# Automatic Object Detection ----YOLO & SSD

**Jianping Fan**
**Dept of Computer Science**
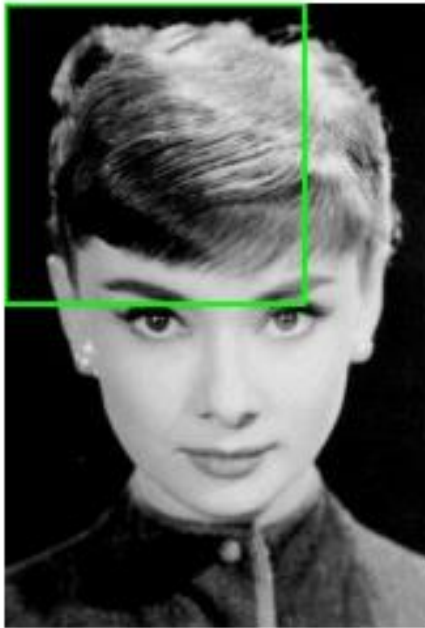**UNC-Charlotte**

# Before Deep Learning

Sliding windows.
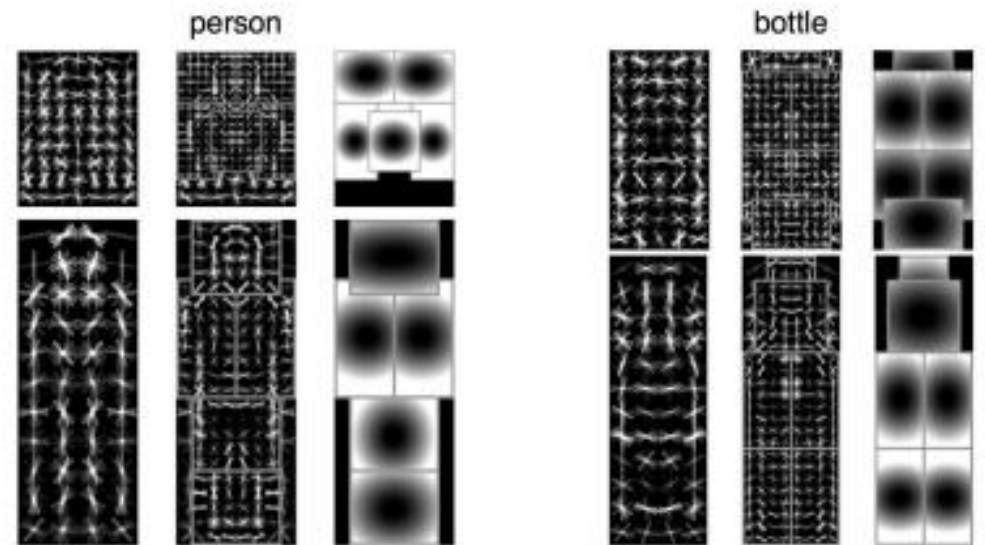
- Score every subwindow.

Deformable part models (DPM)

- Uses HOG features
- Very fast

person                    bottle



https://cs.brown.edu/~pff/papers/lsvm-pami.pdf

# Techniques for Key Components

## Selective Search



http://www.huppelen.nl/publications/selectiveSearchDraft.pdf

Uijlings, J. R., van de Sande, K. E., Gevers, T., & Smeulders, A. W. (2013). Selective search for object recognition

**(a) Two-stage approaches**

**(b) One-stage approaches**

## Hard Negative Mining

Imbalance between positive and negative examples.

Use negative examples with higher confidence score.

## Non Maximum Suppression

If output boxes overlap, only consider the most confident.

## Bounding Box Regression

Regression used to find bounding box parameters

Applied in one of two ways

- Bounding box refinement
- Complete object detection

# You Only Look Once:
# Unified, Real-Time Object Detection

Joseph Redmon[*], Santosh Divvala[*†], Ross Girshick[¶], Ali Farhadi[*†]

University of Washington[*], Allen Institute for AI[†], Facebook AI Research[¶]

http://pjreddie.com/yolo/

## Abstract

*We present YOLO, a new approach to object detection. Prior work on object detection repurposes classifiers to perform detection. Instead, we frame object detection as a regression problem to spatially separated bounding boxes and associated class probabilities. A single neural network predicts bounding boxes and class probabilities directly from full images in one evaluation. Since the whole detection pipeline is a single network, it can be optimized end-to-end directly on detection performance.*

*Our unified architecture is extremely fast. Our base YOLO model processes images in real-time at 45 frames per second. A smaller version of the network, Fast YOLO, processes an astounding 155 frames per second while still achieving double the mAP of other real-time detectors. Compared to state-of-the-art detection systems, YOLO makes more localization errors but is less likely to predict false positives on background. Finally, YOLO learns very general representations of objects. It outperforms other detection methods, including DPM and R-CNN, when generalizing from natural images to other domains like artwork.*
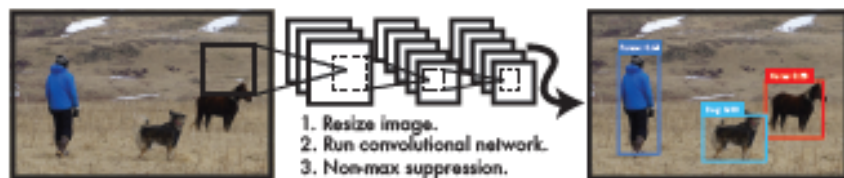
**Figure 1: The YOLO Detection System.** Processing images with YOLO is simple and straightforward. Our system (1) resizes the input image to $448 \times 448$, (2) runs a single convolutional network on the image, and (3) thresholds the resulting detections by the model's confidence.

methods to first generate potential bounding boxes in an image and then run a classifier on these proposed boxes. After classification, post-processing is used to refine the bounding boxes, eliminate duplicate detections, and rescore the boxes based on other objects in the scene [13]. These complex pipelines are slow and hard to optimize because each individual component must be trained separately.

We reframe object detection as a single regression problem, straight from image pixels to bounding box coordinates and class probabilities. Using our system, you only

# Shortcoming:

1. Slow, impossible for real-time detection

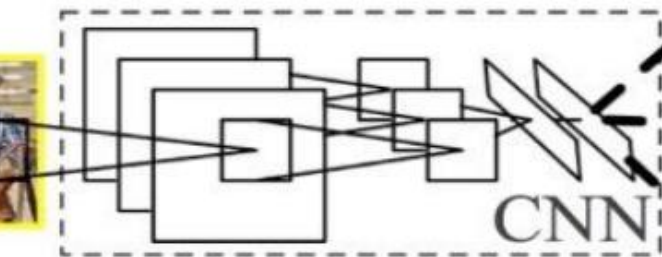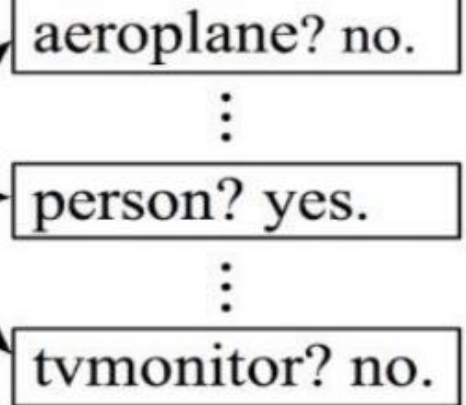2. Hard to optimize



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

warped region

aeroplane? no.
person? yes.
tvmonitor? no.

CNN

**1.** Input image
**2.** Extract region proposals (~2k)
**3.** Compute CNN features
**4.** Classify regions

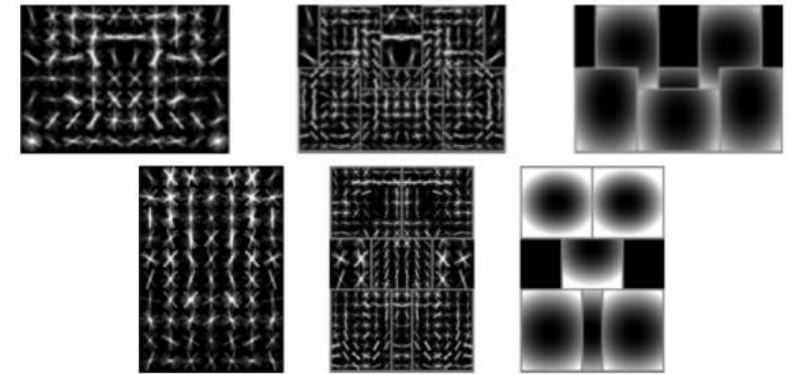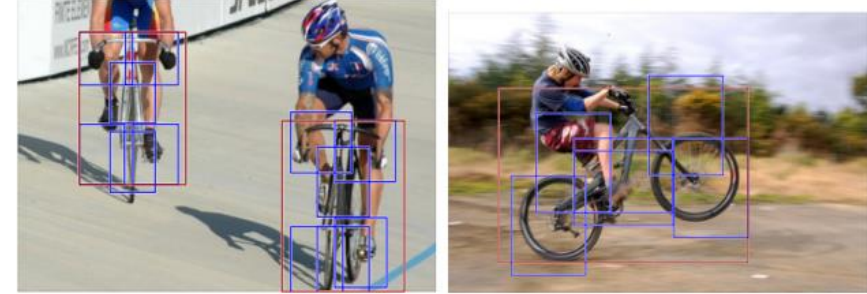# Object Detection as Regression Problem

- **Previous:** Repurpose classifiers to perform detection

  - <u>Deformable Parts Models (DPM)</u>
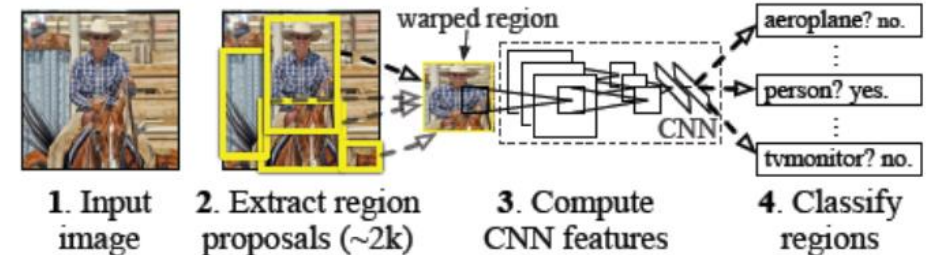
    - Sliding window

  - <u>R-CNN based methods</u>

    - 1) generate potential bounding boxes.

    - 2) run classifiers on these proposed boxes

    - 3) post-processing (refinement, elimination, rescore)



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

1. Input image
2. Extract region proposals (~2k)
3. Compute CNN features
4. Classify regions

aeroplane? no.
person? yes.
tvmonitor? no.

# What's new?  Regression

## YOLO Features :

1. Extremely fast (45 frames per second)

2. Reason Globally on the Entire Image

3. Learn Generalizable Representations

# Object Detection as Regression Problem

- **YOLO:** Single Regression Problem

  - Image → bounding box coordinate and class probability.

  - Extremely Fast
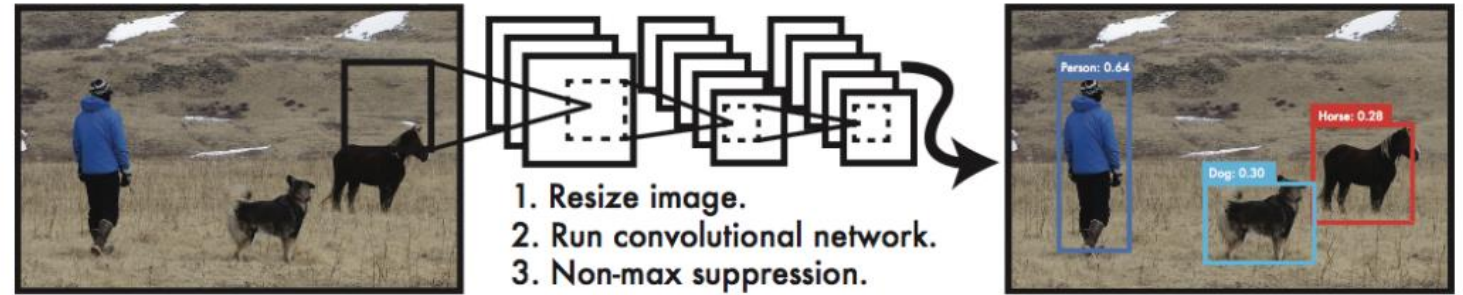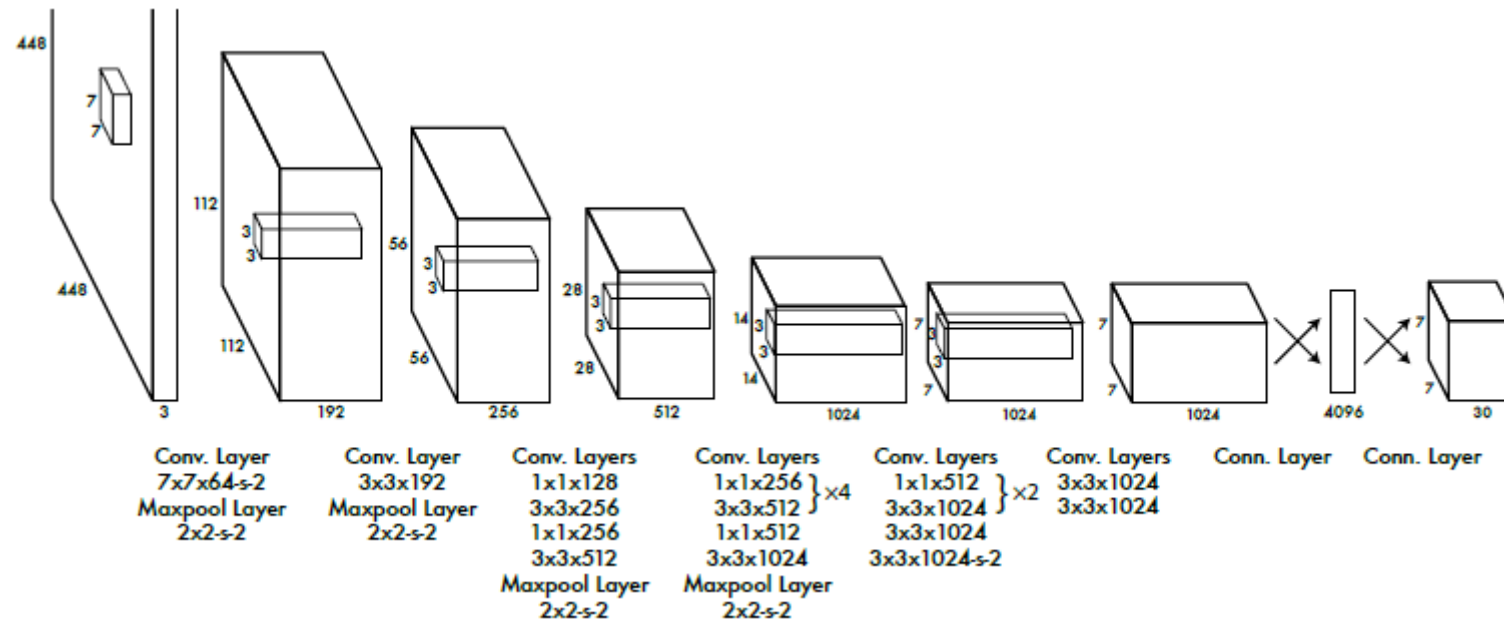
  - Global reasoning

  - Generalizable representation



**Figure 1: The YOLO Detection System.** Processing images with YOLO is simple and straightforward. Our system (1) resizes the input image to 448 × 448, (2) runs a single convolutional network on the image, and (3) thresholds the resulting detections by the model's confidence.

# YOLO
## You Only Look Once: Unified, Real-Time Object Detection

? Tool?-- A single neural network, unified architecture

? Framework?– Darknet

? Technology background?--related methods are slow , not real-time and devoid of generalization ability

• Solution and Advantages:

@Simpler structure of network

@much more faster, even with real-time property: 150fps: able to process streaming video in real-time with less than 25 milliseconds of latency

@Maintaining  a proper accuracy range

# How YOLO gets its goal?



- For speed: <span style="color:red">Also! structure advantage!</span>

➢ No bounding box proposals and subsequent pixel or feature resampling stage

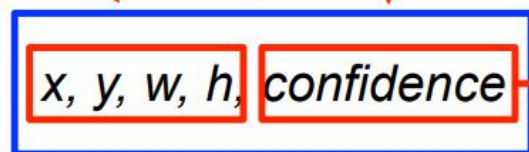✓ A neural network predicts bounding boxes and class probabilities directly from full images in one evaluation

# Unified Detection

- All BBox, All classes

1) Image → **S x S** grids

2) grid cell
→ **B**: BBoxes and Confidence score

$x, y, w, h,$ *confidence* → $\mathrm{Pr}(\mathrm{Object}) * \mathrm{IOU}_{\mathrm{pred}}^{\mathrm{truth}}$

→ **C**: class probabilities w.r.t #classes

$\mathrm{Pr}(\mathrm{Class}_i | \mathrm{Object})$



S × S grid on input

Bounding boxes + confidence
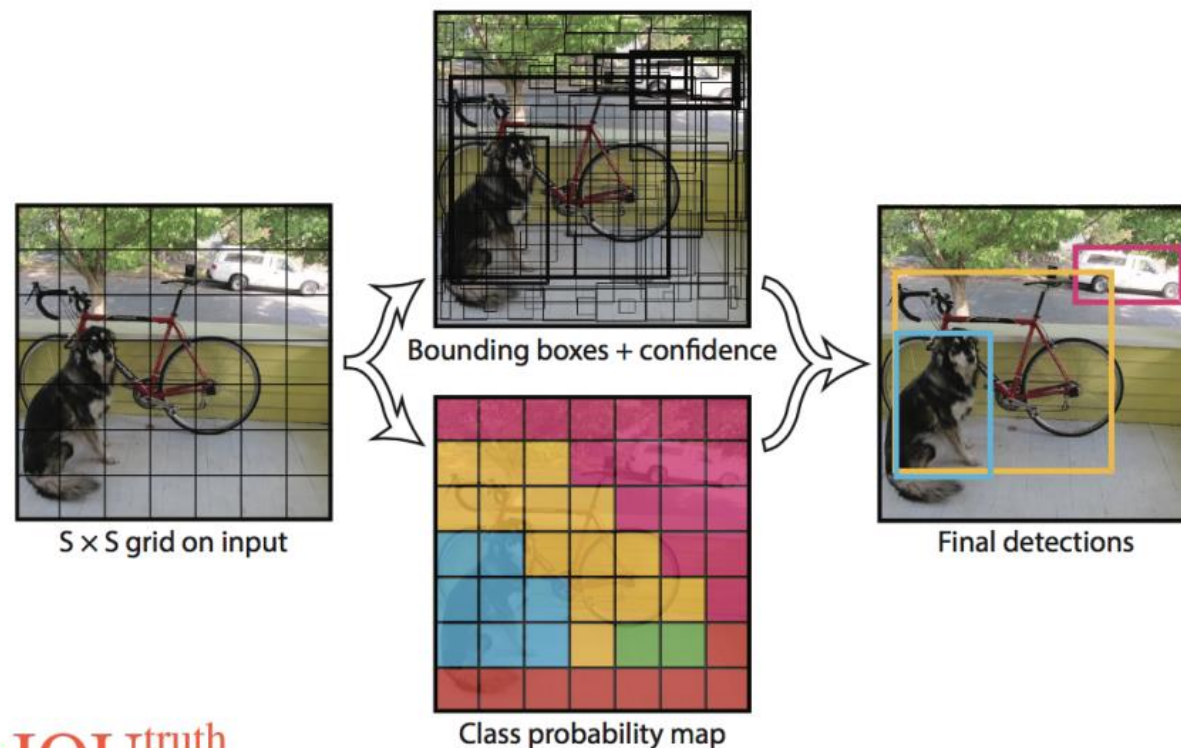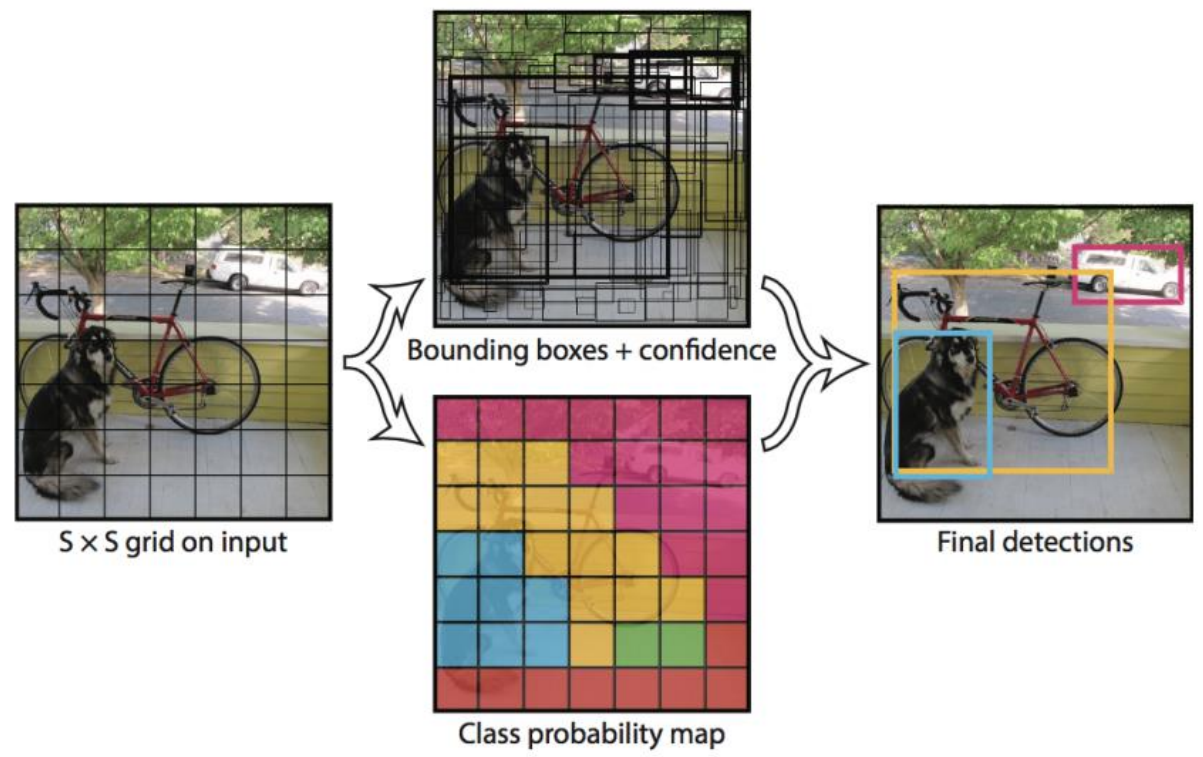
Class probability map
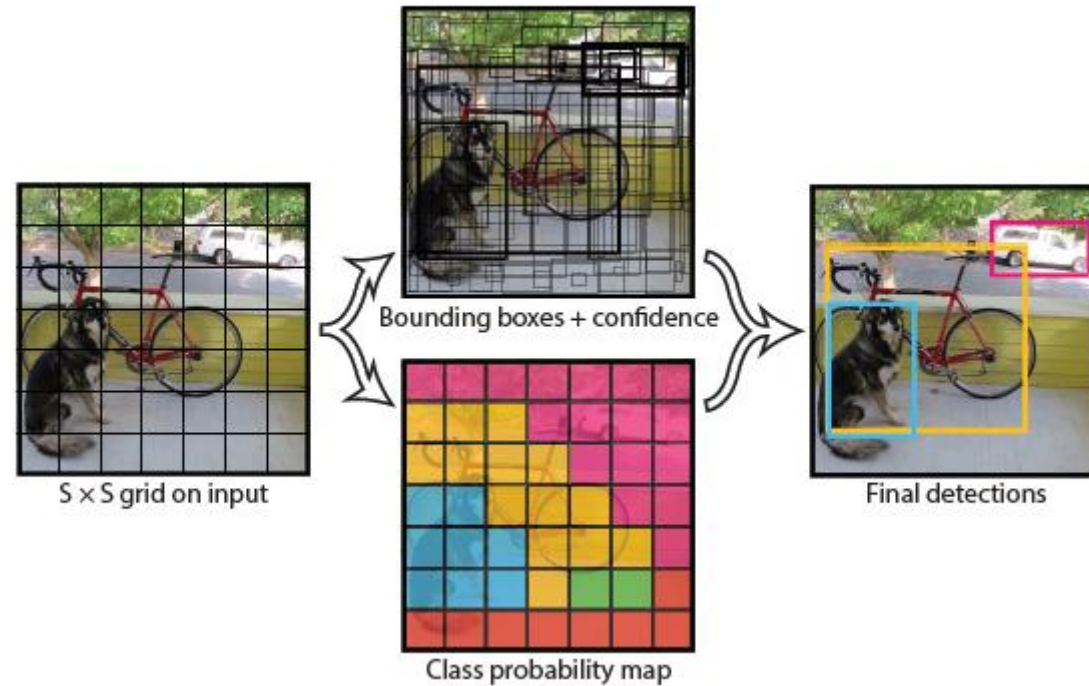
Final detections

**Figure 2: The Model.** Our system models detection as a regression problem. It divides the image into an $S \times S$ grid and for each grid cell predicts $B$ bounding boxes, confidence for those boxes, and $C$ class probabilities. These predictions are encoded as an $S \times S \times (B * 5 + C)$ tensor.

# Unified Detection

- Predict **one set of class probabilities** per grid cell, **regardless of the number of boxes B**.

- At test time, individual box confidence prediction

$$\text{Pr}(\text{Class}_i\,|\text{Object})*\text{Pr}(\text{Object})*\text{IOU}_{\text{pred}}^{\text{truth}}$$

$$=\text{Pr}(\text{Class}_i)*\text{IOU}_{\text{pred}}^{\text{truth}}$$



S × S grid on input

Bounding boxes + confidence
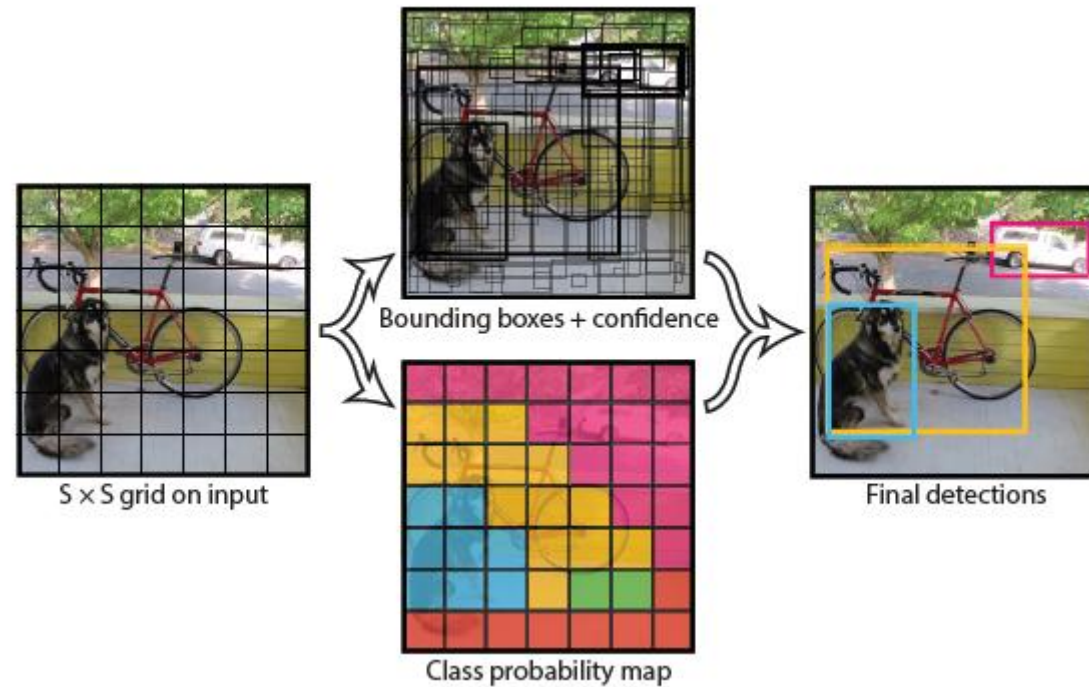
Class probability map

Final detections

**Figure 2: The Model.** Our system models detection as a regression problem. It divides the image into an $S \times S$ grid and for each grid cell predicts $B$ bounding boxes, confidence for those boxes, and $C$ class probabilities. These predictions are encoded as an $S \times S \times (B * 5 + C)$ tensor.

