

Extra Class

YOLO – V1

Nguyen Quoc Thai

CONTENT

(1) – Object Detection

(2) – YOLO V1

(3) – VOC Dataset

1 – Object Detection



Object Detection

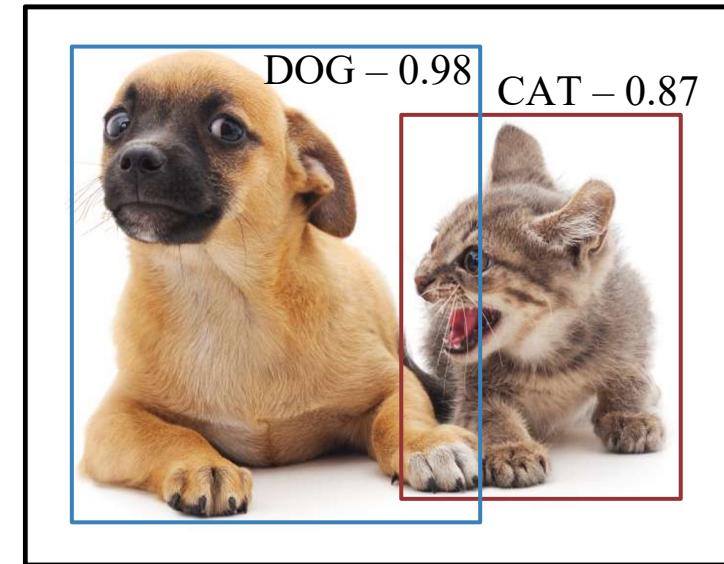
Image Segmentation

0	0	0	0	0	0	0	0
0	1	1	0	0	0	0	0
0	1	1	1	0	2	2	0
0	1	1	1	0	2	2	0
0	1	1	1	2	2	2	0
0	1	1	1	1	2	2	0
1	1	1	1	1	2	2	0
0	0	0	0	0	0	0	0

DOG

CAT

Object Detection

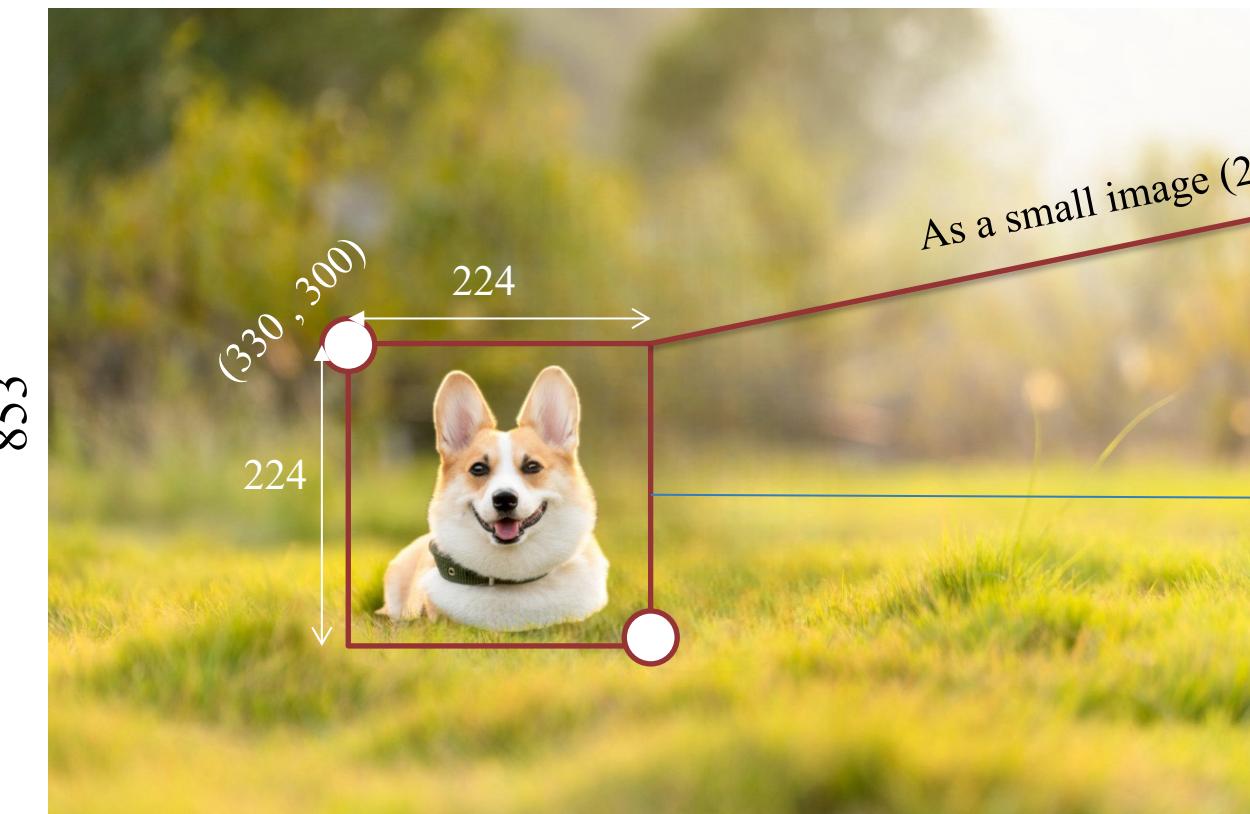


Assign labels, bounding boxes
to objects in the image

1 – Object Detection



Naïve Object Detection



As a small image (224x224)



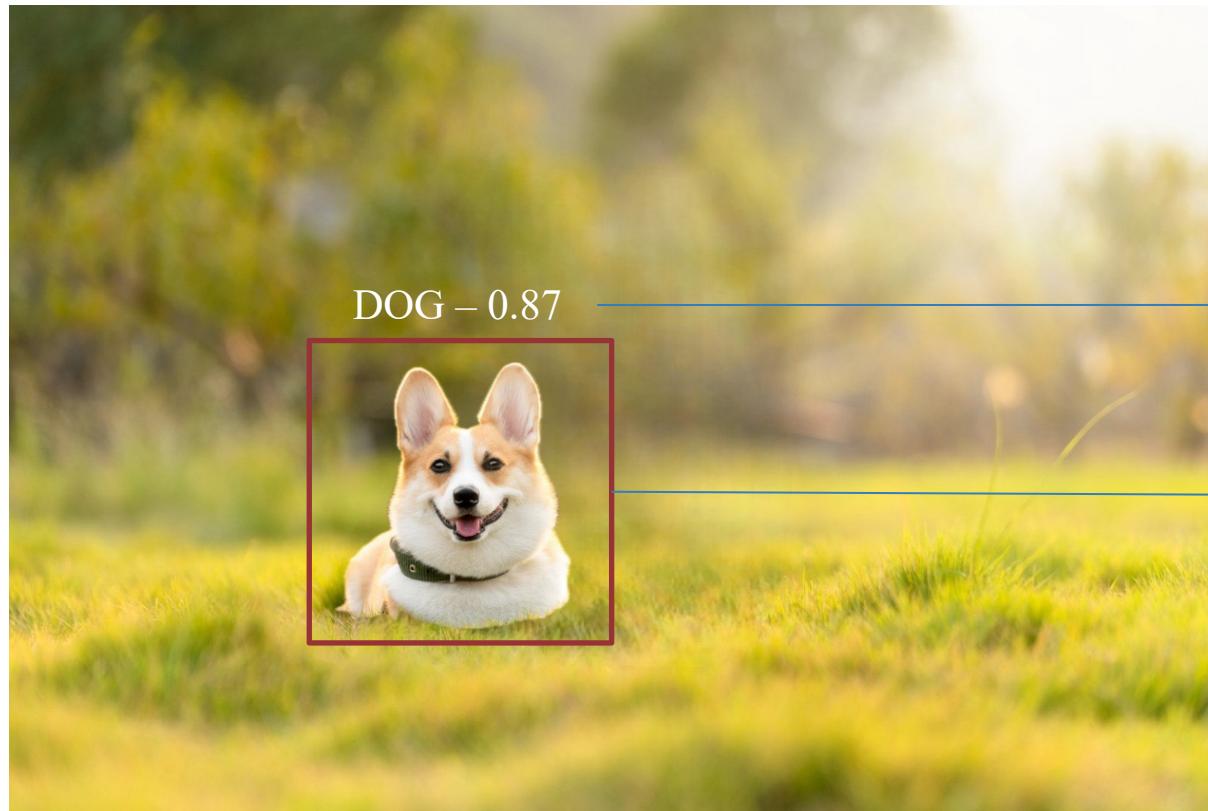
Bounding Box

A rectangular (Describe the spatial location of an object)

1 – Object Detection

!

