# Semantic Segmentation
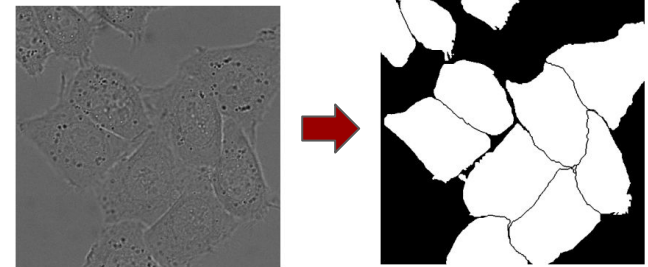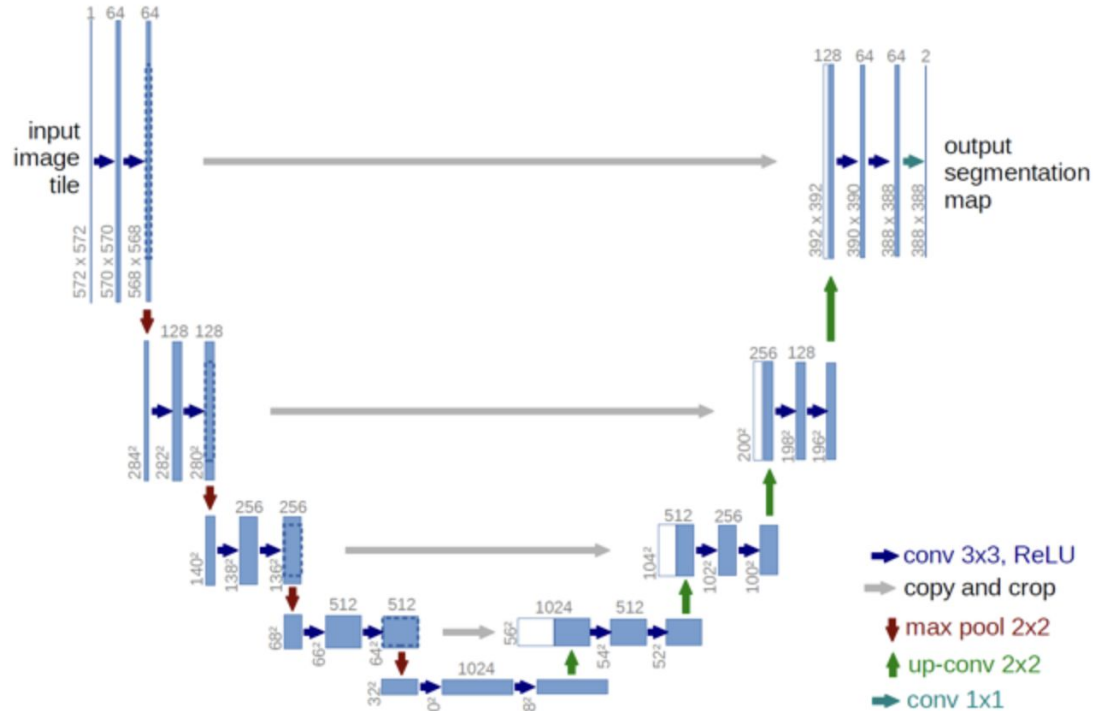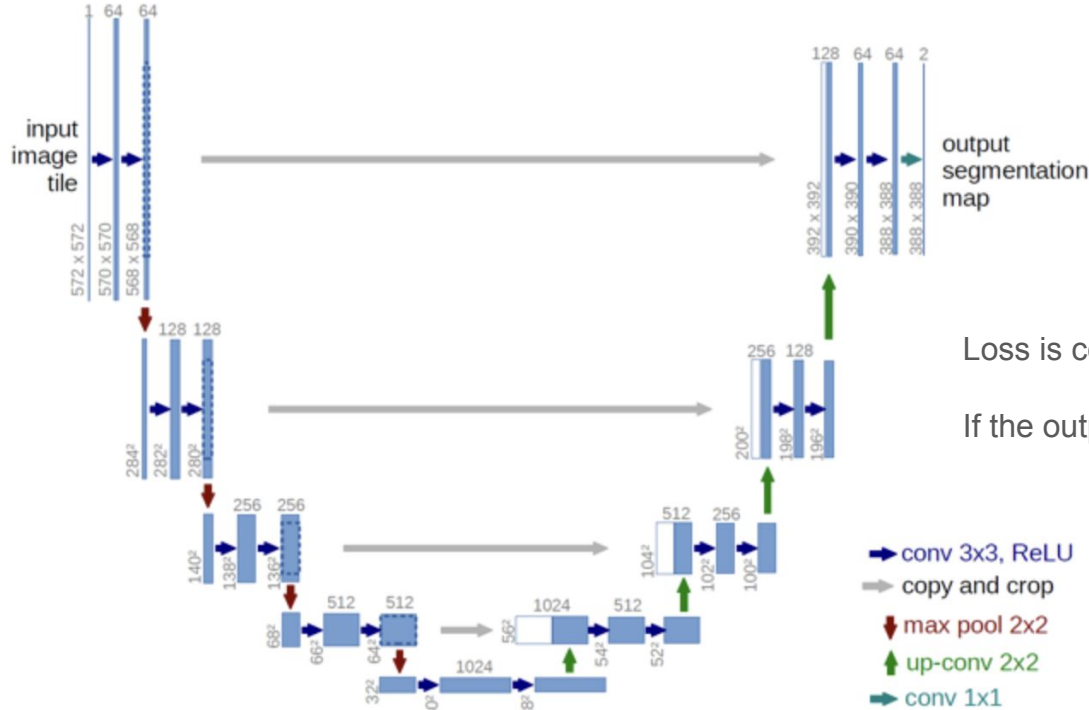
# Consider the binary Image Segmentation Problem
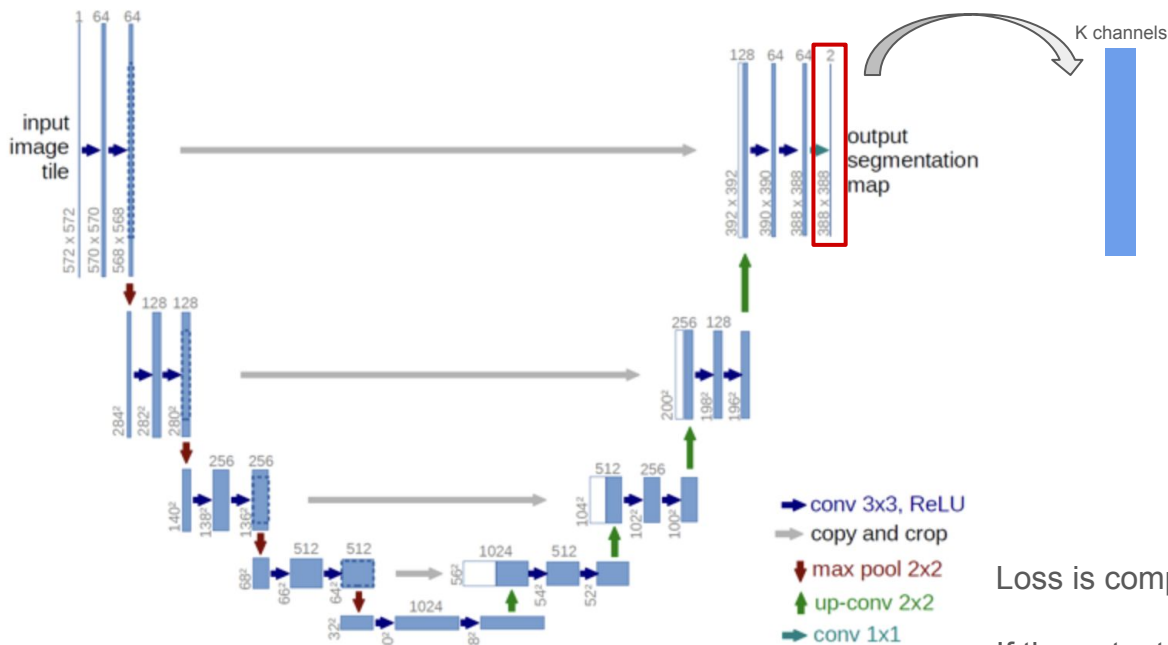