# How unified detection works?



Bounding boxes + confidence

S × S grid on input

Class probability map

Final detections

- uses features from the entire image to predict each bounding box
- predicts all bounding boxes across all classes for an image simultaneously?
- ➤ divides the input image into an s*s grid. If the center of an object falls into a grid cell, the cell is responsible for detecting that object.
- ➤ each grid cell predicts B bounding boxes and confidence scores for those boxes
- ➤ Each grid also predicts C conditional(conditioned on the grid cell containing an object) class probabilities
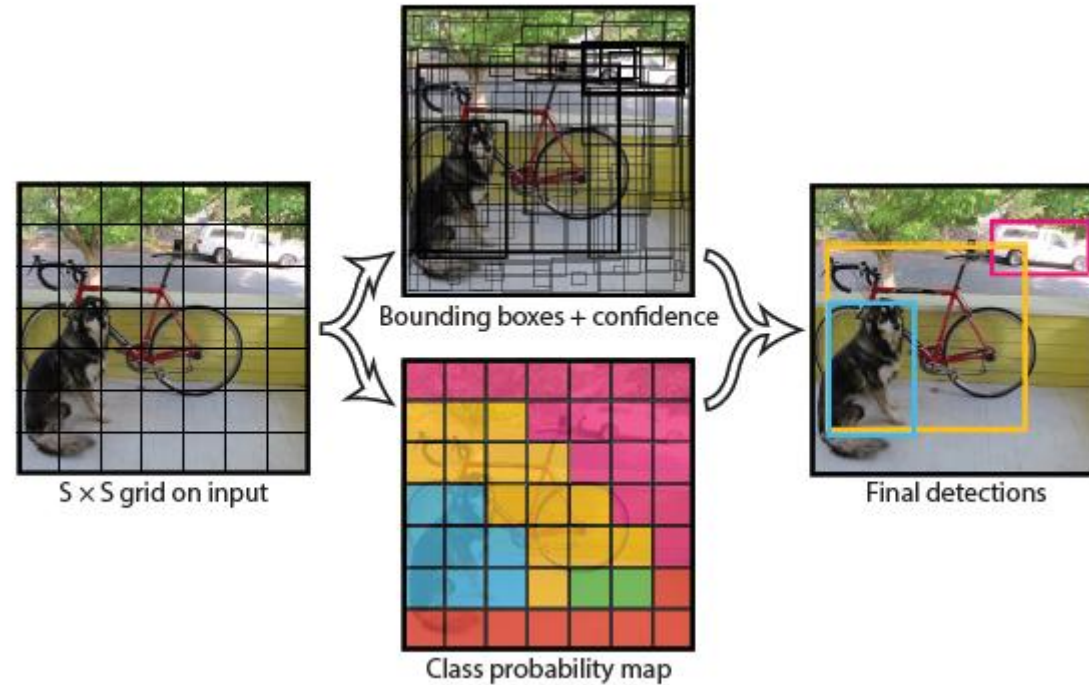
# How unified detection works?



S × S grid on input

Bounding boxes + confidence

Class probability map

Final detections

**confidence scores**: reflect how confident is that the box contains an object+how accurate the box is .

$$Pr(Object) * IOU_{pred}^{truth}$$

**conditional class probabilities:** conditioned on the grid cell containing an object
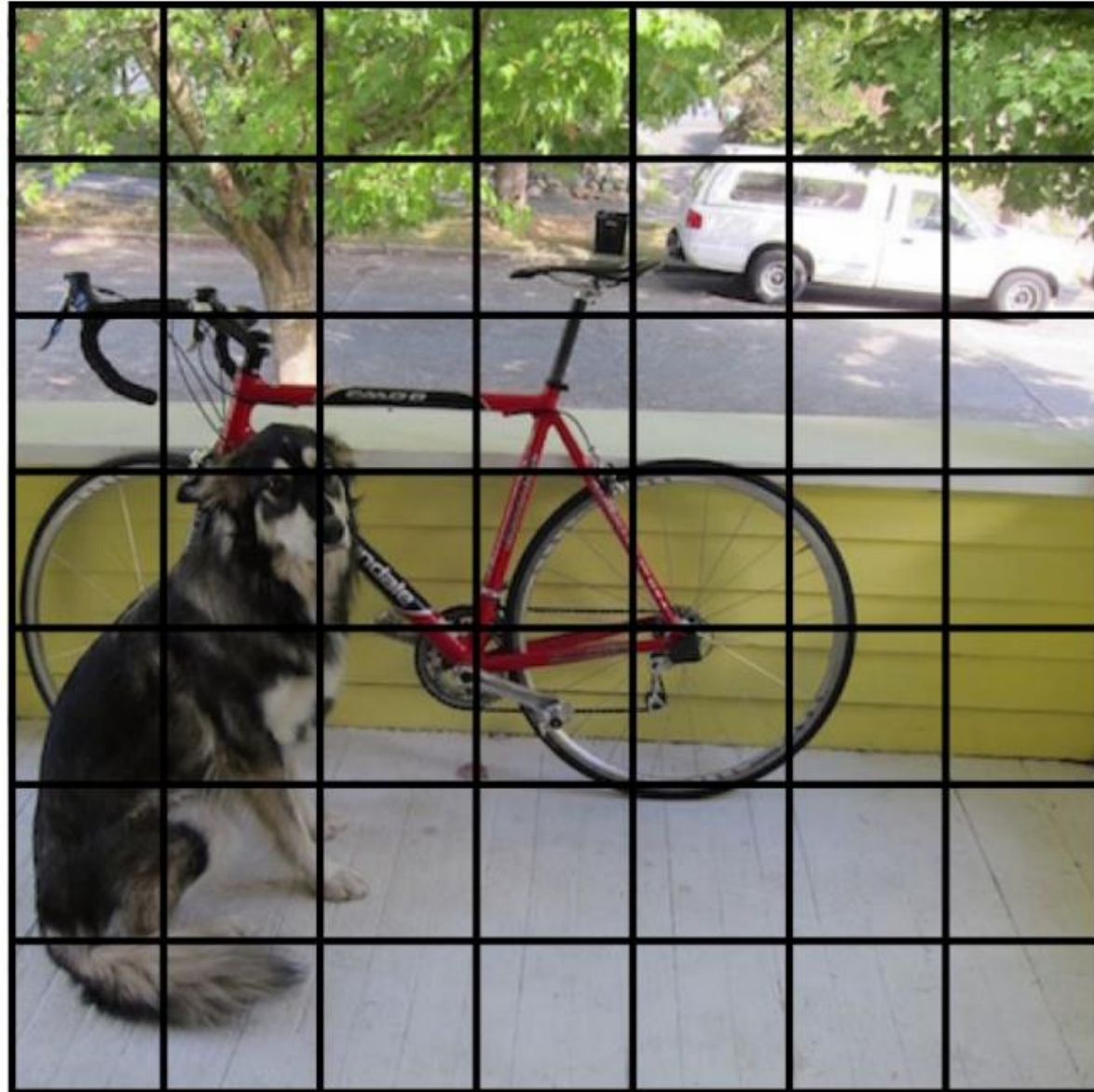
$$Pr(Class_i|Object)$$

# How unified detection works?



$$\text{Pr}(\text{Class}_i | \text{Object}) * \text{Pr}(\text{Object}) * \text{IOU}_{\text{pred}}^{\text{truth}} = \text{Pr}(\text{Class}_i) * \text{IOU}_{\text{pred}}^{\text{truth}}$$

- At test time, multiply the conditional class probabilities and the individual box confidence predictions
- giving class-specific confidence scores for each box
- Showing both the probability of that class appearing in the box and how well the predicted box fits the object

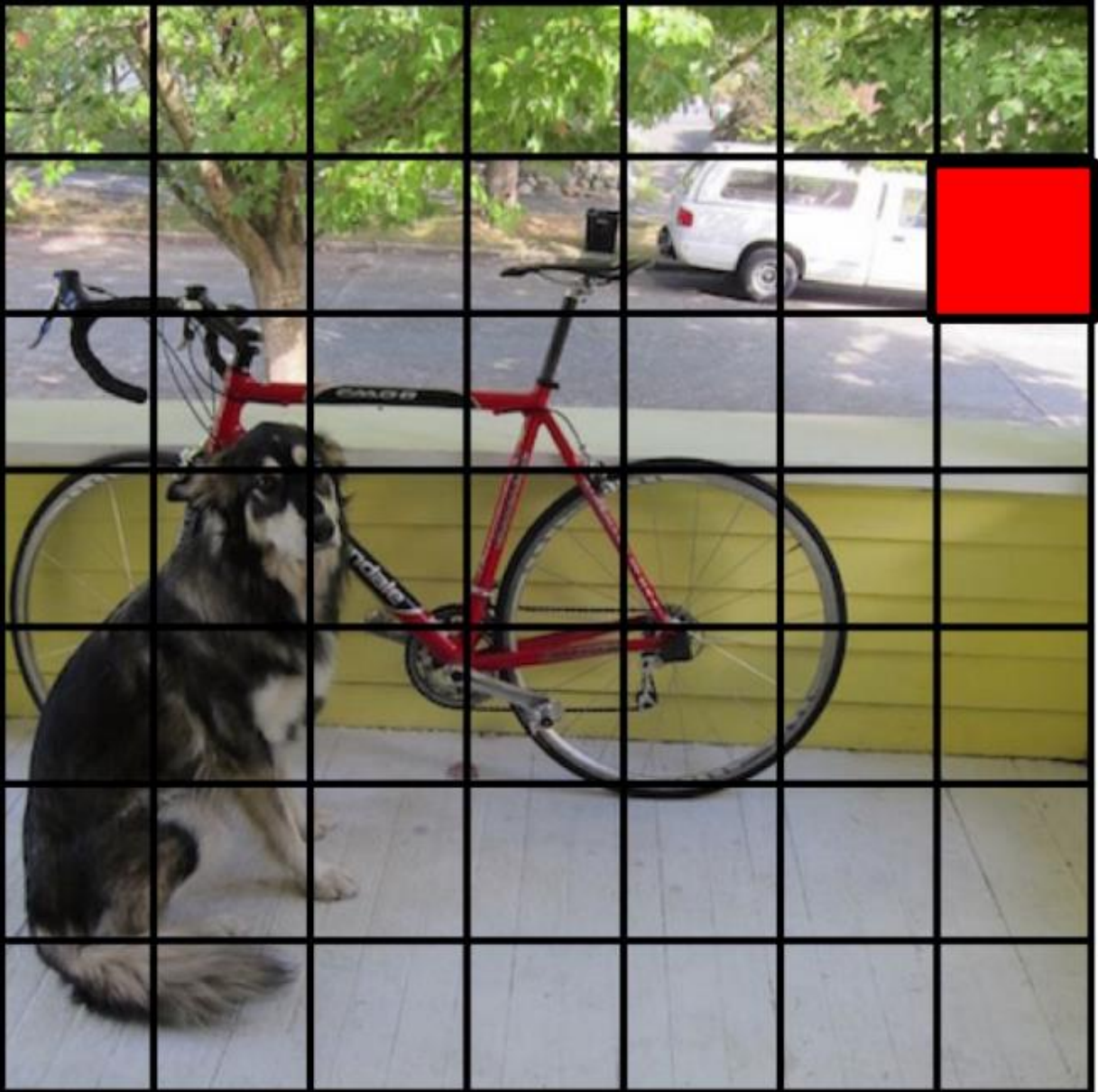# Automatic Object Detection from Image

# We split the image into an S*S grid



**7*7 grids**

# Each cell predicts B boxes(x,y,w,h) and confidences of each box: P(Object)

# Each cell predicts B boxes(x,y,w,h) and confidences of each box: P(Object)
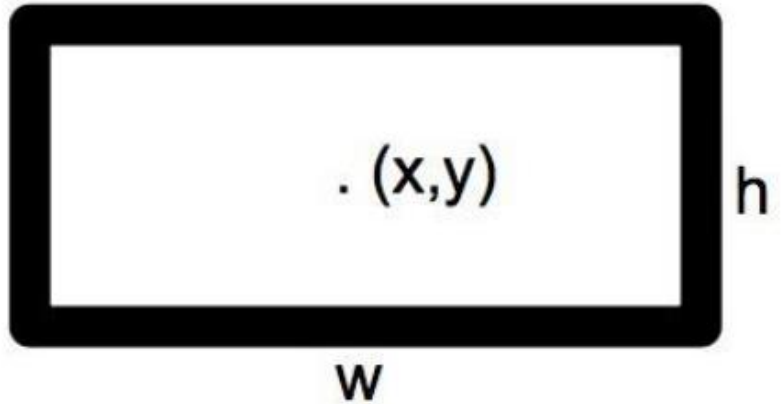
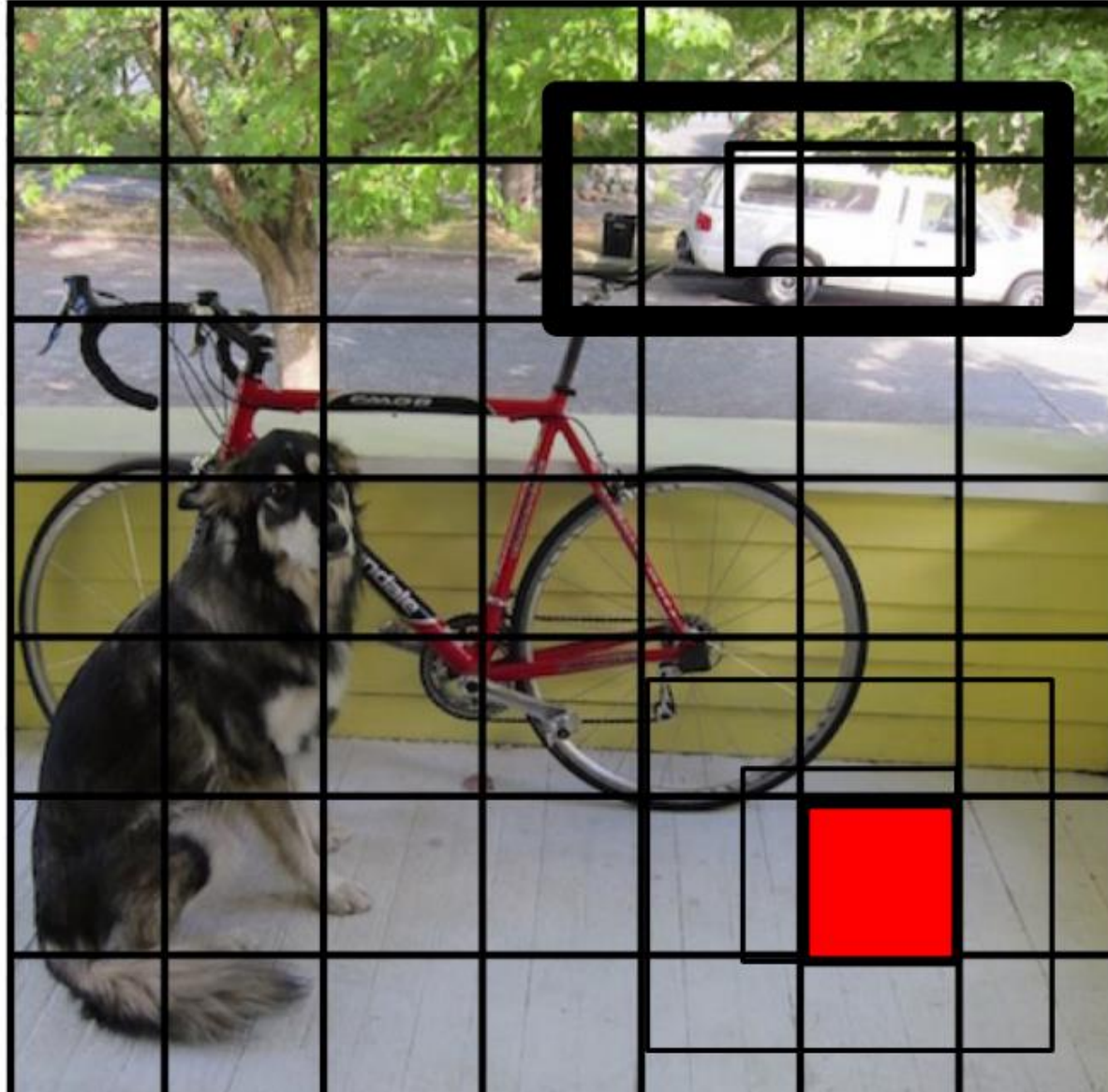# Each cell predicts B boxes(x,y,w,h) and confidences of each box: P(Object)

B = 2

each box predict:

. (x,y)

h

w

P(Object): probability that the box contains an object

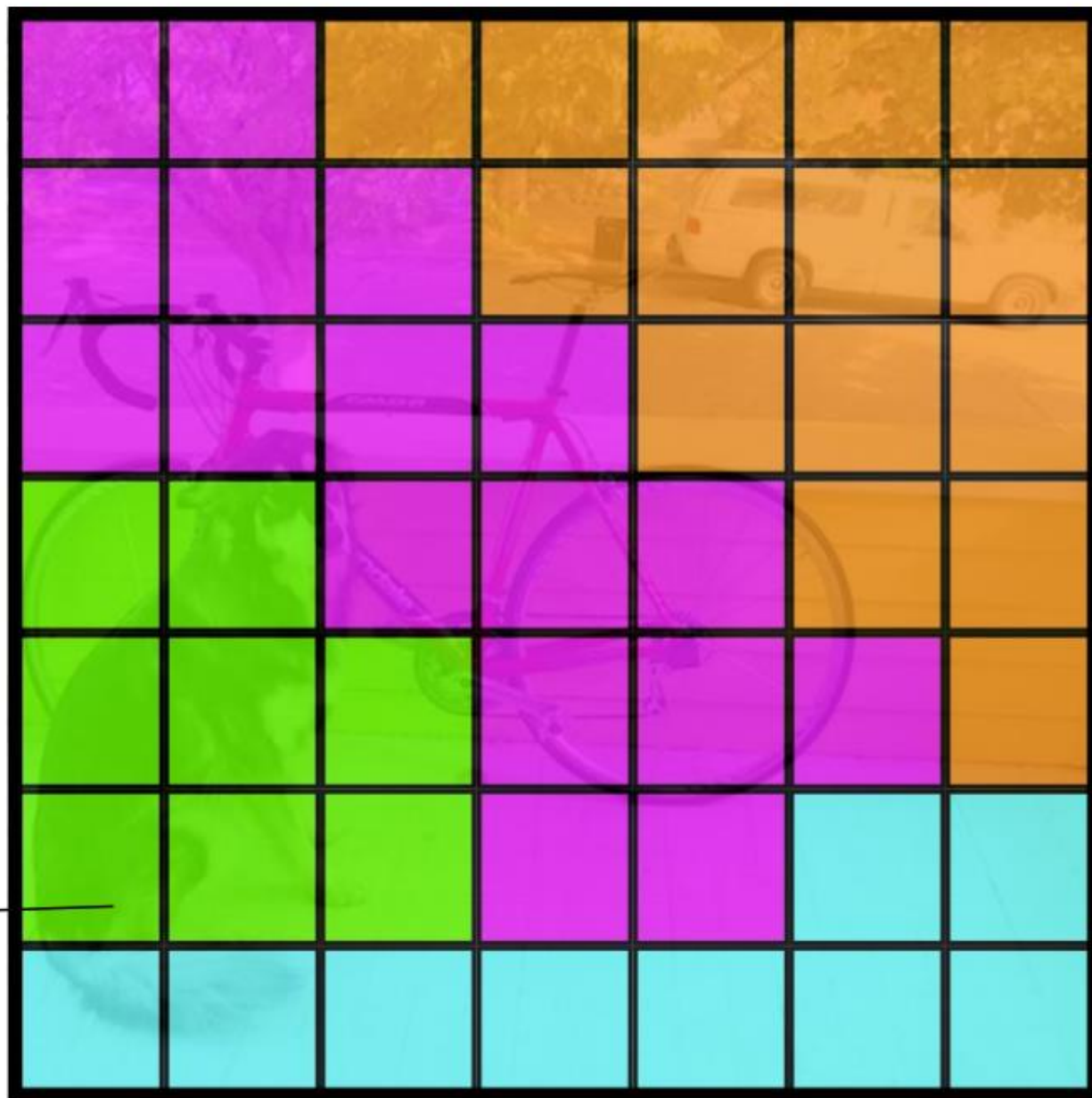# Each cell predicts B boxes(x,y,w,h) and confidences of each box: P(Object)

# Each cell predicts boxes and confidences: P(Object)

# Each cell also predicts a class probability.



Bicycle

Car

Dog

Dining Table

# Conditioned on object: P(Car | Object)



Bicycle

Car

Dog

**Eg.**

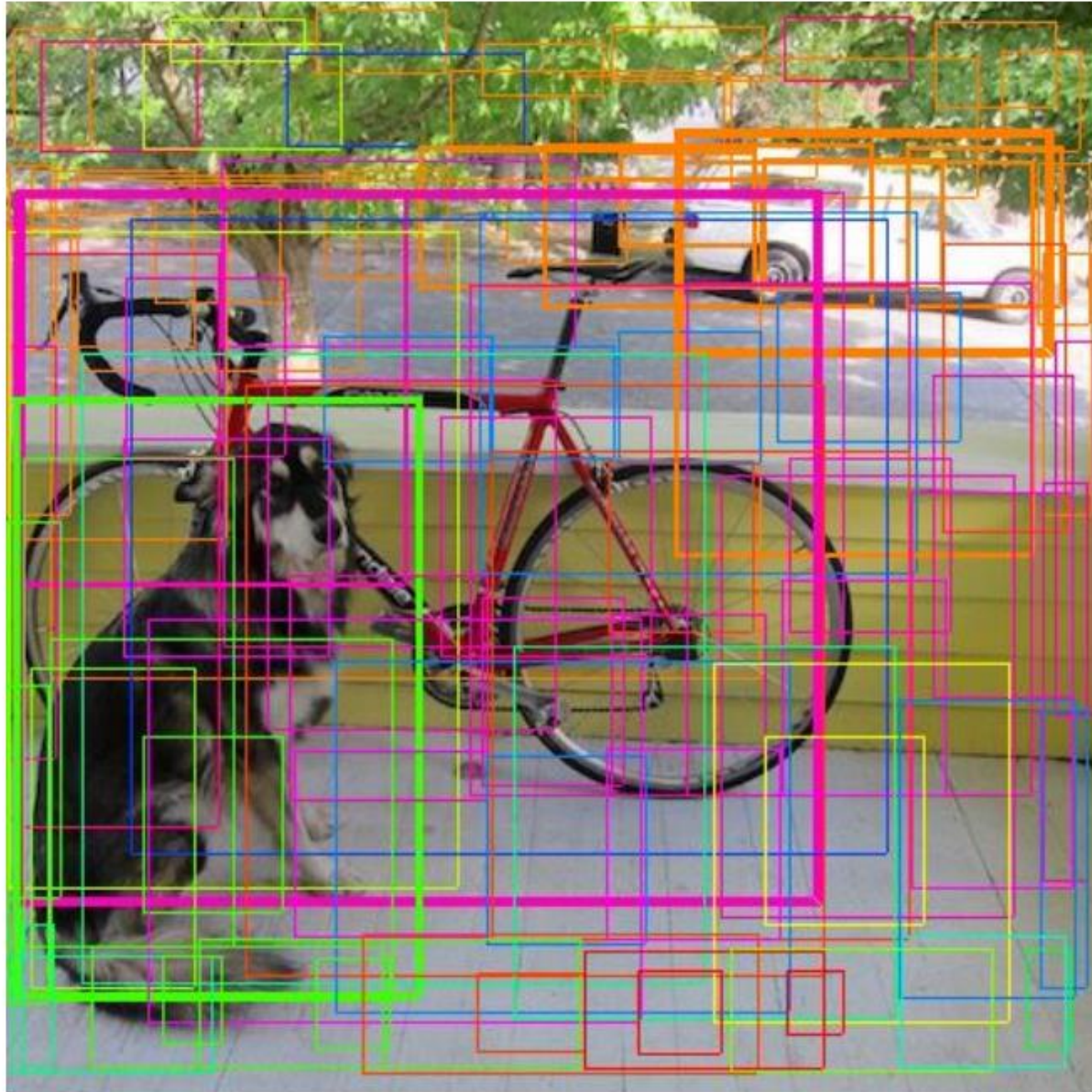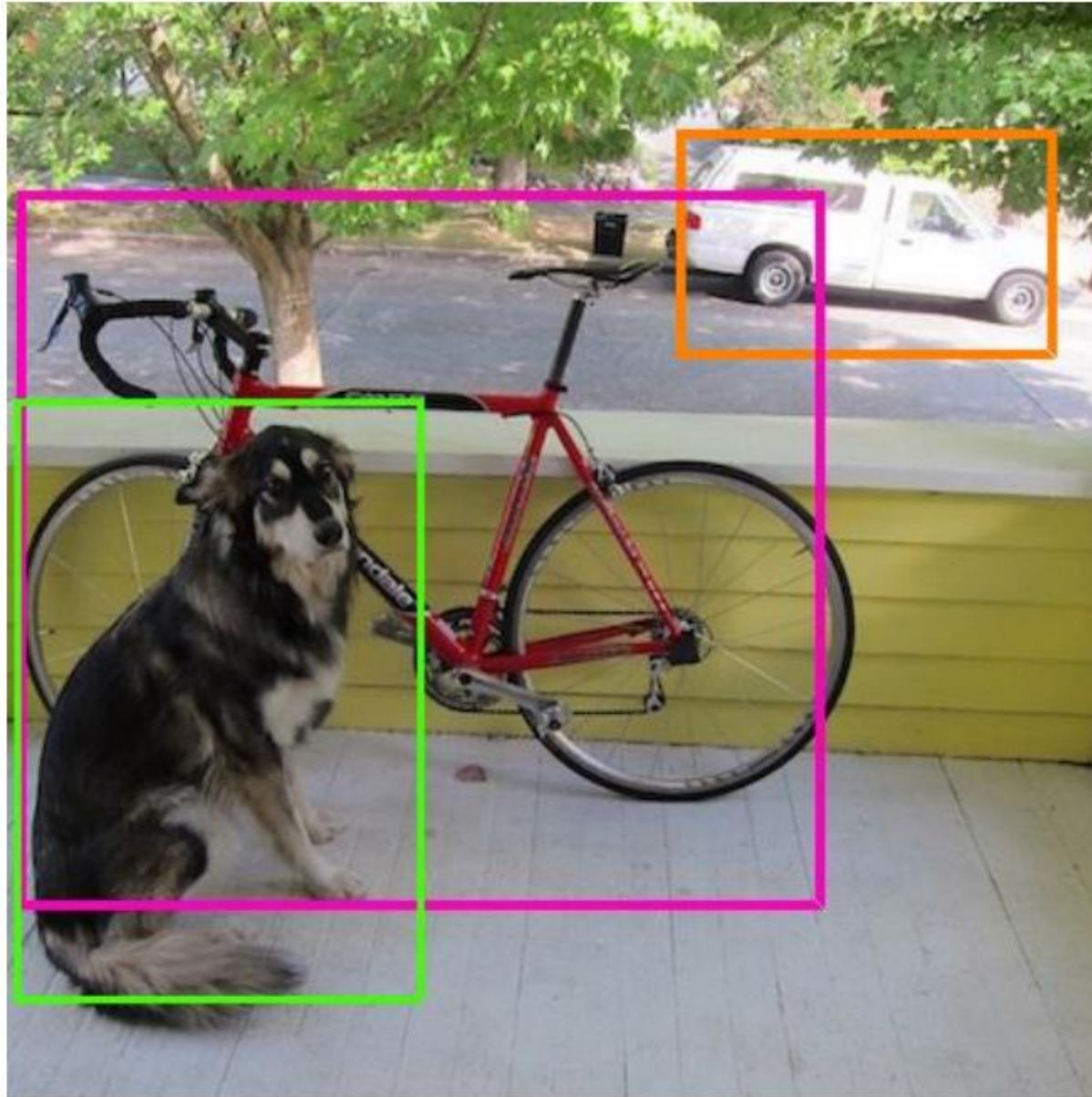**Dog = 0.8**

Cat = 0

Bike = 0

Dining
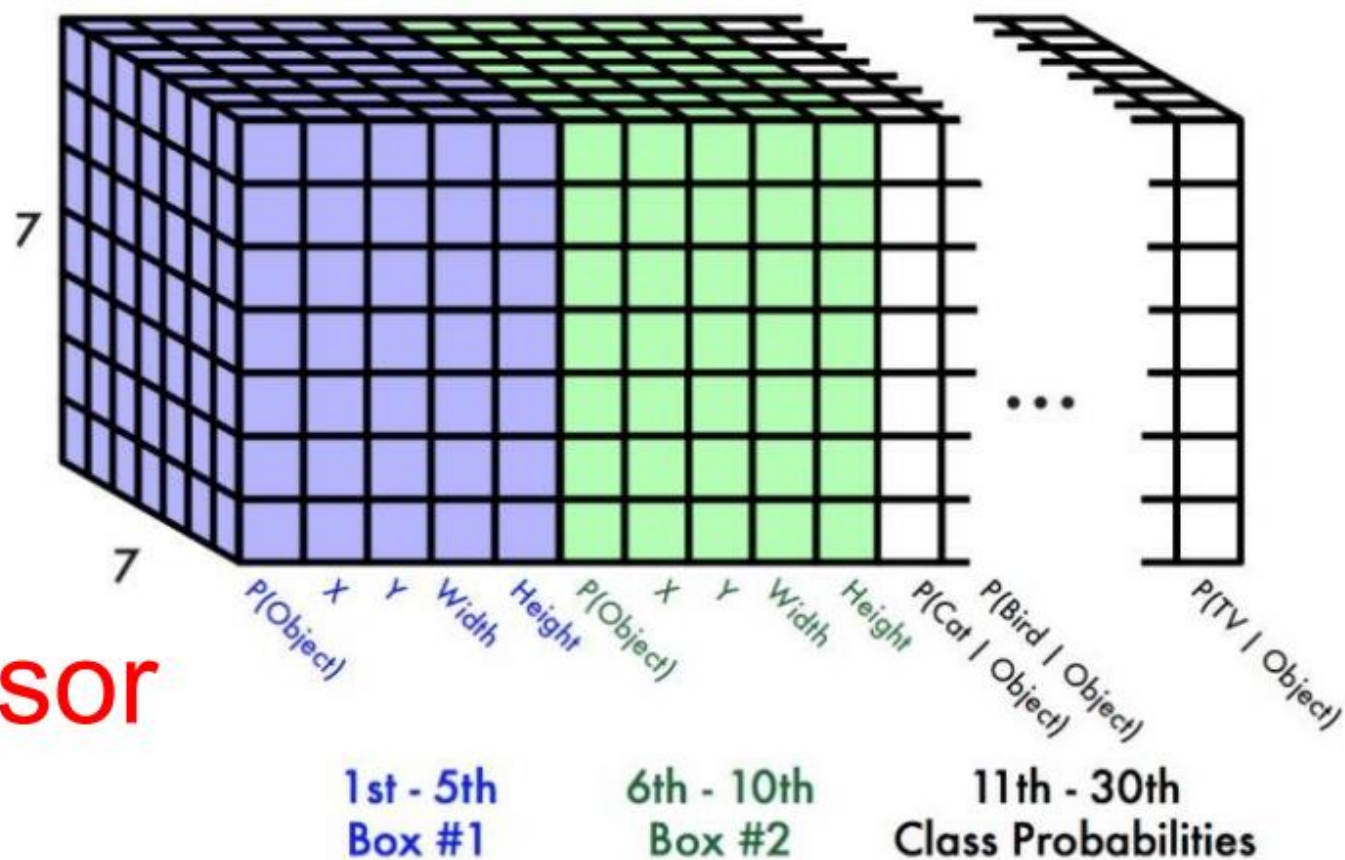Table

# Then we combine the box and class predictions.



P(class|Object) * P(Object)
=P(class)