Naïve Object Detection



Stage 3: Post-Processing

Label

Stage 2: Classification

Bounding Box

Stage 1: Region Proposal

2 – YOLO V1



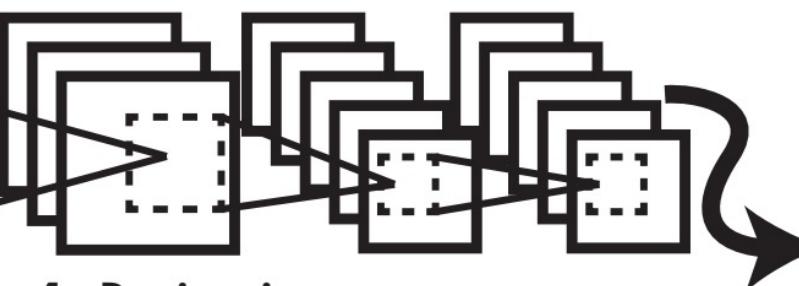
YOLO – V1

You Only Look Once Unified, Real-Time Object Detection

(1) A single neural network for localization and for classification

(2) Need to inference only once

(3) Looks at the entire image each time leading to less false positives



1. Resize image.
2. Run convolutional network.
3. Non-max suppression.

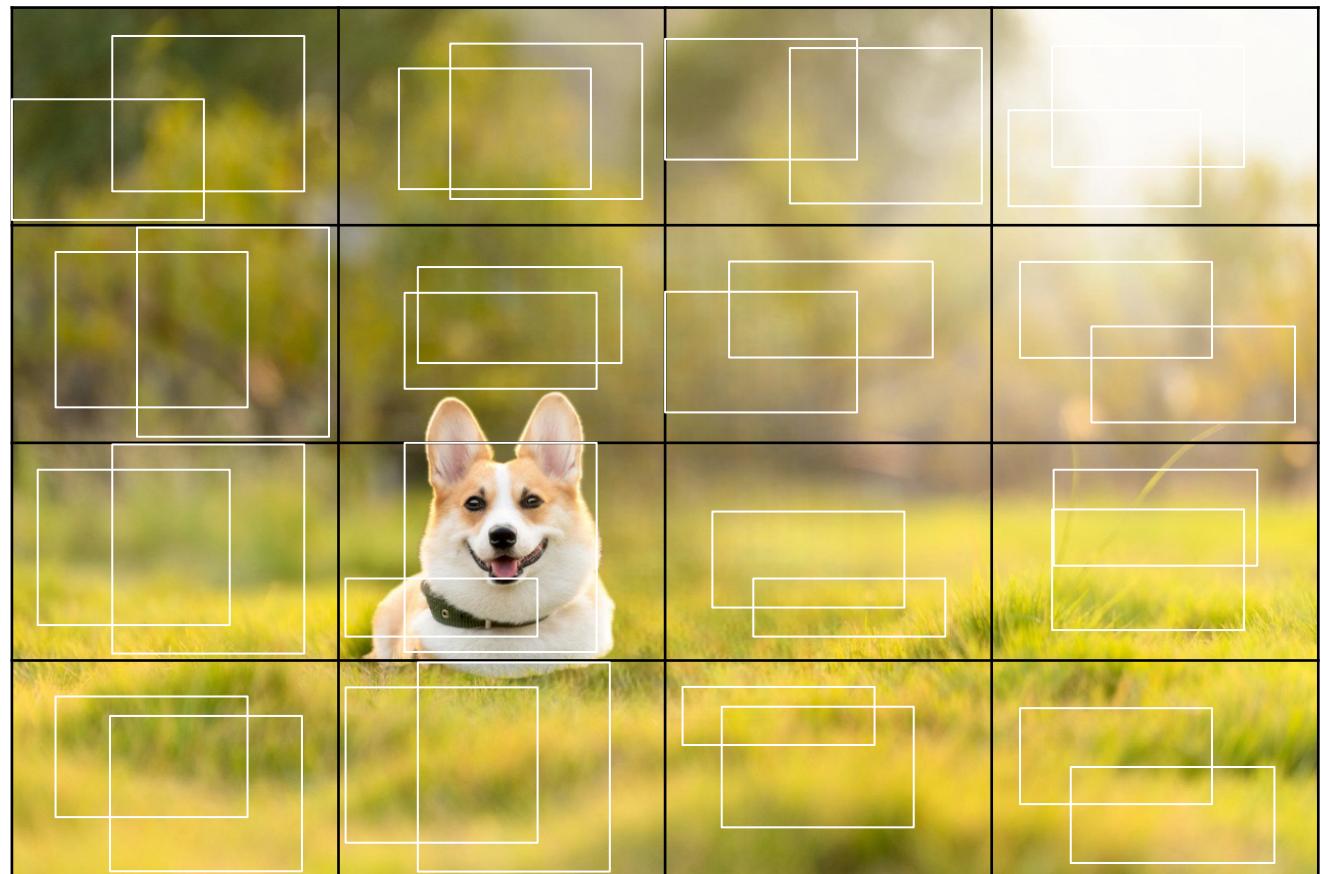


2 – YOLO V1



Unified Detection

- ❖ Split into a $S \times S$ grid
- ❖ For each grid square, generate B bounding boxes
- ❖ For each bounding box, 5 predictions:
(x, y, w, h , confidence)



2 – YOLO V1



Unified Detection

- ❖ Split into a $S \times S$ grid
- ❖ For each grid square, generate B bounding boxes
- ❖ For each bounding box, 5 predictions: $(x, y, w, h, \text{confidence})$

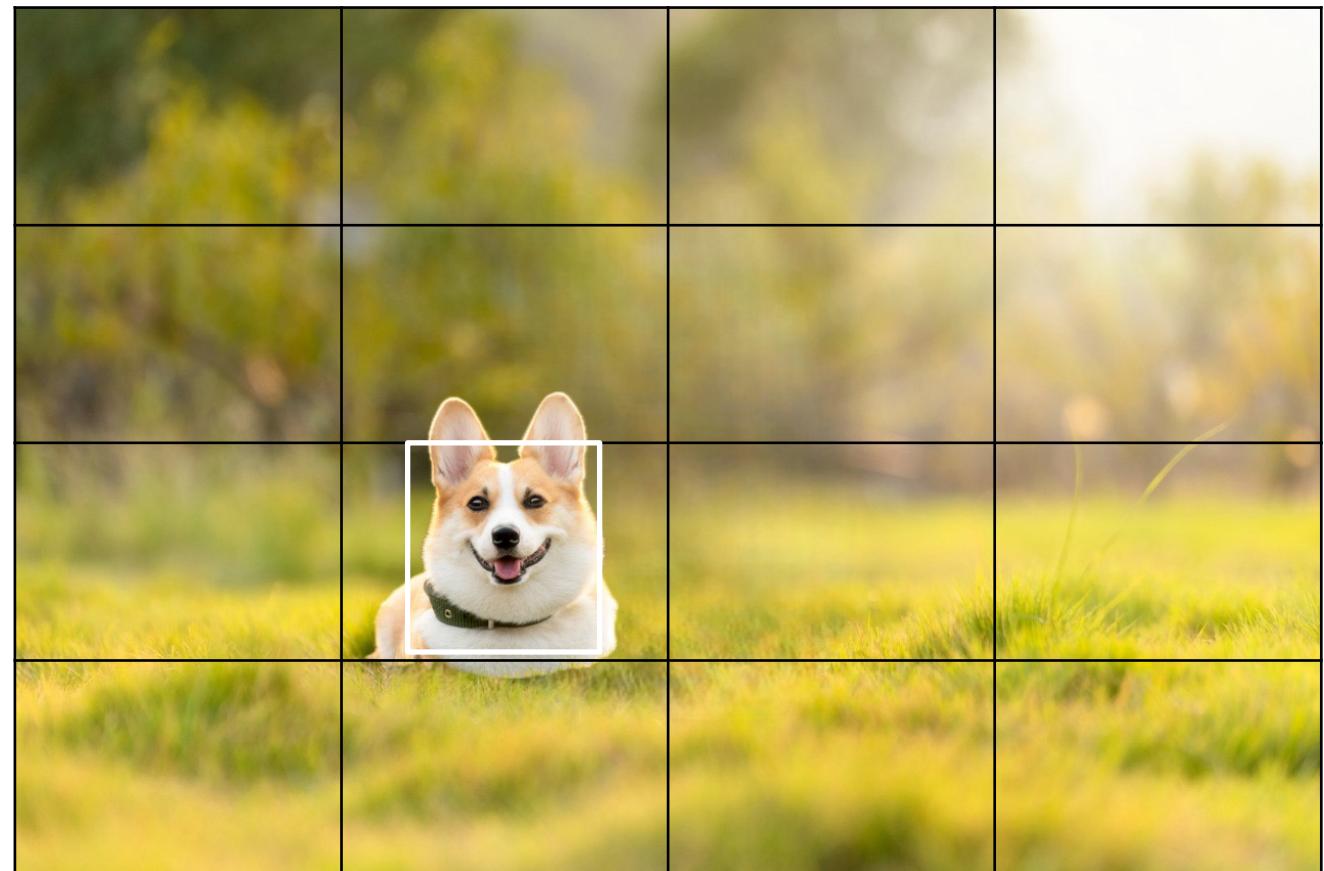
- (x, y) : the center of the box relative to the bounds of the grid

$$- h = \frac{h_{box}}{H_{image}}$$

$$- w = \frac{w_{box}}{W_{image}}$$

- confidence score

$$\text{Pr(Object)} * \text{IOU}_{\text{pred}}^{\text{truth}}$$



2 – YOLO V1



Unified Detection – Demo

- ❖ Split into a $S \times S$ grid
- ❖ For each grid square, generate B bounding boxes
- ❖ For each bounding box, 5 predictions: $(x, y, w, h, \text{confidence})$

