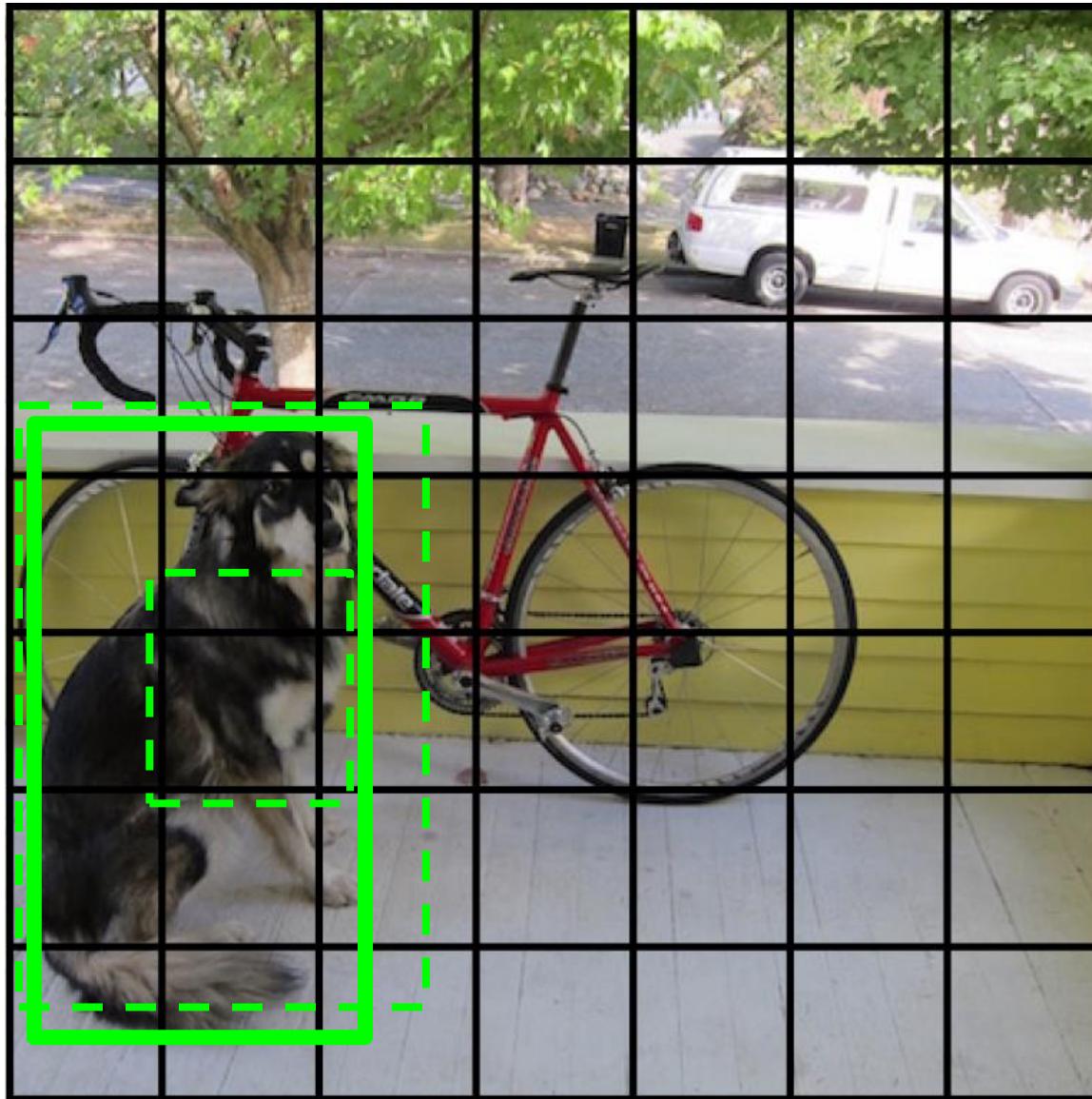
...



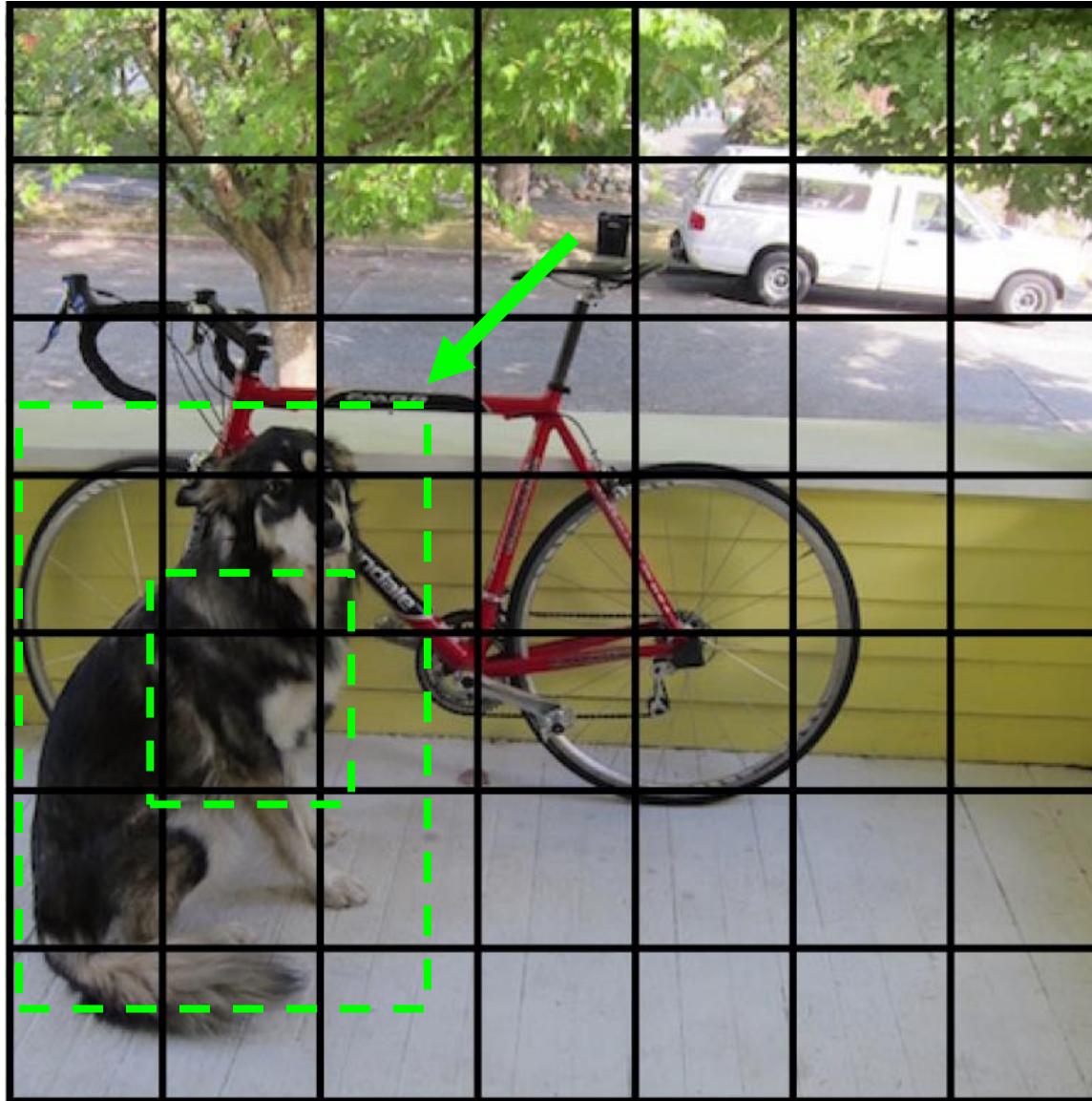
Look at that cell's predicted boxes



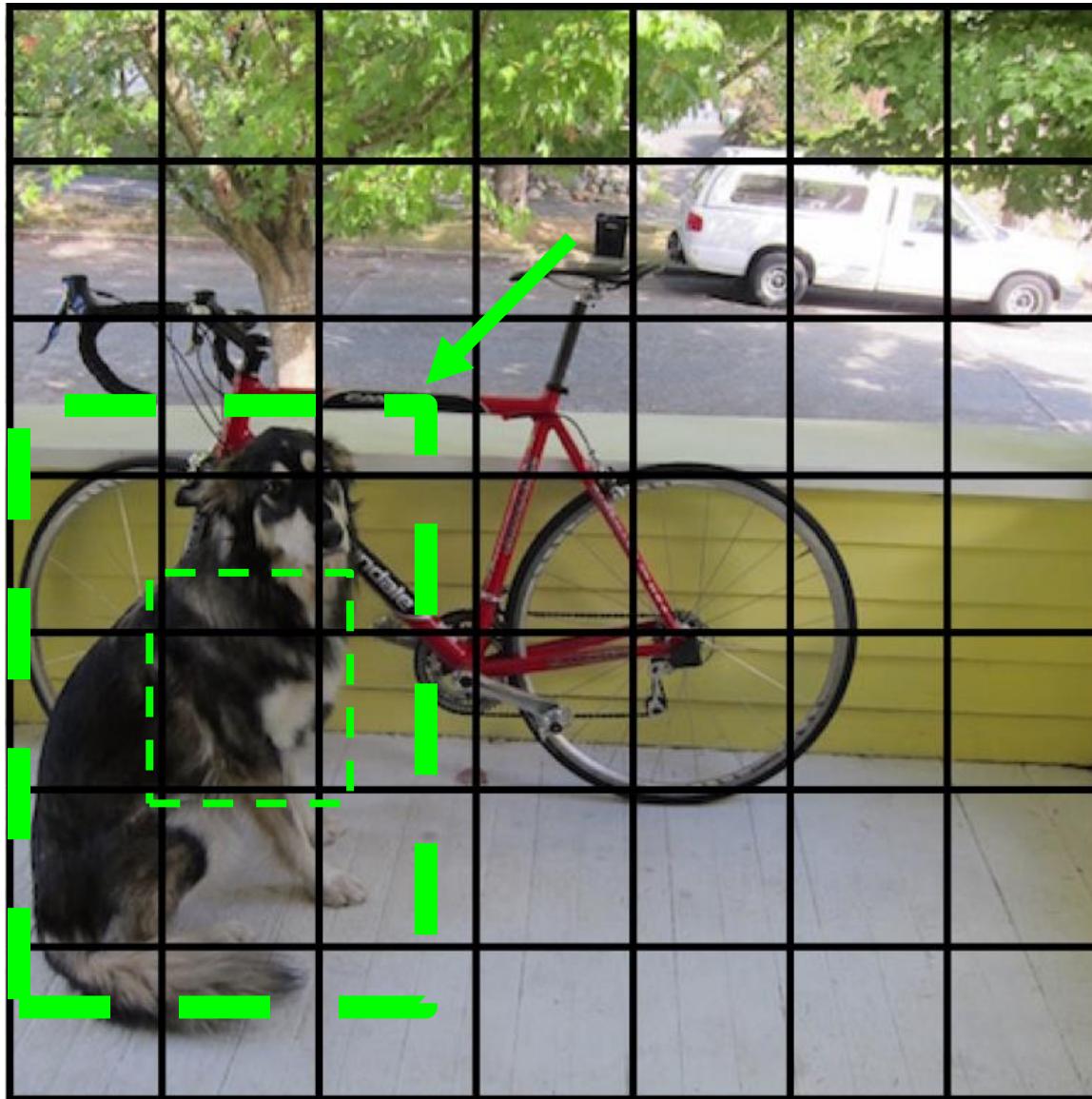
Find the best one, adjust it, increase the confidence



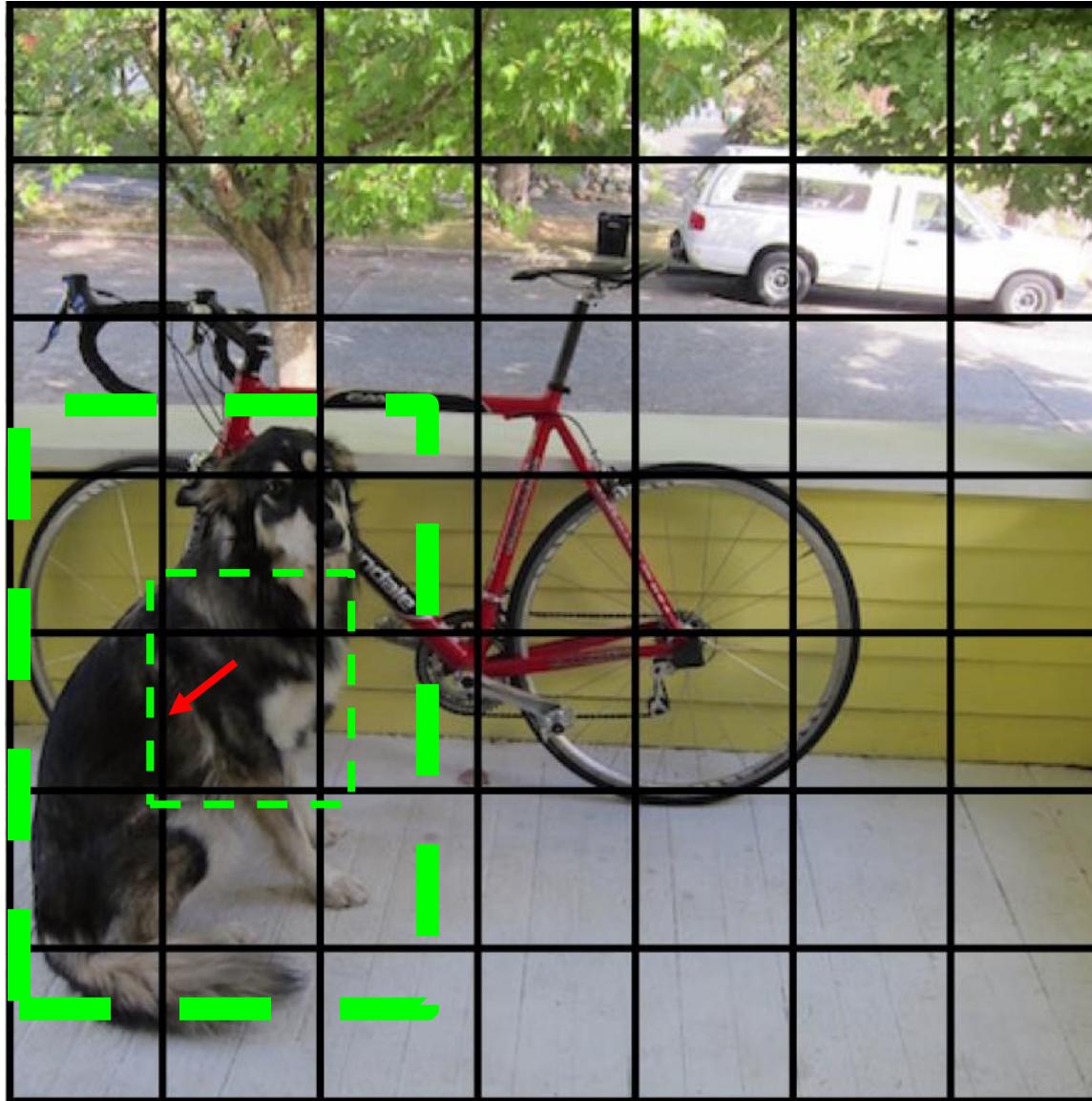
Find the best one, adjust it, increase the confidence