U-Net

# Consider the binary Image Segmentation Problem



$$\sum_{i,j} \mathcal{L}_{CE}(y_{(i,j)}, \hat{y}_{(i,j)})$$

Loss is computed over each pixel

If the output tensor after softmax is y', each y'[i,j, :] is logits

K channels

$$\sum_{i,j} \mathcal{L}_{CE}(y_{(i,j)}, \hat{y}_{(i,j)})$$

→ conv 3x3, ReLU
→ copy and crop
↓ max pool 2x2
↑ up-conv 2x2
→ conv 1x1

Loss is computed over each pixel

If the output tensor after softmax is y', each y'[i,j, :] is logits

Input

Prediction

# IOU score (½ Dice score)

$$\frac{1}{K} \sum_{k=1...K} \frac{\sum_{i,j} \mathbb{I}_{y_{(i,j)}=\hat{y}_{(i,j)}}}{\sum_{i,j} \mathbb{I}_{y_{(i,j)}=k} + \mathbb{I}_{\hat{y}_{(i,j)}=k} - \mathbb{I}_{y_{(i,j)}=\hat{y}_{(i,j)}}}$$
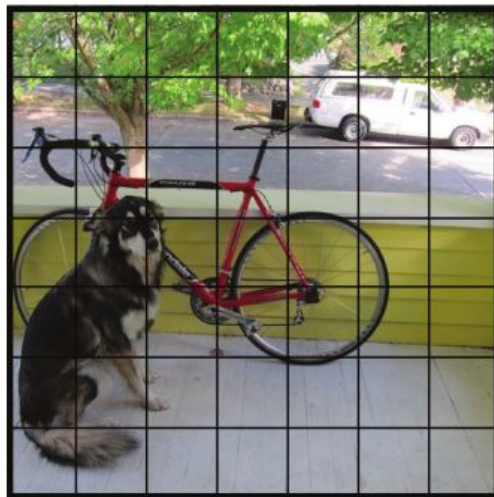
# You Only Look Once:
# Unified, Real-Time Object Detection