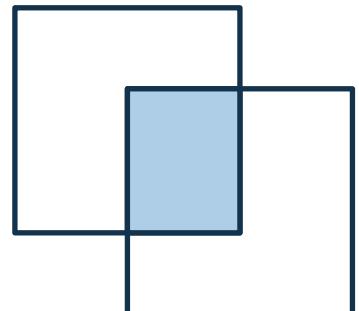
- (x, y) : the center of the box relative to the bounds of the grid

$$- h = \frac{h_{box}}{H_{image}}$$

$$- w = \frac{w_{box}}{W_{image}}$$

- confidence score

$$\text{Pr(Object)} * \text{IOU}_{\text{pred}}^{\text{truth}}$$



$$\text{IoU} = \frac{\text{Area of Overlap}}{\text{Area of Union}} =$$

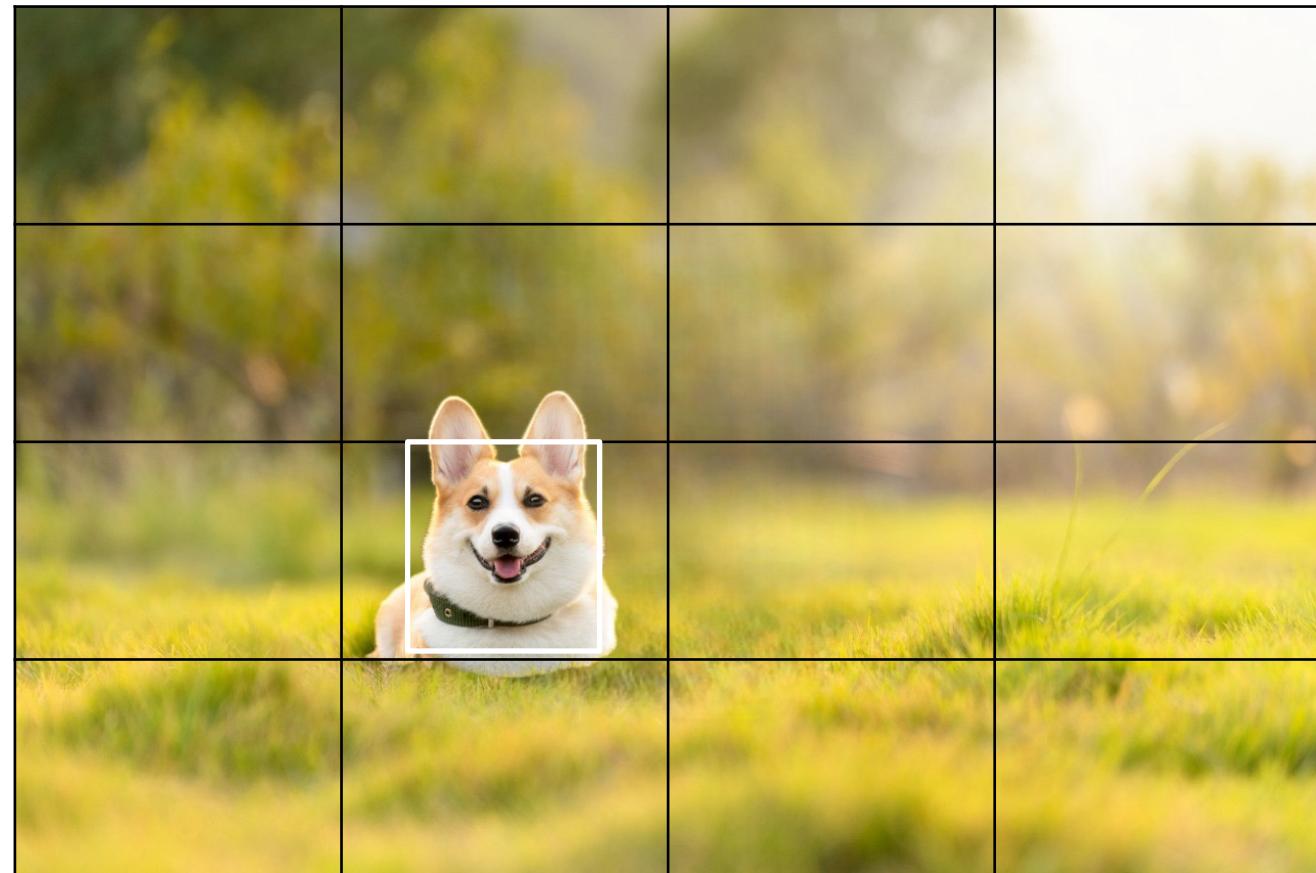


2 – YOLO V1



Unified Detection

- ❖ Split into a $S \times S$ grid
- ❖ For each grid square, generate B bounding boxes
- ❖ For each bounding box, 5 predictions:
 $(x, y, w, h, \text{confidence})$
- ❖ C : conditional class probabilities:
 $\text{Pr}(\text{Class}_i | \text{Object})$



2 – YOLO V1



Unified Detection

- ❖ Split into a $S \times S$ grid
- ❖ For each grid square, generate B bounding boxes
- ❖ For each bounding box, 5 predictions: $(x, y, w, h, \text{confidence})$
- ❖ C : conditional class probabilities: $\Pr(\text{Class}_i | \text{Object})$