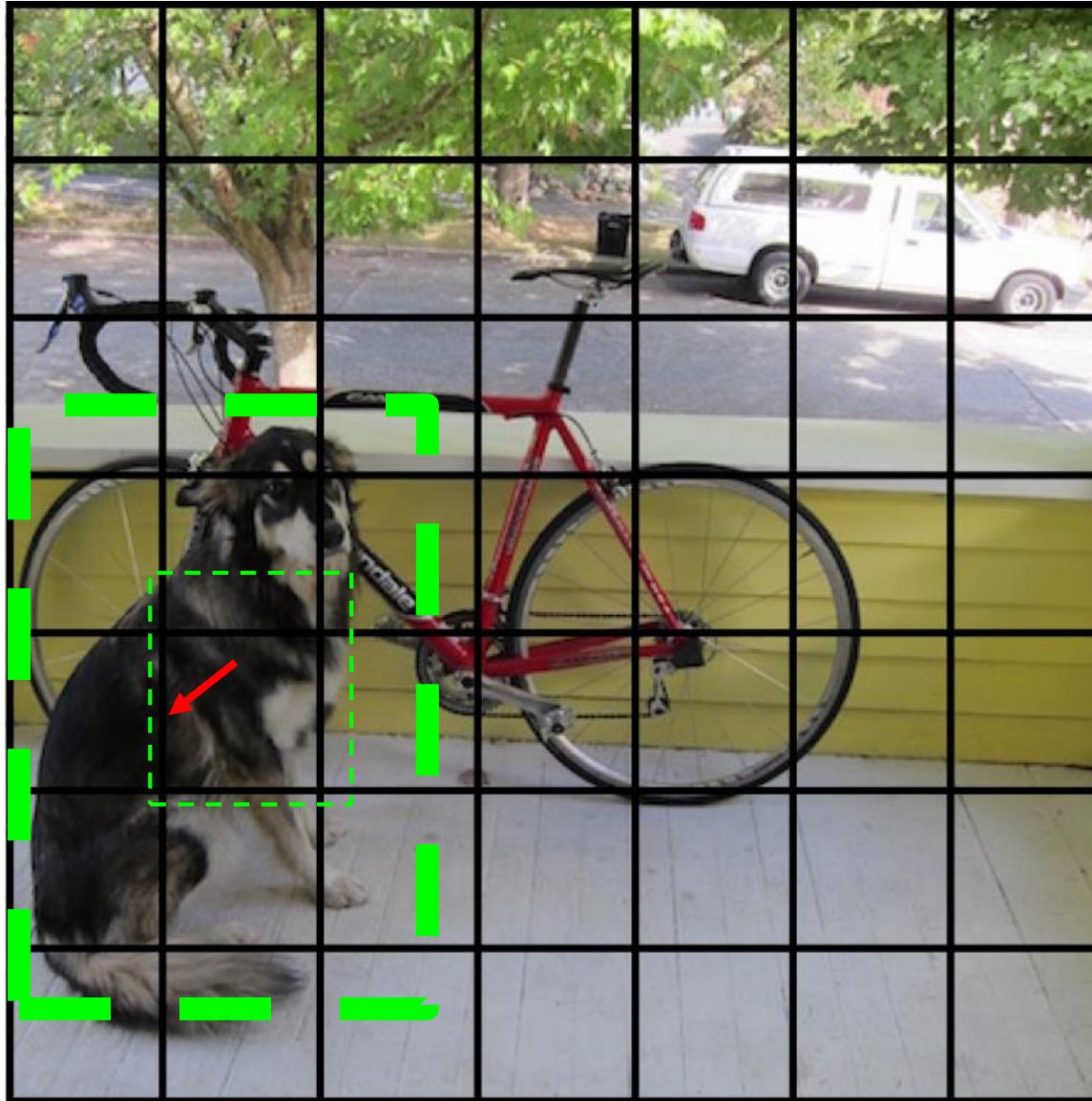
Find the best one, adjust it, increase the confidence



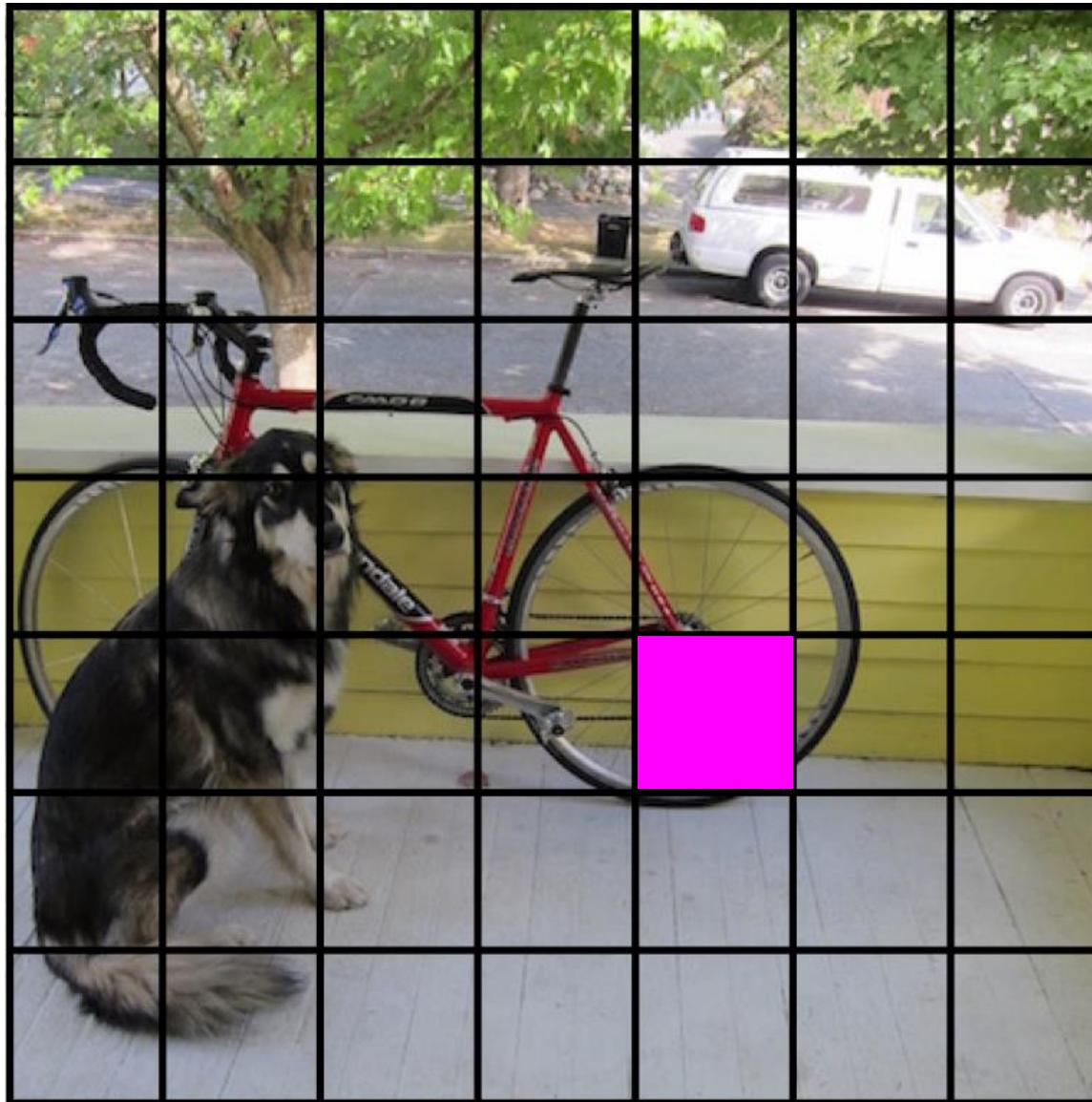
Decrease the confidence of other boxes



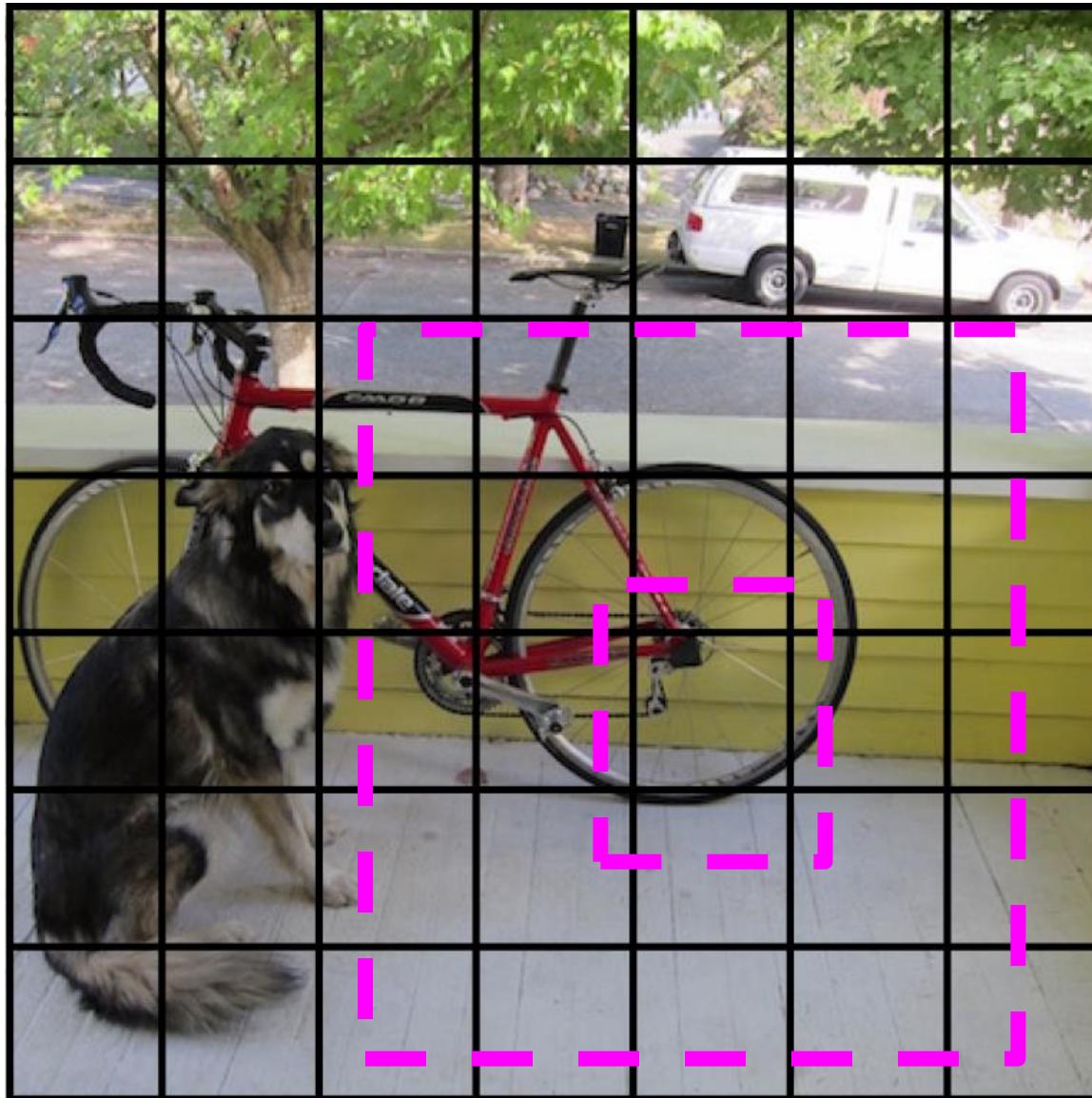
Decrease the confidence of other boxes



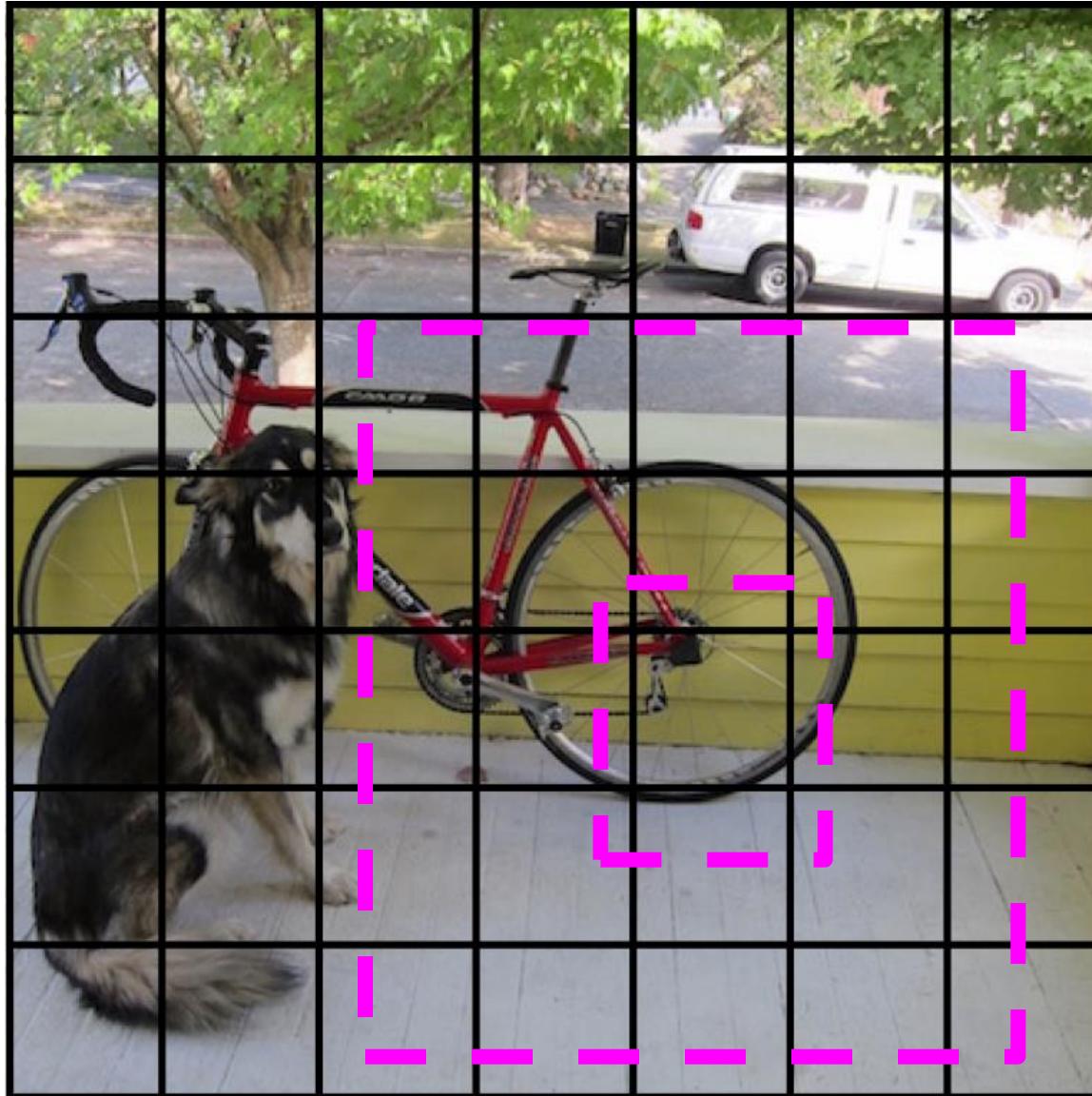
Some cells don't have any ground truth detections!