# Finally we do threshold detections and NMS

Each cell predicts:

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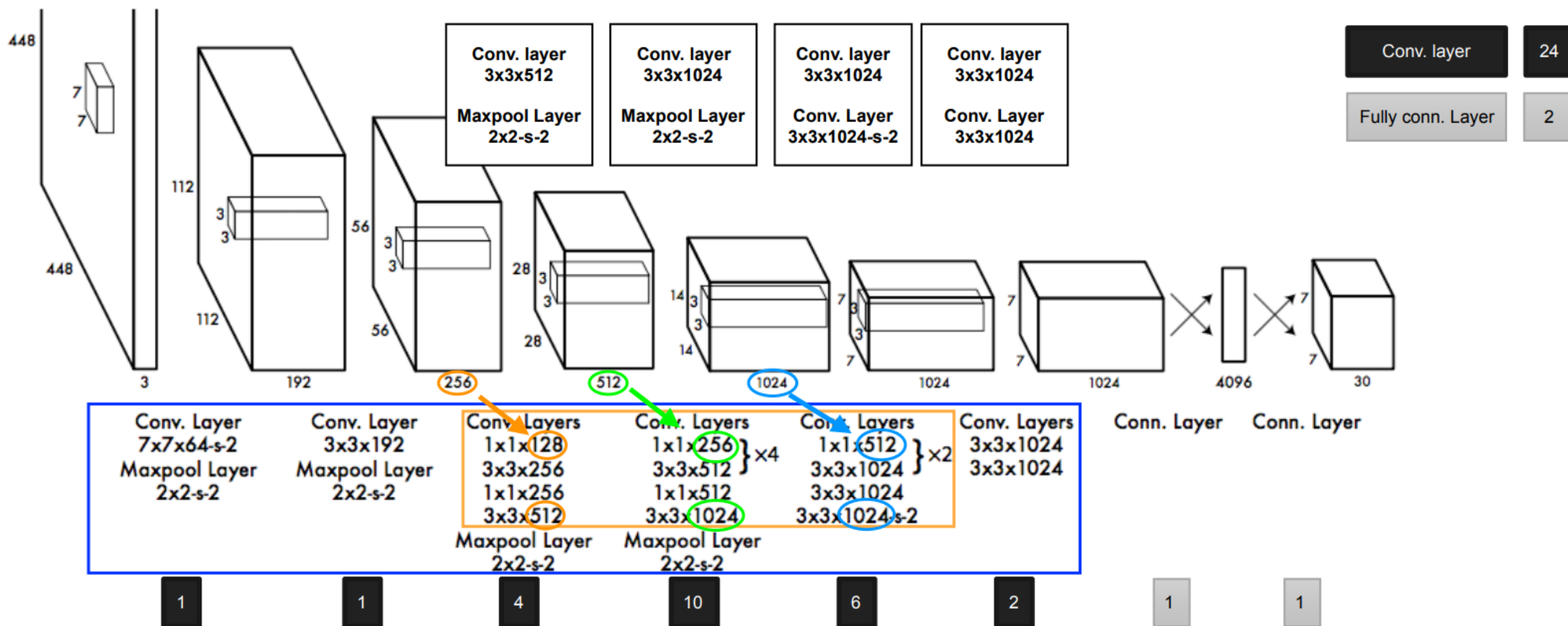    - 1 confidence value
- Some number of class probabilities

S * S * (B * 5 + C) tensor

1st - 5th Box #1    6th - 10th Box #2    11th - 30th Class Probabilities
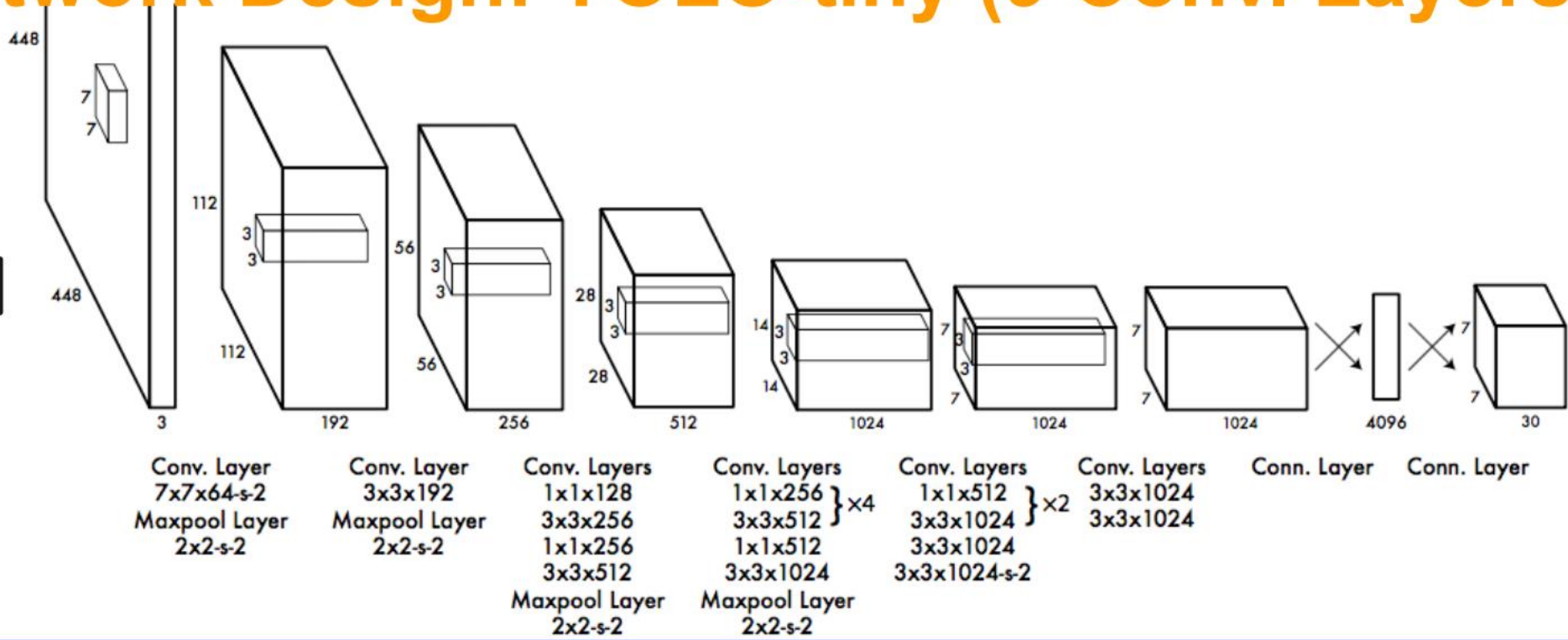
# Network Design: YOLO

- Modified GoogLeNet

- 1x1 reduction layer ("Network in Network")

Our network architecture is inspired by the GoogLeNet model for image classification [34]. Our network has 24 convolutional layers followed by 2 fully connected layers. Instead of the inception modules used by GoogLeNet, we simply use 1 × 1 reduction layers followed by 3 × 3 convolutional layers, similar to Lin et al [22]. The full network is shown in Figure 3.
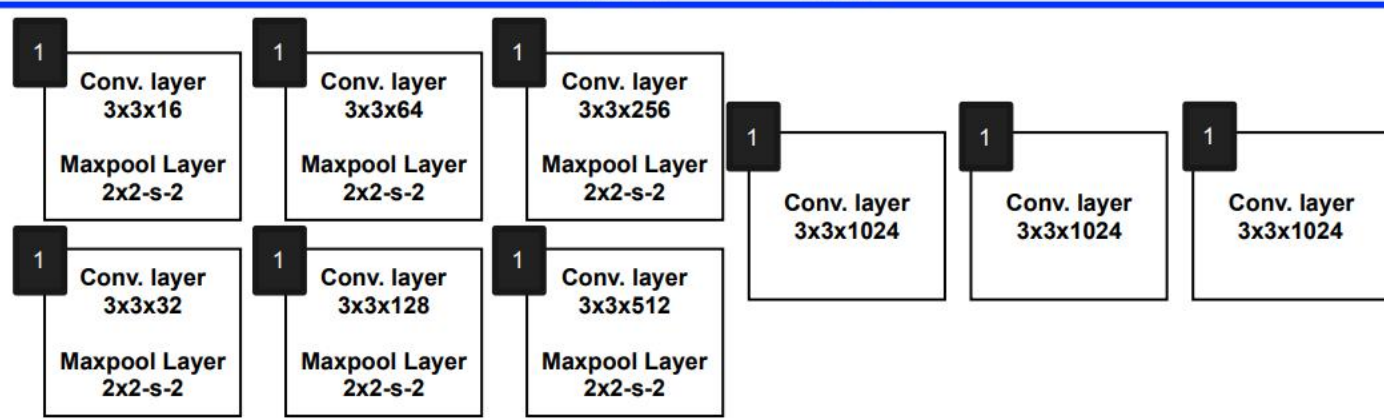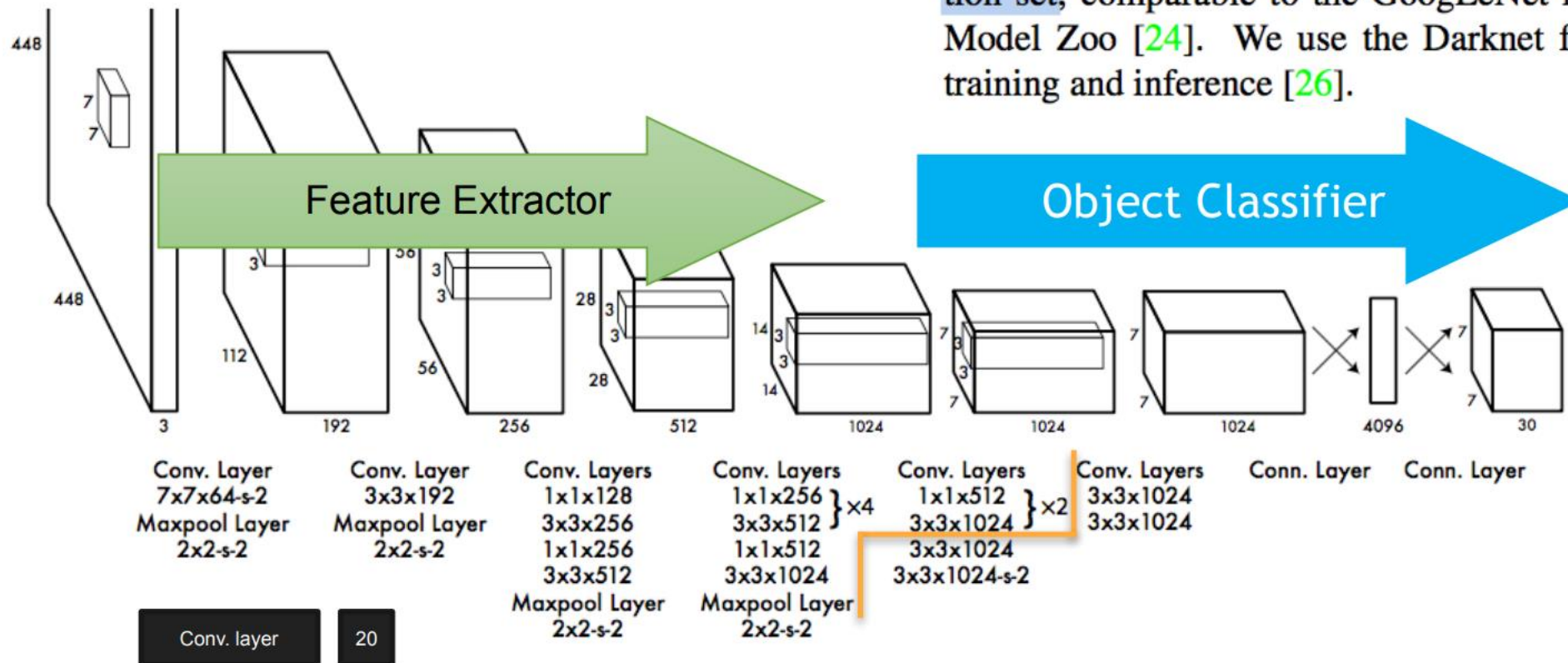
# Network Design: YOLO-tiny (9 Conv. Layers)



**YOLO**

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| Conv. Layer<br>7x7x64-s-2<br>Maxpool Layer<br>2x2-s-2 | Conv. Layer<br>3x3x192<br>Maxpool Layer<br>2x2-s-2 | Conv. Layers<br>1x1x128<br>3x3x256<br>1x1x256<br>3x3x512<br>Maxpool Layer<br>2x2-s-2 | Conv. Layers<br>1x1x256<br>3x3x512 }×4<br>1x1x512<br>3x3x1024<br>Maxpool Layer<br>2x2-s-2 | Conv. Layers<br>1x1x512<br>3x3x1024 }×2<br>3x3x1024<br>3x3x1024-s-2 | Conv. Layers<br>3x3x1024<br>3x3x1024 | Conn. Layer | Conn. Layer |

**YOLO-Tiny**

| 1 Conv. layer<br>3x3x16<br>Maxpool Layer<br>2x2-s-2 | 1 Conv. layer<br>3x3x64<br>Maxpool Layer<br>2x2-s-2 | 1 Conv. layer<br>3x3x256<br>Maxpool Layer<br>2x2-s-2 | 1 Conv. layer<br>3x3x1024 | 1 Conv. layer<br>3x3x1024 | 1 Conv. layer<br>3x3x1024 |
|---|---|---|---|---|---|
| 1 Conv. layer<br>3x3x32<br>Maxpool Layer<br>2x2-s-2 | 1 Conv. layer<br>3x3x128<br>Maxpool Layer<br>2x2-s-2 | 1 Conv. layer<br>3x3x512<br>Maxpool Layer<br>2x2-s-2 | | | |

| Conv. layer | 9 |
|---|---|
| Fully conn. Layer | 2 |

Slide credit to Taegyun Jeon

# Training

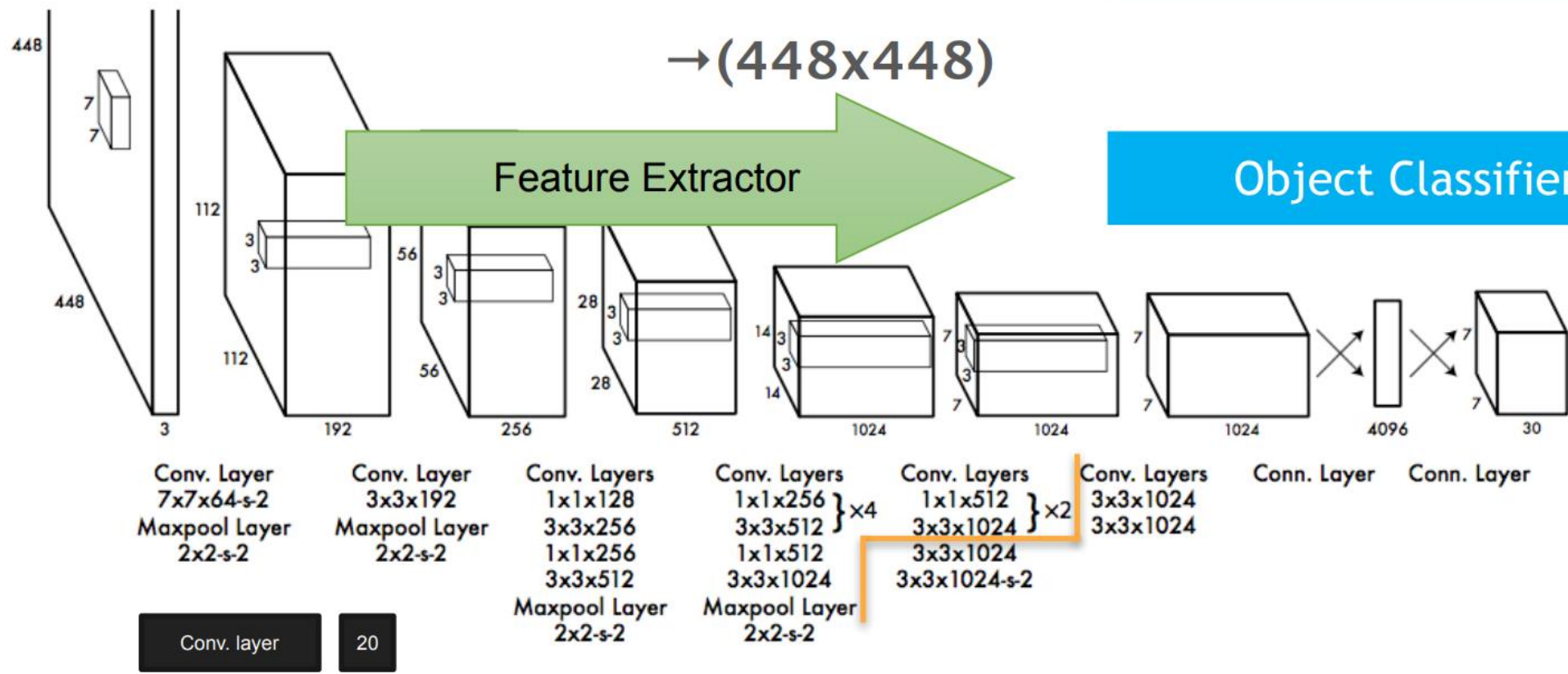1) Pretrain with ImageNet 1000-class competition dataset

We pretrain our convolutional layers on the ImageNet 1000-class competition dataset [30]. For pretraining we use the first 20 convolutional layers from Figure 3 followed by a average-pooling layer and a fully connected layer. We train this network for approximately a week and achieve a single crop top-5 accuracy of 88% on the ImageNet 2012 validation set, comparable to the GoogLeNet models in Caffe's Model Zoo [24]. We use the Darknet framework for all training and inference [26].
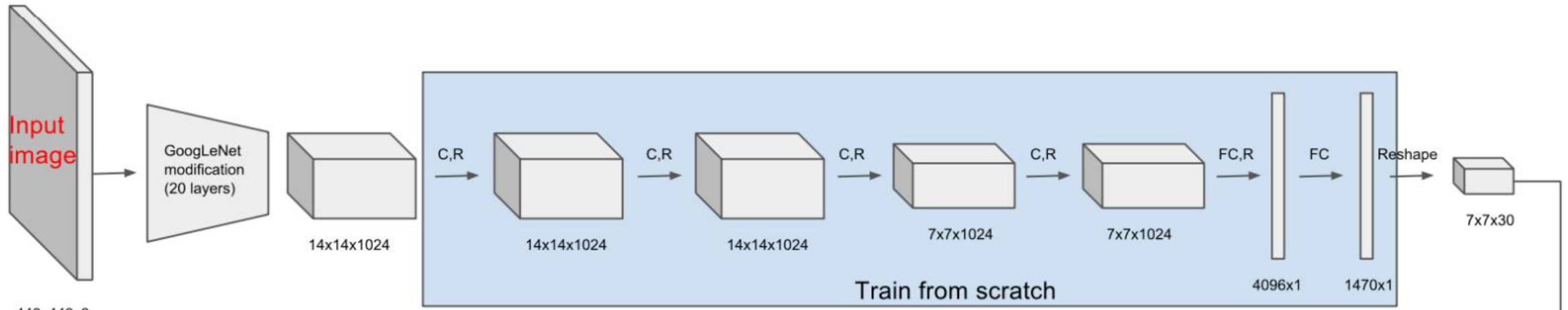
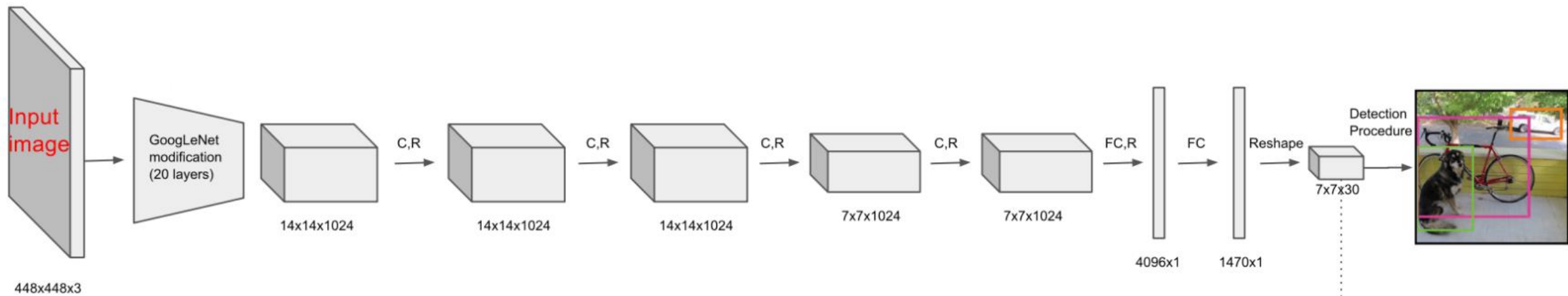# Training

2) "Network on Convolutional Feature Maps"

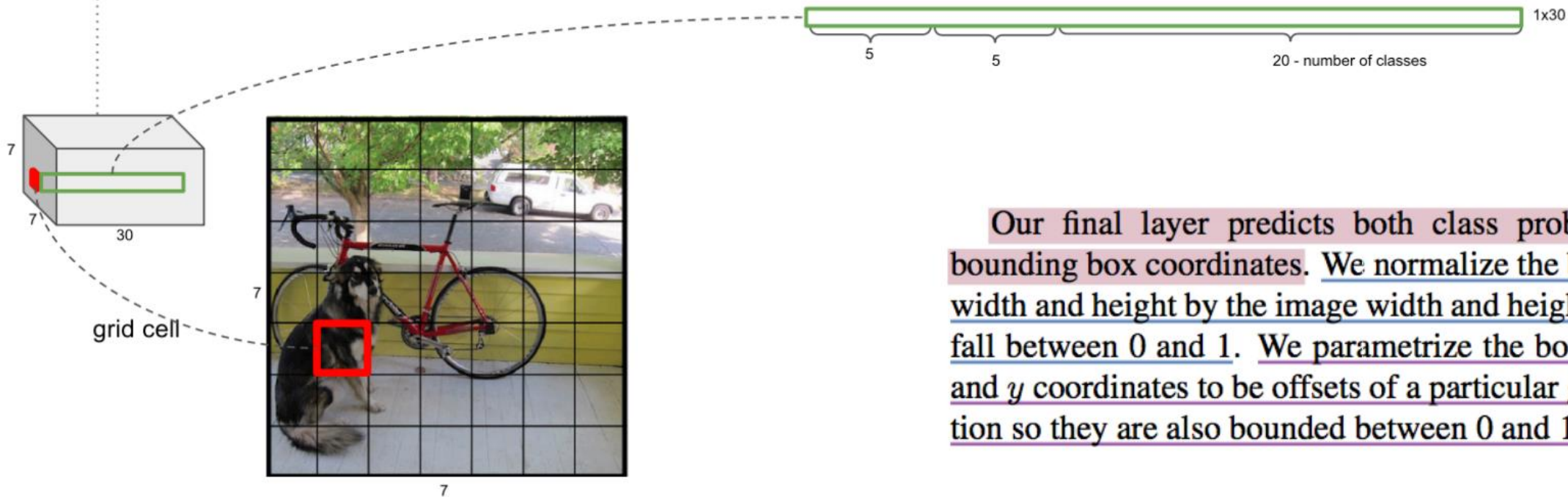Increased input resolution (224x224)

We then convert the model to perform detection. Ren et al. show that adding both convolutional and connected layers to pretrained networks can improve performance [29]. Following their example, we add four convolutional layers and two fully connected layers with randomly initialized weights. Detection often requires fine-grained visual information so we increase the input resolution of the network from 224 × 224 to 448 × 448.
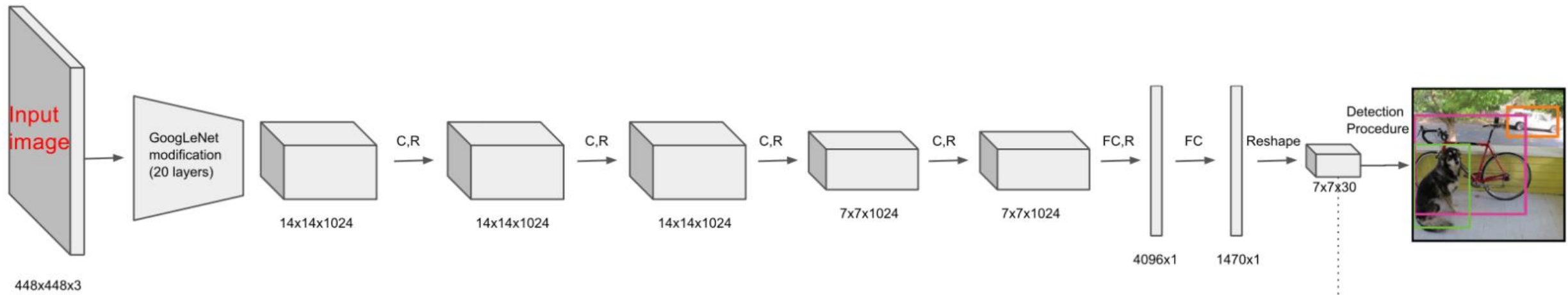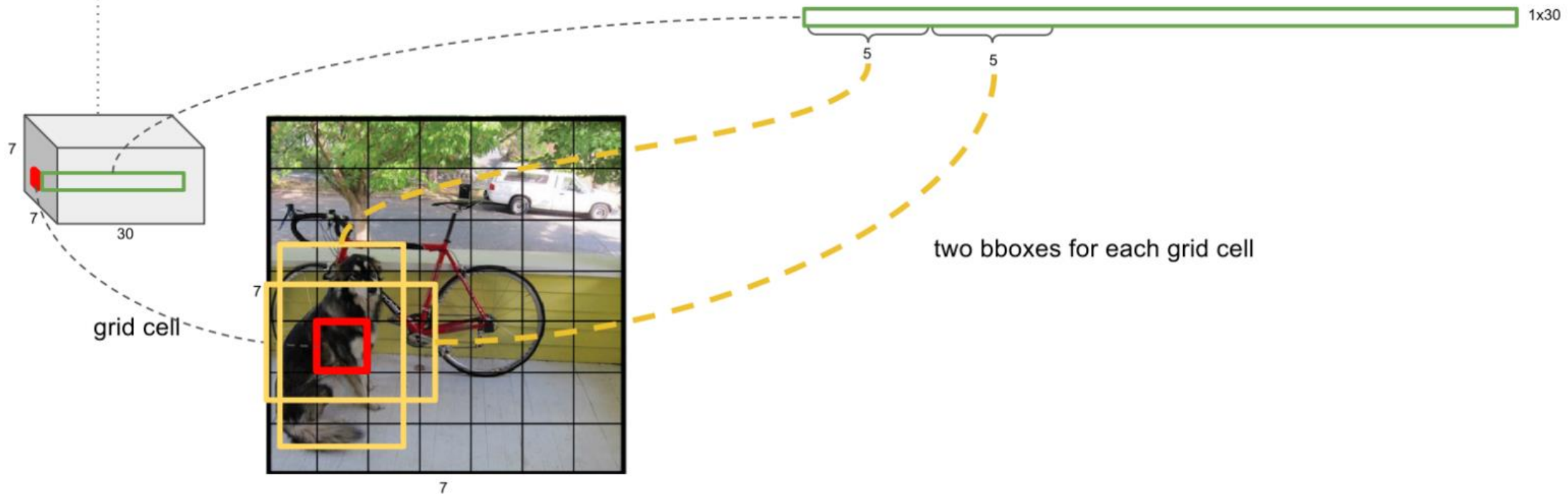
→(448x448)



| Conv. layer | 20 |

| Conv. layer | 4 |
| Fully conn. Layer | 2 |

Input image

448x448x3

GoogLeNet modification (20 layers)