YOLO is a regression algorithm

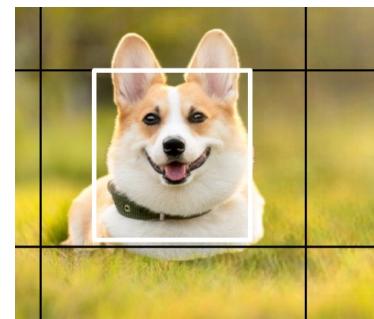
X



$W * H * 3$

Width
Height
RGB Values

Y



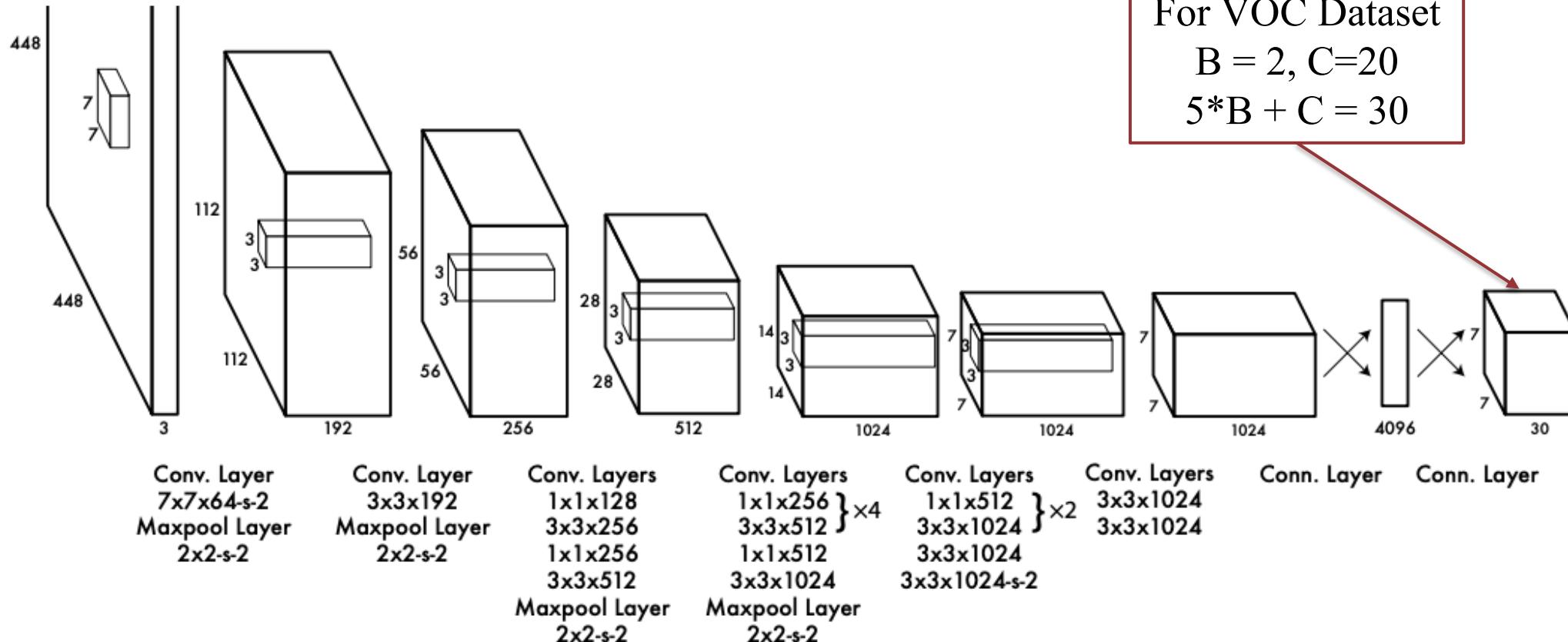
$S * S * (B * 5 + C)$

$(B * 5 + C)$
predictions + class
predicted distribution for a
grid block

2 – YOLO V1



Network Design – Demo

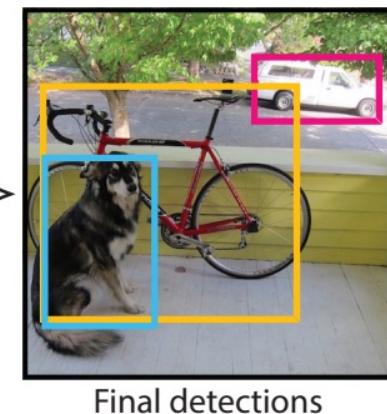
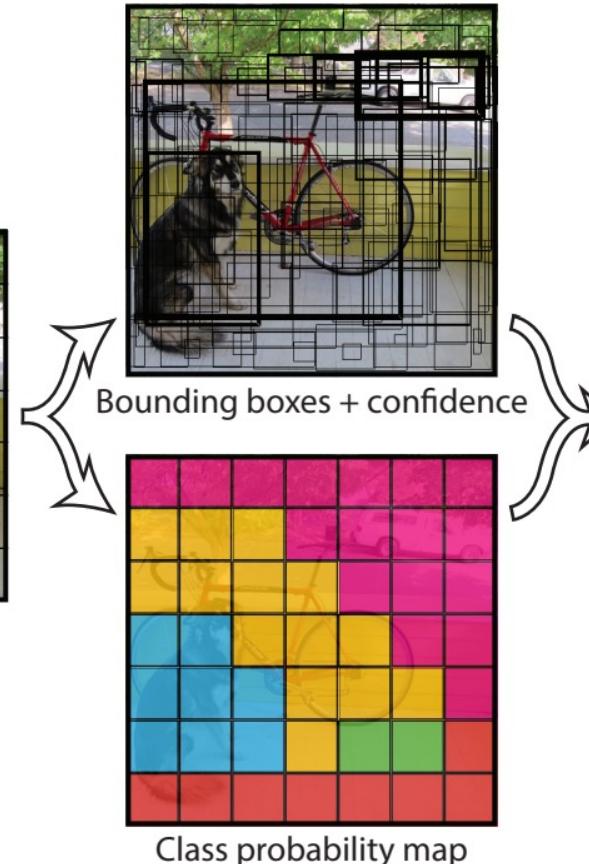
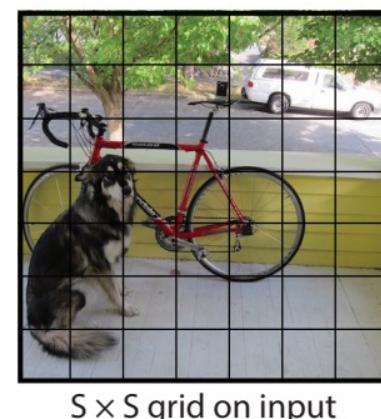


2 – YOLO V1



Inference

- ❖ Use a threshold to filter out bounding boxes with low $P(\text{Object})$
- ❖ Multiply the conditional class probabilities and the individual box confidence predictions



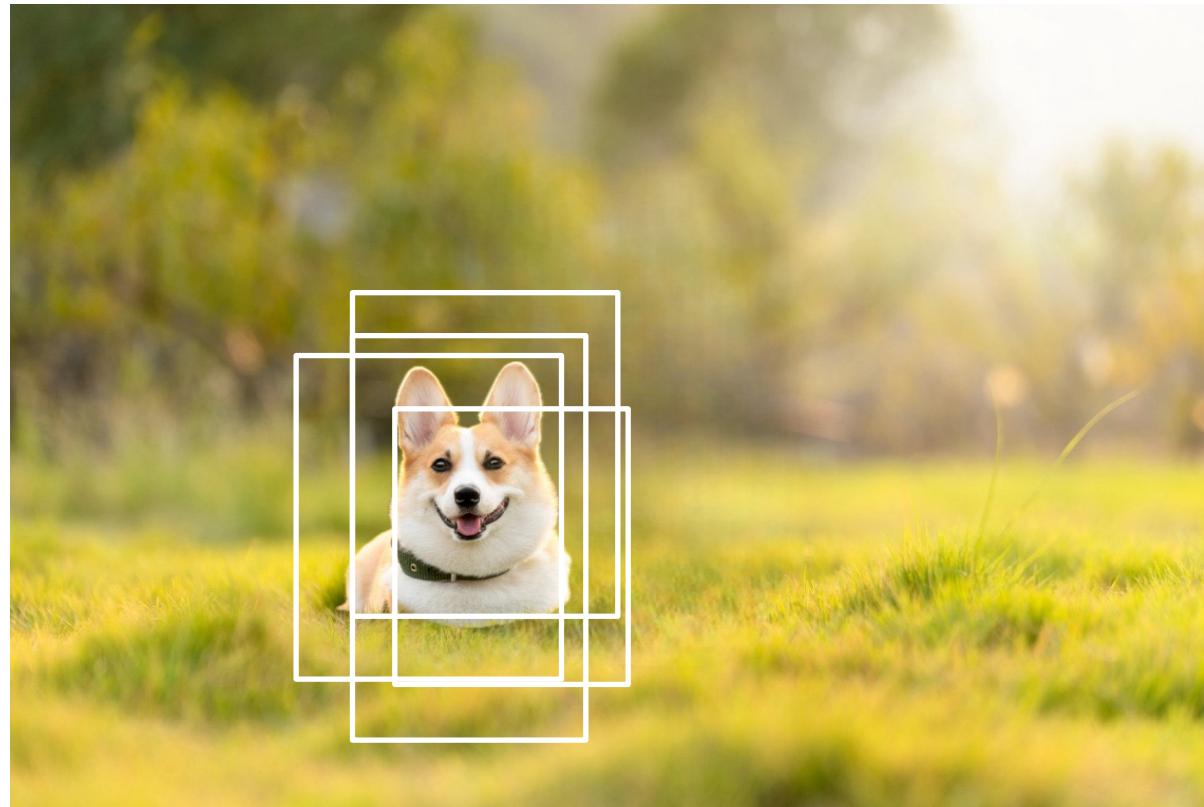
$$\Pr(\text{Class}_i | \text{Object}) * \Pr(\text{Object}) * \text{IOU}_{\text{pred}}^{\text{truth}} = \Pr(\text{Class}_i) * \text{IOU}_{\text{pred}}^{\text{truth}}$$

2 – YOLO V1

!

Non-maximal Supresion

- ❖ Discard bounding box with high overlap (keep the bounding box with highest confidence)