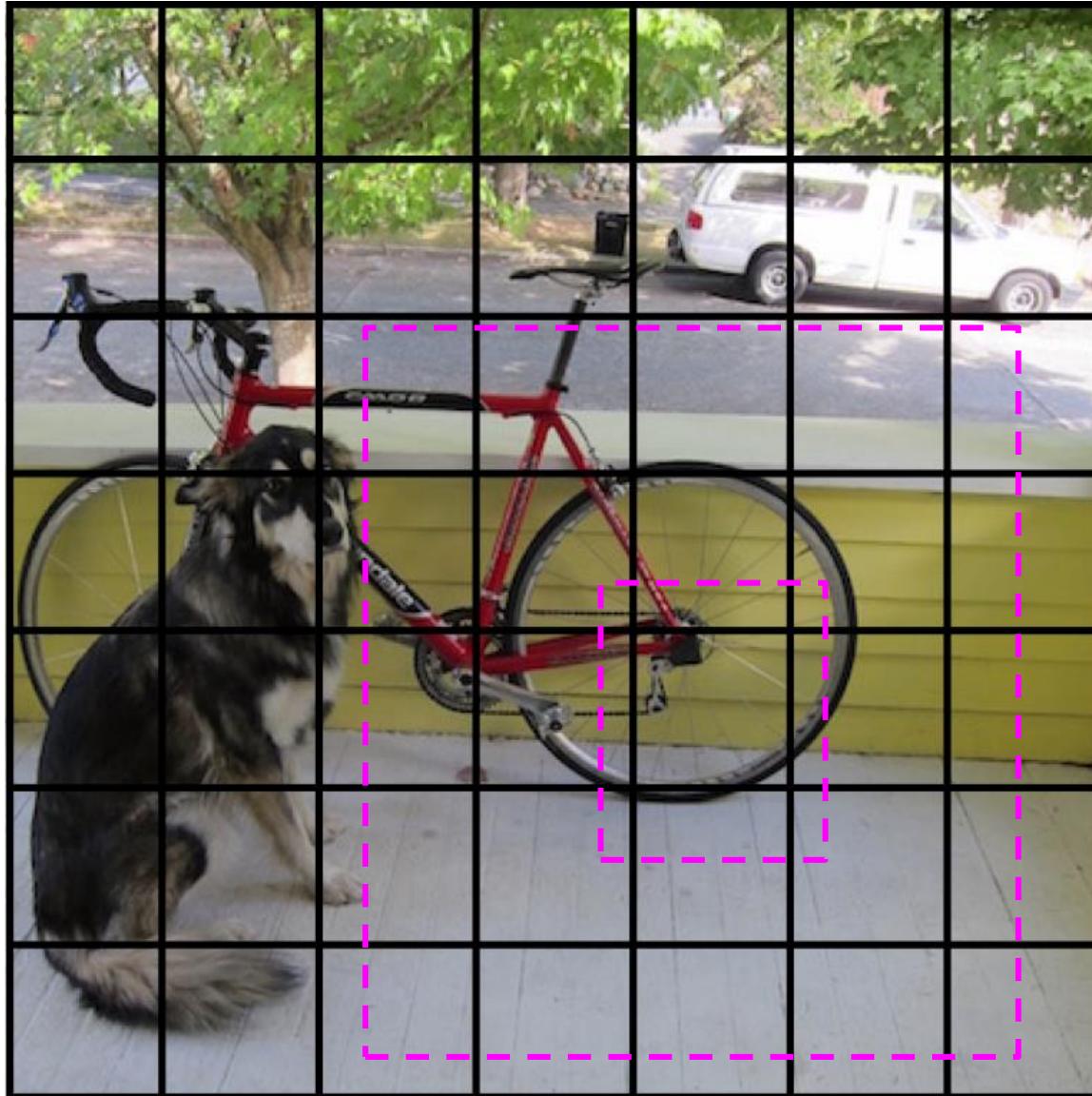
Some cells don't have any ground truth detections!



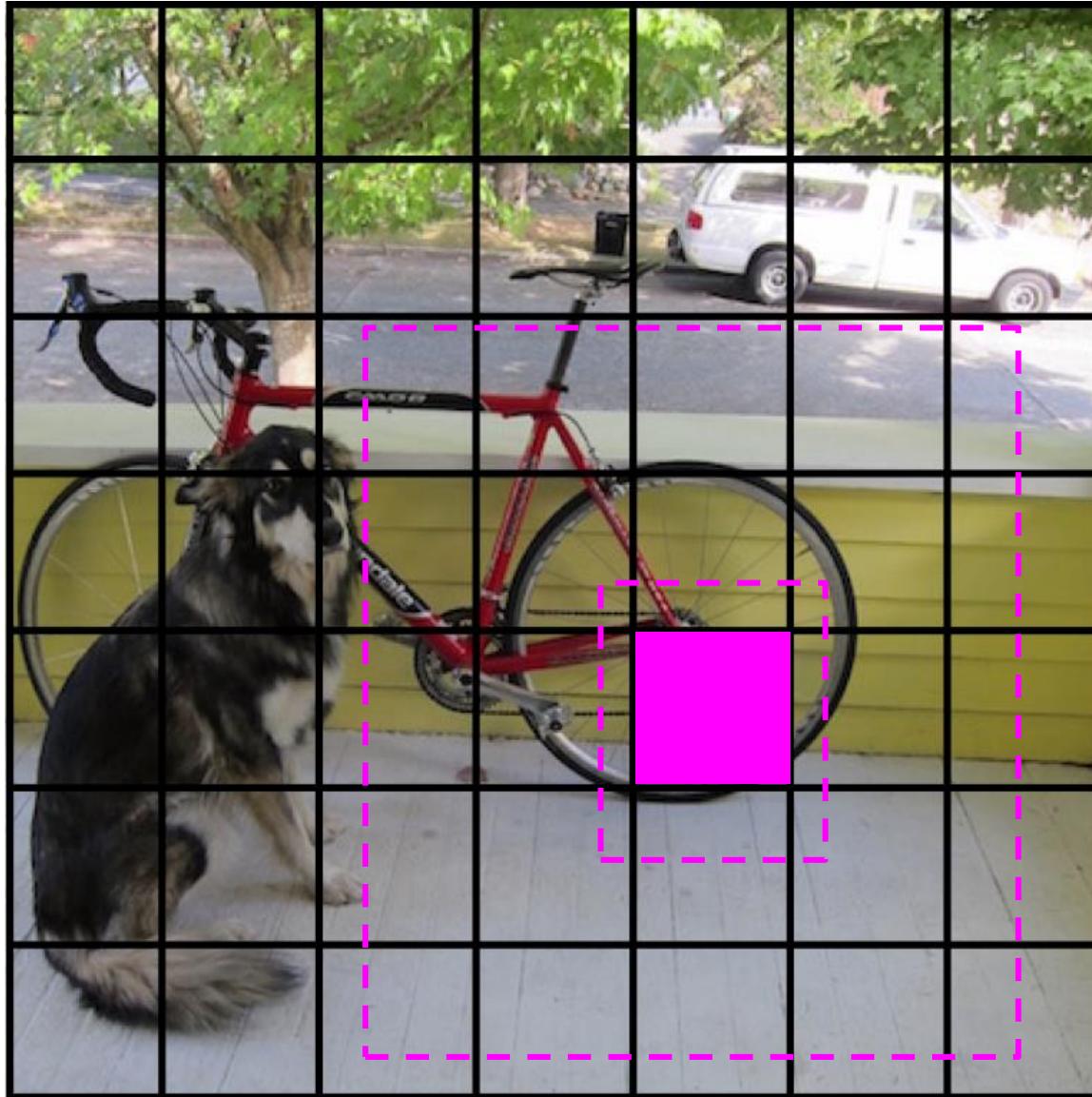
Decrease the confidence of these boxes



Decrease the confidence of these boxes

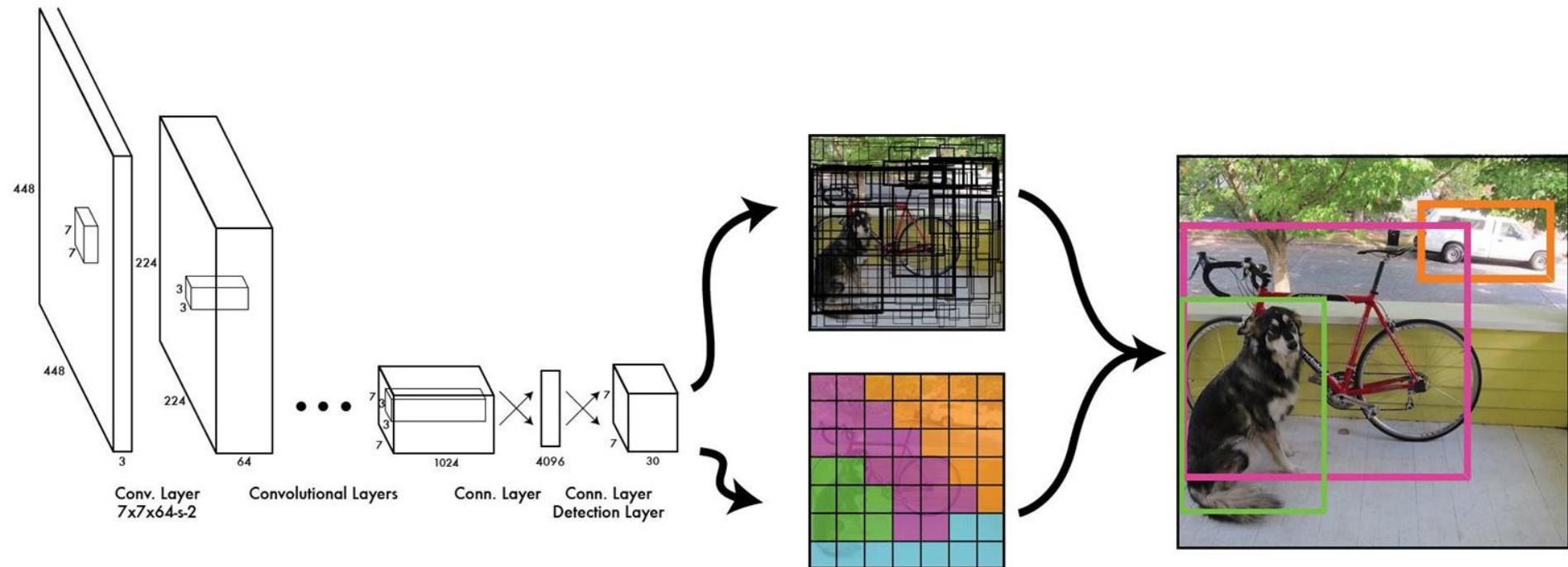


Do not adjust the class probabilities or coordinates

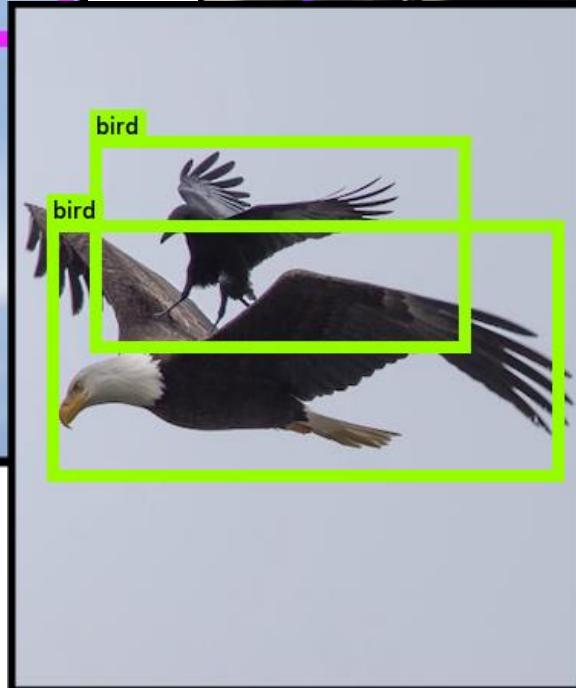
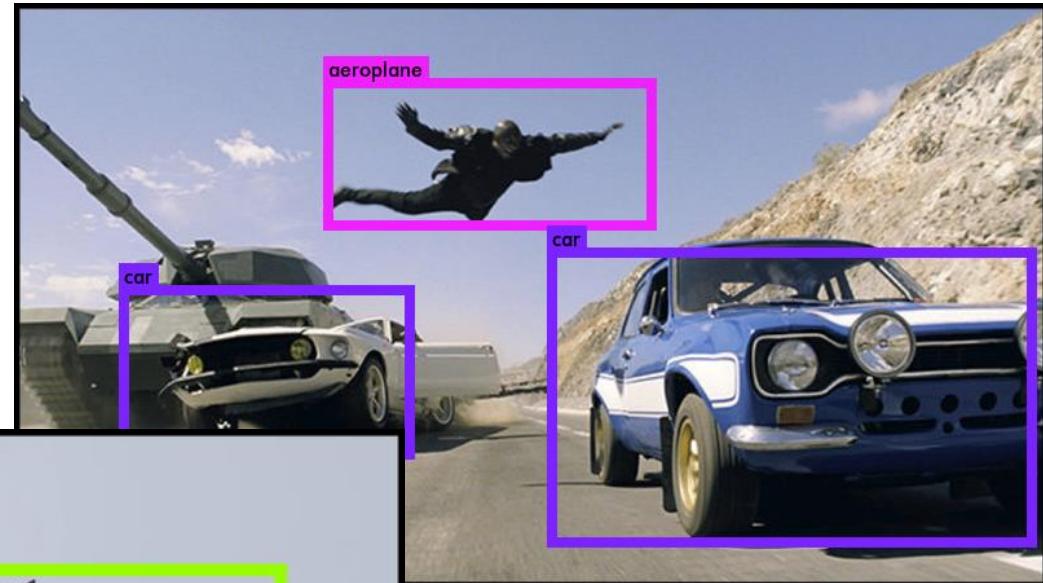
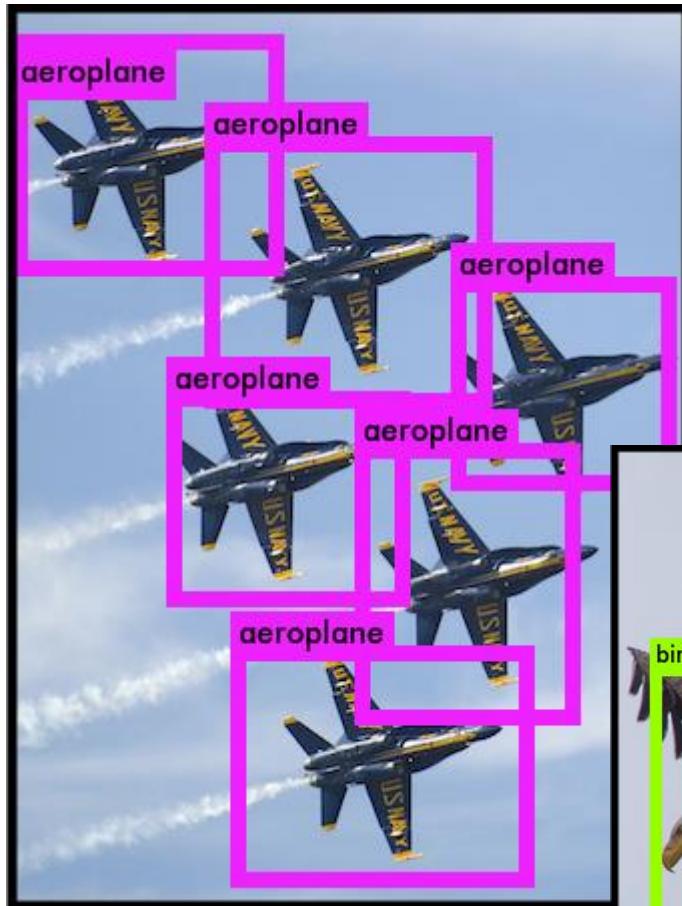