14x14x1024

C,R

14x14x1024

C,R

14x14x1024

C,R

7x7x1024

C,R

7x7x1024

FC,R

4096x1

FC

1470x1

Reshape

7x7x30

Train from scratch

Detection Procedure

Slide credit to Taegyun Jeon

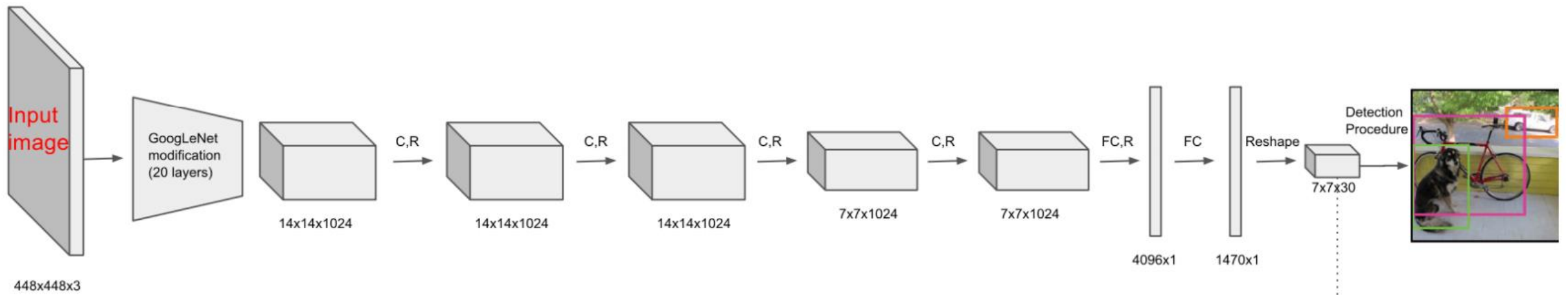Tensor values interpretation

5     5     20 - number of classes

1x30

grid cell

Our final layer predicts both class probabilities and bounding box coordinates. We normalize the bounding box width and height by the image width and height so that they fall between 0 and 1. We parametrize the bounding box $x$ and $y$ coordinates to be offsets of a particular grid cell location so they are also bounded between 0 and 1.
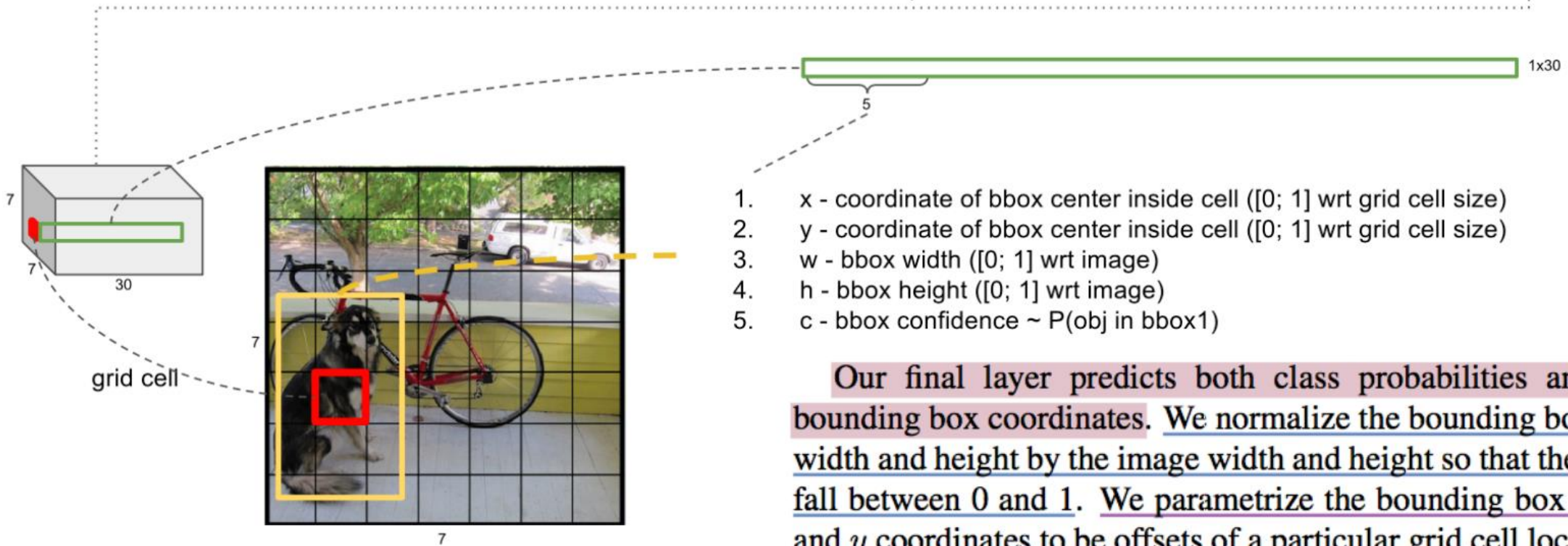
Input image

448x448x3

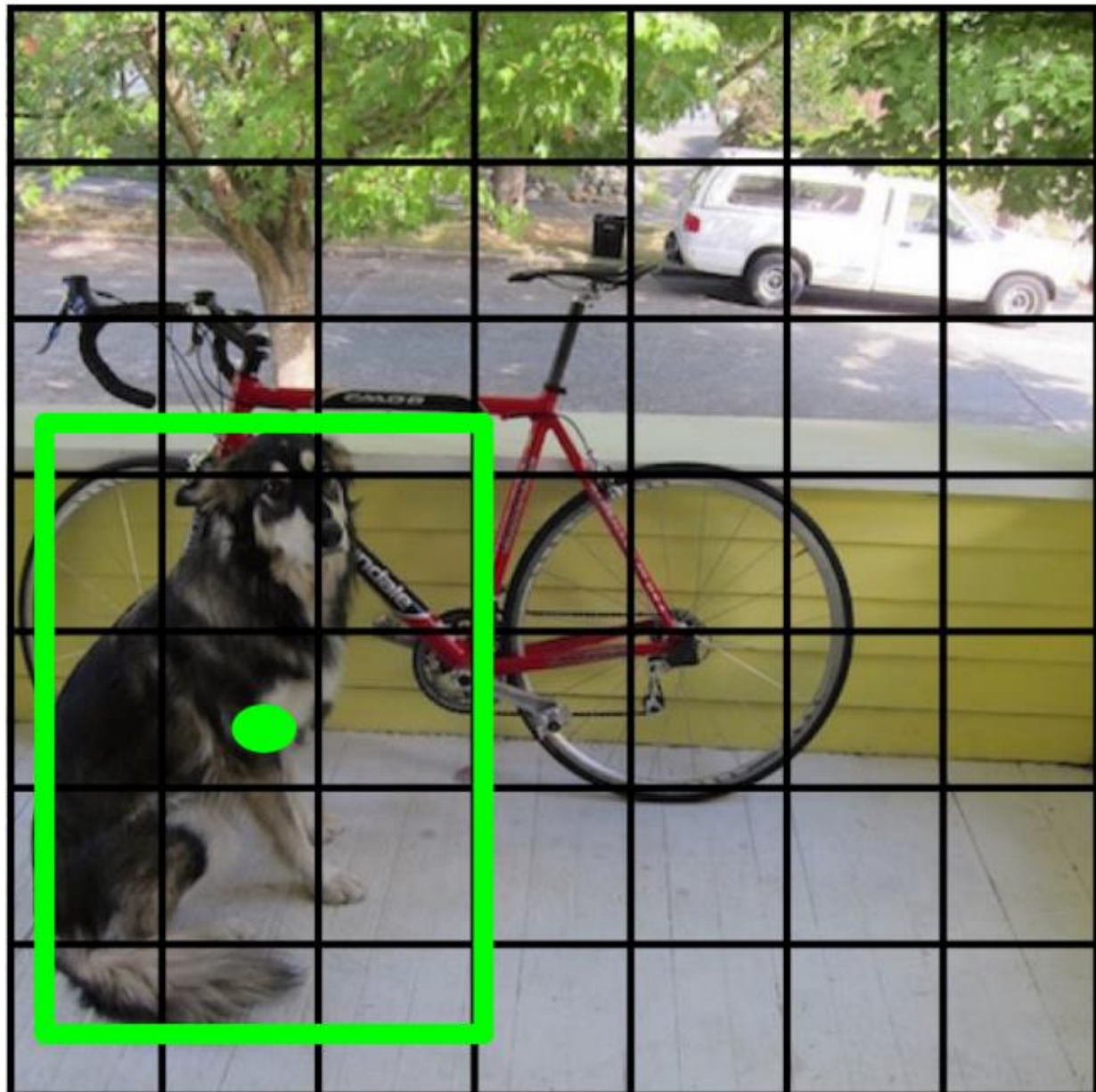GoogLeNet modification (20 layers)

14x14x1024    C,R    14x14x1024    C,R    14x14x1024    C,R    7x7x1024    C,R    7x7x1024    FC,R    FC    Reshape

4096x1    1470x1

7x7x30

Detection Procedure

Tensor values interpretation

1x30

5    5

two bboxes for each grid cell

7    7    30

7    7

grid cell

7

Slide credit to Taegyun Jeon

Detection Procedure

448x448x3 → GoogLeNet modification (20 layers) → 14x14x1024 → C,R → 14x14x1024 → C,R → 14x14x1024 → C,R → 7x7x1024 → C,R → 7x7x1024 → FC,R → 4096x1 → FC → 1470x1 → Reshape → 7x7x30

Tensor values interpretation

1x30

5

grid cell

1.  x - coordinate of bbox center inside cell ([0; 1] wrt grid cell size)
2.  y - coordinate of bbox center inside cell ([0; 1] wrt grid cell size)
3.  w - bbox width ([0; 1] wrt image)
4.  h - bbox height ([0; 1] wrt image)
5.  c - bbox confidence ~ P(obj in bbox1)

Our final layer predicts both class probabilities and bounding box coordinates. We normalize the bounding box width and height by the image width and height so that they fall between 0 and 1. We parametrize the bounding box $x$ and $y$ coordinates to be offsets of a particular grid cell location so they are also bounded between 0 and 1.

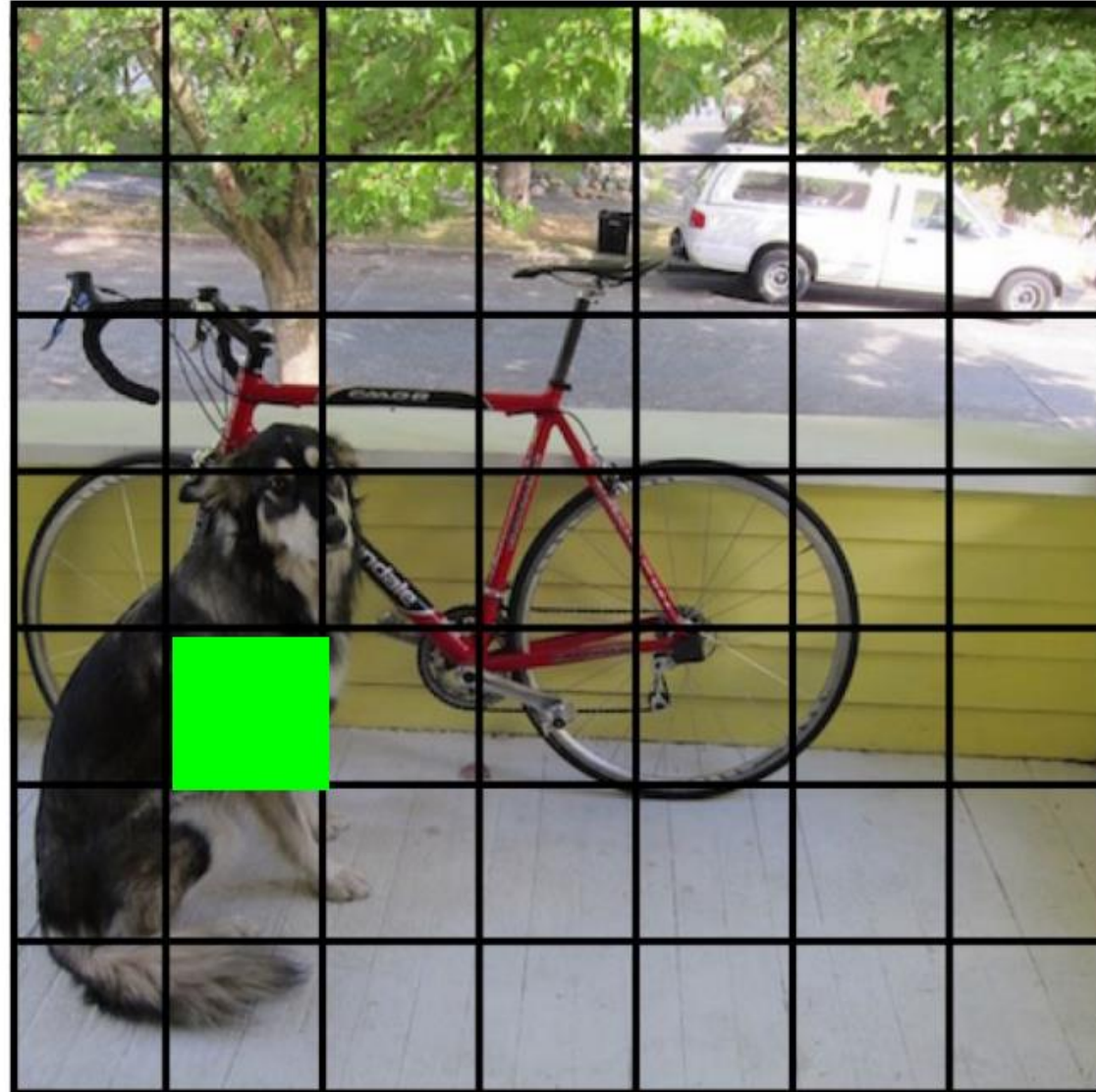# During training, match example to the right cell

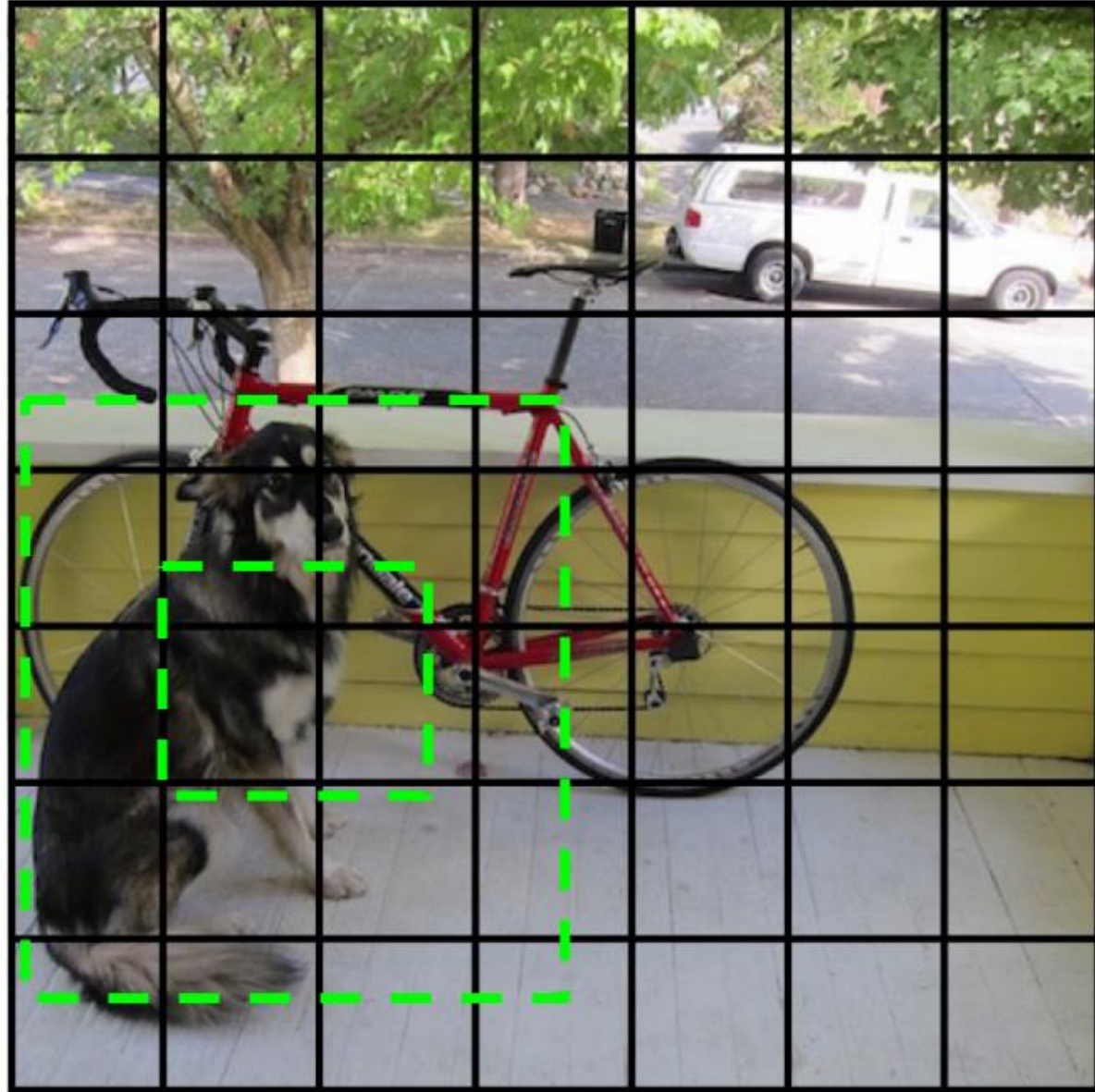# During training, match example to the right cell
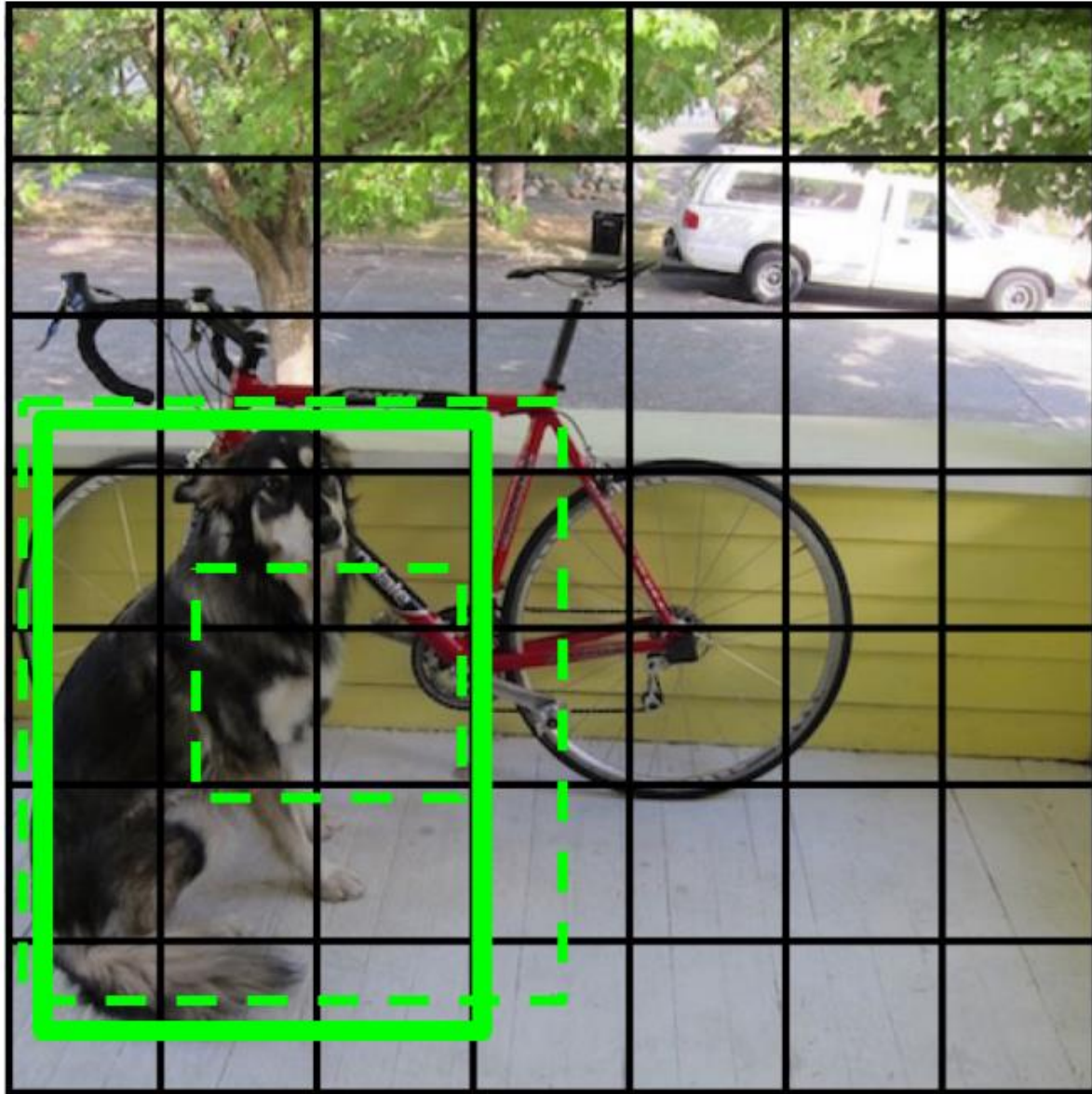
# Adjust that cell's class prediction



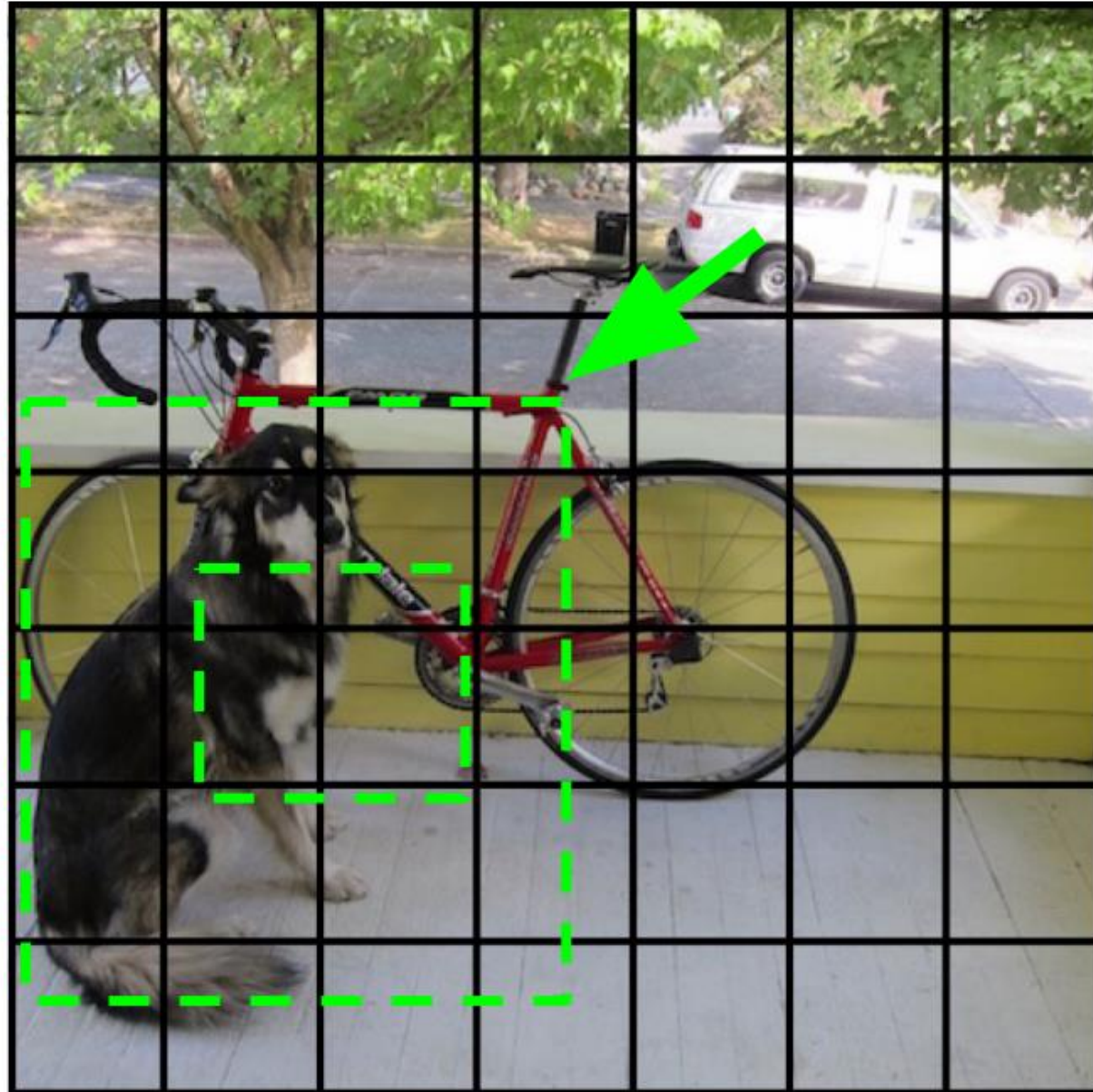**Dog = 1**
Cat = 0
Bike = 0
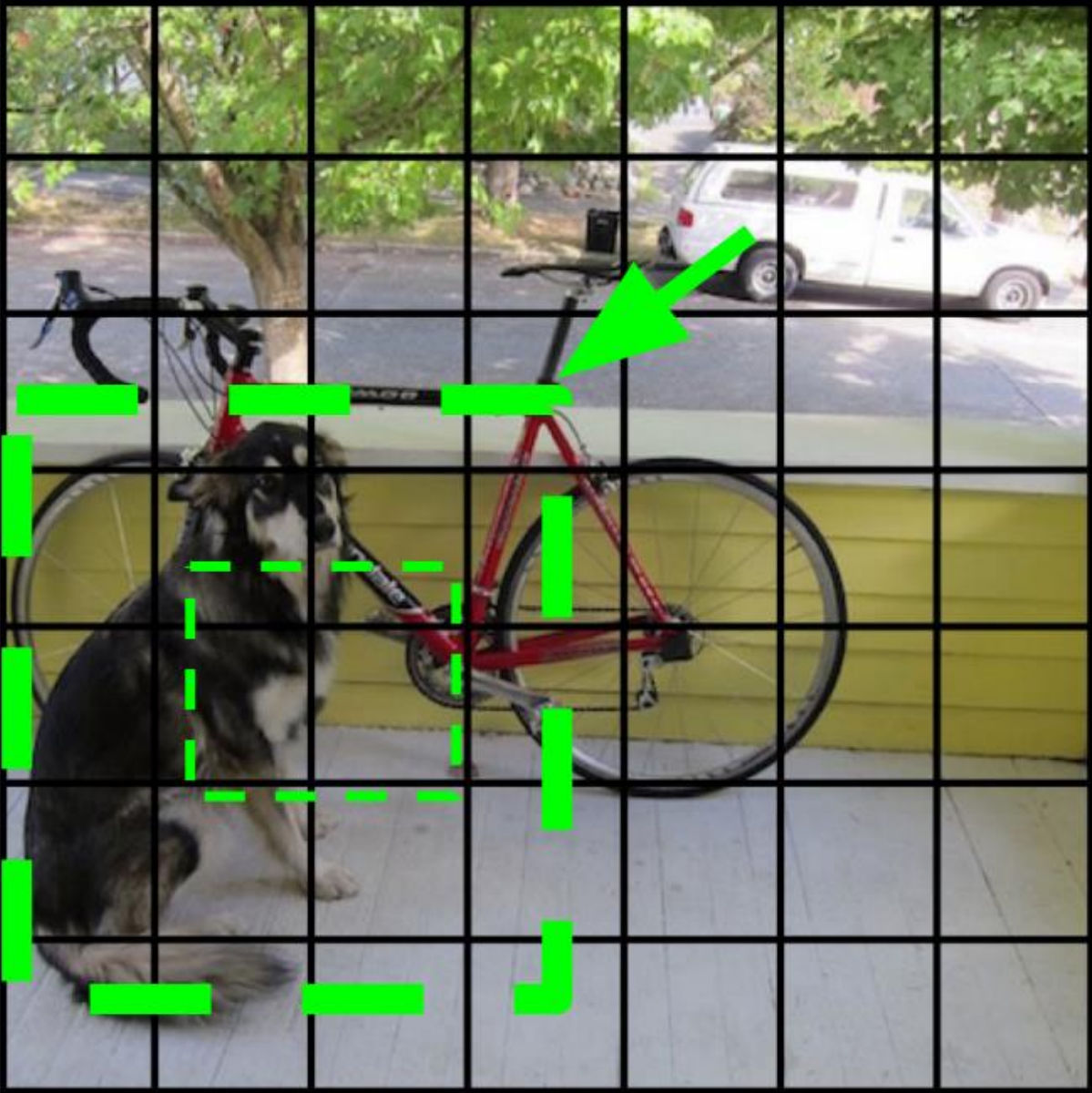...