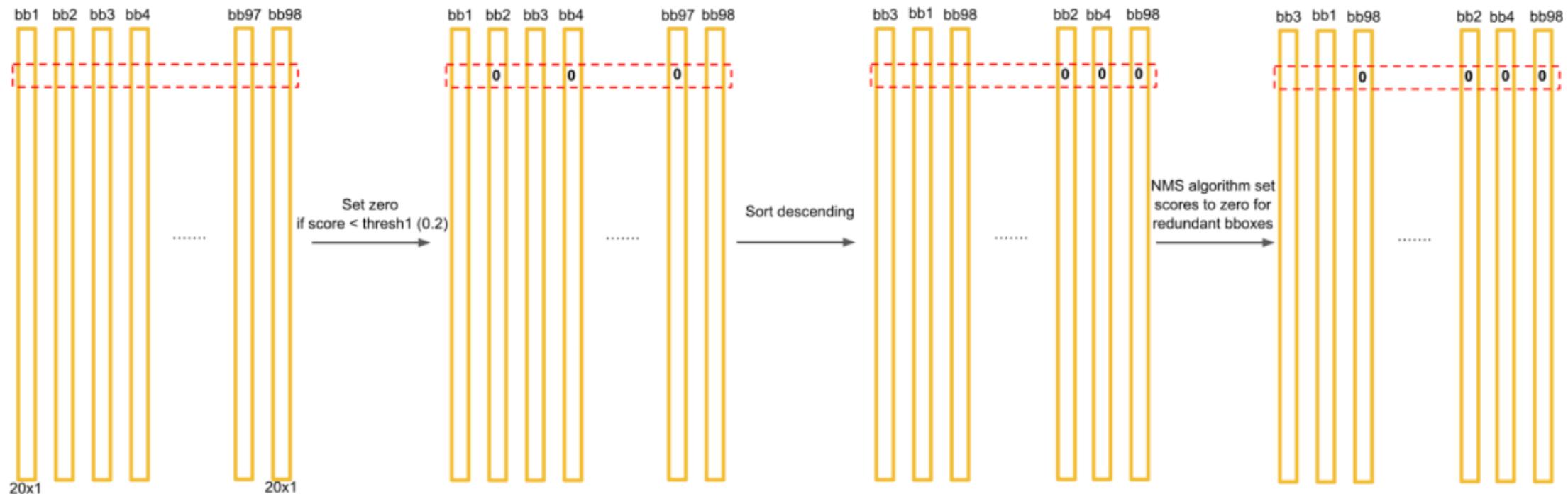


2 – YOLO V1



Non-maximal Suppresion – Demo

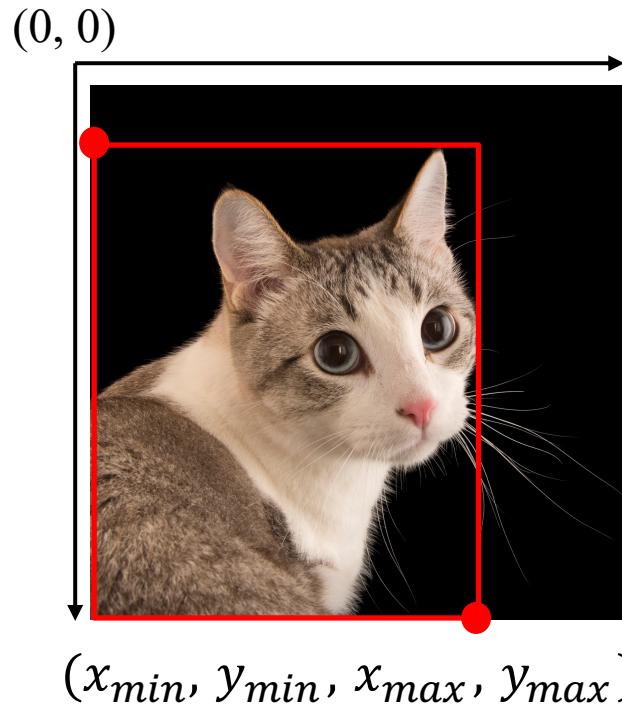
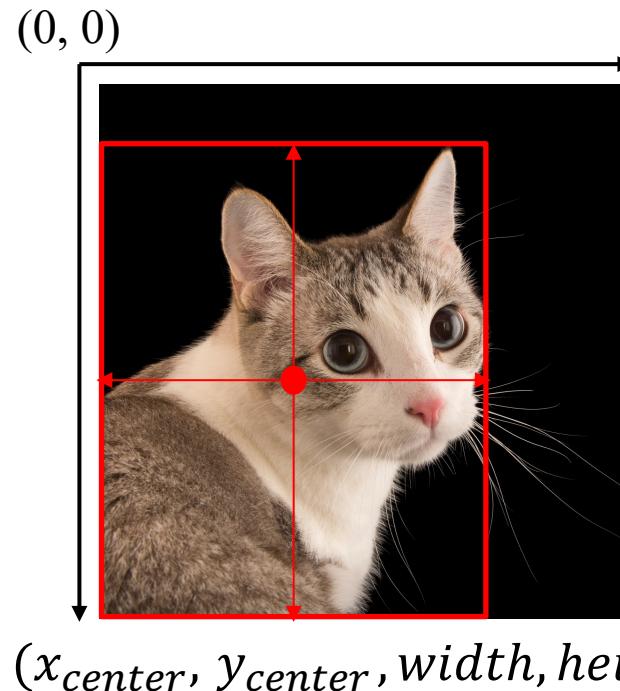
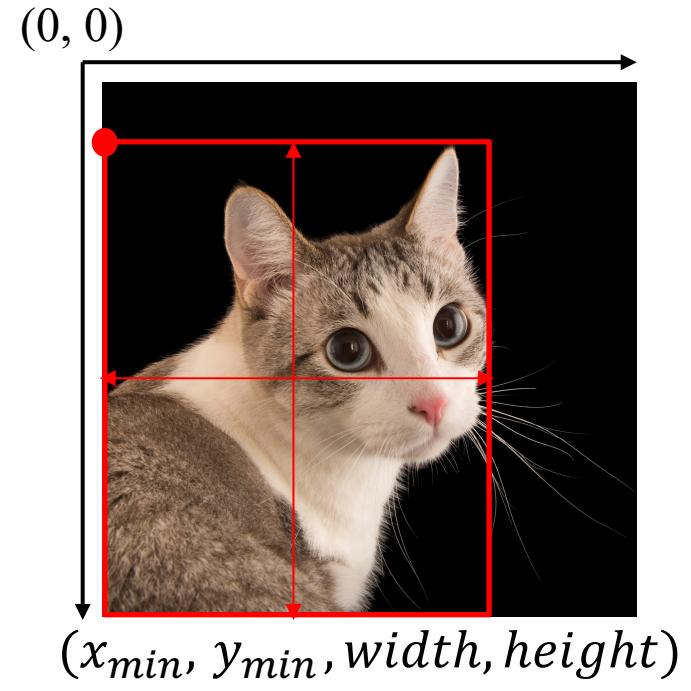
- ❖ Discard bounding box with high overlap (keep the bounding box with highest confidence)
- ❖ $\text{Pr}(\text{Class}_i) * \text{IOU}_{\text{pred}}^{\text{truth}}$



3 – VOC Dataset



Object Detection Dataset Format

VOC**YOLO****COCO**

3 – VOC Dataset

!

Pascal VOC Detection Dataset

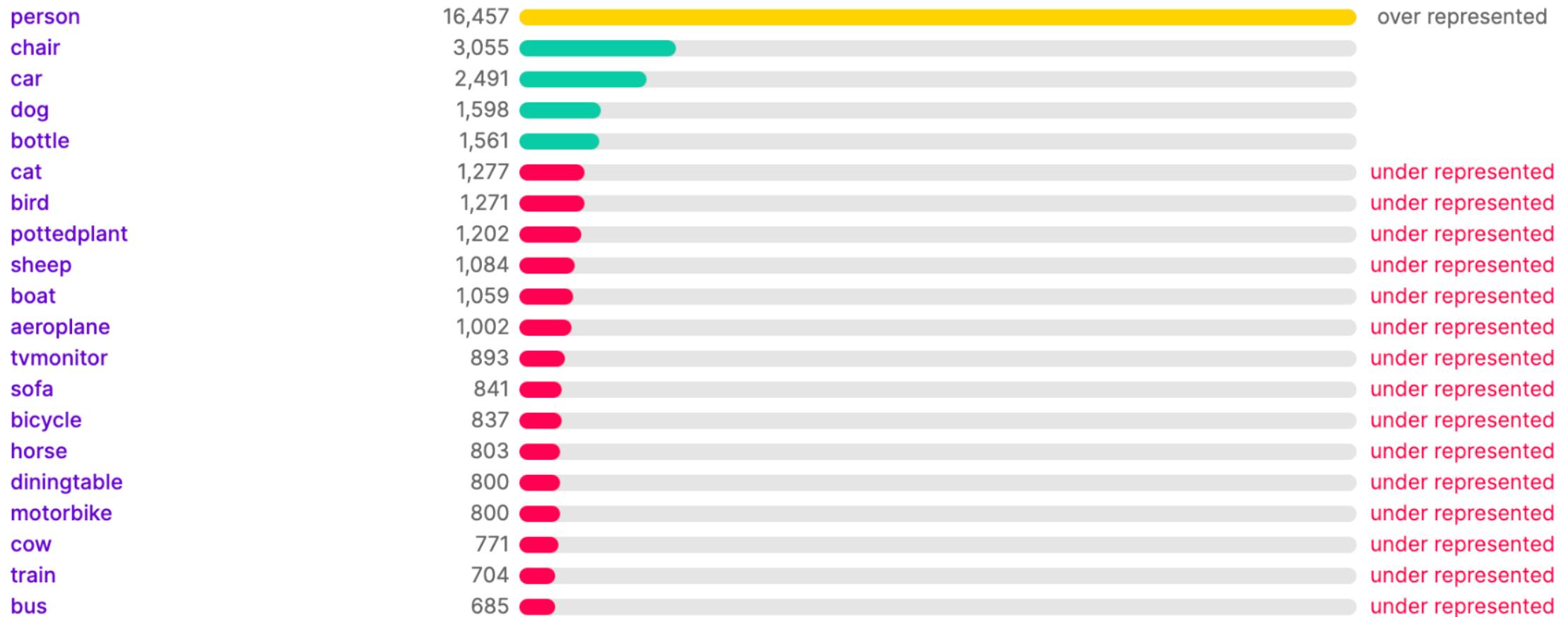


- **Person:** person
- **Animal:** bird, cat, cow, dog, horse, sheep
- **Vehicle:** aeroplane, bicycle, boat, bus, car, motorbike, train
- **Indoor:** bottle, chair, dining table, potted plant, sofa, tv/monitor

3 – VOC Dataset



Pascal VOC Detection Dataset



3 – VOC Dataset



Pascal VOC Detection Dataset

VOCDETECTION

```
CLASS torchvision.datasets.VOCDetection(root: str, year: str = '2012', image_set: str =
    'train', download: bool = False, transform: Optional[Callable] = None,
    target_transform: Optional[Callable] = None, transforms:
    Optional[Callable] = None) [SOURCE]
```

Pascal VOC Detection Dataset.

Parameters:

- **root** (string) – Root directory of the VOC Dataset.
- **year** (string, optional) – The dataset year, supports years "2007" to "2012".
- **image_set** (string, optional) – Select the image_set to use, "train", "trainval" or "val". If `year=="2007"`, can also be "test".
- **download** (bool, optional) – If true, downloads the dataset from the internet and puts it in root directory. If dataset is already downloaded, it is not downloaded again. (default: alphabetic indexing of VOC's 20 classes).
- **transform** (callable, optional) – A function/transform that takes in an PIL image and returns a transformed version. E.g. `transforms.RandomCrop`
- **target_transform** (callable, required) – A function/transform that takes in the target and transforms it.
- **transforms** (callable, optional) – A function/transform that takes input sample and its target as entry and returns a transformed version.