# We train with standard tricks:

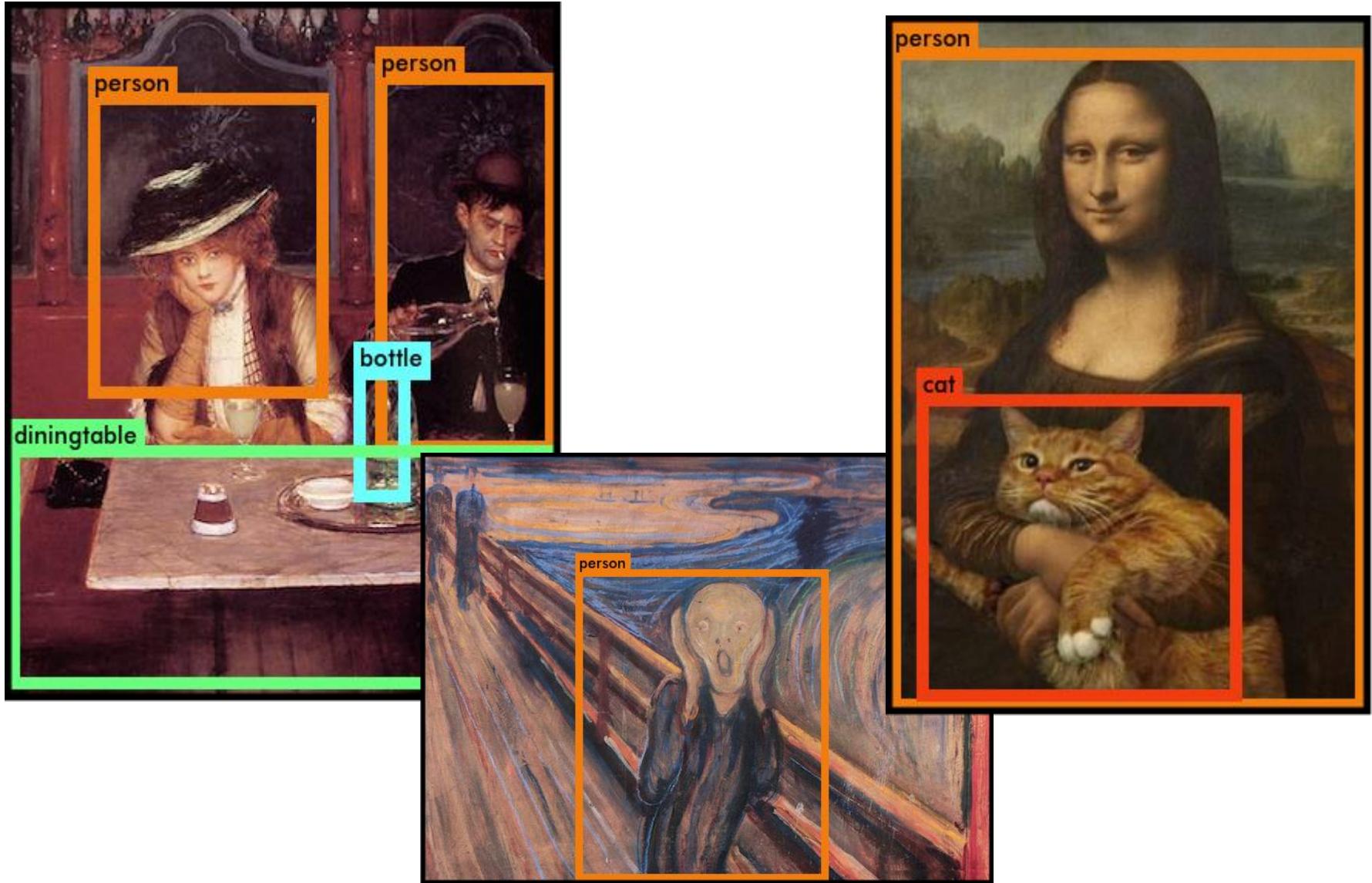
- Pretraining on Imagenet
- SGD with decreasing learning rate
- Extensive data augmentation
- For details, see the paper



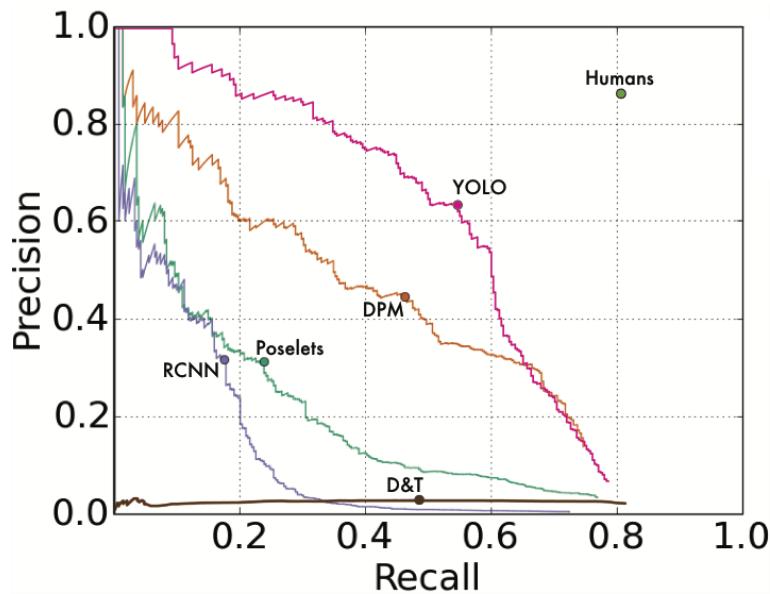
# YOLO works across a variety of natural images



It also generalizes well to new domains (like art)



YOLO outperforms methods like DPM and R-CNN when generalizing to person detection in artwork



	VOC 2007	Picasso		People-Art
	AP	AP	Best $F_1$	AP
<b>YOLO</b>	<b>59.2</b>	<b>53.3</b>	<b>0.590</b>	<b>45</b>
R-CNN	54.2	10.4	0.226	26
DPM	43.2	37.8	0.458	32

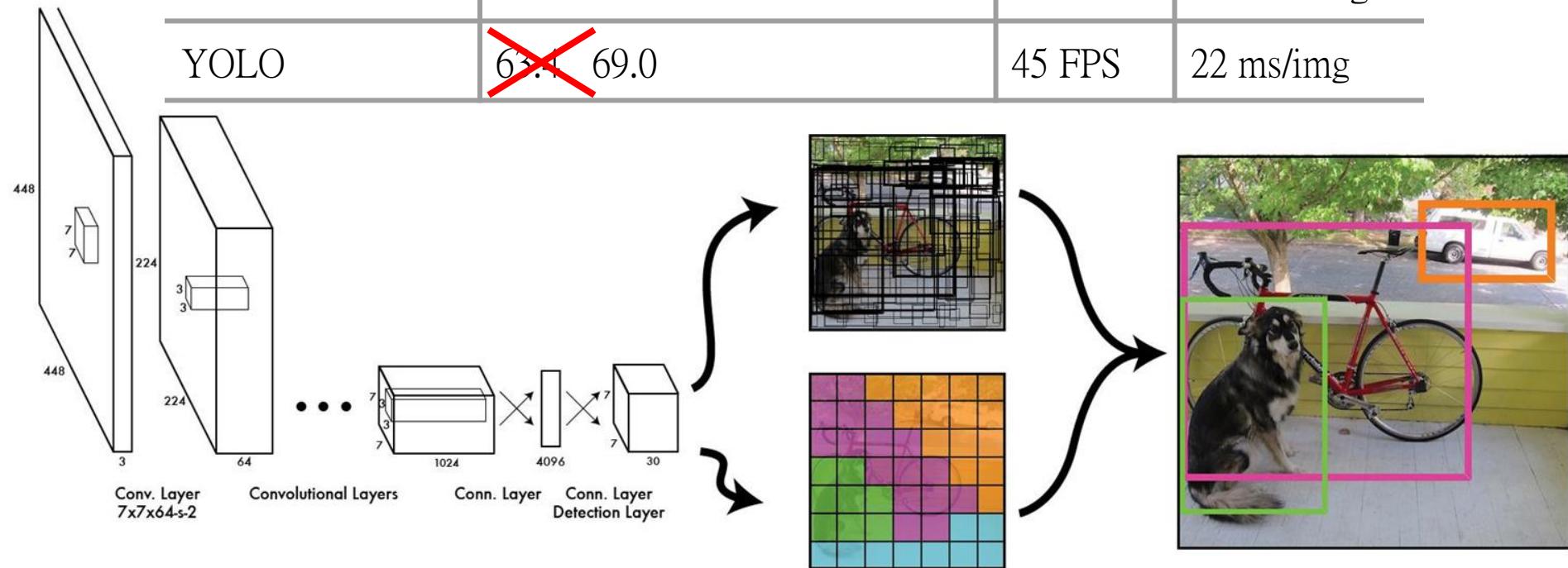
S. Ginosar, D. Haas, T. Brown, and J. Malik. Detecting people in cubist art. In Computer Vision-ECCV 2014 Workshops, pages 101–116. Springer, 2014.

H. Cai, Q. Wu, T. Corradi, and P. Hall. The cross-depiction problem: Computer vision algorithms for recognising objects in artwork and in photographs.

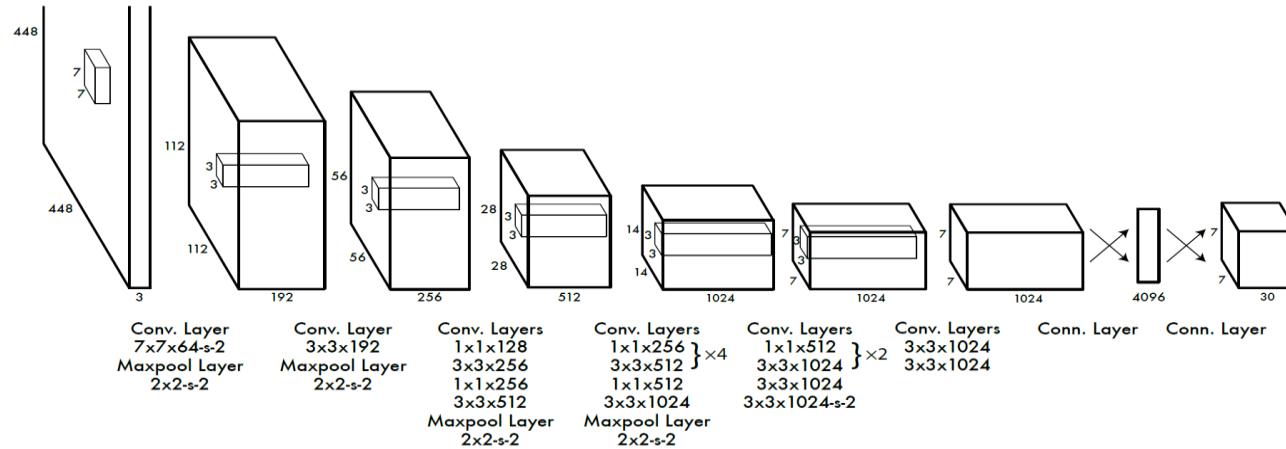
Code available! [pjreddie.com/yolo](http://pjreddie.com/yolo)



	Pascal 2007 mAP	Speed	
DPM v5	33.7	.07 FPS	14 s/img
R-CNN	66.0	.05 FPS	20 s/img
Fast R-CNN	70.0	.5 FPS	2 s/img
Faster R-CNN	73.2	7 FPS	140 ms/img
YOLO	69.0	45 FPS	22 ms/img



# YOLO- You Only Look Once



Divide the image into 7x7 cells.

Each cell trains a detector.