# Look at that cell's predicted boxes

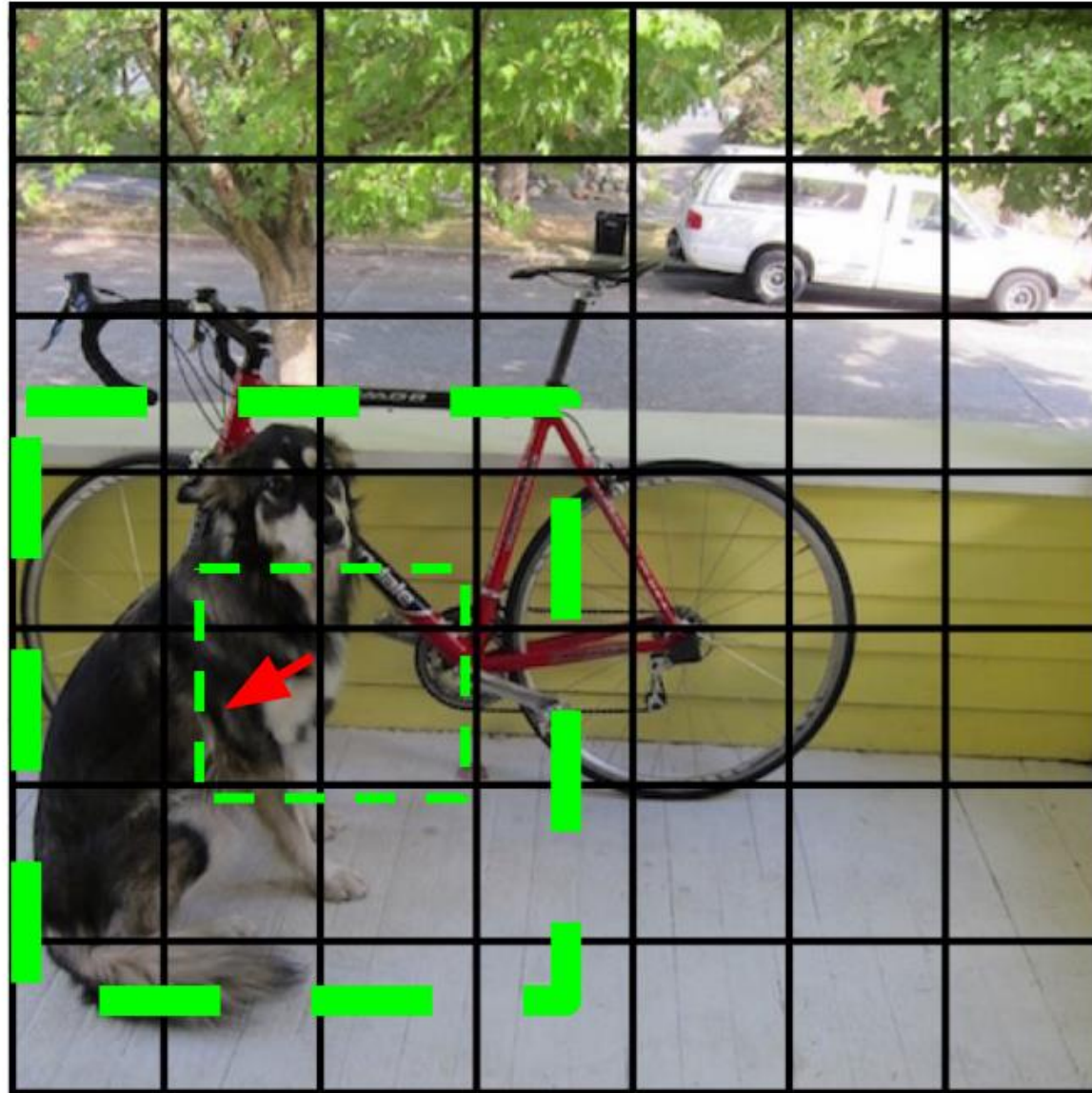# Find the best one, adjust it, increase the confidence

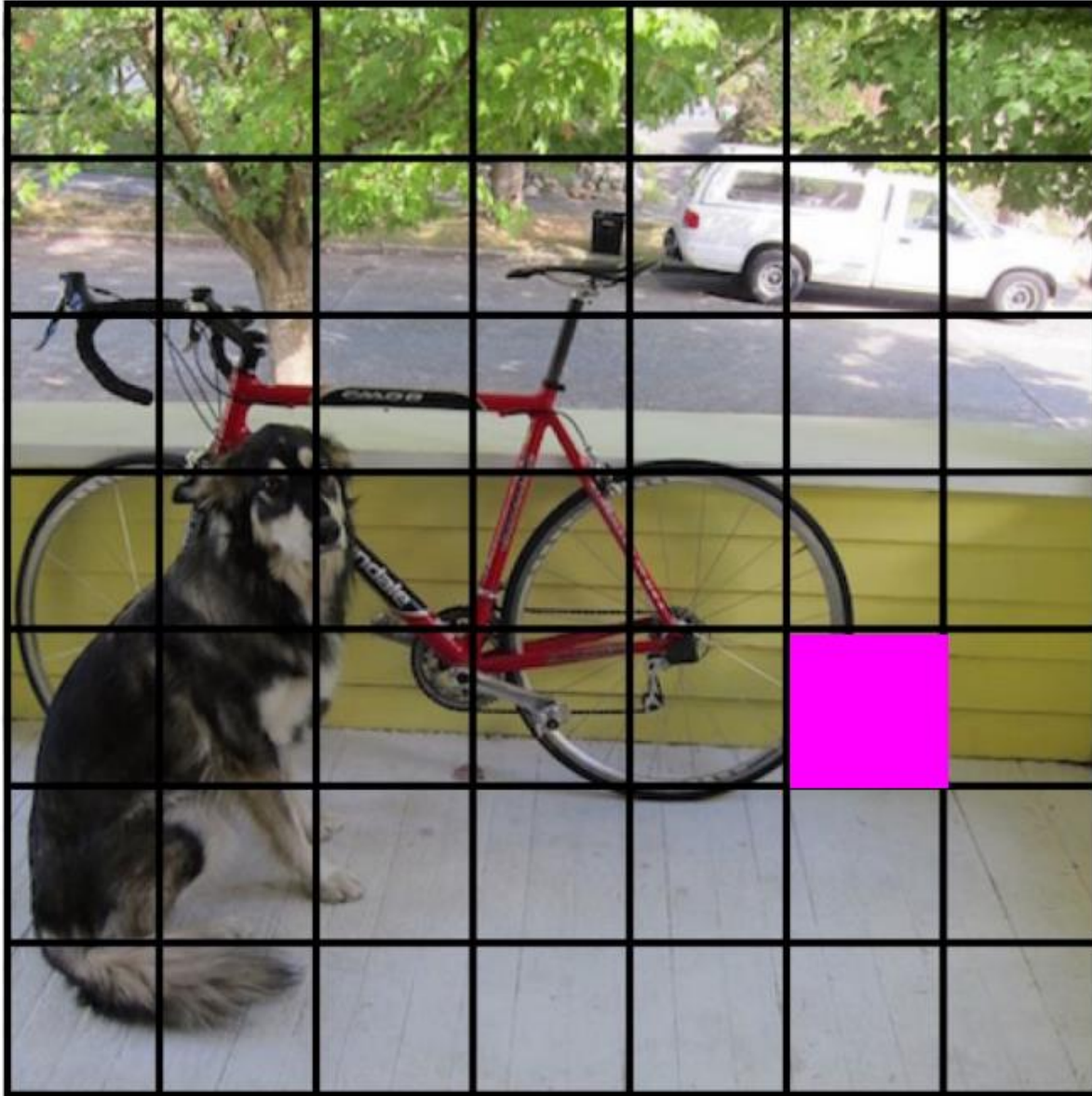# Find the best one, adjust it, increase the confidence
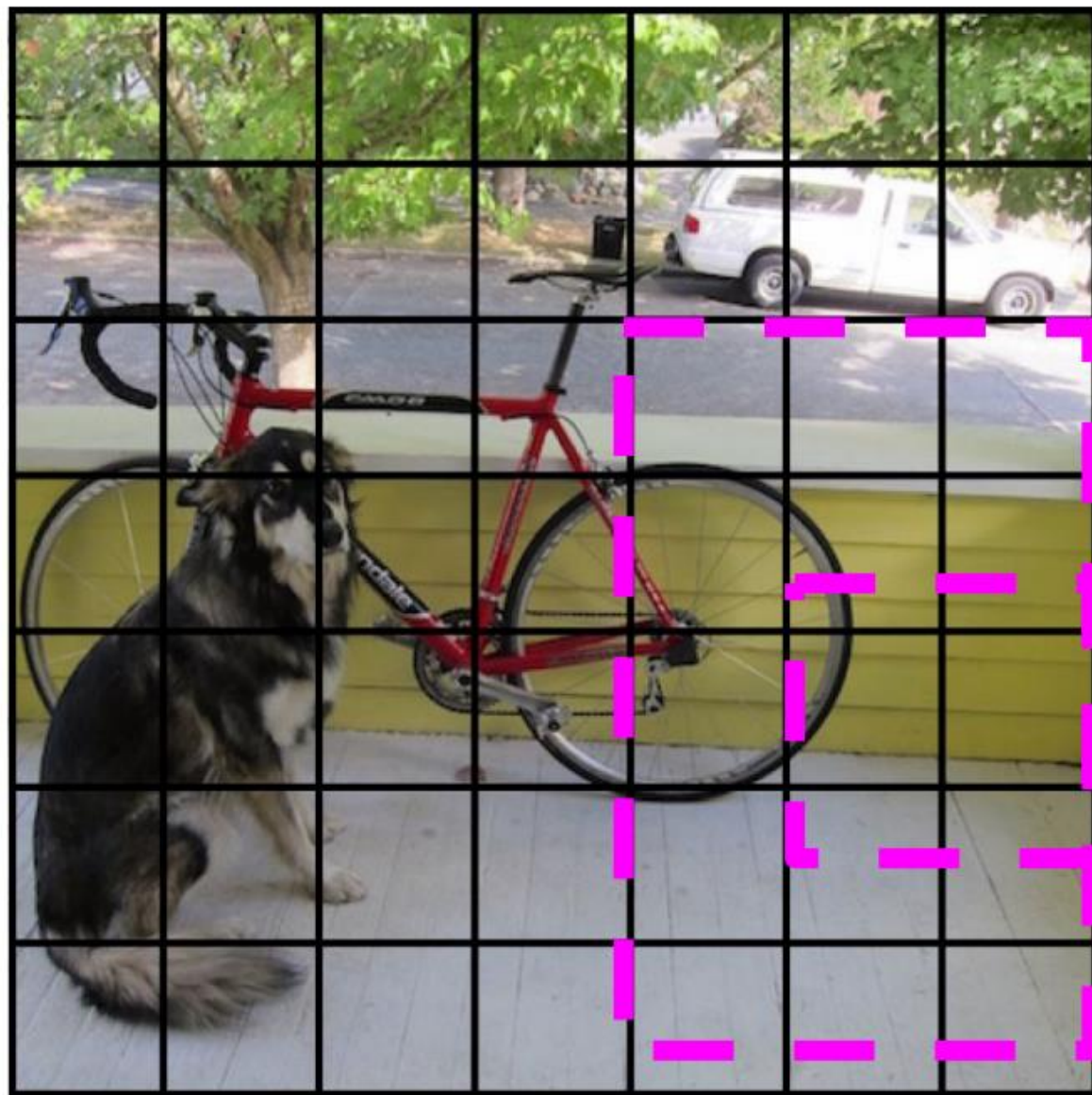
# Find the best one, adjust it, increase the confidence

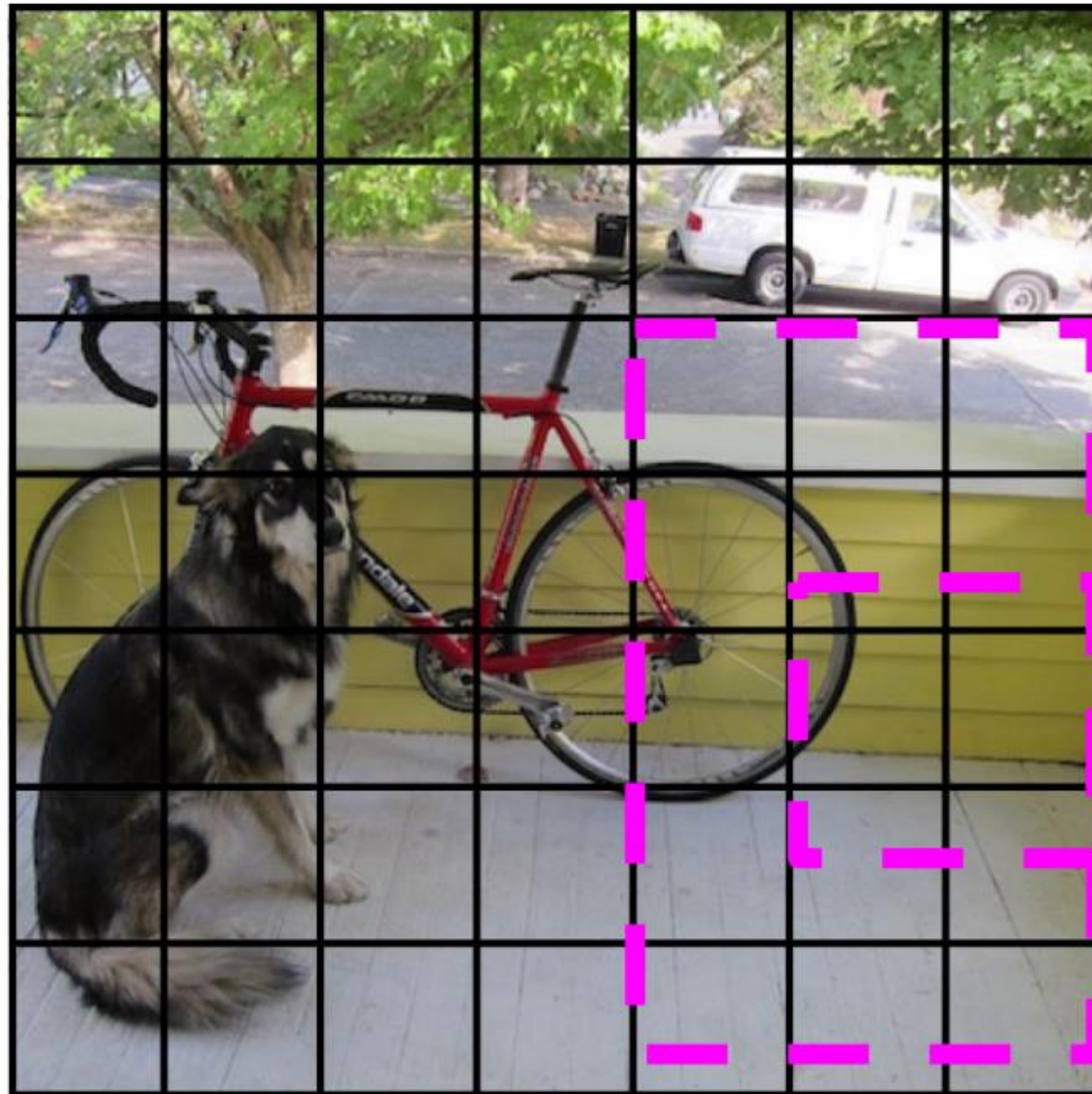# Decrease the confidence of the other box

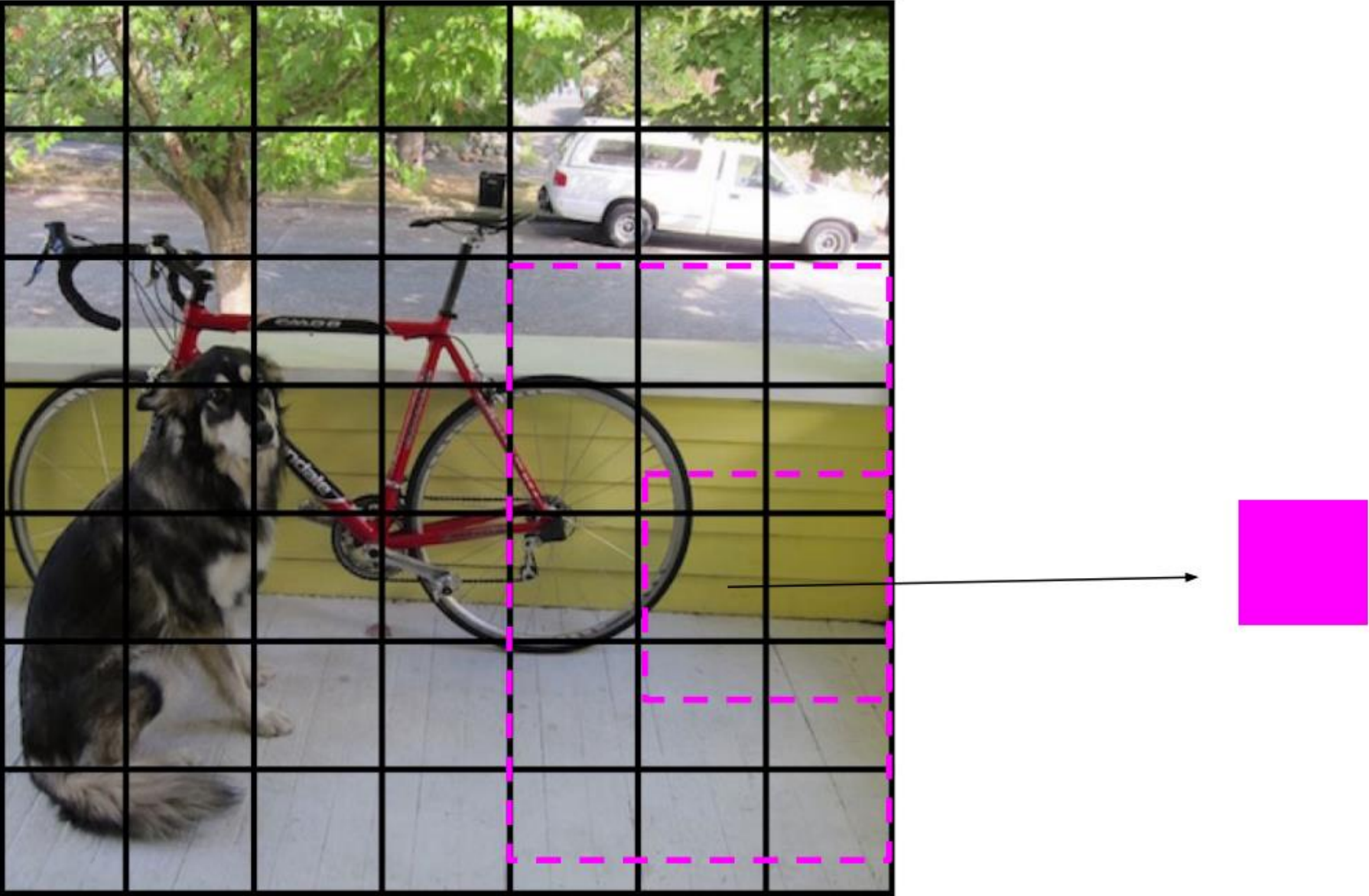# Some cells don't have any ground truth detections!

# Some cells don't have any ground truth detections!

# Decrease the confidence of boxes boxes

# Don't adjust the class probabilities or coordinates

# Loss Function (sum-squared error)

loss function:

$$\lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^{B} \mathbb{1}_{ij}^{\text{obj}} \left[ (x_i - \hat{x}_i)^2 + (y_i - \hat{y}_i)^2 \right]$$

$$+ \lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^{B} \mathbb{1}_{ij}^{\text{obj}} \left[ \left( \sqrt{w_i} - \sqrt{\hat{w}_i} \right)^2 + \left( \sqrt{h_i} - \sqrt{\hat{h}_i} \right)^2 \right]$$

$$+ \sum_{i=0}^{S^2} \sum_{j=0}^{B} \mathbb{1}_{ij}^{\text{obj}} \left( C_i - \hat{C}_i \right)^2$$

$$+ \lambda_{\text{noobj}} \sum_{i=0}^{S^2} \sum_{j=0}^{B} \mathbb{1}_{ij}^{\text{noobj}} \left( C_i - \hat{C}_i \right)^2$$

$$+ \sum_{i=0}^{S^2} \mathbb{1}_{i}^{\text{obj}} \sum_{c \in \text{classes}} (p_i(c) - \hat{p}_i(c))^2 \qquad (3)$$

model. We use sum-squared error because it is easy to optimize, however it does not perfectly align with our goal of maximizing average precision. It weights localization error equally with classification error which may not be ideal. Also, in every image many grid cells do not contain any object. This pushes the "confidence" scores of those cells towards zero, often overpowering the gradient from cells that do contain objects. This can lead to model instability, causing training to diverge early on.

To remedy this, we increase the loss from bounding box coordinate predictions and decrease the loss from confidence predictions for boxes that don't contain objects. We use two parameters, $\lambda_{\text{coord}}$ and $\lambda_{\text{noobj}}$ to accomplish this. We set $\lambda_{\text{coord}} = 5$ and $\lambda_{\text{noobj}} = .5$.
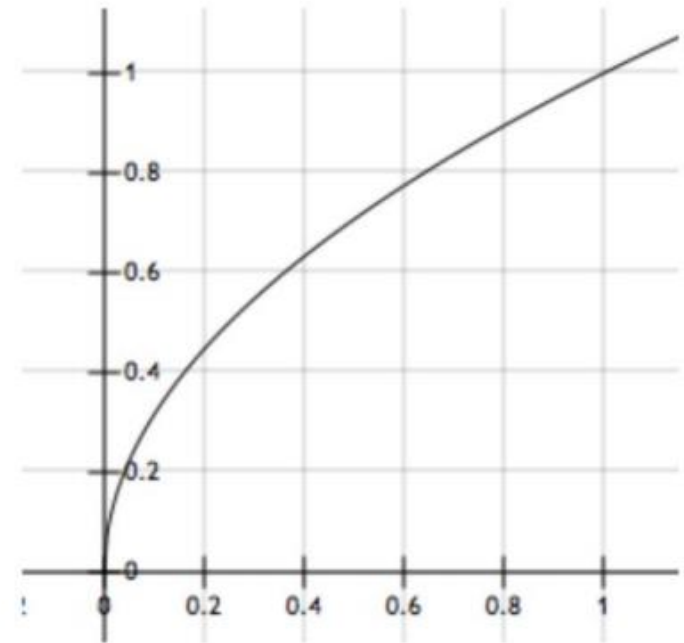
$$\lambda_{\text{coord}} = 5, \quad \lambda_{\text{noobj}} = 0.5$$

# Loss Function (sum-squared error)

loss function:

$$
\lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^{B} \mathbb{1}_{ij}^{\text{obj}} \left[ (x_i - \hat{x}_i)^2 + (y_i - \hat{y}_i)^2 \right]
$$

$$
+ \lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^{B} \mathbb{1}_{ij}^{\text{obj}} \left[ \left( \sqrt{w_i} - \sqrt{\hat{w}_i} \right)^2 + \left( \sqrt{h_i} - \sqrt{\hat{h}_i} \right)^2 \right]
$$

$$
+ \sum_{i=0}^{S^2} \sum_{j=0}^{B} \mathbb{1}_{ij}^{\text{obj}} \left( C_i - \hat{C}_i \right)^2
$$

$$
+ \lambda_{\text{noobj}} \sum_{i=0}^{S^2} \sum_{j=0}^{B} \mathbb{1}_{ij}^{\text{noobj}} \left( C_i - \hat{C}_i \right)^2
$$

$$
+ \sum_{i=0}^{S^2} \mathbb{1}_{i}^{\text{obj}} \sum_{c \in \text{classes}} (p_i(c) - \hat{p}_i(c))^2 \quad (3)
$$

Sum-squared error also equally weights errors in large boxes and small boxes. Our error metric should reflect that small deviations in large boxes matter less than in small boxes. To partially address this we predict the square root of the bounding box width and height instead of the width and height directly.

# Loss Function (sum-squared error)

loss function:

$$\lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^{B} \mathbb{1}_{ij}^{obj} \left[ (x_i - \hat{x}_i)^2 + (y_i - \hat{y}_i)^2 \right]$$

$$+ \lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^{B} \mathbb{1}_{ij}^{obj} \left[ \left( \sqrt{w_i} - \sqrt{\hat{w}_i} \right)^2 + \left( \sqrt{h_i} - \sqrt{\hat{h}_i} \right)^2 \right]$$

$$+ \sum_{i=0}^{S^2} \sum_{j=0}^{B} \mathbb{1}_{ij}^{obj} \left( C_i - \hat{C}_i \right)^2$$

$$+ \lambda_{\text{noobj}} \sum_{i=0}^{S^2} \sum_{j=0}^{B} \mathbb{1}_{ij}^{noobj} \left( C_i - \hat{C}_i \right)^2$$

$$+ \sum_{i=0}^{S^2} \mathbb{1}_{i}^{obj} \sum_{c \in \text{classes}} (p_i(c) - \hat{p}_i(c))^2 \qquad (3)$$

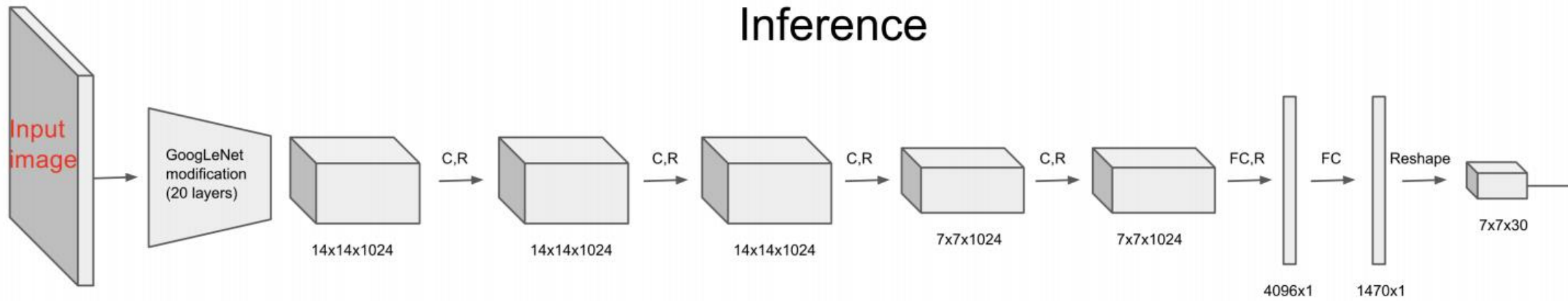$\mathbb{1}_{ij}^{obj}$ **The $j$th bbox predictor** in *cell i* is "responsible" for that prediction

$\mathbb{1}_{ij}^{noobj}$
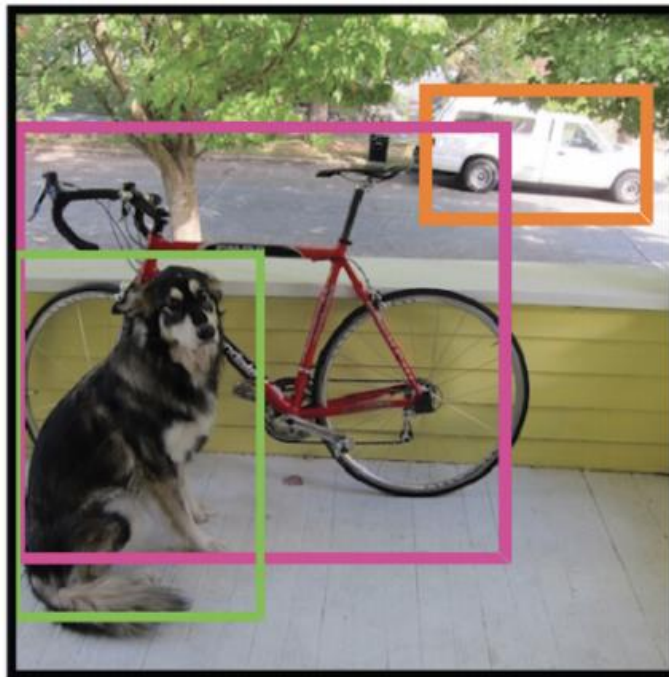
$\mathbb{1}_{i}^{obj}$ If object appears in *cell i*

Note that the loss function only penalizes classification error if an object is present in that grid cell (hence the conditional class probability discussed earlier). It also only penalizes bounding box coordinate error if that predictor is "responsible" for the ground truth box (i.e. has the highest IOU of any predictor in that grid cell).
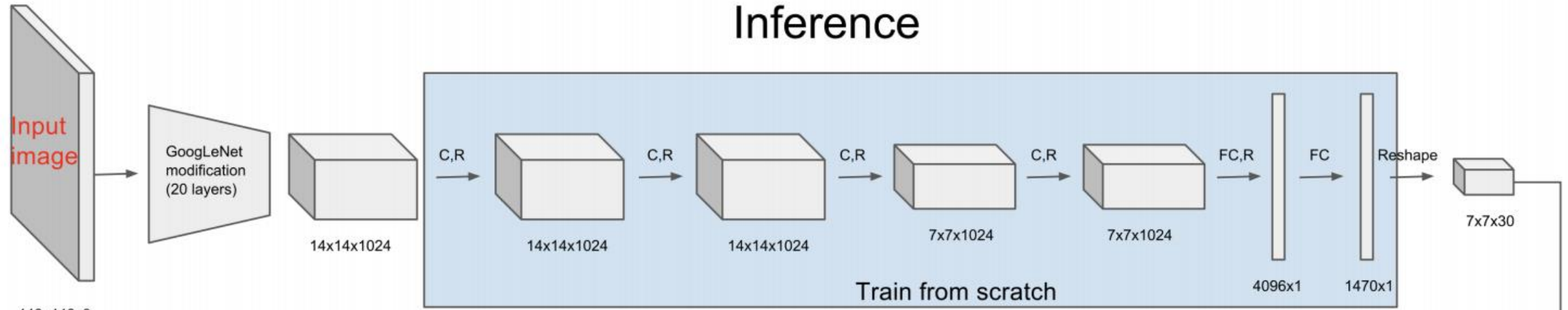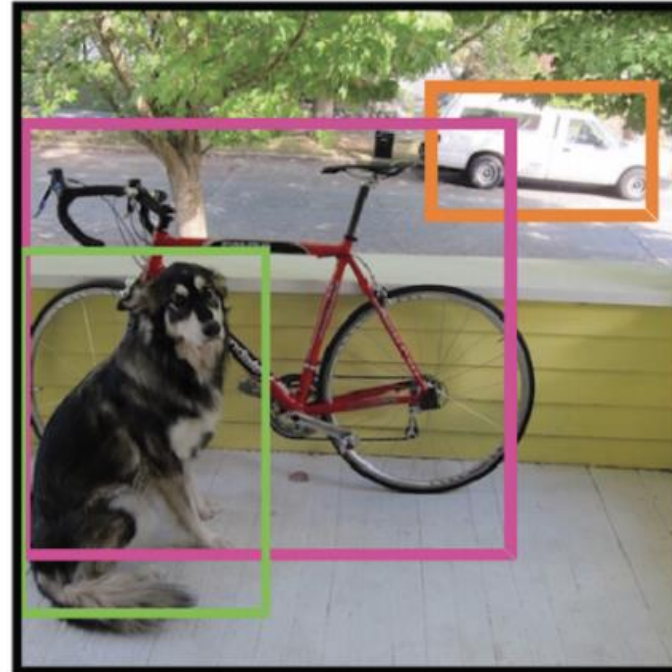
# Inference