3 – VOC Dataset



Pascal VOC Detection Dataset

Image with 1 bounding box

```
<annotation>
  <folder>VOC2012</folder>
  <filename>2008_000616.jpg</filename>
  <source>
    <database>The VOC2008 Database</database>
    <annotation>PASCAL VOC2008</annotation>
    <image>flickr</image>
  </source>
  <size>
    <width>231</width>
    <height>256</height>
    <depth>3</depth>
  </size>
  <segmented>0</segmented>
  <object>
    <name>person</name>
    <bndbox>
      <xmin>51</xmin>
      <ymin>62</ymin>
      <xmax>209</xmax>
      <ymax>256</ymax>
    </bndbox>
  </object>
</annotation>
```

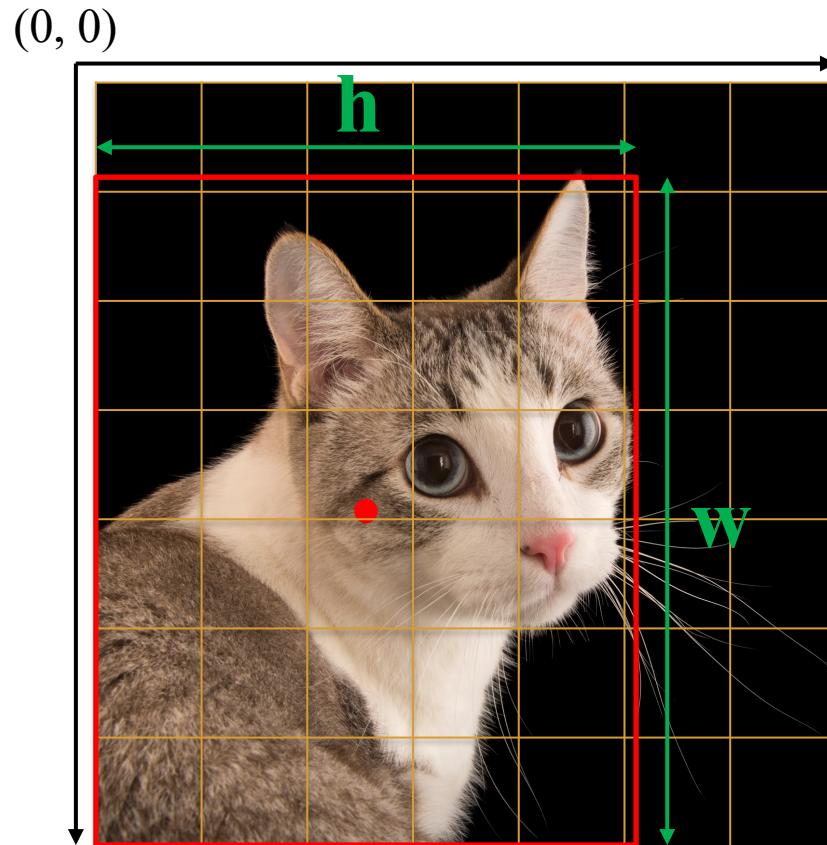
Image with 2 bounding boxes

```
<annotation>
  <folder>VOC2012</folder>
  <filename>2008_000200.jpg</filename>
  <source>
    <database>The VOC2008 Database</database>
    <annotation>PASCAL VOC2008</annotation>
    <image>flickr</image>
  </source>
  <size>
    <width>500</width>
    <height>375</height>
    <depth>3</depth>
  </size>
  <segmented>0</segmented>
  <object>
    <name>person</name>
    <bndbox>
      <xmin>119</xmin>
      <ymin>76</ymin>
      <xmax>184</xmax>
      <ymax>311</ymax>
    </bndbox>
  </object>
  <object>
    <name>person</name>
    <bndbox>
      <xmin>266</xmin>
      <ymin>43</ymin>
      <xmax>338</xmax>
      <ymax>323</ymax>
    </bndbox>
  </object>
</annotation>
```

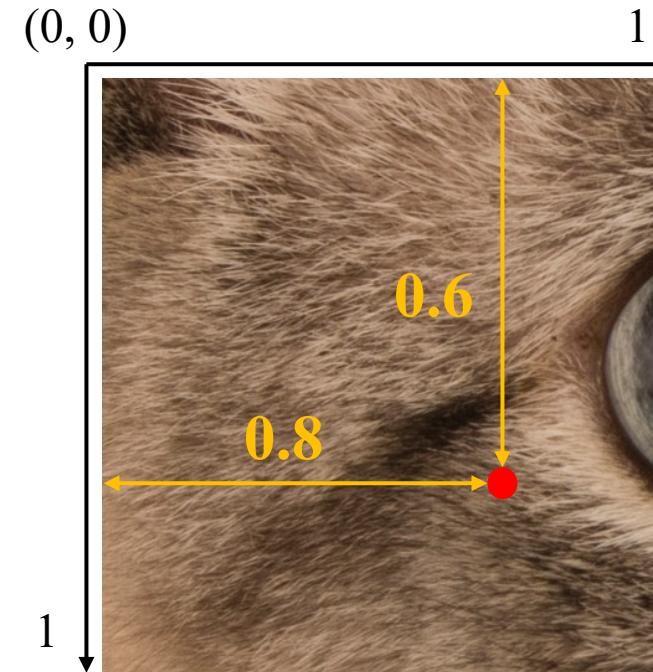
3 – VOC Dataset



Rescale Bounding Box



rescale

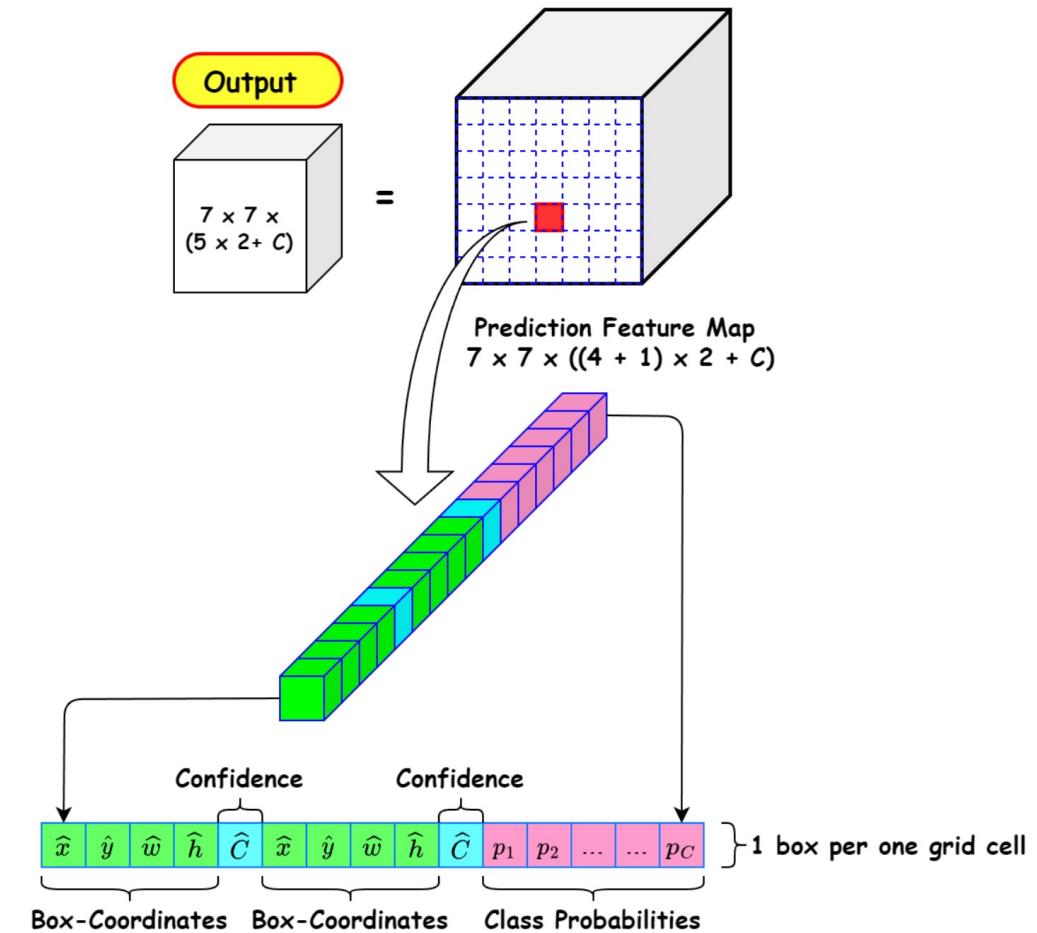
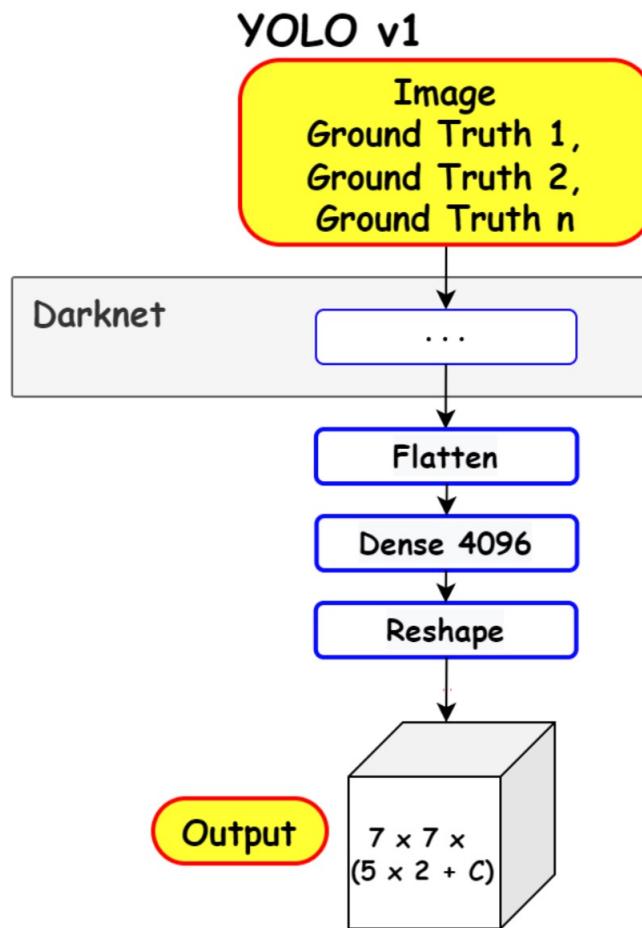