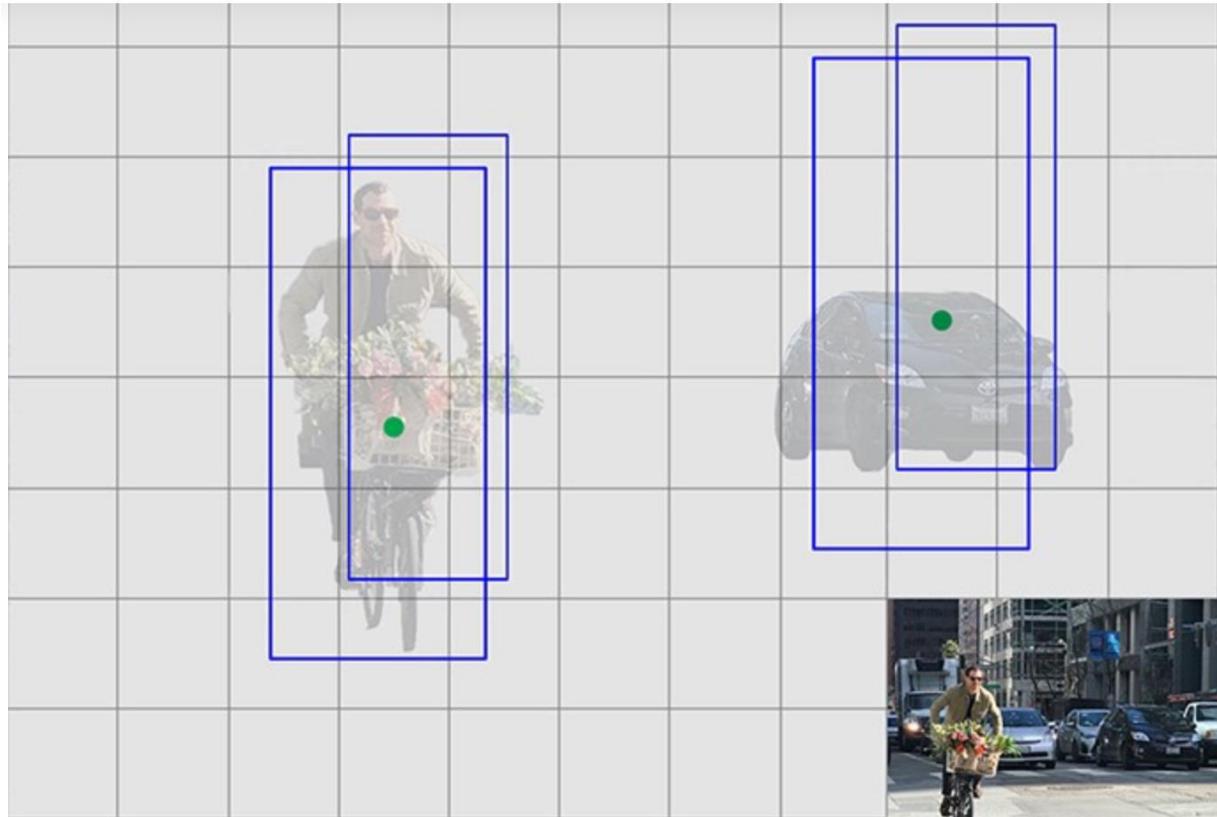
The detector needs to predict the object's class distributions.

The detector has 2 bounding-box predictors to predict bounding-boxes and confidence scores.

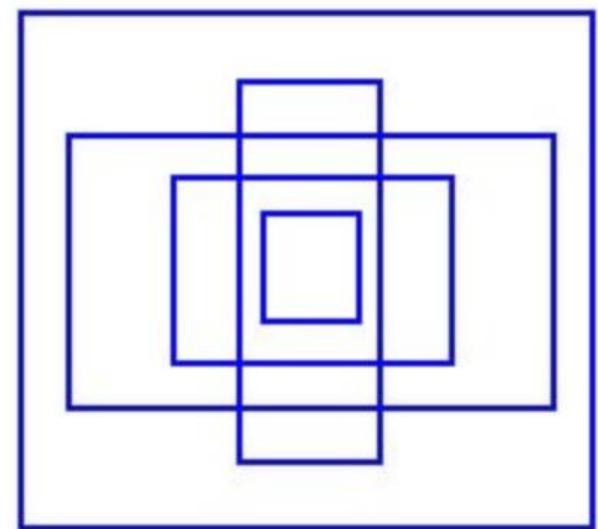
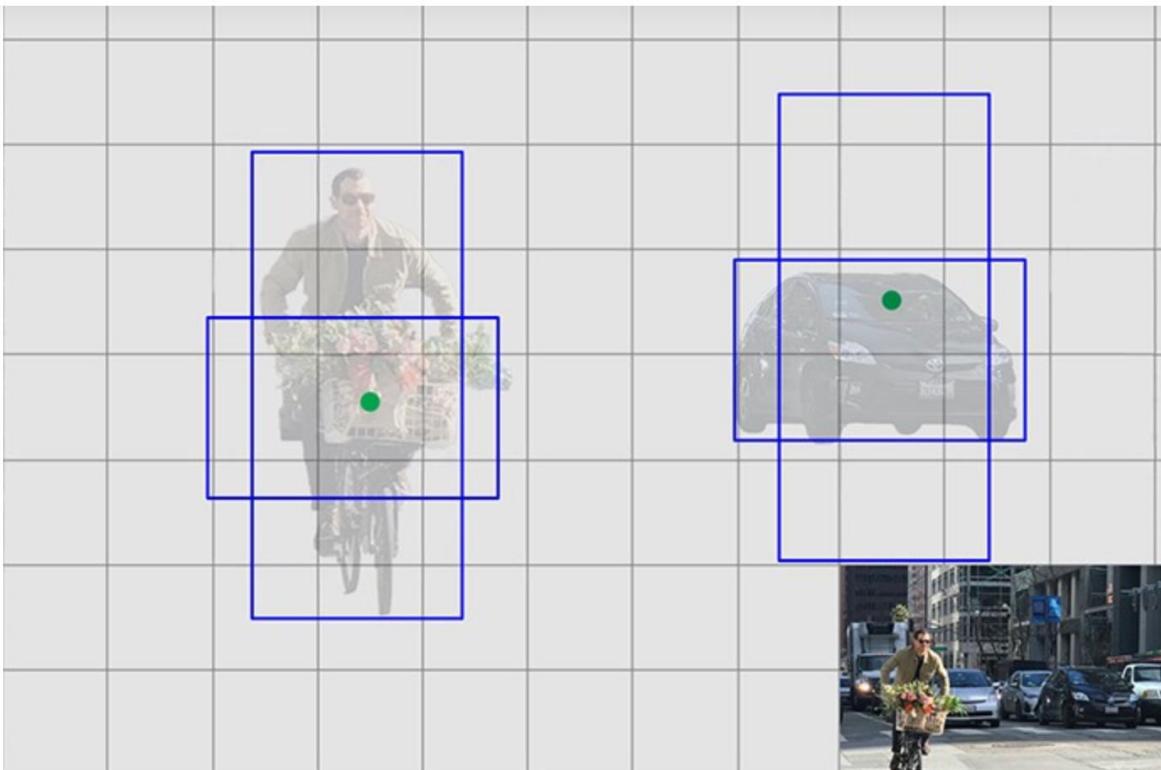
<https://arxiv.org/abs/1506.02640>

Redmon et al. CVPR 2016.

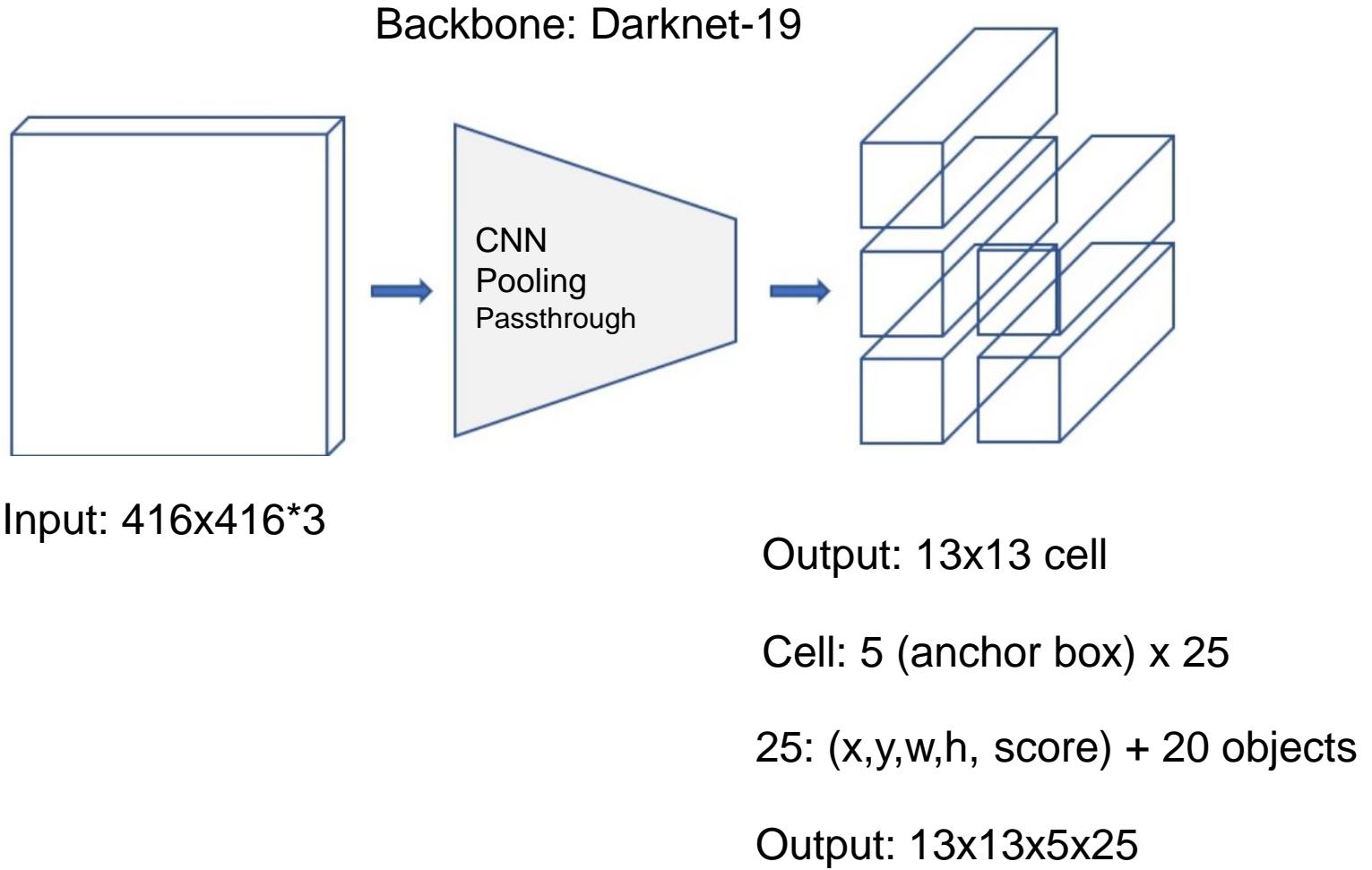
# YOLO- You Only Look Once



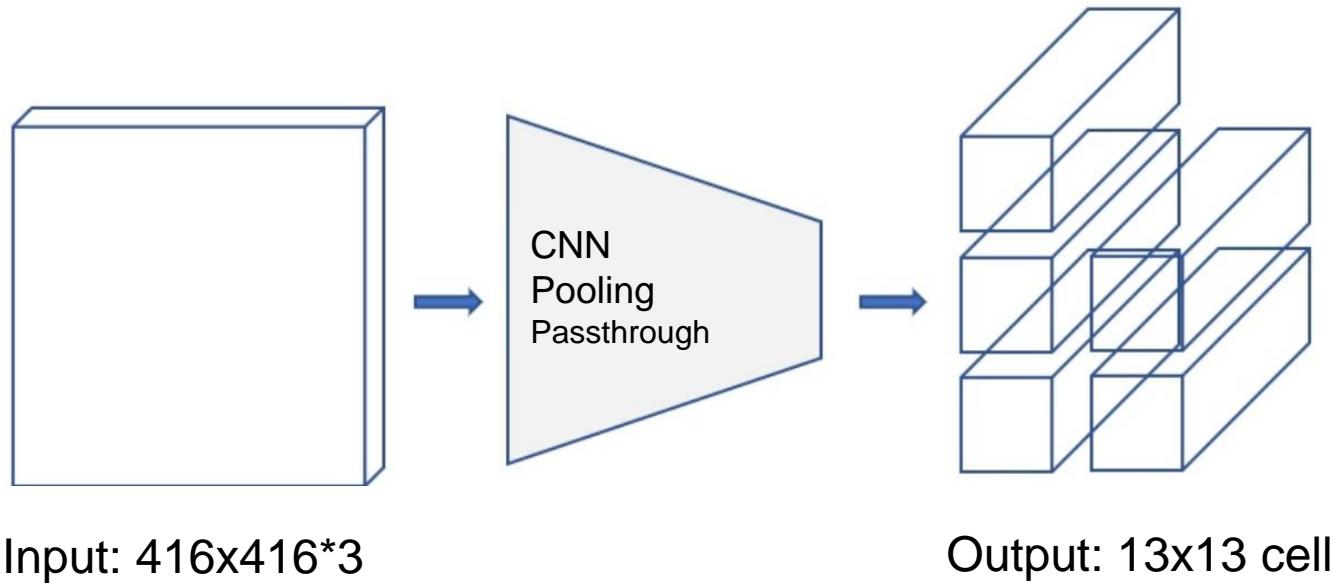
# YOLO 2



# YOLO 2



# YOLO 2

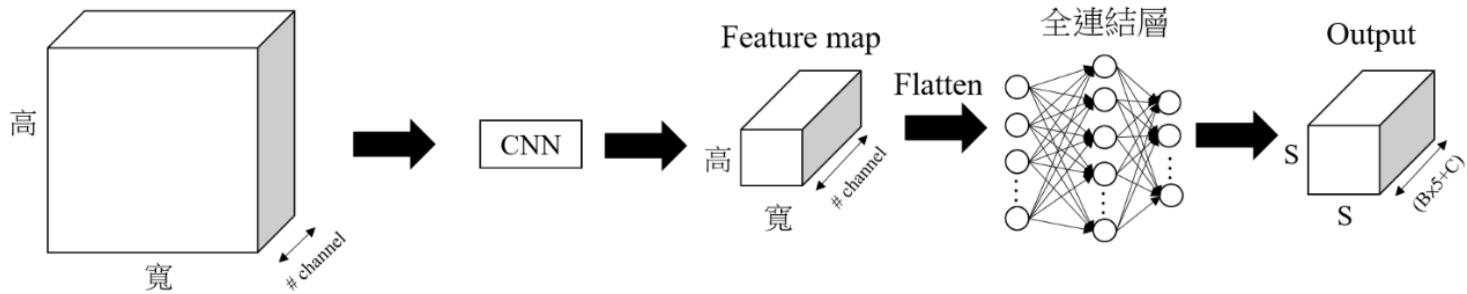


$$414/32 = 13, \text{ 降了 } 2 \times 2 \times 2 \times 2 \times 2 \text{ (5層)}$$

Cell: 5 (anchor box) 如何選呢 ?