Slide credit to Taegyun Jeon

# Inference



Detection Procedure

# Inference

use new additional conv layers => better performance



448x448x3

14x14x1024    14x14x1024    14x14x1024    7x7x1024    7x7x1024    FC,R  4096x1    FC  1470x1    Reshape    7x7x30

C,R    C,R    C,R    C,R

Detection Procedure

Slide credit to Taegyun Jeon

# Inference



Slide credit to Taegyun Jeon

# Inference

# Inference



Input image
448x448x3

GoogLeNet modification (20 layers)

14x14x1024 — C,R → 14x14x1024 — C,R → 14x14x1024 — C,R → 7x7x1024 — C,R → 7x7x1024 — FC,R → 4096x1 — FC → 1470x1 — Reshape → 7x7x30 → Detection Procedure

Tensor values interpretation

7  7  30

grid cell

7  7

# Inference



448x448x3

GoogLeNet modification (20 layers)

14x14x1024 → C,R → 14x14x1024 → C,R → 14x14x1024 → C,R → 7x7x1024 → C,R → 7x7x1024 → FC,R → 4096x1 → FC → 1470x1 → Reshape → 7x7x30 → Detection Procedure

Tensor values interpretation

7 · 7 · 30

grid cell

7 · 7

# Inference



Tensor values interpretation

grid cell

# Inference



Tensor values interpretation

1. x - coordinate of bbox center inside cell ([0; 1] wrt grid cell size)
2. y - coordinate of bbox center inside cell ([0; 1] wrt grid cell size)
3. w - bbox width ([0; 1] wrt image)
4. h - bbox height ([0; 1] wrt image)
5. c - bbox confidence ~ P(obj in bbox1)

grid cell

# Inference



Tensor values interpretation

1. x - coordinate of bbox center inside cell ([0; 1] wrt grid cell size)
2. y - coordinate of bbox center inside cell ([0; 1] wrt grid cell size)
3. w - bbox width ([0; 1] wrt image)
4. h - bbox height ([0; 1] wrt image)
5. c - bbox confidence ~ P(obj in bbox2)

# Inference



448x448x3

GoogLeNet modification (20 layers)

14x14x1024 → C,R → 14x14x1024 → C,R → 14x14x1024 → C,R → 7x7x1024 → C,R → 7x7x1024 → FC,R → 4096x1 → FC → 1470x1 → Reshape → 7x7x30 → Detection Procedure

Tensor values interpretation

1x30

5    5

two bboxes for each grid cell

grid cell

7    7    30

7    7

7

# Inference



448x448x3    GoogLeNet modification (20 layers)    14x14x1024    C,R    14x14x1024    C,R    14x14x1024    C,R    7x7x1024    C,R    7x7x1024    FC,R    4096x1    FC    1470x1    Reshape    7x7x30    Detection Procedure

Tensor values interpretation

1x30

5    5    20 - number of classes

7    7    30    7    7

grid cell

# Inference



Input image
448x448x3

GoogLeNet modification (20 layers)

14x14x1024 → C,R → 14x14x1024 → C,R → 14x14x1024 → C,R → 7x7x1024 → C,R → 7x7x1024 → FC,R → 4096x1 → FC → 1470x1 → Reshape → 7x7x30

Detection Procedure

## Tensor values interpretation

i (from 1 to 20)

1x30

5    5    20 - number of classes

Class score ~ P(obj is class_i | obj in box)

7    7    30

grid cell

# Inference



Tensor values interpretation

grid cell

bb1 confidence

1x30

5    5    20

MUL

Class scores for bb1

20x1

Input image

GoogLeNet modification (20 layers)

448x448x3

14x14x1024    14x14x1024    14x14x1024    7x7x1024    7x7x1024

C,R    C,R    C,R    C,R    FC,R    FC    Reshape

4096x1    1470x1

7x7x30

Detection Procedure

7    30    7

7    7

# Inference



Tensor values interpretation

grid cell

Class scores for bb2

Slide credit to Taegyun Jeon

# Inference



Input image
448x448x3

GoogLeNet modification (20 layers)

14x14x1024 → C,R → 14x14x1024 → C,R → 14x14x1024 → C,R → 7x7x1024 → C,R → 7x7x1024 → FC,R → 4096x1 → FC → 1470x1 → Reshape → 7x7x30

Detection Procedure

Tensor values interpretation

bb2 confidence

1x30

5    5    20

MUL

Class scores for bb2

20x1

grid cell

7 7 30

7

7

Do this operation for each bbox in each grid cell

# Inference

Tensor values interpretation

2 bboxes for first cell (1, 1)

bb1   bb2

448x448x3

14x14x1024   14x14x1024   14x14x1024   7x7x1024   7x7x1024

C,R   C,R   C,R   C,R   FC,R   FC   Reshape

4096x1   1470x1   7x7x30

Detection Procedure

grid cell (1, 1)

7   7   30   7   7

20x1   20x1

Slide credit to Taegyun Jeon

# Inference



Input image

GoogLeNet modification (20 layers)

448x448x3

14x14x1024    C,R    14x14x1024    C,R    14x14x1024    C,R    7x7x1024    C,R    7x7x1024    FC,R    4096x1    FC    1470x1    Reshape    7x7x30

Detection Procedure

## Tensor values interpretation

2 bboxes for second cell (1, 2)

bb1    bb2    bb3    bb4

7    30

grid cell (1, 2)

7    7

7

20x1    20x1    20x1    20x1

# Inference



448x448x3 · 14x14x1024 · C,R · 14x14x1024 · C,R · 14x14x1024 · C,R · 7x7x1024 · C,R · 7x7x1024 · FC,R · 4096x1 · FC · 1470x1 · Reshape · 7x7x30

GoogLeNet modification (20 layers)

Detection Procedure

Tensor values interpretation

2 bboxes for last cell (7, 7)

bb1  bb2  bb3  bb4 · bb97  bb98

grid cell (7, 7)

20x1 20x1 20x1 20x1 20x1 20x1 20x1 20x1 20x1 · 20x1 20x1 20x1 20x1 20x1 20x1

# Inference



448x448x3

GoogLeNet modification (20 layers)

14x14x1024 → C,R → 14x14x1024 → C,R → 14x14x1024 → C,R → 7x7x1024 → C,R → 7x7x1024 → FC,R → 4096x1 → FC → 1470x1 → Reshape → 7x7x30 → Detection Procedure

Tensor values interpretation

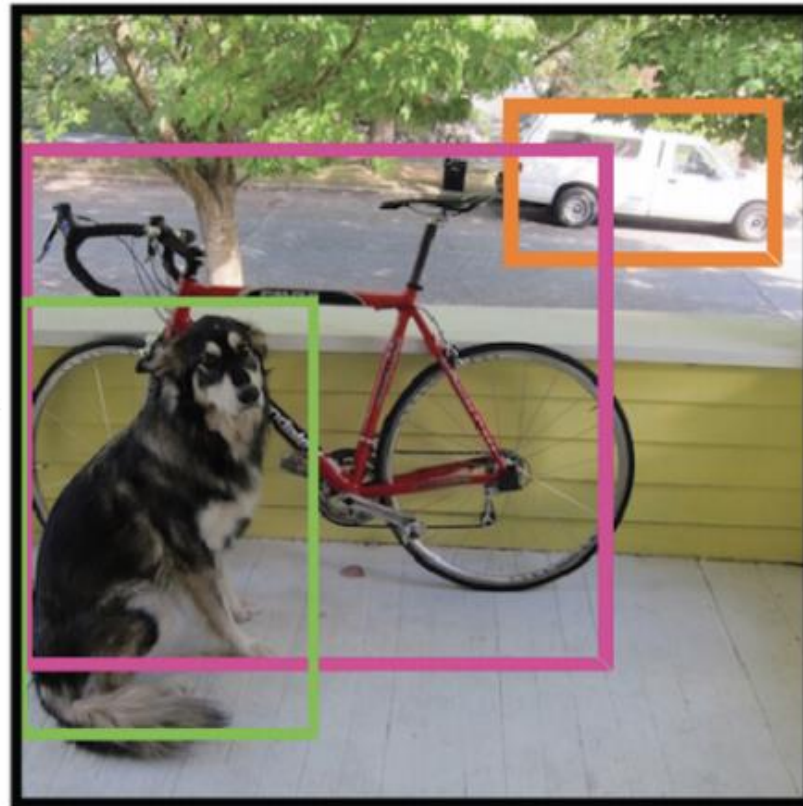Total 7*7*2 = 98 bboxes

bb1  bb2  bb3  bb4  ........  bb97  bb98

7

30

grid cell (7, 7)

7

7

# Look at detection procedure



7x7x30

Detection
Procedure

Get first class scores for each bbox

Dog scores

bb1 bb2 bb3 bb4 bb97 bb98

20x1    20x1

Set zero
if score < thresh1 (0.2)

bb1 bb2 bb3 bb4 bb97 bb98
    0       0       0

Sort descending

bb3 bb1 bb98 bb2 bb4 bb98
            0   0   0

NMS algorithm set
scores to zero for
redundant bboxes

bb3 bb1 bb98 bb2 bb4 bb98
            0   0   0   0

How it works

# Non-Maximum Suppression: intuition

class (dog) scores for each bbox

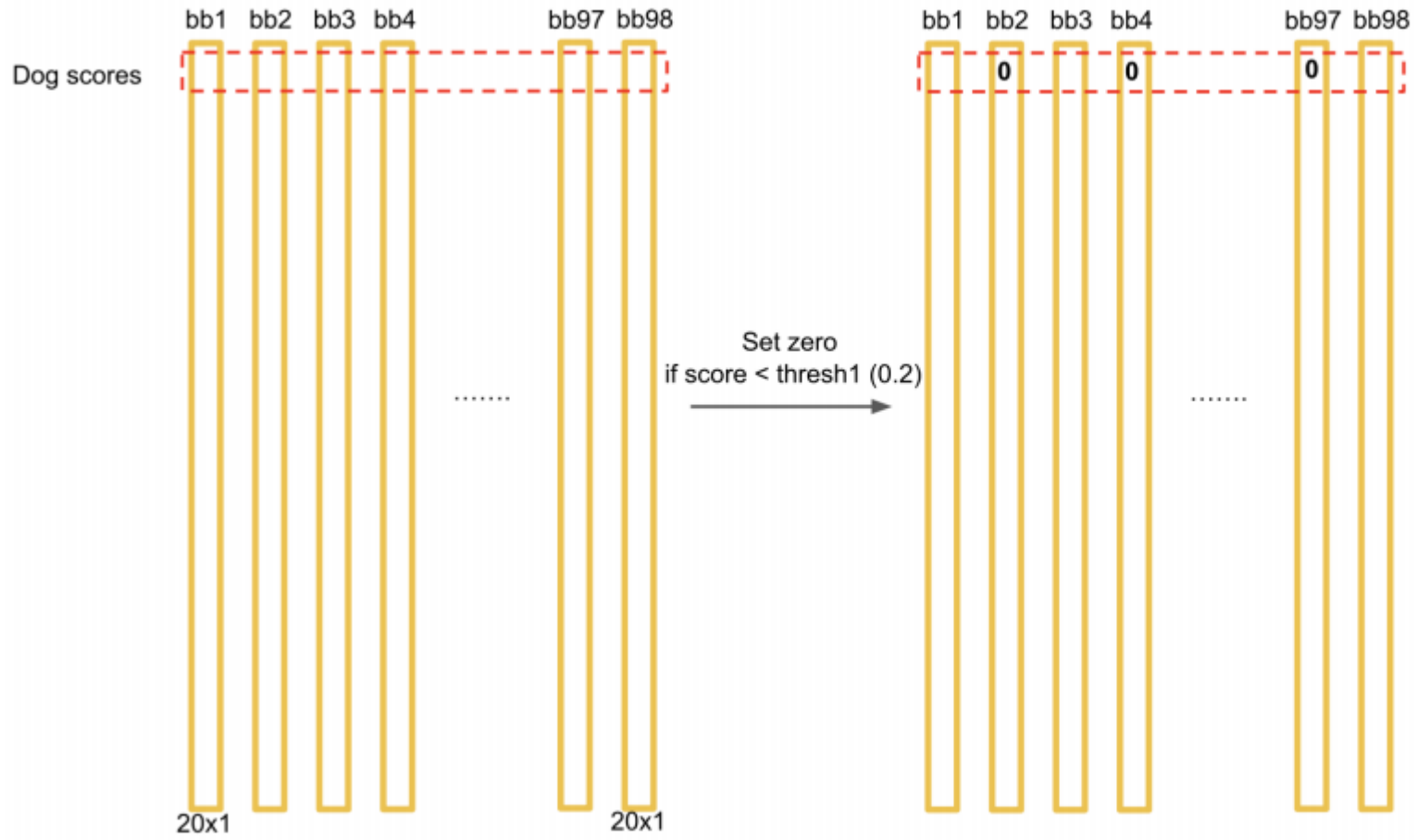| | bb47 | bb20 | bb15 | bb7 | | | | | | | | | bb1 | bb4 | bb8 | bb98 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| class: dog | 0.5 | 0.3 | 0.2 | 0.1 | | | | | | | | | 0 | 0 | 0 | 0 | 1x98 |

# Non-Maximum Suppression: intuition

class (dog) scores for each bbox

| | bb47 | bb20 | bb15 | bb7 | | | | | | | | bb1 | bb4 | bb8 | bb98 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| class: dog | 0.5 | 0.3 | 0.2 | 0.1 | | | | | | | | 0 | 0 | 0 | 0 | 1x98 |

# Non-Maximum Suppression: intuition

# Non-Maximum Suppression: intuition

class (dog) scores for each bbox

| | bb47 | bb20 | bb15 | bb7 | | | | | | | bb1 | bb4 | bb8 | bb98 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| class: dog | 0.5 | 0.3 | 0.2 | 0.1 | | | | | | | 0 | 0 | 0 | 0 | 1x98 |

# Non-Maximum Suppression: intuition

class (dog) scores for each bbox

| | bb47 | bb20 | bb15 | bb7 | | | | | | | | bb1 | bb4 | bb8 | bb98 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| class: dog | 0.5 | 0.3 | 0.2 | 0.1 | | | | | | | | 0 | 0 | 0 | 0 | 1x98 |

# Non-Maximum Suppression: intuition

class (dog) scores for each bbox

| | bb47 | bb20 | bb15 | bb7 | | | | | | | | bb1 | bb4 | bb8 | bb98 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| class: dog | 0.5 | 0.3 | 0.2 | 0.1 | | | | | | | | 0 | 0 | 0 | 0 | 1x98 |

Get bbox with max score. Let's denote it "bbox_max"

# Non-Maximum Suppression: intuition

class (dog) scores for each bbox

|  | bb47 | bb20 | bb15 | bb7 |  |  |  |  |  |  | bb1 | bb4 | bb8 | bb98 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| class: dog | 0.5 | 0.3 | 0.2 | 0.1 |  |  |  |  |  |  | 0 | 0 | 0 | 0 |

1x98

Compare "bbox_max" with others less score (non-zero!) bboxes. Let's denote it "bbox_cur"

# Non-Maximum Suppression: intuition

class (dog) scores for each bbox



If IoU(bbox_max, bbox_cur) > 0.5 then set 0 score to bbox_cur.

# Non-Maximum Suppression: intuition

class (dog) scores for each bbox

class: dog

| bb47 | bb20 | bb15 | bb7 | | | | | | | bb1 | bb4 | bb8 | bb98 | 1x98 |
|------|------|------|-----|--|--|--|--|--|--|-----|-----|-----|------|------|
| 0.5 | 0 | 0.2 | 0.1 | | | | | | | 0 | 0 | 0 | 0 | |

If IoU(bbox_max, bbox_cur) > 0.5 then set 0 score to bbox_cur.

In this case: set to 0.

# Non-Maximum Suppression: intuition

class (dog) scores for each bbox

| | bb47 | bb20 | bb15 | bb7 | | | | | | bb1 | bb4 | bb8 | bb98 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| class: dog | 0.5 | 0 | 0.2 | 0.1 | | | | | | 0 | 0 | 0 | 0 | 1x98 |

Go to next bbox_cur.

# Non-Maximum Suppression: intuition



Go to next bbox_cur.

If IoU(bbox_max, bbox_cur) > 0.5 then set 0 score to bbox_cur.

Slide credit to Taegyun Jeon

# Non-Maximum Suppression: intuition

class (dog) scores for each bbox

| | bb47 | bb20 | bb15 | bb7 | | | | | | bb1 | bb4 | bb8 | bb98 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| class: dog | 0.5 | **0** | 0.2 | 0.1 | | | | | | 0 | 0 | 0 | 0 | 1x98 |



Go to next bbox_cur.

If IoU(bbox_max, bbox_cur) > 0.5 then set 0 score to bbox_cur.

In this case: continue.

# Non-Maximum Suppression: intuition

class (dog) scores for each bbox

| | bb47 | bb20 | bb15 | bb7 | | | | | | | | | bb1 | bb4 | bb8 | bb98 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| class: dog | 0.5 | 0 | 0.2 | 0.1 | | | | | | | | | 0 | 0 | 0 | 0 | 1x98 |

Go to next bbox_cur.

# Non-Maximum Suppression: intuition

class (dog) scores for each bbox

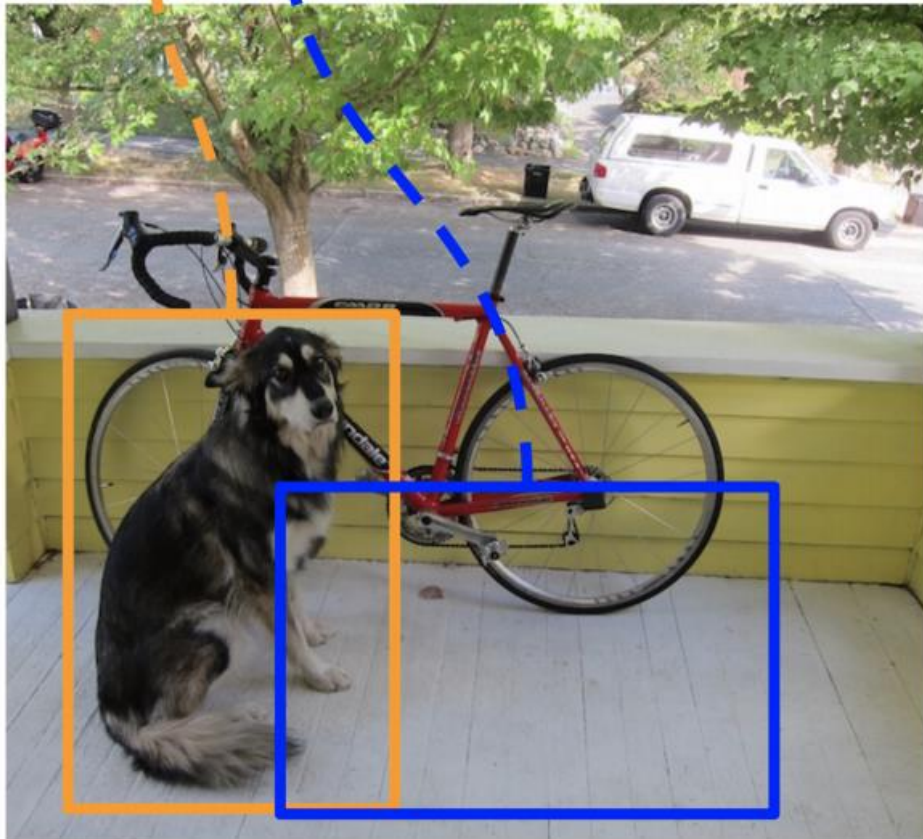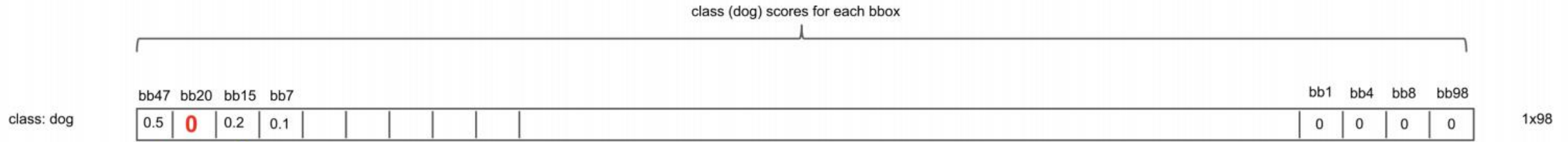| | bb47 | bb20 | bb15 | bb7 | | | | | | | | | | | bb1 | bb4 | bb8 | bb98 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| class: dog | 0.5 | **0** | 0.2 | 0.1 | | | | | | | | | | | 0 | 0 | 0 | 0 | 1x98 |

Go to next bbox_cur.

If IoU(bbox_max, bbox_cur) > 0.5 then set 0 score to bbox_cur.

# Non-Maximum Suppression: intuition

class (dog) scores for each bbox



Go to next bbox_cur.

If IoU(bbox_max, bbox_cur) > 0.5 then set 0 score to bbox_cur.

In this case: continue.

# Non-Maximum Suppression: intuition

class (dog) scores for each bbox

| | bb47 | bb20 | bb15 | bb7 | | | | | | bb1 | bb4 | bb8 | bb98 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| class: dog | 0.5 | 0 | 0.2 | 0.1 | | | | | | 0 | 0 | 0 | 0 | 1x98 |

Go to next bbox_cur.

If IoU(bbox_max, bbox_cur) > 0.5 then set 0 score to bbox_cur.

In this case: continue.

Do this procedure for other "bbox_cur". After that ...

# Non-Maximum Suppression: intuition

class (dog) scores for each bbox

class: dog

| | bb47 | bb20 | bb15 | bb7 | | | | | | | | bb1 | bb4 | bb8 | bb98 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 0.5 | 0 | 0.2 | 0.1 | | | | | | | | 0 | 0 | 0 | 0 | 1x98 |

Go to next bbox with big score.
Let's denote it "bbox_max"

Slide credit to Taegyun Jeon

# Non-Maximum Suppression: intuition

class (dog) scores for each bbox

| | bb47 | bb20 | bb15 | bb7 | | | | | | | | bb1 | bb4 | bb8 | bb98 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| class: dog | 0.5 | 0 | 0.2 | **0** | | | | | | | | 0 | 0 | 0 | 0 | 1x98 |



Go to next bbox_cur.

If IoU(bbox_max, bbox_cur) > 0.5 then set 0 score to bbox_cur.

In this case: set to 0.

Do this procedure for other "bbox_max" and for other corresponding "bbox_cur".

# Non-Maximum Suppression: intuition

class (dog) scores for each bbox

| | bb47 | bb20 | bb15 | bb7 | | | | | | | | | | bb1 | bb4 | bb8 | bb98 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| class: dog | 0.5 | 0 | 0.2 | 0 | | | | | | | | | | 0 | 0 | 0 | 0 | 1x98 |



After comparison almost all pairs of bboxes the only two bboxes left with non-zero class score value.

Do this procedure for all classes

After this procedure - a lot of zeros

Select bboxes to draw by class score values

bb1 bb2 bb3 bb4 ... bb97 bb98

Set zero if score < thresh1 (0.2)

bb1 bb2 bb3 bb4 ... bb97 bb98

20x1 20x1

Sort descending

bb3 bb1 bb98 ... bb2 bb4 bb97

NMS algorithm set scores to zero for redundant bboxes

bb3 bb1 bb98 ... bb2 bb4 bb97

class = max_index(scores for bb3);
score = max(scores for bb3);

Score > 0

no → skip bbox

yes → draw bbox with class color

bb1 bb2 bb3 bb4    bb97 bb98

bb1 bb2 bb3 bb4    bb97 bb98

bb3 bb1 bb98    bb2 bb4 bb97

bb3 bb1 bb98    bb2 bb4 bb97

Set zero
if score < thresh1 (0.2)

Sort descending

NMS algorithm set
scores to zero for
redundant bboxes

20x1    20x1

0    0    0

0

0

0

0

0

0    0    0

0

0

0    0    0

class = max_index(scores for bb1);
score = max(scores for bb1);

no    skip bbox

Score > 0

yes

draw bbox with class color

bb1 bb2 bb3 bb4    bb97 bb98

20x1    20x1

Set zero
if score < thresh1 (0.2)

bb1 bb2 bb3 bb4    bb97 bb98

0    0    0

Sort descending

bb3 bb1 bb98    bb2 bb4 bb97

0    0    0    0

NMS algorithm set
scores to zero for
redundant bboxes

bb3 bb1 bb98    bb2 bb4 bb97

0
0

0

0

0

0    0    0

0

0    0    0

class = max_index(scores for bb98);
score = max(scores for bb98);

no    Score > 0

skip bbox

yes

draw bbox with class color

| bb1 | bb2 | bb3 | bb4 | ... | bb97 | bb98 |

Set zero
if score < thresh1 (0.2)

| bb1 | bb2 | bb3 | bb4 | ... | bb97 | bb98 |

Sort descending

| bb3 | bb1 | bb98 | ... | bb2 | bb4 | bb97 |

NMS algorithm set
scores to zero for
redundant bboxes

| bb3 | bb1 | bb98 | ... | bb2 | bb4 | bb97 |

20x1

class = max_index(scores for bb97);
score = max(scores for bb97);

Score > 0

no → skip bbox

yes → draw bbox with class color

Slide credit to Taegyun Jeon

# Strengths and Weaknesses

- Strengths:
  - Fast: 45fps, smaller version 155fps
  - End2end training
  - Background error is low

- Weaknesses:
  - Performance is lower than state-of-art
  - Makes more localization errors

# Open Questions

- How to determine the number of cell, bounding box and the size of the box
- Why normalization x,y,w,h even all the input images have the same resolution?

-

# Limitation of YOLO

- Group of small objects

- Unusual aspect ratios

- Coarse feature

- Localization error of bounding box

## 2.4. Limitations of YOLO

YOLO imposes strong spatial constraints on bounding box predictions since each grid cell only predicts two boxes and can only have one class. This spatial constraint limits the number of nearby objects that our model can predict. Our model struggles with small objects that appear in groups, such as flocks of birds.

Since our model learns to predict bounding boxes from data, it struggles to generalize to objects in new or unusual aspect ratios or configurations. Our model also uses relatively coarse features for predicting bounding boxes since our architecture has multiple downsampling layers from the input image.

Finally, while we train on a loss function that approximates detection performance, our loss function treats errors the same in small bounding boxes versus large bounding boxes. A small error in a large box is generally benign but a small error in a small box has a much greater effect on IOU. Our main source of error is incorrect localizations.

# Extension Part

YOLOv2!

| 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 $288 \times 288$ | 2007+2012 | 69.0 | 91 |
| YOLOv2 $352 \times 352$ | 2007+2012 | 73.7 | 81 |
| YOLOv2 $416 \times 416$ | 2007+2012 | 76.8 | 67 |
| YOLOv2 $480 \times 480$ | 2007+2012 | 77.8 | 59 |
| YOLOv2 $544 \times 544$ | 2007+2012 | **78.6** | 40 |

Table 3: Detection frameworks on PASCAL VOC 2007. YOLOv2 is faster and more accurate than prior detection methods. It can also run at different resolutions for an easy tradeoff between speed and accuracy. Each YOLOv2 entry is actually the same trained model with the same weights, just evaluated at a different size. All timing information is on a Geforce GTX Titan X (original, not Pascal model).