x	y	w	h
0.8	0.6	5.1	6.2

3 – VOC Dataset



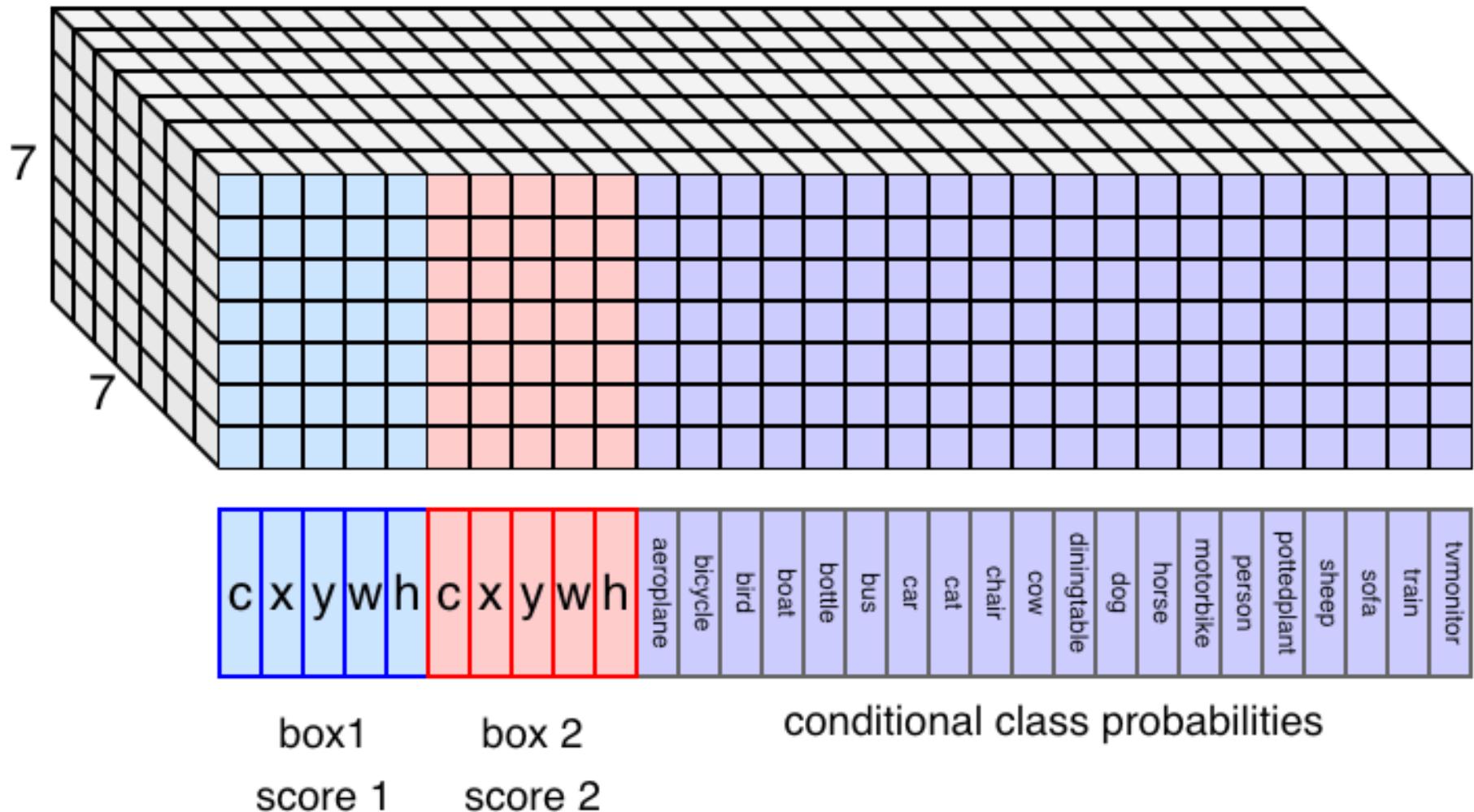
Output Format



3 – VOC Dataset

!

Output Shape

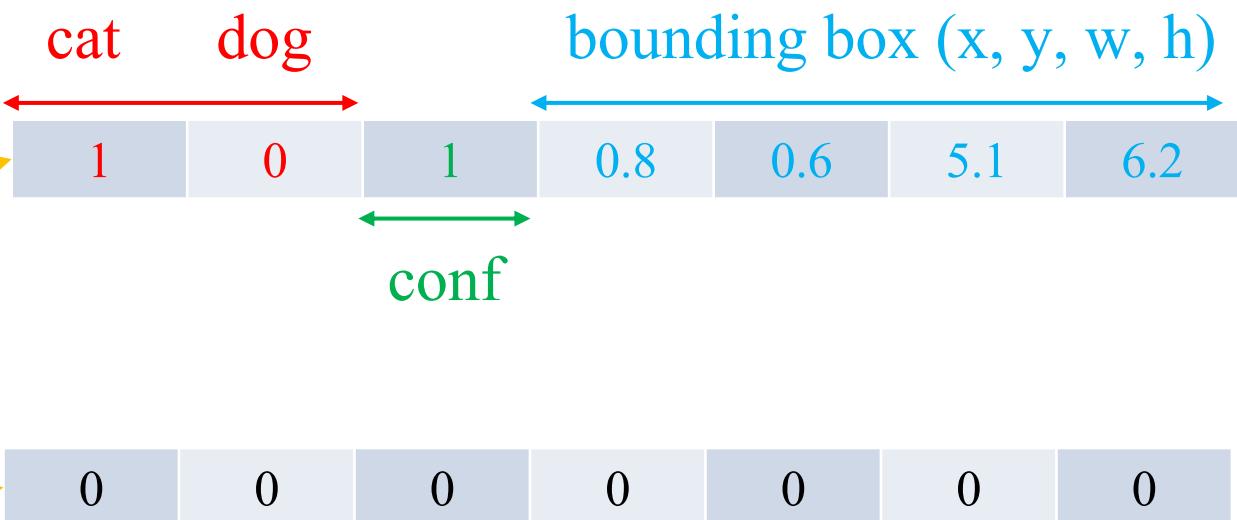
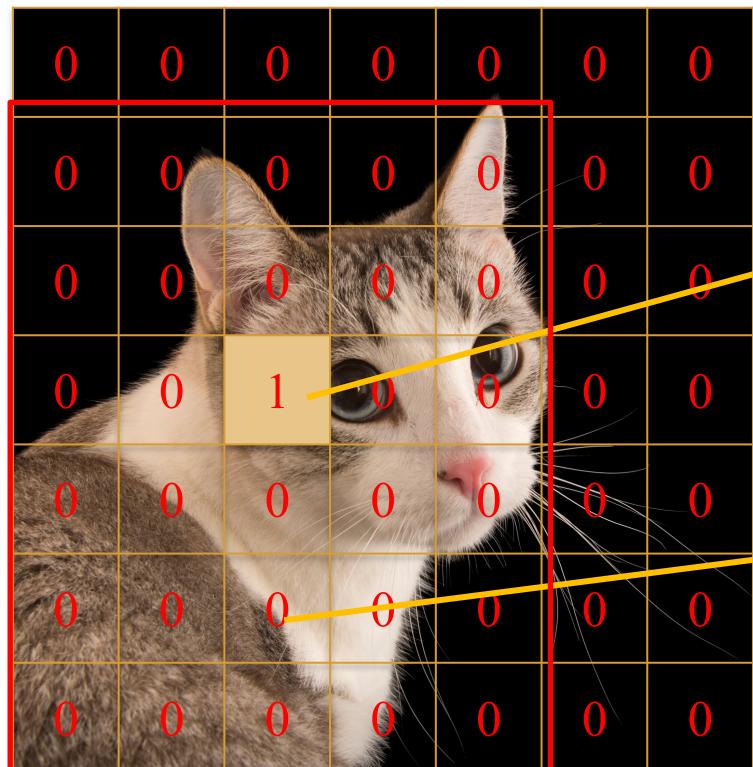


3 – VOC Dataset



Target for Training

For example, we have 2 classes: *cat* and *dog*.