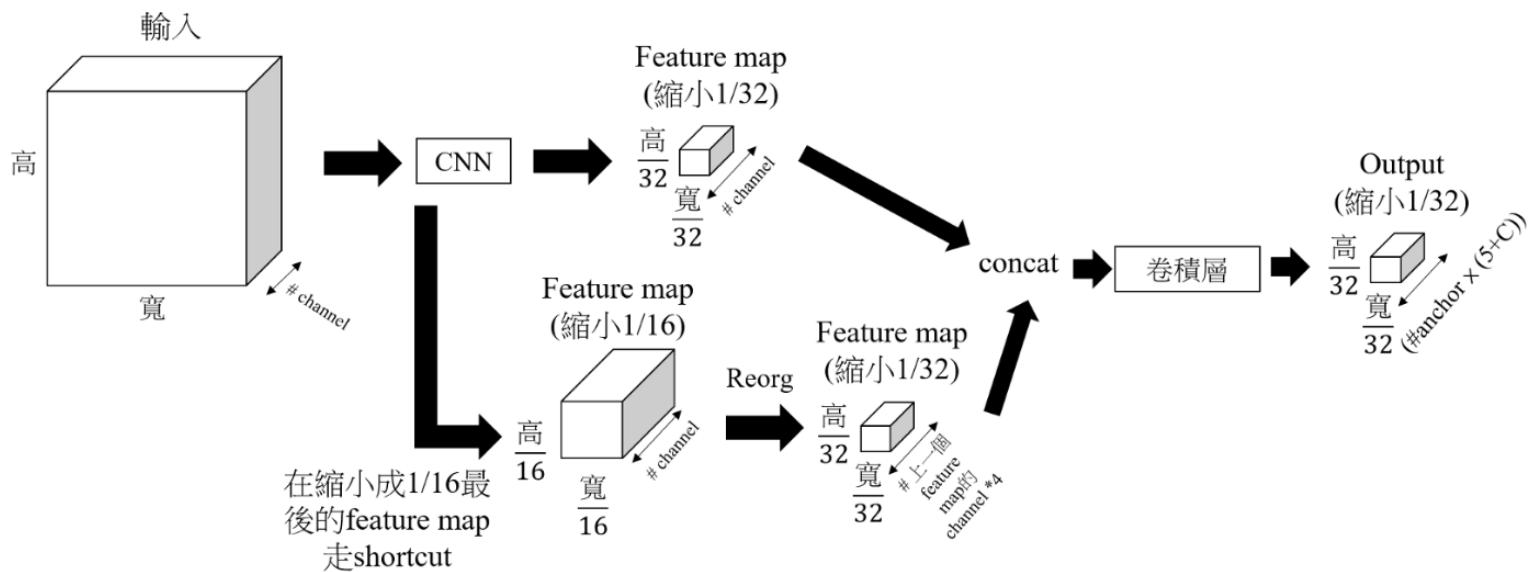
K-means

# YOLO 1



# YOLO 2

## YOLOv2



# Yolo v2 – Redmon et al. 2017

- Batch normalization
- Pre-train on higher resolution ImageNet
- Use and improve anchor box idea from Faster RCNN
- Train at multiple resolutions
- Very good accuracy, very fast

	YOLO	YOLOv2							
batch norm?	✓	✓	✓	✓	✓	✓	✓	✓	
hi-res classifier?	✓	✓	✓	✓	✓	✓	✓	✓	
convolutional?		✓	✓	✓	✓	✓	✓	✓	
anchor boxes?		✓	✓						
new network?			✓	✓	✓	✓	✓	✓	
dimension priors?				✓	✓	✓	✓	✓	
location prediction?					✓	✓	✓	✓	
passthrough?						✓	✓	✓	
multi-scale?							✓	✓	
hi-res detector?								✓	
VOC2007 mAP	63.4	65.8	69.5	69.2	69.6	74.4	75.4	76.8	<b>78.6</b>

Detection Frameworks	Train	mAP	FPS
Fast R-CNN [5]	2007+2012	70.0	0.5
Faster R-CNN VGG-16[15]	2007+2012	73.2	7
Faster R-CNN ResNet[6]	2007+2012	76.4	5
YOLO [14]	2007+2012	63.4	45
SSD300 [11]	2007+2012	74.3	46
SSD500 [11]	2007+2012	76.8	19
YOLOv2 288 × 288	2007+2012	69.0	91
YOLOv2 352 × 352	2007+2012	73.7	81
YOLOv2 416 × 416	2007+2012	76.8	67
YOLOv2 480 × 480	2007+2012	77.8	59
YOLOv2 544 × 544	2007+2012	<b>78.6</b>	40

<https://youtu.be/VOC3huqHrss>

# YOLOv3

## YOLOv3: An Incremental Improvement

Joseph Redmon, Ali Farhadi

University of Washington

### Abstract

We present some updates to YOLO! We made a bunch of little design changes to make it better. We also trained this new network that's pretty swell. It's a little bigger than last time but more accurate. It's still fast though, don't worry. At  $320 \times 320$  YOLOv3 runs in 22 ms at 28.2 mAP, as accurate as SSD but three times faster. When we look at the old .5 IOU mAP detection metric YOLOv3 is quite good. It achieves 57.9 AP<sub>50</sub> in 51 ms on a Titan X, compared to 57.5 AP<sub>50</sub> in 198 ms by RetinaNet, similar performance but 3.8× faster. As always, all the code is online at <https://pjreddie.com/yolo/>.

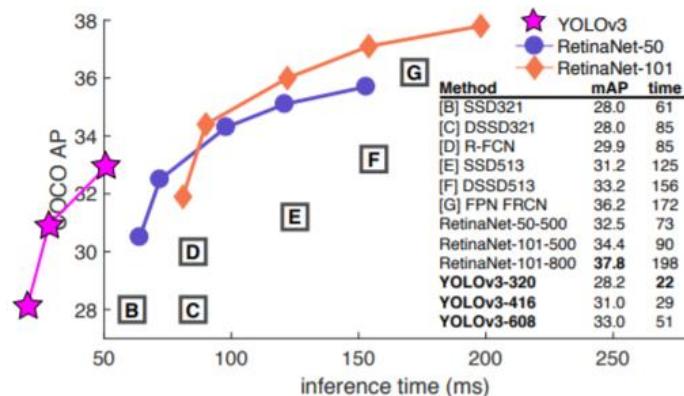
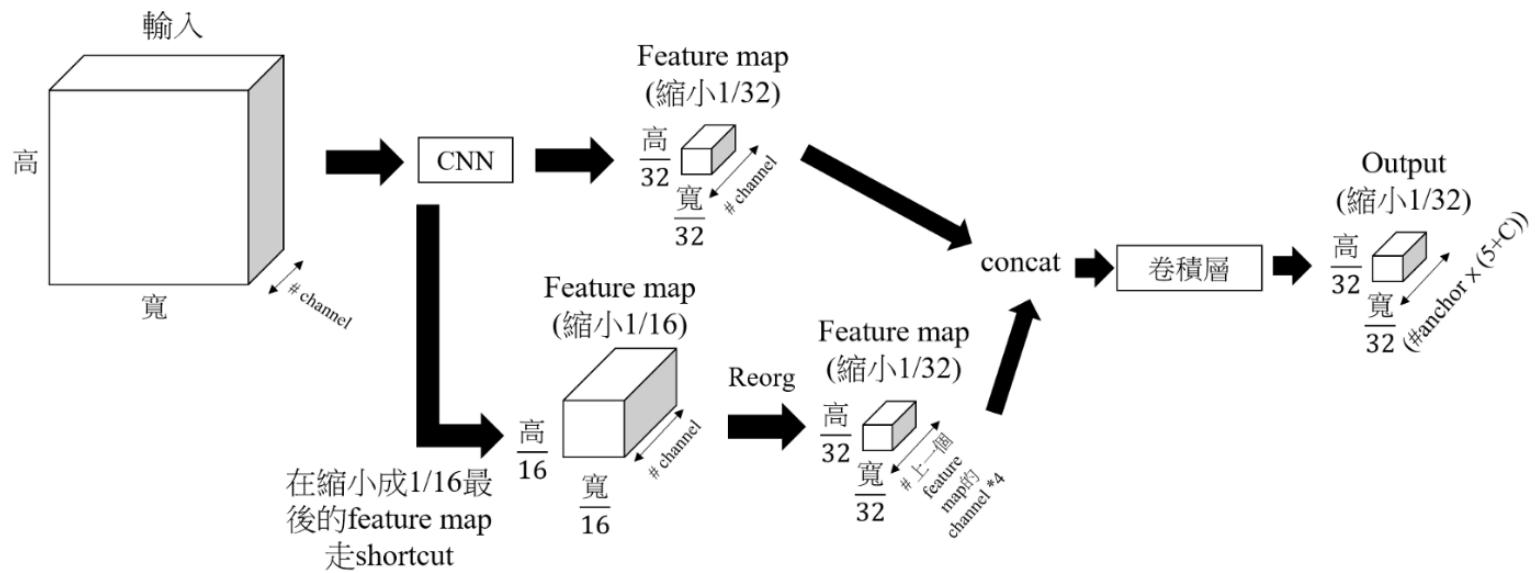


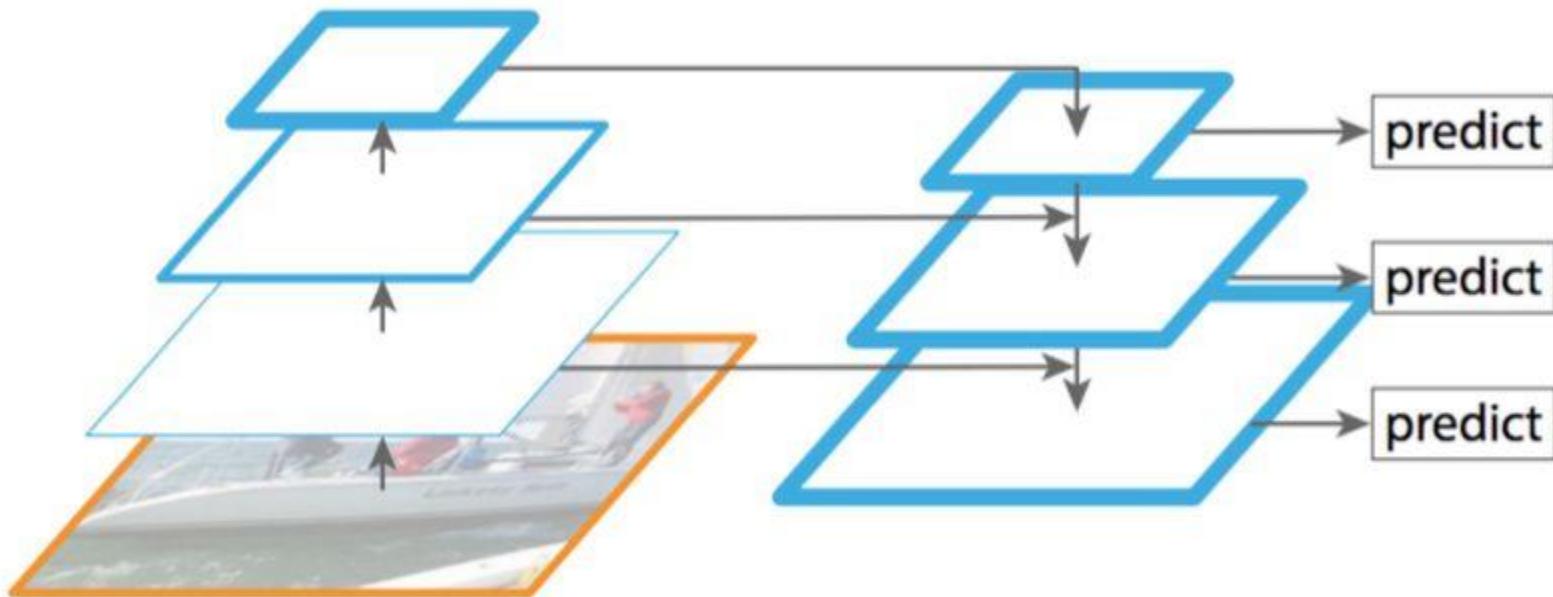
Figure 1. We adapt this figure from the Focal Loss paper [9]. YOLOv3 runs significantly faster than other detection methods