# Improvement

- **YOLO**: (1) make each cell predict more bounding boxes
  (2) also put the idea of multi-scale into it in order to process
  small objects

Features

- **SSD**: simplify its structure and speed up

# Comparison to Other Detection System

## Big picture



mAP

**DPM**
FPS: 0.5
mAP: 34.3

**R-CNN**
FPS: -
mAP: 58.5

**Fast R-CNN**
FPS: 0.5
mAP: 70

**Faster R-CNN**
FPS: 7
mAP: 73.2

**YOLO**
FPS: 45
mAP: 63.4

**SSD**
FPS: 58
mAP: 72.1

Nov 2013        Apr 2015        June 2015        Dec 2015        Time

# Appendix | Implementation

- YOLO (darknet): https://pjreddie.com/darknet/yolov1/

- YOLOv2 (darknet): https://pjreddie.com/darknet/yolo/

- YOLO (caffe): https://github.com/xingwangsfu/caffe-yolo

- YOLO (TensorFlow: Train+Test): https://github.com/thtrieu/darkflow

- YOLO (TensorFlow: Test): https://github.com/gliese581gg/YOLO_tensorflow

# SSD: Single Shot MultiBox Detector

# SSD: Single Shot MultiBox Detector

Wei Liu[1], Dragomir Anguelov[2], Dumitru Erhan[3], Christian Szegedy[3],
Scott Reed[4], Cheng-Yang Fu[1], Alexander C. Berg[1]

[1]UNC Chapel Hill [2]Zoox Inc. [3]Google Inc. [4]University of Michigan, Ann-Arbor
[1]`wliu@cs.unc.edu`, [2]`drago@zoox.com`, [3]`{dumitru,szegedy}@google.com`,
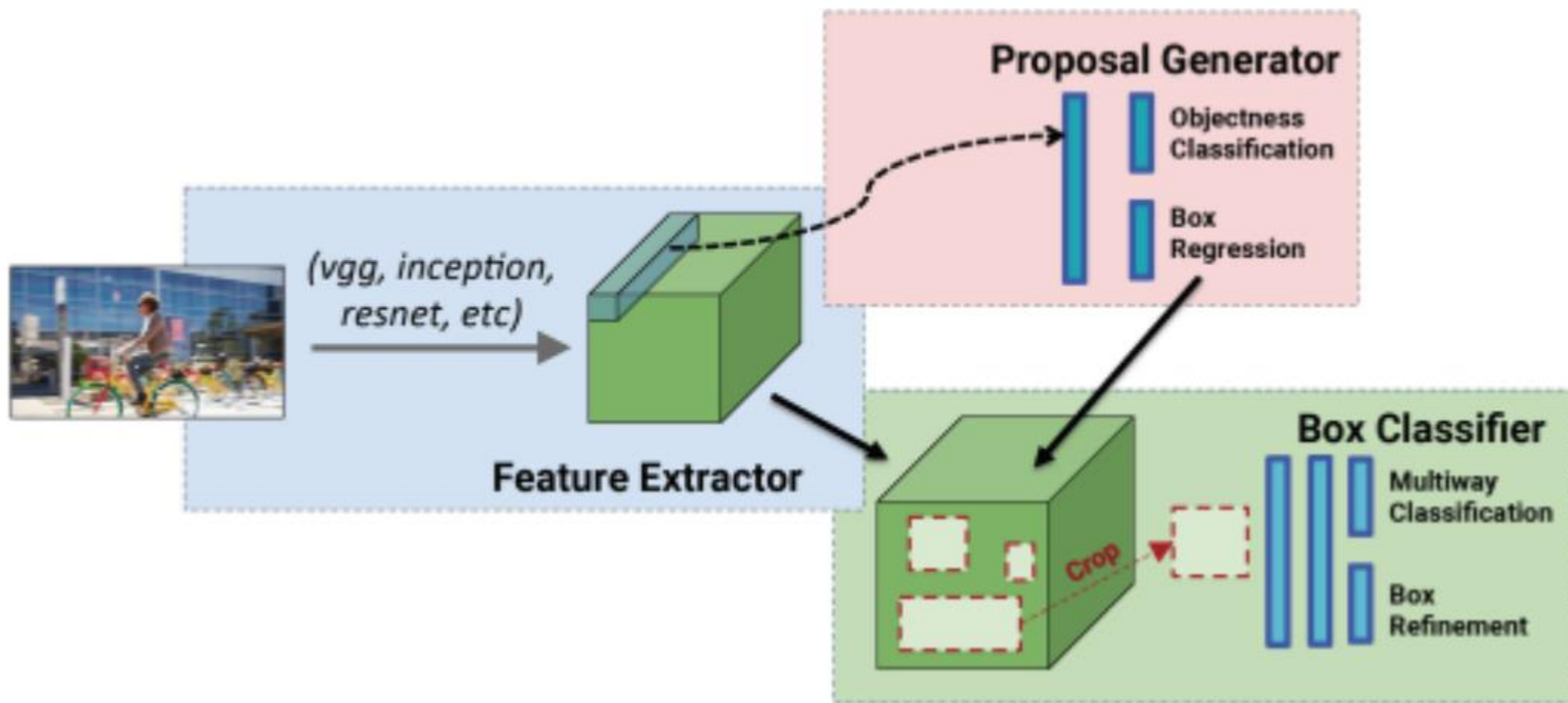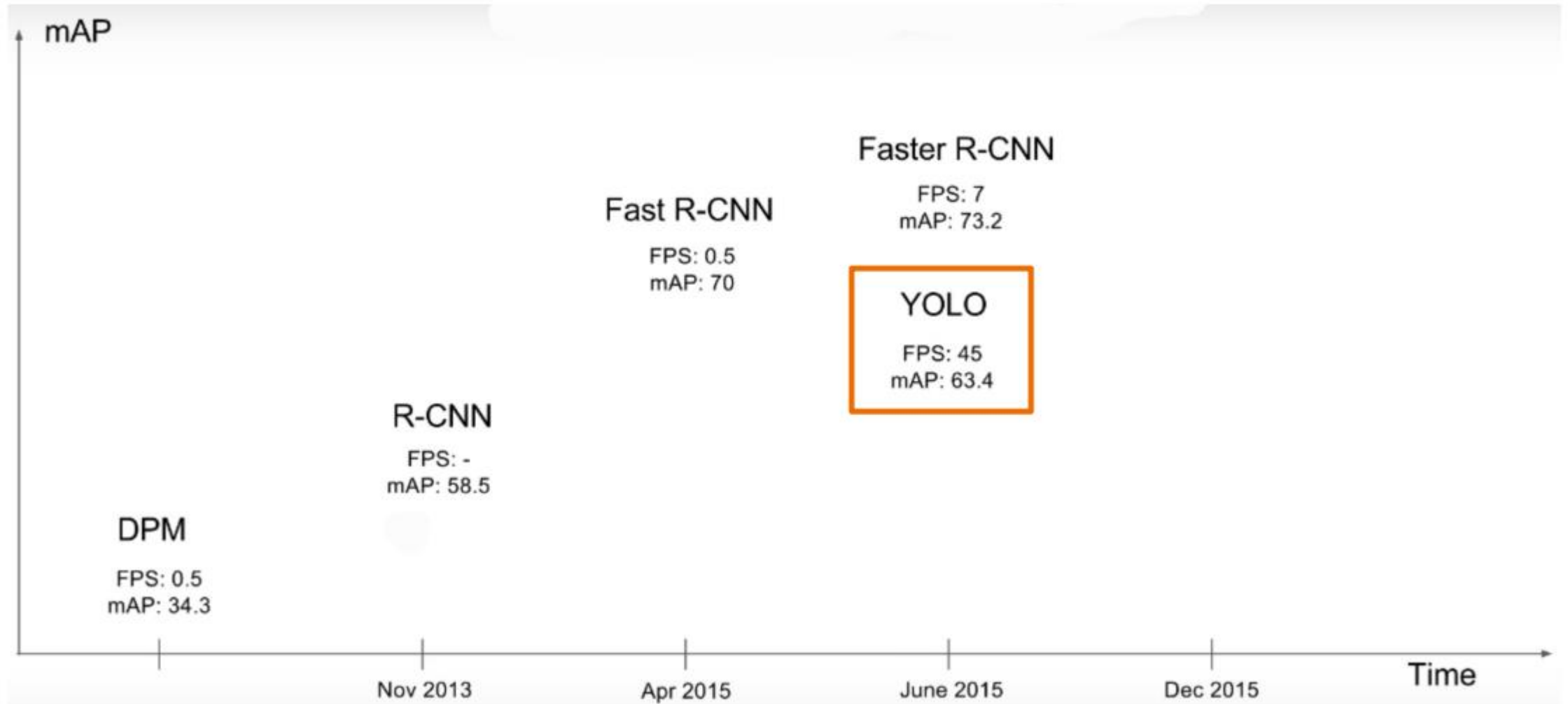[4]`reedscot@umich.edu`, [1]`{cyfu,aberg}@cs.unc.edu`

**Abstract.** We present a method for detecting objects in images using a single deep neural network. Our approach, named SSD, discretizes the output space of bounding boxes into a set of default boxes over different aspect ratios and scales per feature map location. At prediction time, the network generates scores for the presence of each object category in each default box and produces adjustments to the box to better match the object shape. Additionally, the network combines predictions from multiple feature maps with different resolutions to naturally handle objects of various sizes. Our SSD model is simple relative to methods that require object proposals because it completely eliminates proposal generation and subsequent pixel or feature resampling stage and encapsulates all computation in a single network. This makes SSD easy to train and straightforward to integrate into systems that require a detection component. Experimental results on the PASCAL VOC, MS COCO, and ILSVRC datasets confirm that SSD has comparable accuracy to methods that utilize an additional object proposal step and is much faster, while providing a unified framework for both training and inference. Compared to other single stage methods, SSD has much better accuracy, even with a smaller input image size. For $300 \times 300$ input, SSD achieves 72.1% mAP on VOC2007 `test` at 58 FPS on a Nvidia Titan X and for $500 \times 500$ input, SSD achieves 75.1% mAP, outperforming a comparable state of the art Faster R-CNN model. Code is available at `https://github.com/weiliu89/caffe/tree/ssd`.

# Faster R-CNN: Box Classification and Regression are being done 2 times.



Two stage object detection is time-consuming. Faster R-CNN is faster but not fast enough.

# YOLO: Fast but not accurate enough.

Object detection needs a good tradeoff between accuracy and speed.

Slide courtesy of DeepSystem.io "YOLO: You only look once Review"

# SSD: Single Shot MultiBox Detection (ECCV 2016)



1. Achieves competitive mAP (72.1) as faster R-CNN (73.2).
2. Much faster (58 fps) than faster R-CNN (7fps) and YOLO (45fps), making accurate real-time detection possible.
3. Makes predictions on multiple feature maps with different resolutions to handle objects of different sizes.

# SSD: Single Shot MultiBox Detector
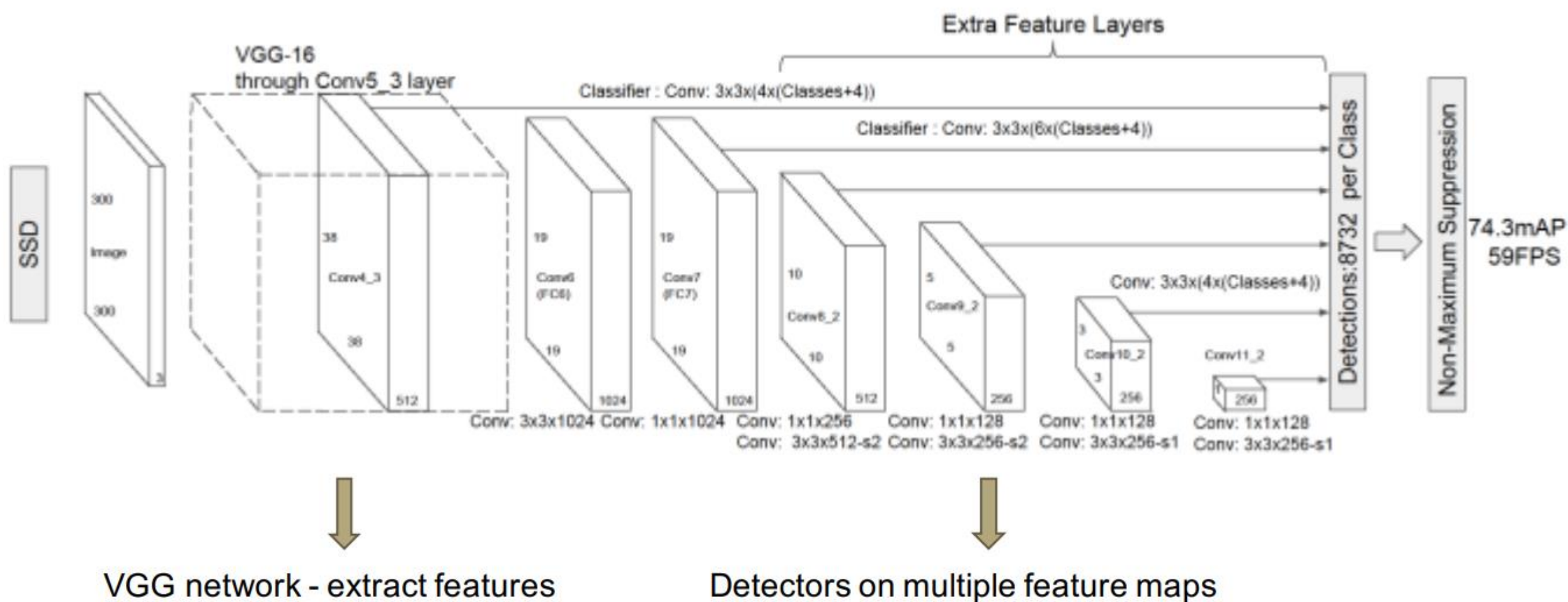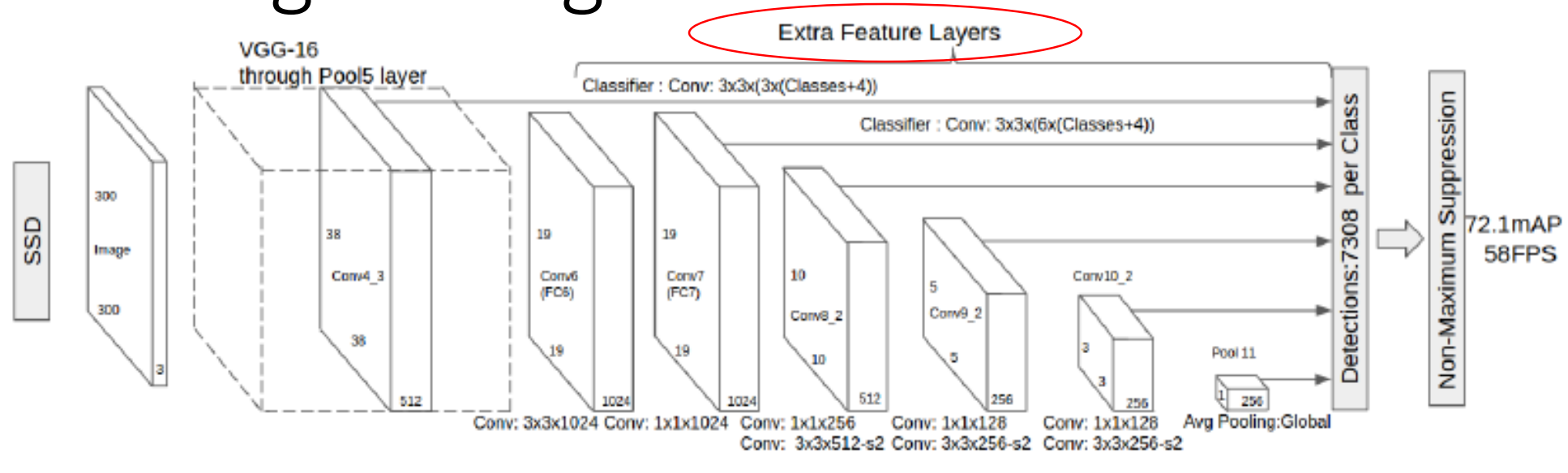
? Tool?-- A single deep neural network

? Framework?– Caffe

? Technology background?--related methods are structure-complicated and hard to bring high speed and good accuracy

• Solution and Advantages:

@Providing a unified framework

@much faster

@better accuracy, even with a smaller input image size

# Network Architecture



VGG-16 through Conv5_3 layer

Extra Feature Layers

SSD

300 × 300 Image × 3

Conv4_3
38 × 38 × 512

Classifier : Conv: 3x3x(4x(Classes+4))

Classifier : Conv: 3x3x(6x(Classes+4))

Conv6 (FC6)
19 × 19 × 1024

Conv7 (FC7)
19 × 19 × 1024

Conv8_2
10 × 10 × 512

Conv9_2
5 × 5 × 256

Conv: 3x3x(4x(Classes+4))

Conv10_2
3 × 3 × 256

Conv11_2
256

Conv: 3x3x1024   Conv: 1x1x1024   Conv: 1x1x256        Conv: 1x1x128        Conv: 1x1x128        Conv: 1x1x128
                                  Conv: 3x3x512-s2     Conv: 3x3x256-s2     Conv: 3x3x256-s1     Conv: 3x3x256-s1

Detections:8732 per Class

Non-Maximum Suppression

74.3mAP
59FPS

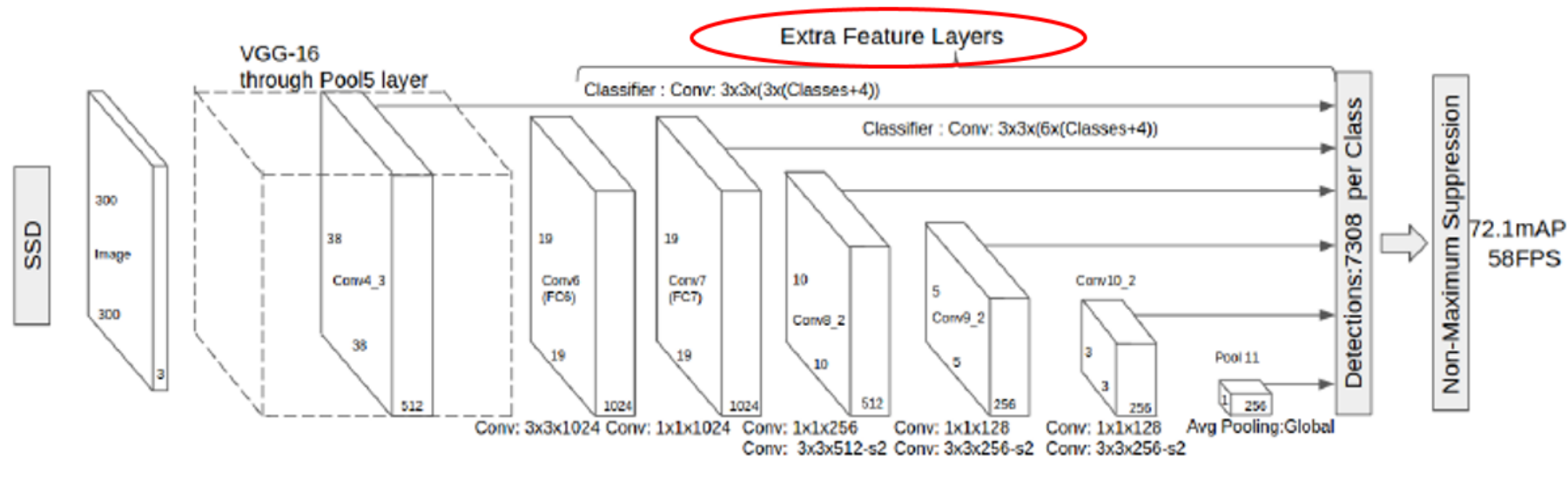VGG network - extract features

Detectors on multiple feature maps

# How SSD gets its goal?



- For speed: structure advantage!

➢Eliminating bounding box proposals and subsequent pixel or feature resampling stage

✓Adding convolution feature layers to the end of the tructed base network to predict detections at multiple scalaes

# How SSD gets its goal?



- For speed: structure advantage!

➢ Eliminating bounding box proposals and subsequent pixel or feature resampling stage

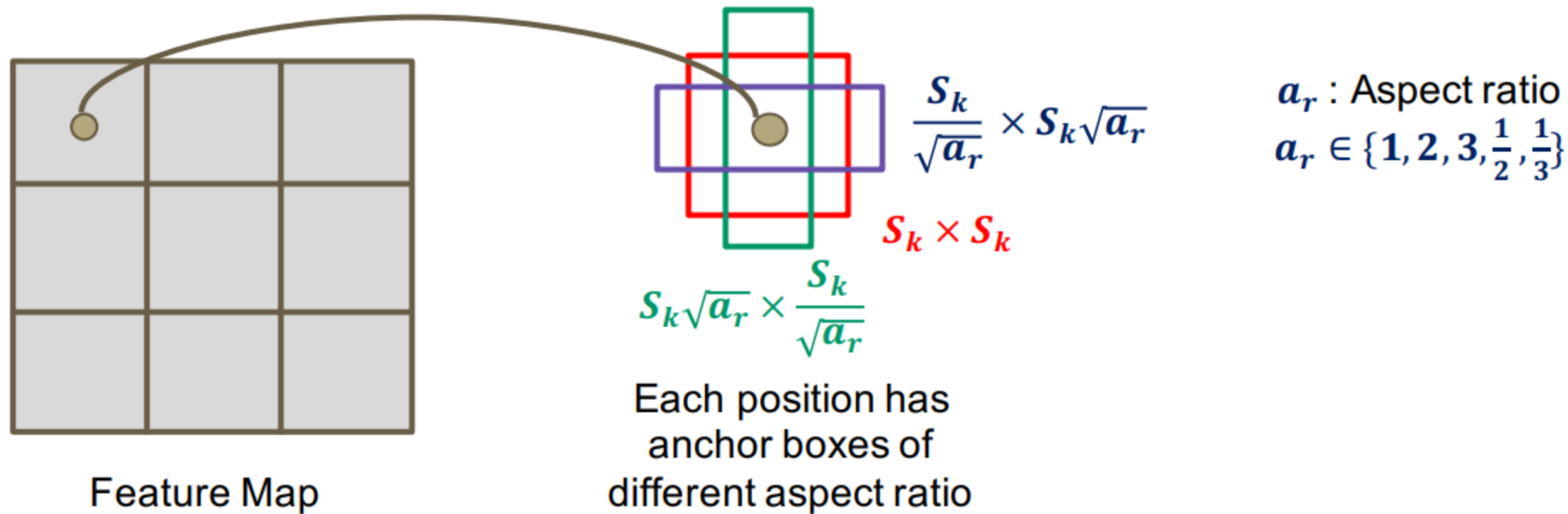✓ Adding convolution feature layers to the end of the tructed base network to predict detections at multiple scalaes

# Default (Anchor) Boxes

- Similar to faster R-CNN, SSD also uses anchor boxes.
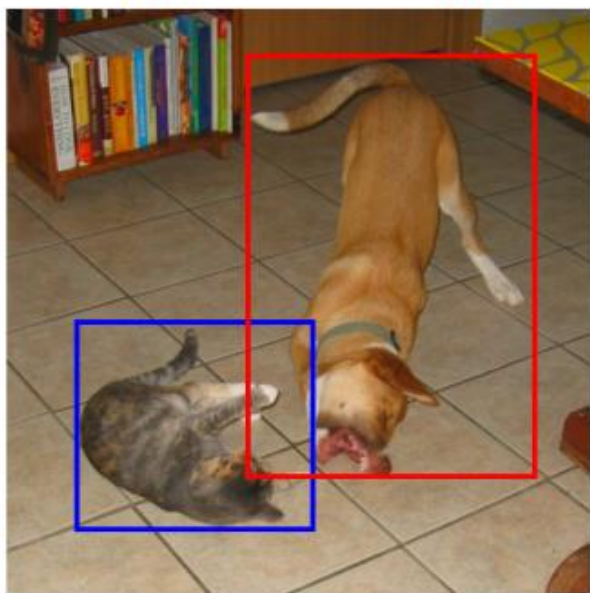- At each feature map position, anchor boxes have different aspect ratios.



$$\frac{S_k}{\sqrt{a_r}} \times S_k\sqrt{a_r}$$

$$S_k \times S_k$$

$$S_k\sqrt{a_r} \times \frac{S_k}{\sqrt{a_r}}$$

$a_r$ : Aspect ratio

$a_r \in \{1, 2, 3, \frac{1}{2}, \frac{1}{3}\}$

Feature Map

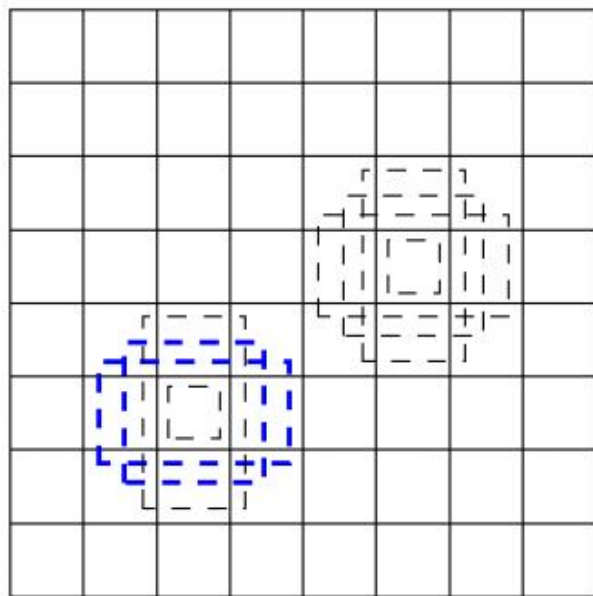Each position has anchor boxes of different aspect ratio

# Default (Anchor) Boxes

- Anchor boxes of different feature maps have unique scales. So different feature map are responsible for objects of different sizes.

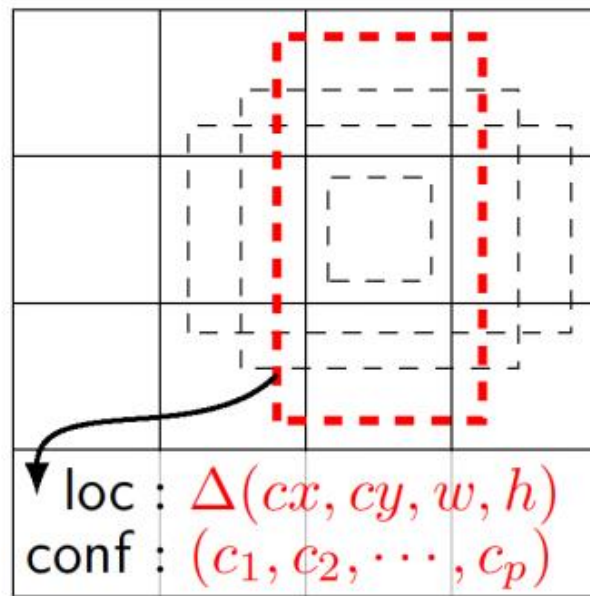$$s_k = s_{\min} + \frac{s_{\max} - s_{\min}}{m - 1}(k - 1), \quad k \in [1, m]$$

- $S_k$ denotes the scale of of the k-th feature map. $m$ is the number of feature maps for prediction. $S_{min} = 0.2, S_{max} = 0.9$.



loc : $\Delta(cx, cy, w, h)$
conf : $(c_1, c_2, \cdots, c_p)$
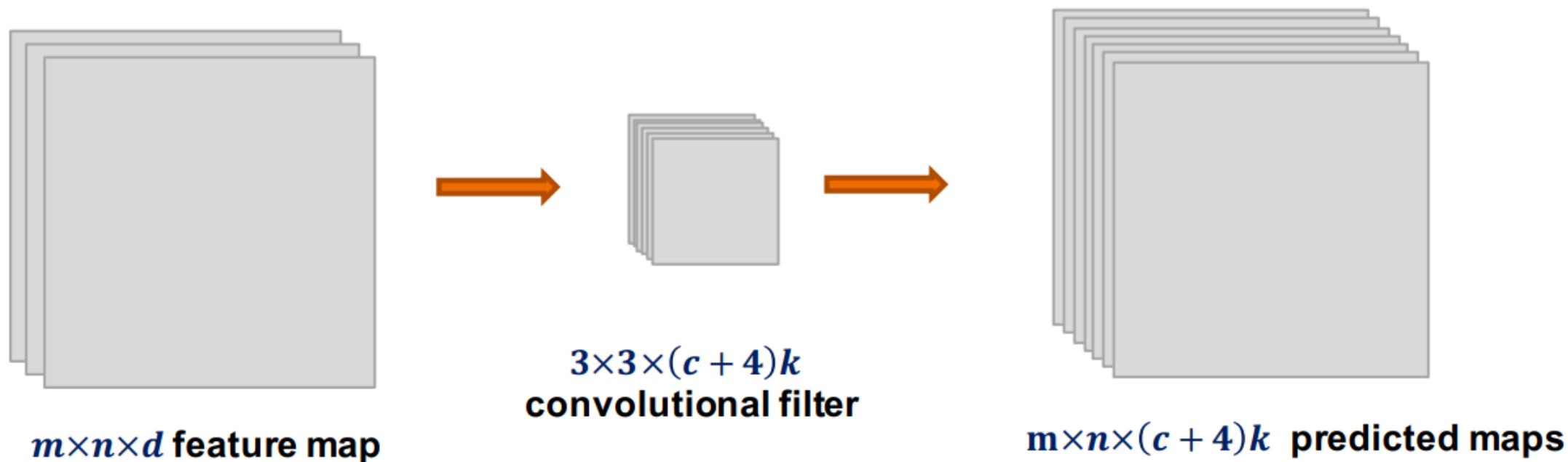
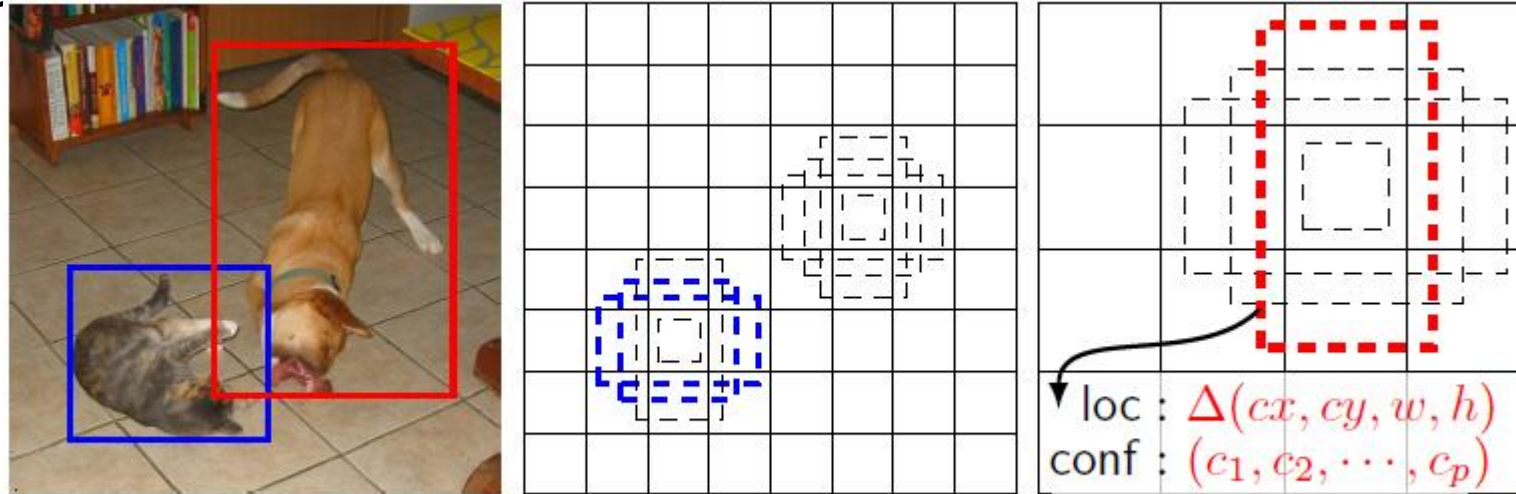(a) Image with GT boxes    (b) $8 \times 8$ feature map    (c) $4 \times 4$ feature map

# Convolutional Filters for Prediction

- On each feature map, two types of convolutional filters will be applied:
  - C filters for category prediction, where C is the number of object categories.
  - 4 filters for bounding box regression, 4 for the coordinates x, y, w, h.