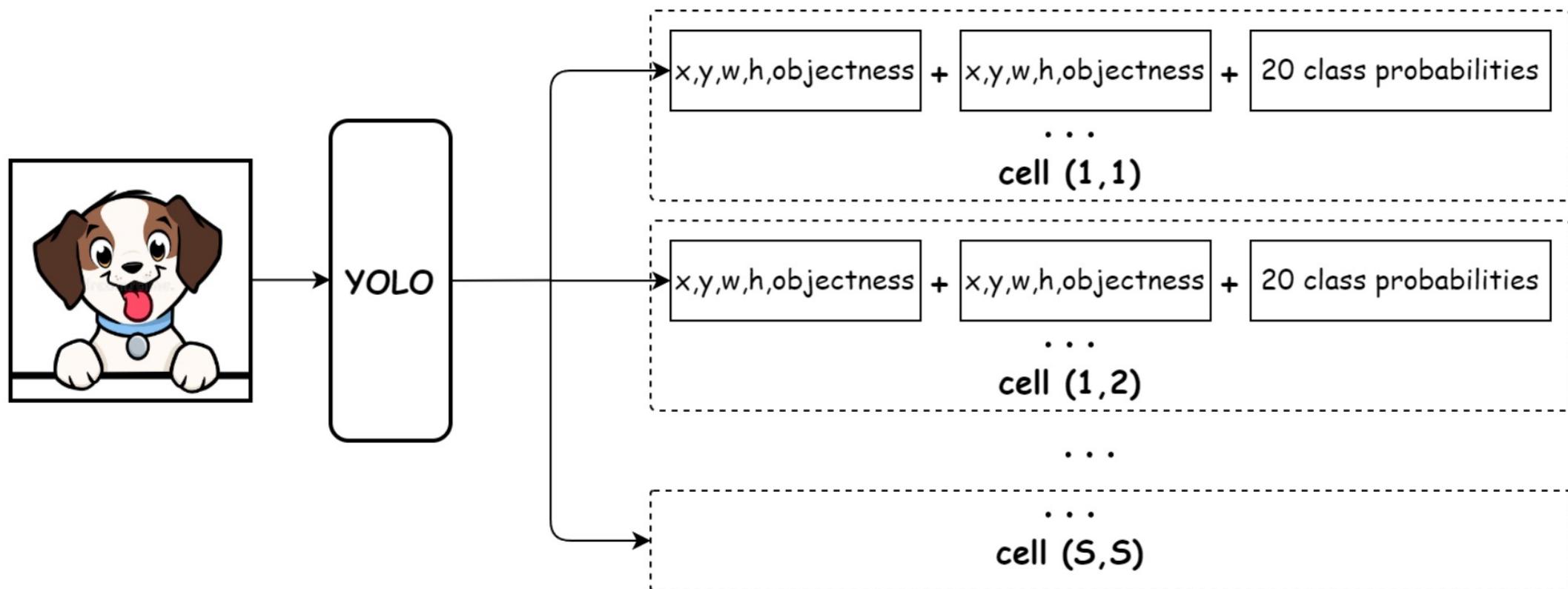


For each object, **only 1 grid cell** is responsible for the prediction.

3 – VOC Dataset



Output Example



3 – VOC Dataset

!

Ground Truth Example

	cat	dog	conf	bounding box (x, y, w, h)			
cell 1	1	0	1	0.8	0.6	5.1	6.2
cell 2	0	0	0	0	0	0	0

3 – VOC Dataset



Localization loss

The localization loss measures the errors in the predicted boundary box locations and sizes. We only count the box responsible for detecting the object.

$$\begin{aligned}\mathcal{L}_{\text{loc}} = & \lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^B 1_{ij}^{\text{obj}} [(x_i - \hat{x}_i)^2 + (y_i - \hat{y}_i)^2 +] \\ & + \lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^B 1_{ij}^{\text{obj}} [(\sqrt{w_i} - \sqrt{\hat{w}_i})^2 + (\sqrt{h_i} - \sqrt{\hat{h}_i})^2]\end{aligned}$$

where

$1_{ij}^{\text{obj}} = 1$ if the j th boundary box in cell i is responsible for detecting the object, otherwise 0.

λ_{coord} increase the weight for the loss in the boundary box coordinates.

$\lambda_{\text{coord}} = 5$

	cat	dog	conf	bounding box (x, y, w, h)			
cell 1	1	0	1	0.8	0.6	5.1	6.2
cell 2	0	0	0	0	0	0	0
cell 1	0.8	0.2	0.8	0.7	0.65	3.1	7.2
cell 2	0.45	0.55	0.2	0.35	0.2	0.5	1.2

$$\begin{aligned}\mathcal{L}_{\text{loc}} = & \lambda_{\text{coord}} [(0.8 - 0.7)^2 + (0.6 - 0.65)^2 \\ & + (\sqrt{5.1} - \sqrt{3.1})^2 + (\sqrt{6.2} - \sqrt{7.2})^2]\end{aligned}$$

3 – VOC Dataset

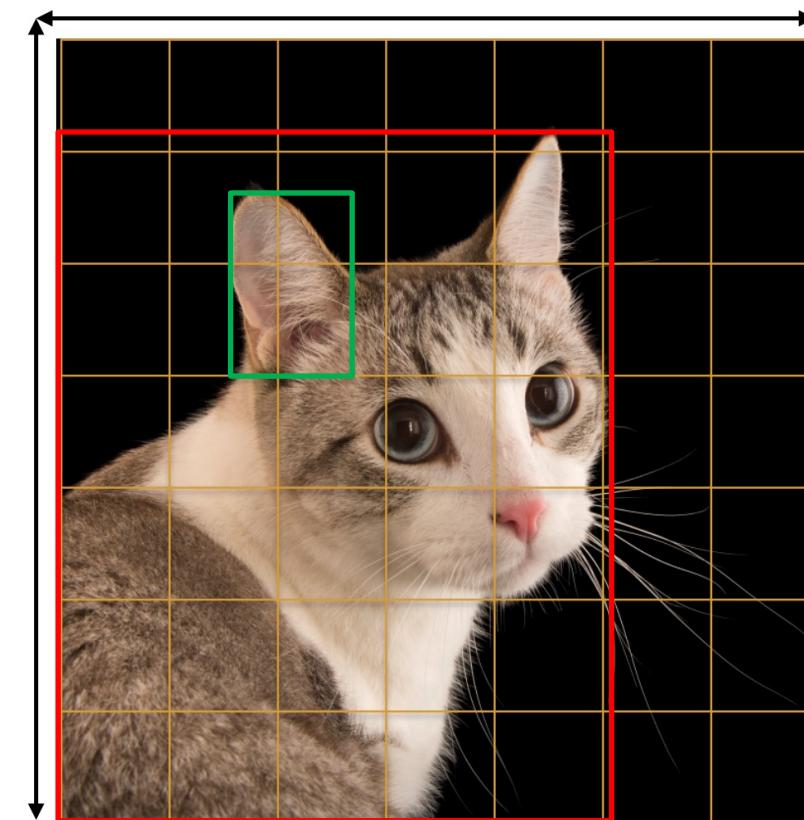


Why we use Square Root for width and height in localization loss

$$\begin{aligned}\mathcal{L}_{\text{loc}} = & \lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^B 1_{ij}^{\text{obj}} [(x_i - \hat{x}_i)^2 + (y_i - \hat{y}_i)^2 +] \\ & + \lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^B 1_{ij}^{\text{obj}} [(\sqrt{w_i} - \sqrt{\hat{w}_i})^2 + (\sqrt{h_i} - \sqrt{\hat{h}_i})^2]\end{aligned}$$

width_gt	width_pred	Error normal	Square Root Error
5.5	5.0	0.5	0.11
1.0	0.5	0.5	0.29

The difference between a small box is much more significant than a large box.



3 – VOC Dataset



Confidence loss

If an object is detected in the box, the confidence loss (measuring the objectness of the box) is:

$$\mathcal{L}_{\text{obj}} = \sum_{i=0}^{S^2} \sum_{j=0}^B 1_{ij}^{\text{obj}} (C_{ij} - \hat{C}_{ij})^2$$

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

where

\hat{C}_i is the box confidence score of the box j in cell i .

$1_{ij}^{\text{obj}} = 1$ if the j th boundary box in cell i is responsible for detecting the object, otherwise 0.

$$\lambda_{\text{noobj}} = 0.5$$

	cat	dog	conf		bounding box		
cell 1	1	0	1	0.8	0.6	5.1	6.2
cell 2	0	0	0	0	0	0	0
cell 1	0.8	0.2	0.8	0.7	0.65	3.1	7.2
cell 2	0.45	0.55	0.2	0.35	0.2	0.5	1.2

$$\mathcal{L}_{\text{obj}} = (1 - 0.8)^2 + \lambda_{\text{noobj}}(0.0 - 0.2)^2$$

3 – VOC Dataset



Classification loss

If an object is detected, the classification loss at each cell is the squared error of the class conditional probabilities for each class:

$$\mathcal{L}_{cls} = \sum_{i=0}^{S^2} 1_i^{obj} \sum_{c \in classes} (p_i(c) - \hat{p}_i(c))^2$$

where

$1_i^{obj} = 1$ if an object appears in cell i , otherwise 0.

$\hat{p}_i(c)$ denotes the conditional class probability for class c in cell i .

	cat	dog	conf	bounding box			
cell 1	1	0	1	0.8	0.6	5.1	6.2
cell 2	0	0	0	0	0	0	0
cell 1	0.8 0.2		0.8	0.7	0.65	3.1	7.2
cell 2	0.45	0.55	0.2	0.35	0.2	0.5	1.2

$$\mathcal{L}_{cls} = (1 - 0.8)^2 + (0 - 0.2)^2$$

3 – VOC Dataset



Total Loss Function

$$\mathcal{L} = \mathcal{L}_{\text{loc}} + \mathcal{L}_{\text{obj}} + \mathcal{L}_{\text{cls}}$$

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

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

$$\mathcal{L}_{\text{obj}} = \sum_{i=0}^{S^2} \sum_{j=0}^B 1_{ij}^{\text{obj}} (C_{ij} - \hat{C}_{ij})^2$$

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

$$\mathcal{L}_{\text{cls}} = \sum_{i=0}^{S^2} 1_i^{\text{obj}} \sum_{c \in \text{classes}} (p_i(c) - \hat{p}_i(c))^2$$

1. penalize bad localization of center of cells
2. penalize the bounding box with inaccurate height and width. The square root is present so that errors in small bounding boxes are more penalizing than errors in big bounding boxes.
3. make confidence score equal to the IoU between the object and the prediction when there is one object
4. make confidence score close to 0 when there are no object in the cell
5. a simple classification loss (not explained in the article)

Thanks!

Any questions?