## YOLOv2

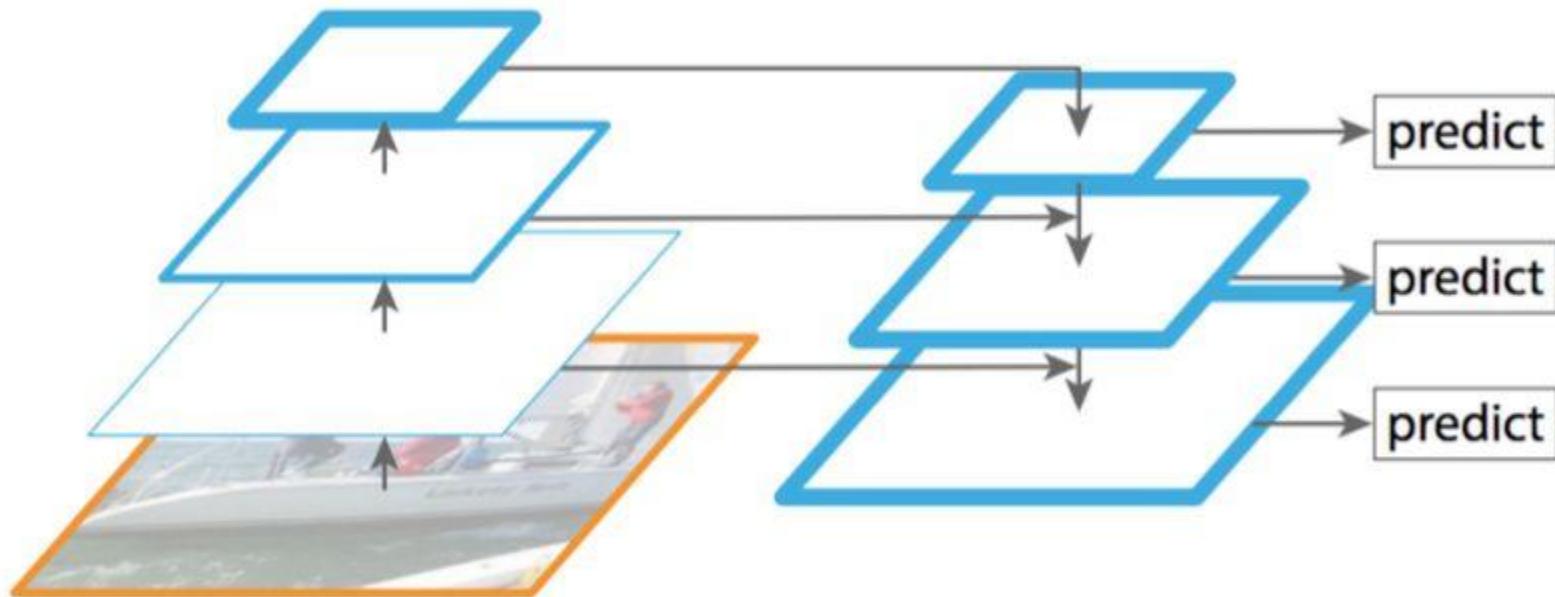


Backbone: Darknet-19

# YOLO3

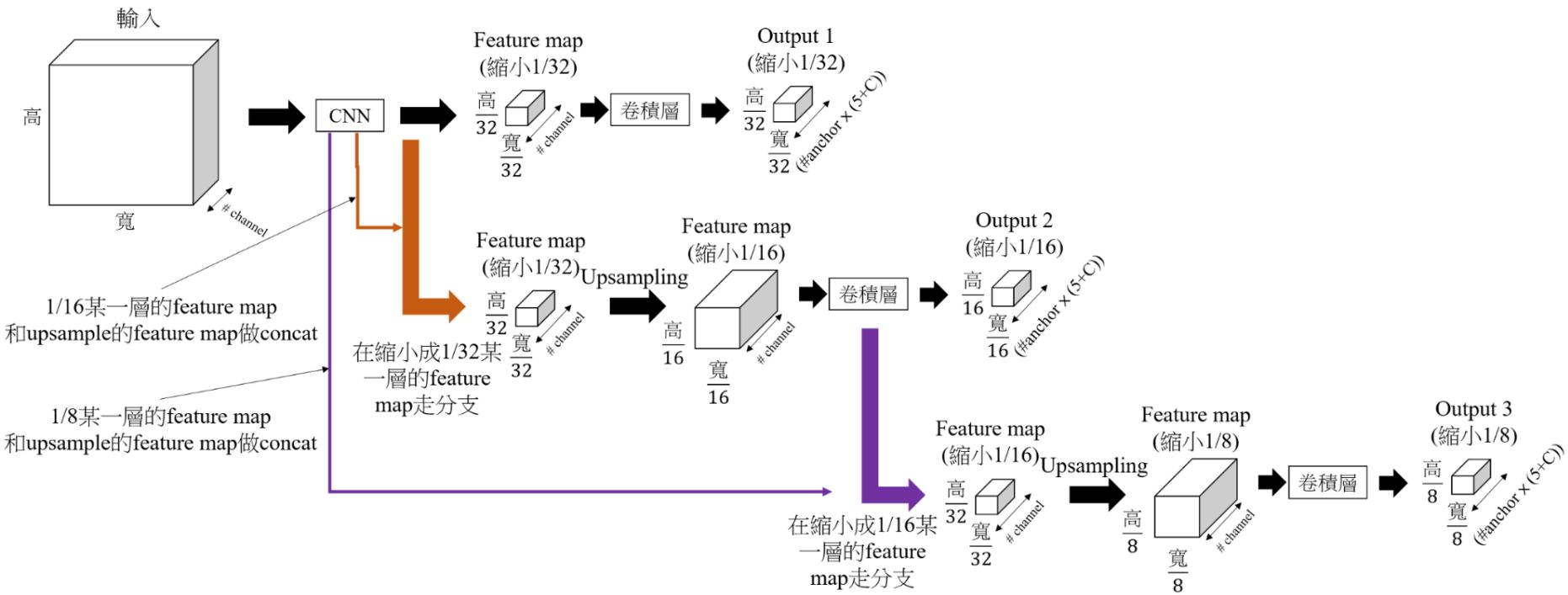


# YOLO3



Backbone: Darknet-53

# YOLO 3



Backbone: Darknet-53

# YOLOV3

YOLO v3中使用53層的卷積網路

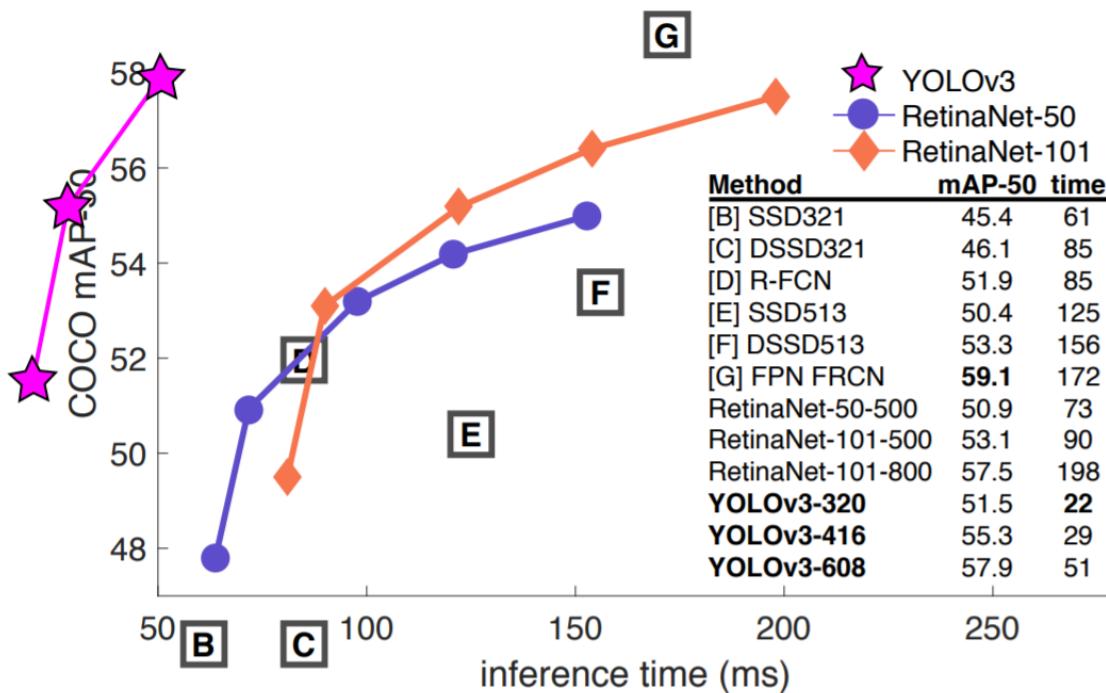
不同於YOLO2，YOLO v3從三種不同尺度的特徵圖譜上進行預測任務

模型不再使用softmax函數作為最終的分類器，使用logistic作為分類器使用 binary cross-entropy作為損失函數

Type	Filters	Size	Output
Convolutional	32	$3 \times 3$	$256 \times 256$
Convolutional	64	$3 \times 3 / 2$	$128 \times 128$
1x	Convolutional	$32$	$1 \times 1$
	Convolutional	$64$	$3 \times 3$
	Residual		$128 \times 128$
2x	Convolutional	$128$	$3 \times 3 / 2$
	Convolutional	$64$	$1 \times 1$
	Convolutional	$128$	$3 \times 3$
8x	Residual		$64 \times 64$
	Convolutional	$256$	$3 \times 3 / 2$
	Convolutional	$128$	$1 \times 1$
8x	Convolutional	$256$	$3 \times 3$
	Residual		$32 \times 32$
	Convolutional	$512$	$3 \times 3 / 2$
8x	Convolutional	$256$	$1 \times 1$
	Convolutional	$512$	$3 \times 3$
	Residual		$16 \times 16$
4x	Convolutional	$1024$	$3 \times 3 / 2$
	Convolutional	$512$	$1 \times 1$
	Convolutional	$1024$	$3 \times 3$
	Residual		$8 \times 8$
	Avgpool		Global
	Connected		1000
	Softmax		

Table 1. Darknet-53.

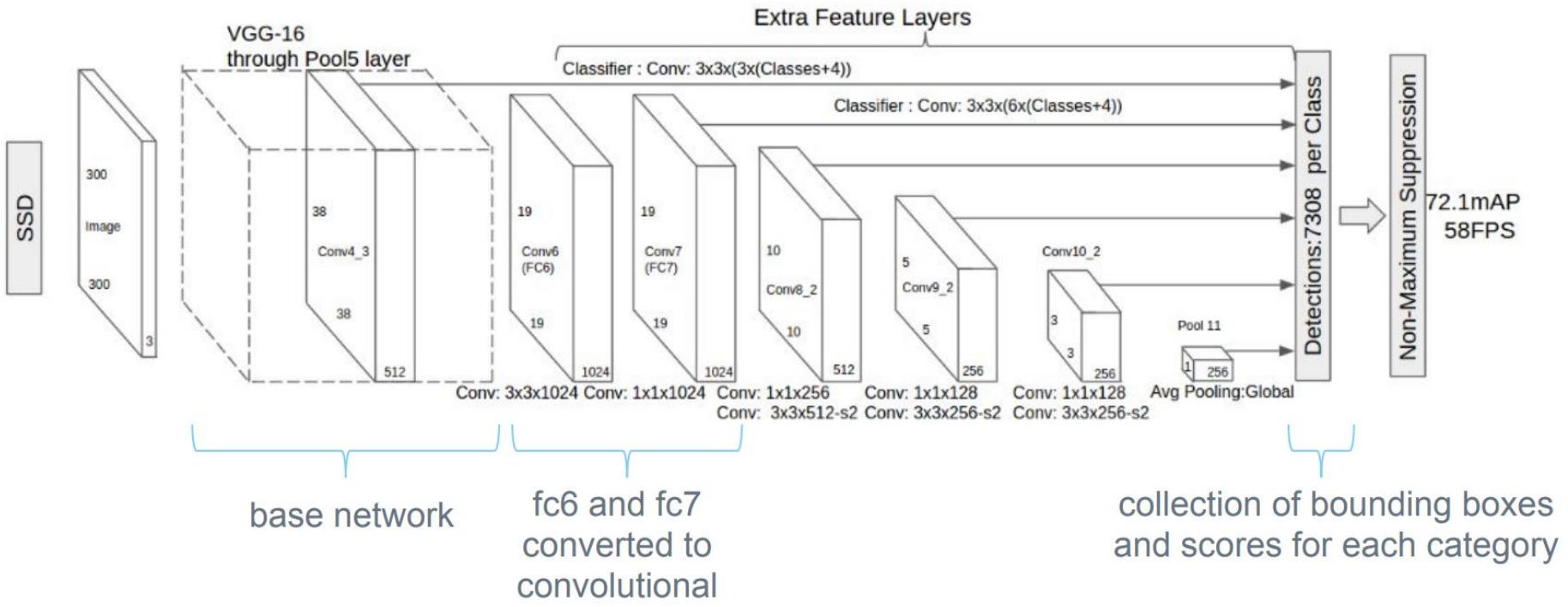
# YOLOv3 –Performance



<https://pjreddie.com/media/files/papers/YOLOv3.pdf>

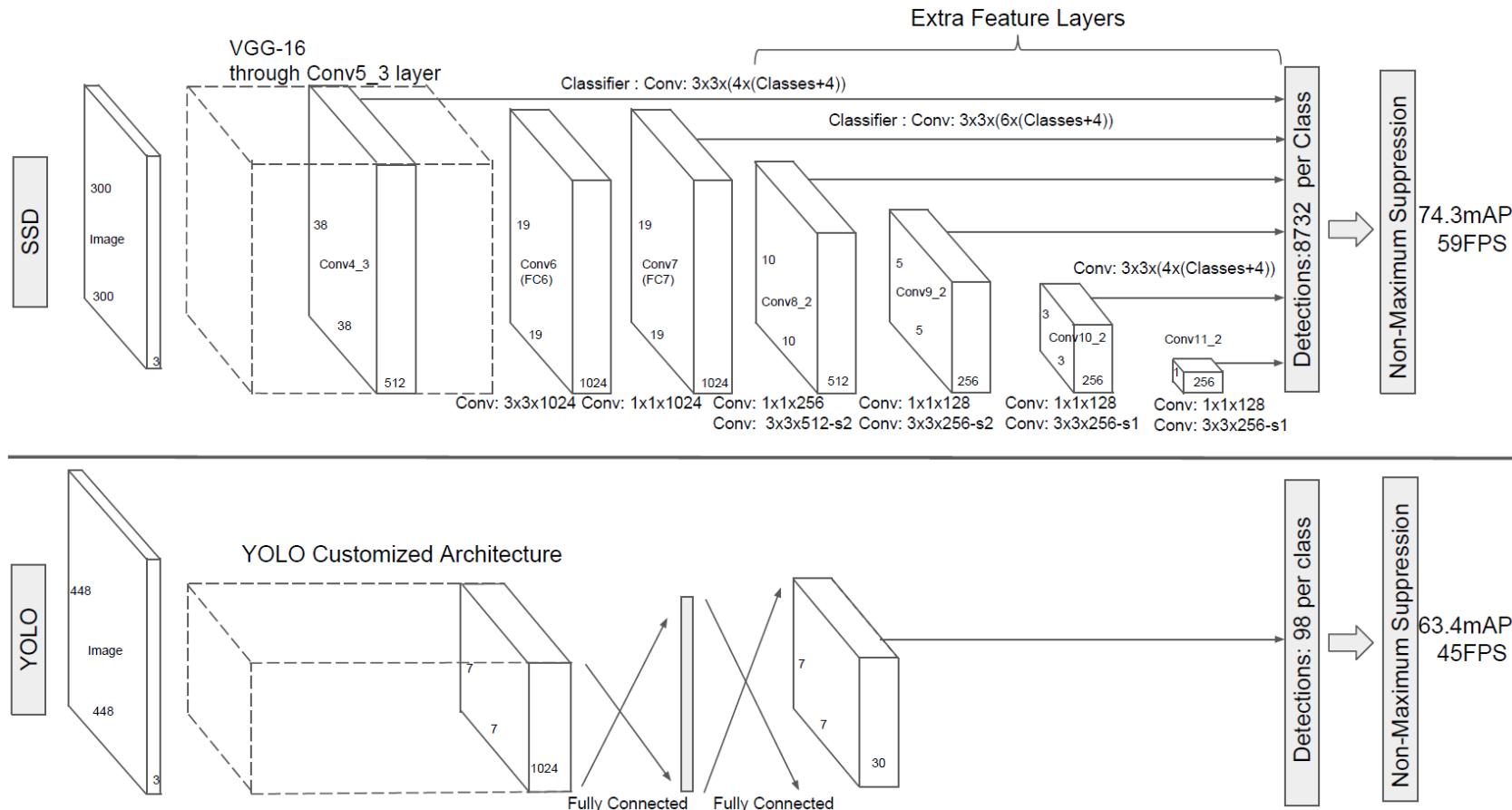
# SSD: Single Shot Detector

Liu et al. ECCV 2016.