- Together there will be $(c + 4)k$ filters for each feature map, $k$ is the number of anchor box types.

- The output for a $m{\times}n$ feature map will be a map of $m{\times}n{\times}(c + 4)k$, indicating the category and coordinates for each bounding box.



$$3{\times}3{\times}(c+4)k$$
**convolutional filter**

$m{\times}n{\times}d$ **feature map**
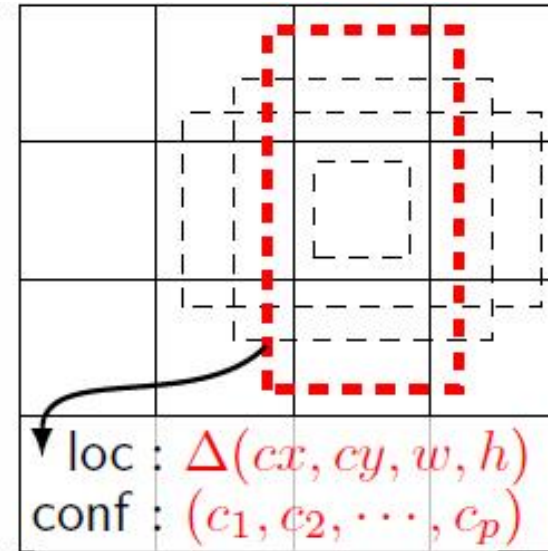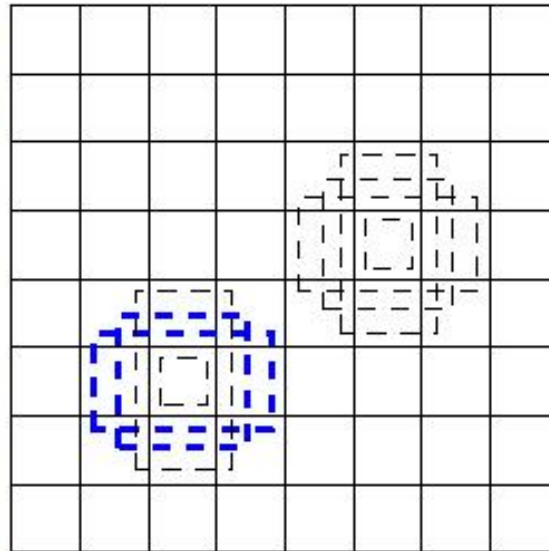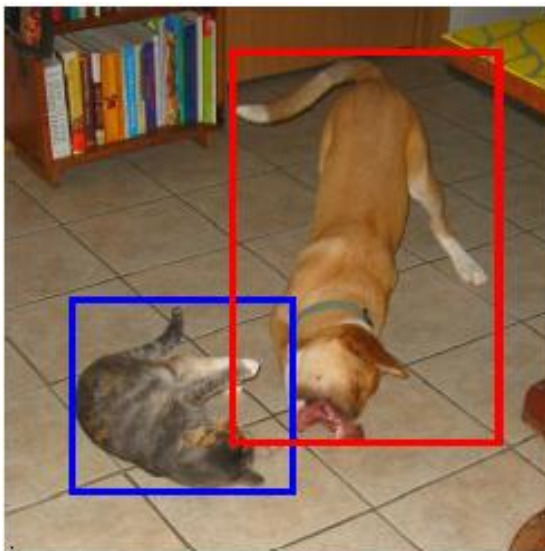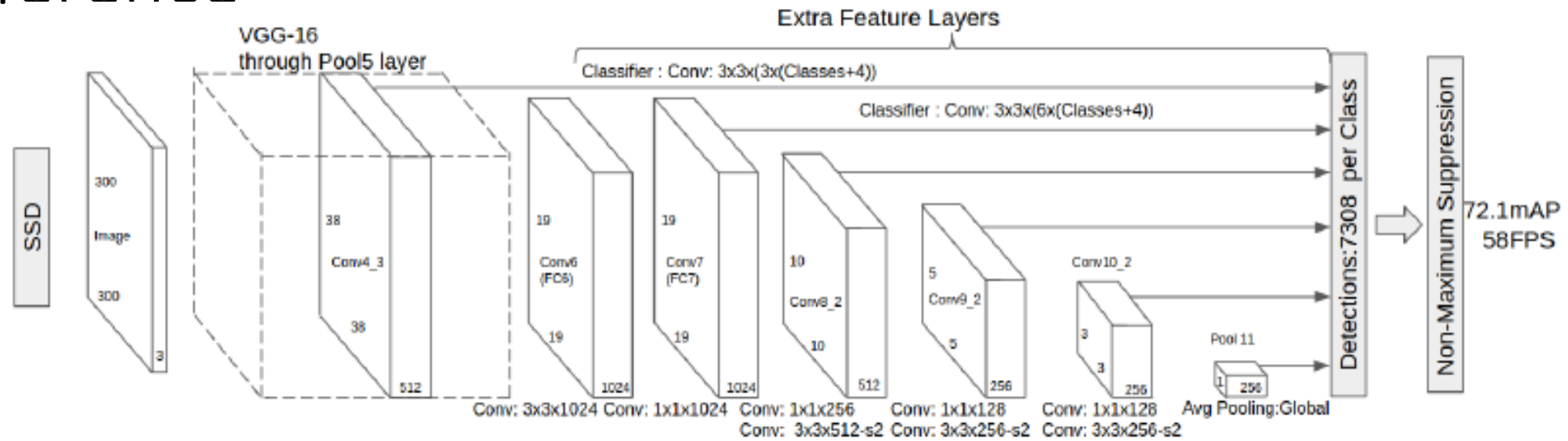
$m{\times}n{\times}(c+4)k$ **predicted maps**

# What do convolution feature layers do?



(a) Image with GT boxes    (b) $8 \times 8$ feature map    (c) $4 \times 4$ feature map

loc : $\Delta(cx, cy, w, h)$
conf : $(c_1, c_2, \cdots, c_p)$

➢ Needs an input image and ground truth boxes for each object during training

➢ Evaluate a small set(e.g.4)of default boxes of different aspect ratios at each locations in several feature maps with different scales.

➢ Filters match these default boxes to the ground truth boxes and predict both the shape offsets and confidence for all object categories for each default box

➢ The feed-forward convolutional network produces a fixed-size collection of bounding boxes and scores for the presence of object class in those boxes

# SSD overcoming the influence of resolution difference



loc : $\Delta(cx, cy, w, h)$
conf : $(c_1, c_2, \cdots, c_p)$

# Training objective

- The loss used in SSD is a combination of confidence loss and localization loss.

$$L(x, c, l, g) = \frac{1}{N}(L_{conf}(x, c) + \alpha L_{loc}(x, l, g))$$

- Confidence loss is the softmax loss over multiple classes confidences.

$$L_{conf}(x, c) = -\sum_{i \in Pos}^{N} x_{ij}^{p} log(\hat{c}_{i}^{p}) - \sum_{i \in Neg} log(\hat{c}_{i}^{0})$$

- Localization loss is the same smooth L1 loss as faster R-CNN.

$$L_{loc}(x, l, g) = \sum_{i \in Pos}^{N} \sum_{m \in \{cx, cy, w, h\}} x_{ij}^{k} \text{smooth}_{L1}(l_{i}^{m} - \hat{g}_{j}^{m})$$

$$\hat{g}_{j}^{cx} = (g_{j}^{cx} - d_{i}^{cx})/d_{i}^{w} \qquad \hat{g}_{j}^{cy} = (g_{j}^{cy} - d_{i}^{cy})/d_{i}^{h}$$

$$\hat{g}_{j}^{w} = \log\left(\frac{g_{j}^{w}}{d_{i}^{w}}\right) \qquad \hat{g}_{j}^{h} = \log\left(\frac{g_{j}^{h}}{d_{i}^{h}}\right)$$

# Matching Strategy

- During training, we need to determine which anchor boxes are corresponding to ground truth boxes.

  1. Match each ground truth box to the anchor box with the best jaccard overlap.
  2. Match default box to any ground truth with jaccard overlap higher than 0.5.

- Question : Why not using predicted box (anchor box after regression) for matching?

- Question : Negative anchor boxes are much more than positive anchor boxes, how to deal with this imbalance?

# Hard Negative Mining

- Number of negative anchor boxes  >>  number of positive anchor boxes
- Calculate the **confidence loss** of each negative anchor boxes.
- Select anchor boxes that have **highest** loss as negative training samples.
- Keep the ratio between negatives and positives **3:1**.
- This leads to faster convergence and stable training.

- Question: Is this approach sufficient for training negative samples?
  - **No**. A recent paper show that negative samples need more elegant loss calculation.
  - (ICCV 2017 best student paper award: *Focal Loss for Dense Object Detection*, Lin et al.)

- Base network: VGG16, pretrained on ILSVRC dataset.
- Add layers on top of VGG Conv5 layer.
- Use dilated convolution in Conv6 layer.

# Results – small objects

- SSD does not perform well on small objects.
- No feature resampling step in SSD.
- Relatively low-level feature maps are responsible for detecting small objects. These feature maps do not have sufficient high-level semantic information.

# Data Augmentation


crop & resize

- Photo-metric distortions.
- Random crop:
  - Use the entire original input image.
  - Crop a patch so that the minimum jaccard overlap with the objects is 0.1, 0.3, 0.5, 0.7 or 0.9.
  - Randomly crop a patch.
- Question: Will data augmentation also benefit this much to Faster R-CNN?
  - No, because faster R-CNN uses a feature pooling step which is relatively robust to object translation.
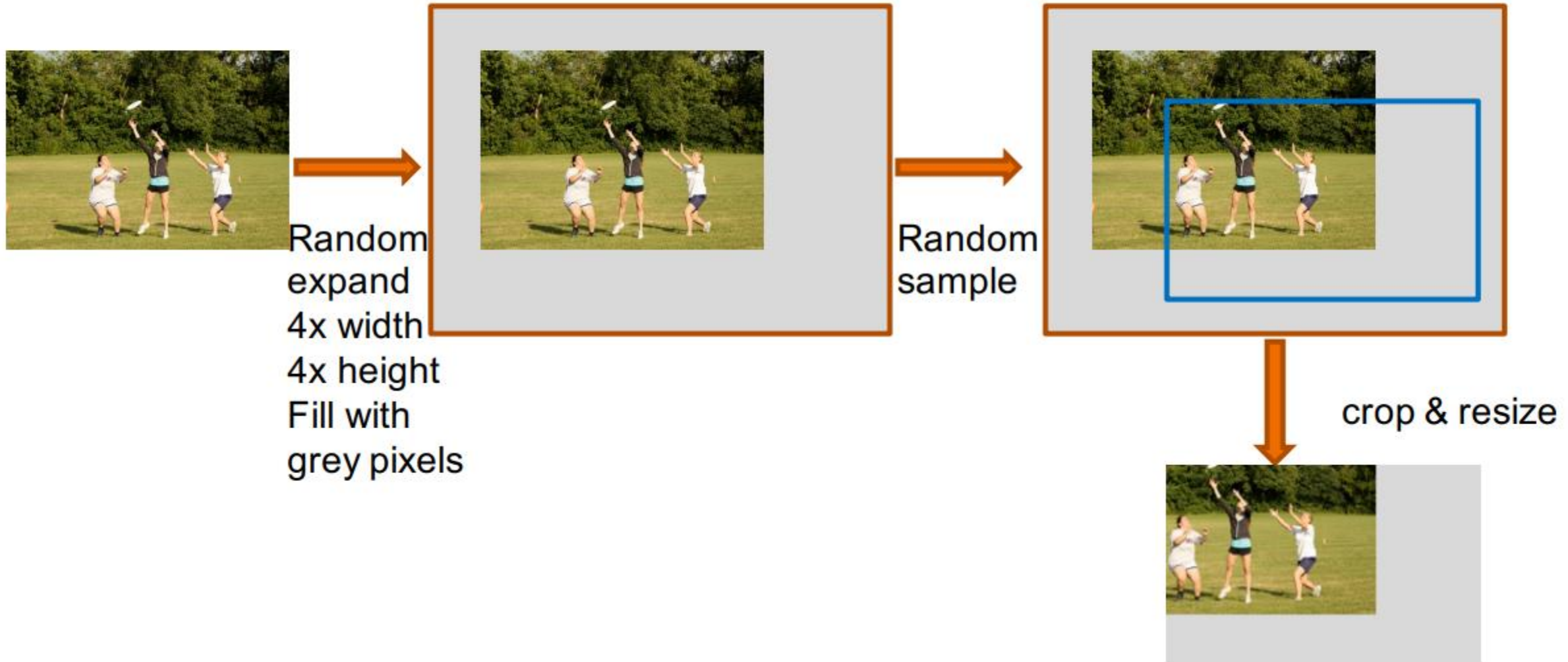
# More data augmentation for small objects

- Keep small objects small.



Random expand
4x width
4x height
Fill with
grey pixels

Random sample

crop & resize

# Quantity choice

- Imposing different aspects ratios for the default boxes, and denote them as :

$$a_r = \left[ 1, 2, 3, \frac{1}{2}, \frac{1}{3} \right]$$

- Instead of using all the negative examples, SSD sorts them using the highest confidence for each default box and pick the top ones so that the ratio between the negatives and positives is at most 3:1——leading to faster optimization and more stable training

# Quantity choice

$$L(x,c,l,g) = \frac{1}{N}(L_{conf}(x,c) + \alpha L_{loc}(x,l,g))$$

--The overall objective loss function is a weighted sum of the localization loss and the confidence loss(conf)
--N: the number of matched default boxes
--l: predicted boxes
--g: the ground truth box
--x=1 denotes some certain default box is matched to a ground truth box

# Choosing scales and aspect ratios for default boxes:

▷ If *m* feature maps are used for prediction, the **scale** of the default boxes for each feature map is:

$$s_k = s_{\min} + \frac{s_{\max} - s_{\min}}{m - 1}(k - 1), \quad k \in [1, m]$$

$s_{\min} = \quad 0.1$

$s_{\max} = \quad 0.95$

# Choosing Scales and Aspect Ratios for Default Boxes

# Choosing Scales and Aspect Ratios for Default Boxes

- At each scale, different aspect ratios are considered



(a) Image with GT boxes    (b) $8 \times 8$ feature map    (c) $4 \times 4$ feature map

$$a_r \in \{1, 2, 3, \tfrac{1}{2}, \tfrac{1}{3}\}$$

$$w_k^a = s_k \sqrt{a_r}$$

$$h_k^a = s_k / \sqrt{a_r}$$

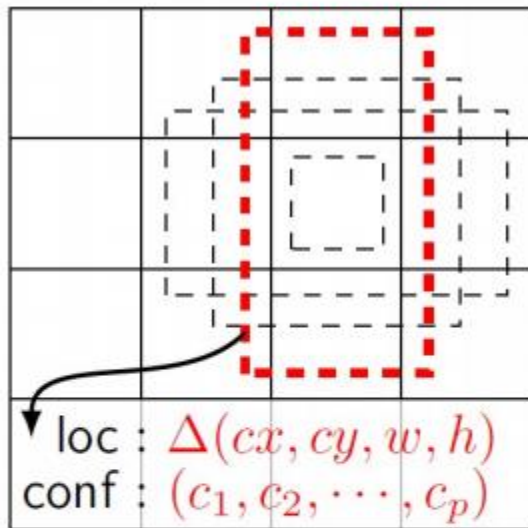For the aspect ratio of 1, one default box is added whose scale is

$$s_k' = \sqrt{s_k s_{k+1}}$$

# Hard Negative Mining

- Significant imbalance between positive and negative training examples.
  - After the matching step, most of the default boxes are negatives, especially when the number of possible default boxes is large.

- Sorting them using the highest confidence loss for each default box.
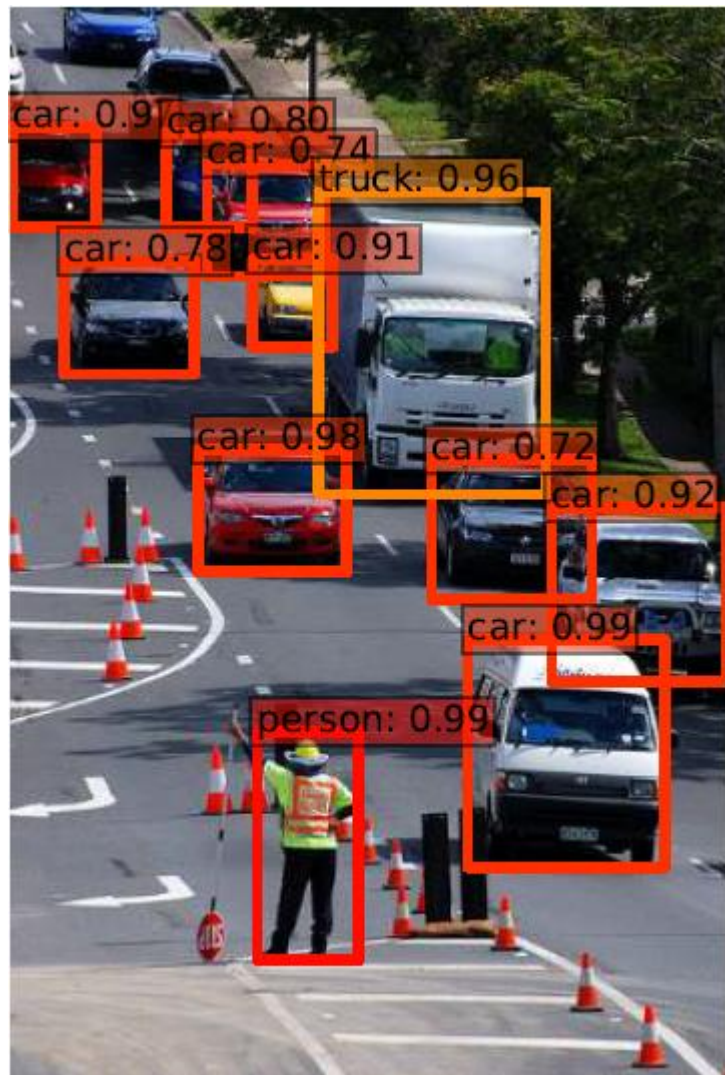- Pick the top ones so that the ratio between the negatives and positives is at most 3:1
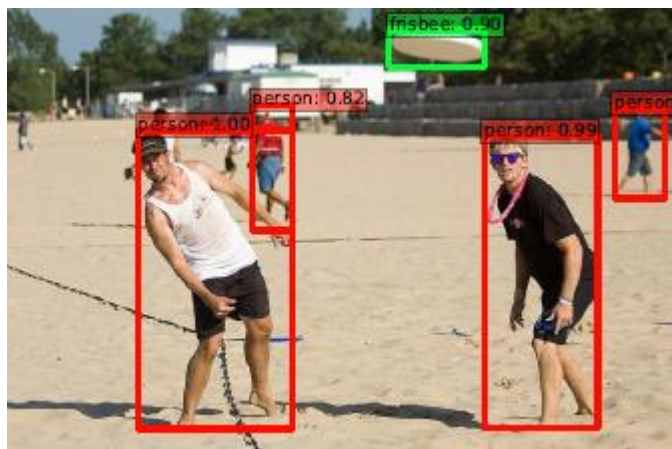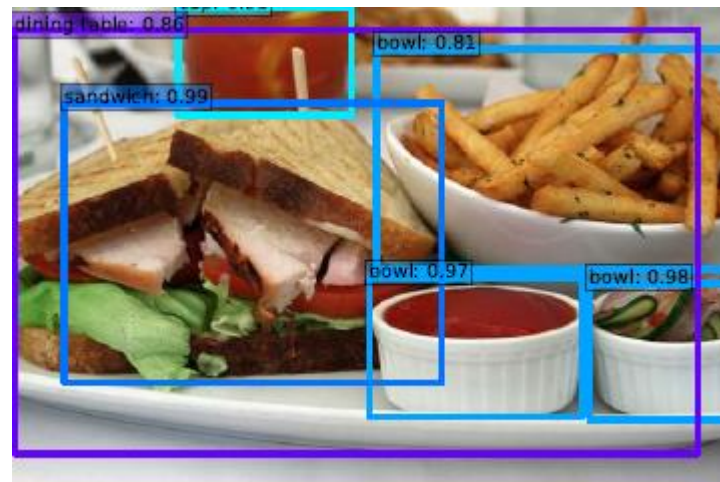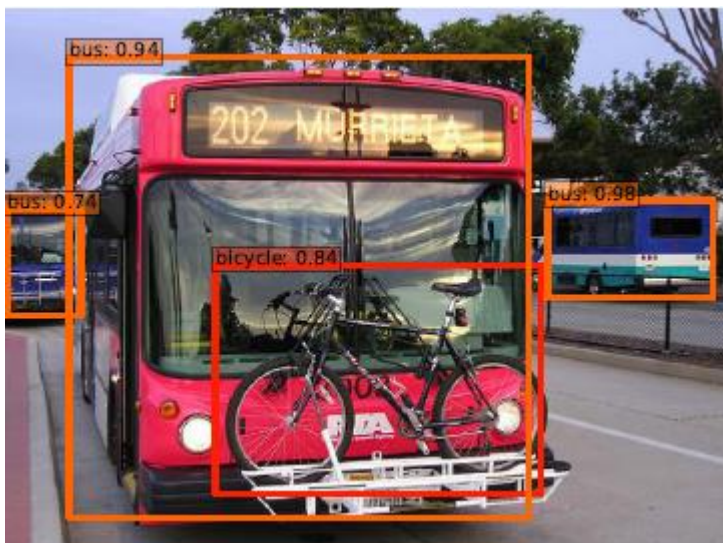
# Data Augmentation

- Use the entire original input image.
- Sample a patch so that the minimum jaccard overlap with the objects is 0.1, 0.3, 0.5, 0.7 or 0.9.
- Randomly sample a patch.
- The size of each sampled patch is [0,1, 1] of the original image size
- Aspect ratio is between ½ and 2.
- Horizontally flipped with probability of 0.5

# Results - comparison with other methods

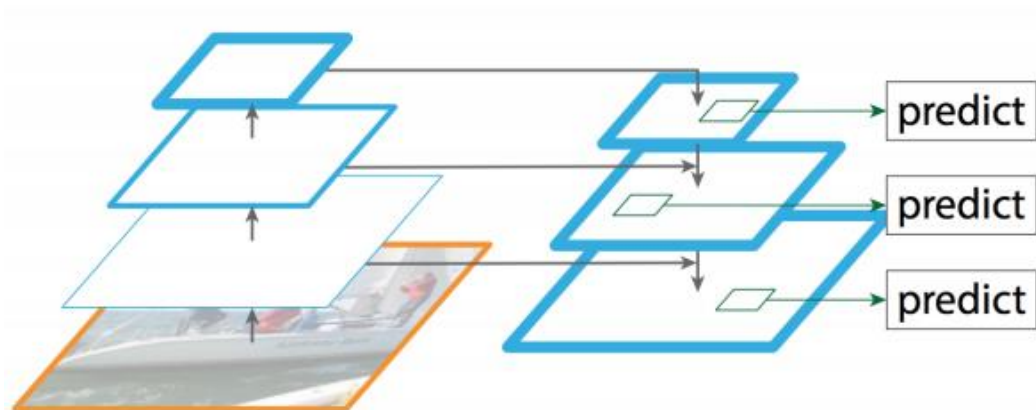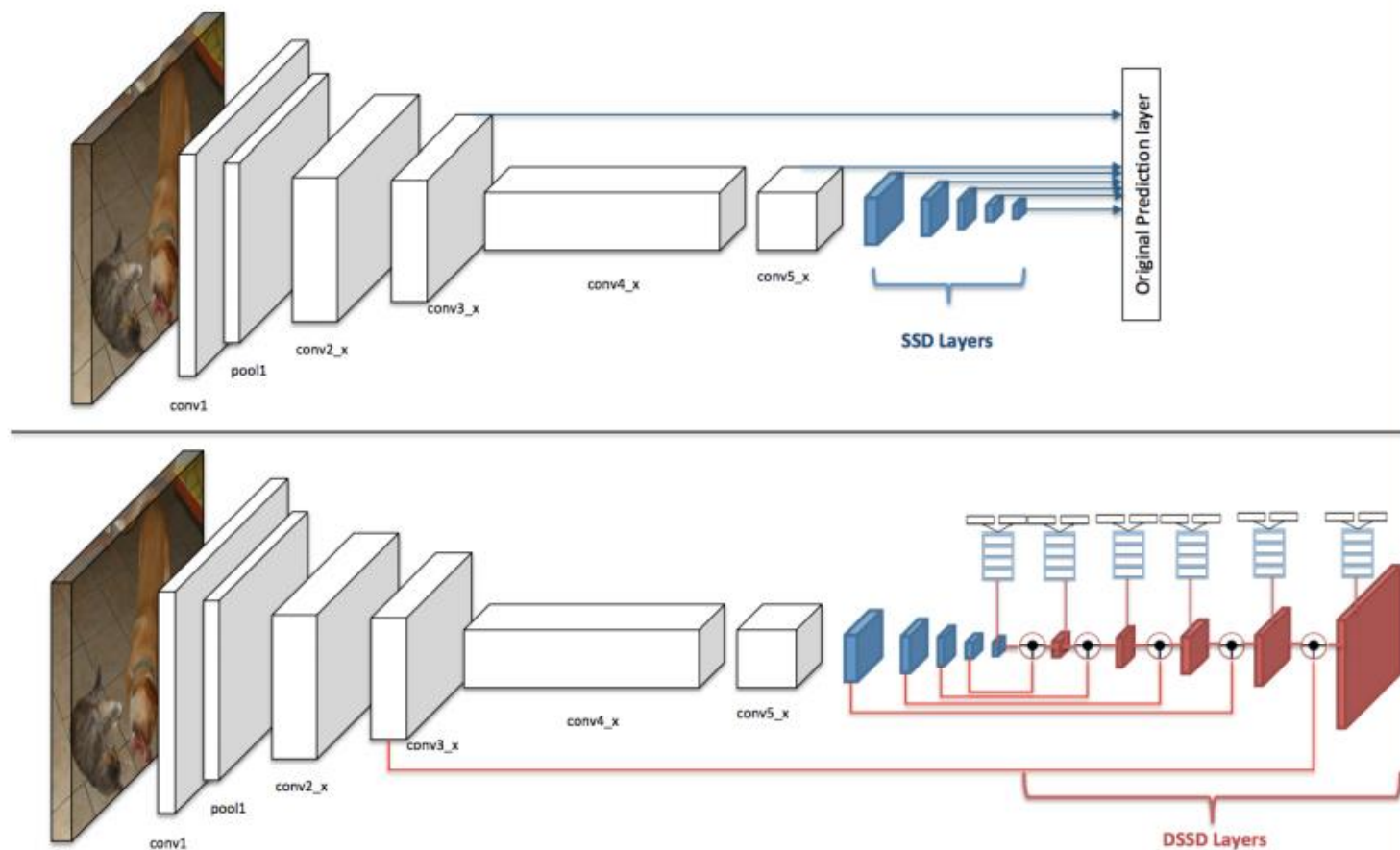- SSD leads in detection speed; it is good at speed / accuracy tradeoff.

# Conclusion

- Introduces a single-stage detector for object detection.
- Uses convolutional predictor on multiple feature maps, each responsible for a unique scale of objects.
- Achieves competitive accuracy and faster speed on various datasets.

# Extensions?

- Use skip connections and deconvolution to integrate low-level location information and high-level semantic information.



*Feature Pyramid Networks for Object Detection (CVPR 2017)*

*DSSD : Deconvolutional Single Shot Detector (2017)*