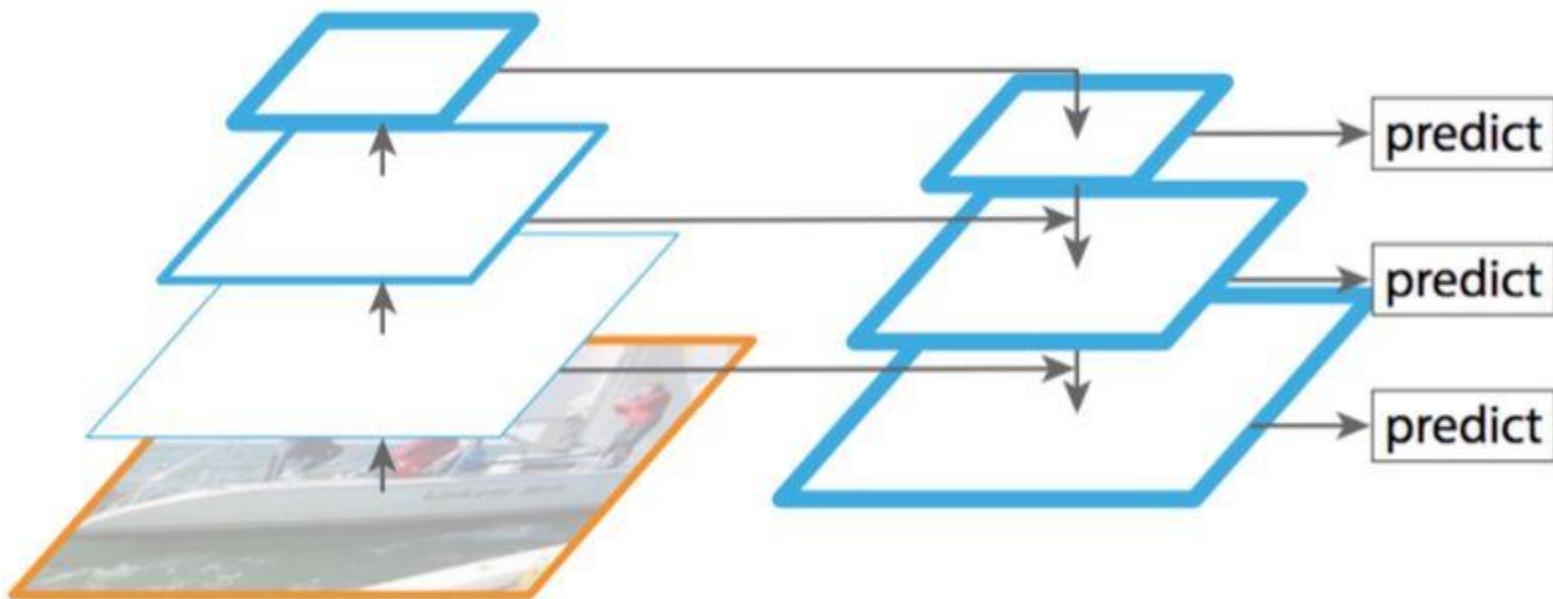


Multi-scale feature maps for detection: observe how conv feature maps decrease in size and allow predictions at multiple scales

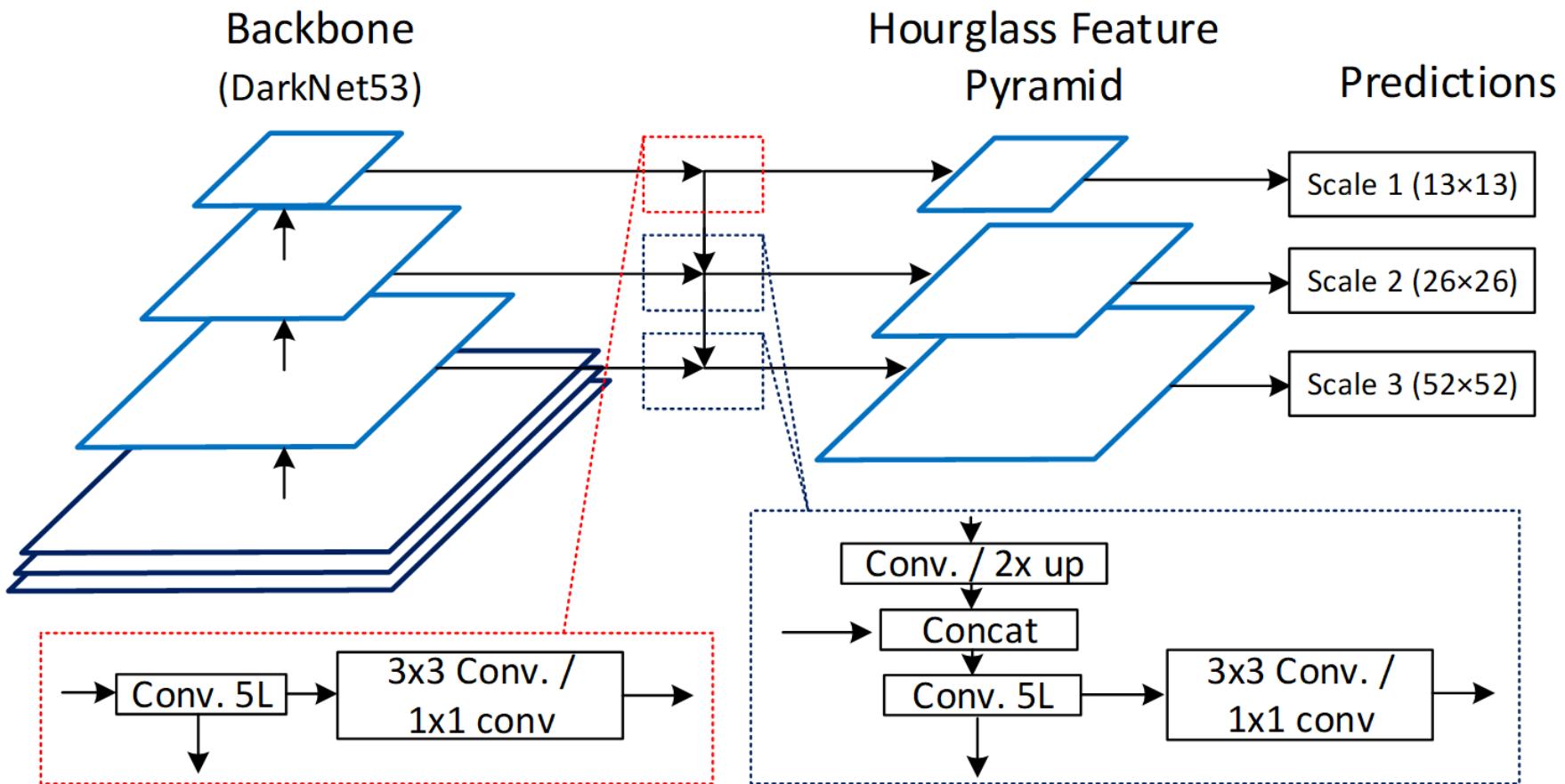
# SSD 與 YOLO 比較



# YOLO3

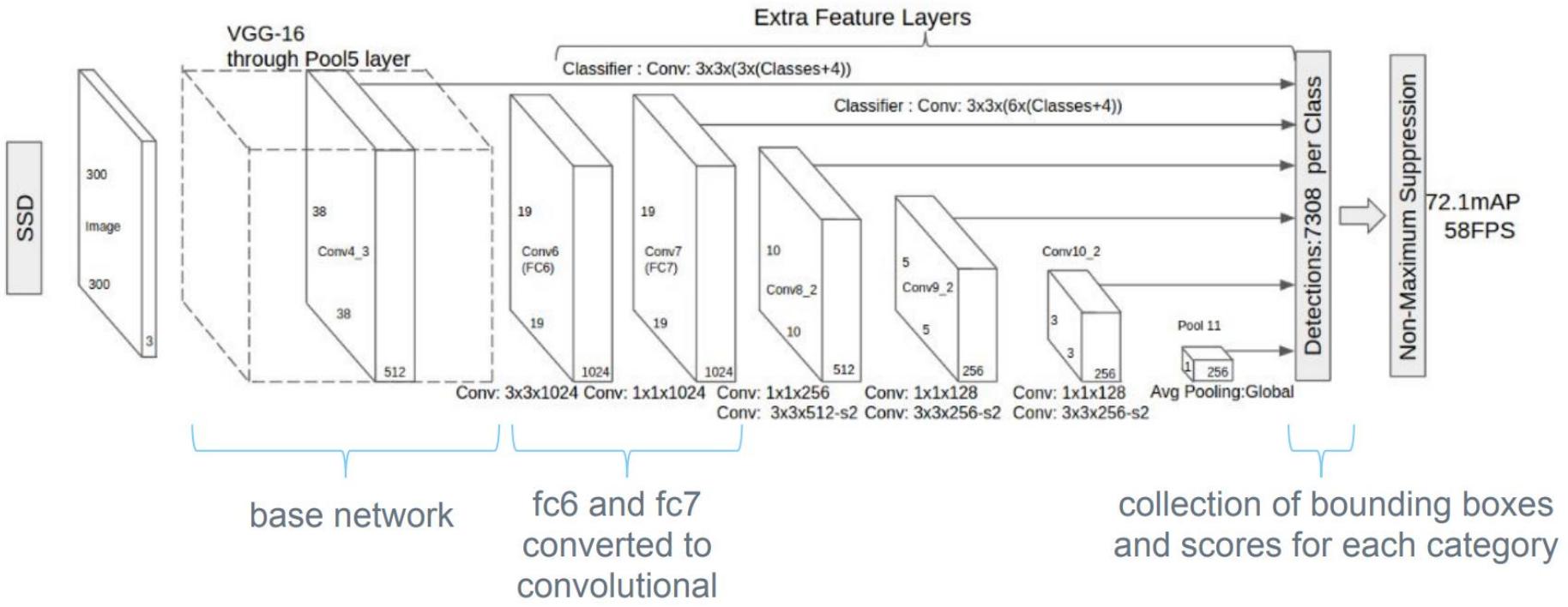


# YOLO3



# SSD: Single Shot Detector

Liu et al. ECCV 2016.



Convolutional predictors for detection: a set of filters that predict detections for different aspect ratios

# SSD: Single Shot Detector

Liu et al. ECCV 2016.

## What is a detection ?



Described by **four parameters** (center bounding box x and y, width and height)

Class category

**For all categories we need for a detection a total of #classes + 4 values**

# SSD: Single Shot Detector

Liu et al. ECCV 2016.

Different “classes” of detections



aspect ratio 2:1  
for cats



aspect ratio 1:2  
for cats

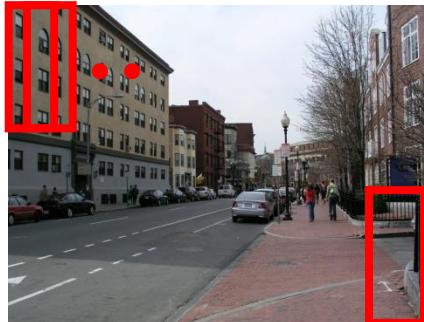
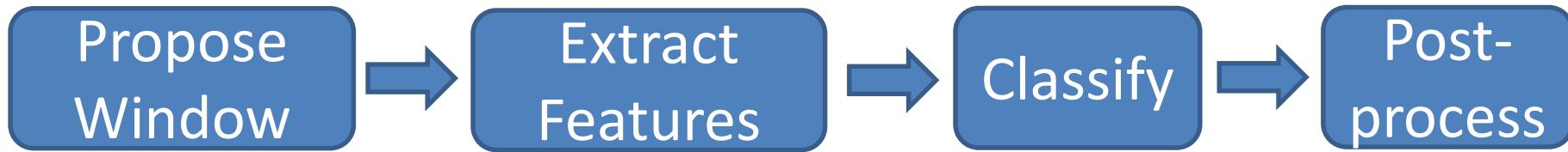


aspect ratio 1:1  
for cats

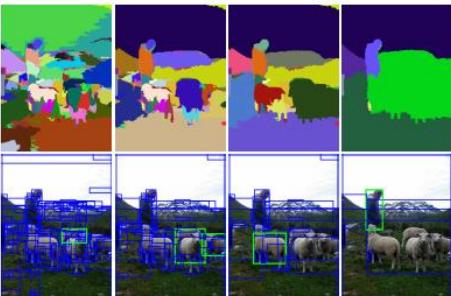
# Influential Works in Detection

- Sung-Poggio (1994, 1998) : ~2412 citations
  - Basic idea of statistical template detection (I think), bootstrapping to get “face-like” negative examples, multiple whole-face prototypes (in 1994)
- Rowley-Baluja-Kanade (1996-1998) : ~4953
  - “Parts” at fixed position, non-maxima suppression, simple cascade, rotation, pretty good accuracy, fast
- Schneiderman-Kanade (1998-2000,2004) : ~2600
  - Careful feature/classifier engineering, excellent results, cascade
- Viola-Jones (2001, 2004) : ~27,000
  - Haar-like features, Adaboost as feature selection, hyper-cascade, very fast, easy to implement
- Dalal-Triggs (2005) : ~18000
  - Careful feature engineering, excellent results, HOG feature, online code
- Felzenszwalb-Huttenlocher (2000): ~2100
  - Efficient way to solve part-based detectors
- Felzenszwalb-McAllester-Ramanan (2008,2010): ~7200
  - Excellent template/parts-based blend
- Girshick-Donahue-Darrell-Malik (2014) : ~4700
  - Region proposals + fine-tuned CNN features (marks significant advance in accuracy over hog-based methods)
- Redmon, Divvala, Girshick, Farhadi (2016): ~210
  - Refine and simplify RCNN++ approach to predict directly from last conv layer

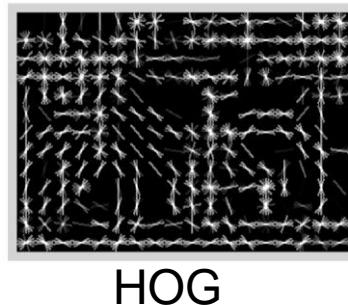
# Summary: statistical templates



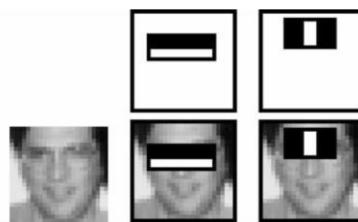
Sliding window: scan image pyramid



Region proposals:  
edge/region-based,  
resize to fixed window



HOG



Fast randomized features

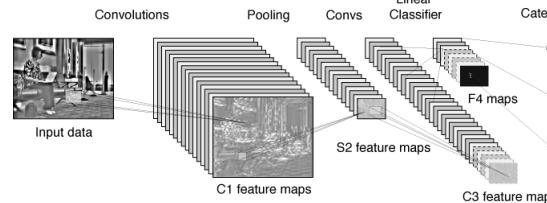
SVM

Boosted stubs

Neural network

Non-max suppression

Segment or refine localization



CNN features

# On the Design of Modern Deep Detectors

**Backbone:** the feature extractor from which the detector draws its discriminative power.

Same backbone: **Region Proposal Network**, VGGNet-16

Faster-RCNN: [Ren et al., 2015]

R-FCN: [Dai et al., 2016]

SSD [Liu et al., 2016]

# On the Design of Modern Deep Detectors

## Common Backbones

1. AlexNet:
2. VGGNet: [Simonyan and Zisserman, 2014]
3. Inception:[Ioffe, 2017, Szegedy et al., 2015, 2016]
4. ResNet: [He et al., 2016]

# On the Design of Modern Deep Detectors

## New Backbones

1. Inception-ResNet: [Szegedy et al., 2017]
2. Xception: group convolution [Chollet, 2017]
3. DarkNet [Redmon and Farhadi, 2017]
4. DenseNet: [Huang et al., 2017b]
5. Hourglass [Newell et al., 2016]
6. Wide-Residual Net: [Lee et al., 2017]
7. ResNeXt [Xie et al., 2017]
8. Dilated Residual Networks: [Yu et al., 2017]
9. DetNet [Li et al., 2018]

For speed

10. MobileNet: [Howard et al., 2017]
11. SqueezeNet [Iandola et al., 2016]
12. ShuffleNet [Zhang et al., 2017]
13. PeLLe [NIPS 2018]