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

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

http://pjreddie.com/yolo/

# Philosophy behind YOLO

Split the entire image into S X S grids, S = 7  will be considered



S × S grid on input

# Philosophy behind YOLO

For each grid cell, predict the centroid coordinates of the bounding box assuming an object is present in its grid cell

# Philosophy behind YOLO

For each grid cell, make B = 2 prediction

# Philosophy behind YOLO

For each grid cell, predict the probability that there is an object present,
We will call it P(Object)

# Philosophy behind YOLO

For each grid cell, assuming an object is present, what is the classification among K classes. We will call it P( . | Object)

# Philosophy behind YOLO

For each grid cell, assuming an object is present, what is the classification among K classes. We will call it P( . | Object)

# YOLO V1

## Architecture



The final output of our network is the $7 \times 7 \times 30$ tensor of predictions.

# YOLO V1

## Architecture



The output tensor y' = 7 X 7 X 30 in shape (h,w,d)

Each grid cell y'[i, j, :] is a 30 dimensional vector representing the predicted:

1. Height, Width, and centroids
2. Objectness score for 2 candidate bounding boxes
3. 20 dimensional classification vector
   (PASCAL VOC 2007 dataset has 20 classes)

The final output of our network is the $7 \times 7 \times 30$ tensor of predictions.

# YOLO V1

## Architecture



The output tensor y' = 7 X 7 X 30 in shape (h,w,d)

Each grid cell y'[i, j, :] is a 30 dimensional vector representing the predicted:
1. Height, Width, and centroids
2. Objectness score for 2 candidate bounding boxes
3. 20 dimensional classification vector
   (PASCAL VOC 2007 dataset has 20 classes)

The final output of our network is the $7 \times 7 \times 30$ tensor of predictions.

# The loss function (YOLO-v1)

Given that a prediction exists in a cell, we include only that prediction in the loss function which has the higher IOU among the proposals.

If the cell has no BBox in the ground truth, we ignore our regression errors for the Bounding box location and size, and the classification loss for class prediction

We include the Objectness score(does this cell have an object to be predicted) is computed for all the cells

# Loss Function of YOLO v1



Bounding boxes + confidence

S × S grid on input

Class probability map

Final detections

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

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

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

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

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

# Loss Function of YOLO v1



S × S grid on input

Bounding boxes + confidence

Class probability map

Final detections

P(Object)  X  Y  Width  Height  P(Object)  X  Y  Width  Height  P(Cat | Object)  P(Bird | Object)  P(TV | Object)

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

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

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

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

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

# Loss Function of YOLO v1



S × S grid on input

Bounding boxes + confidence

Class probability map

Final detections

7

7

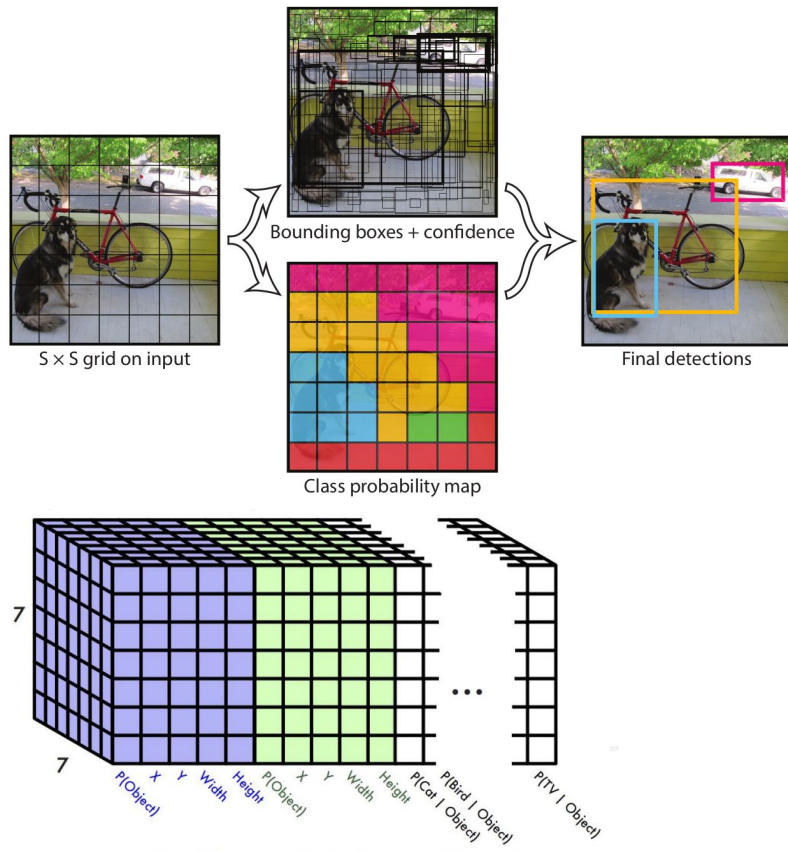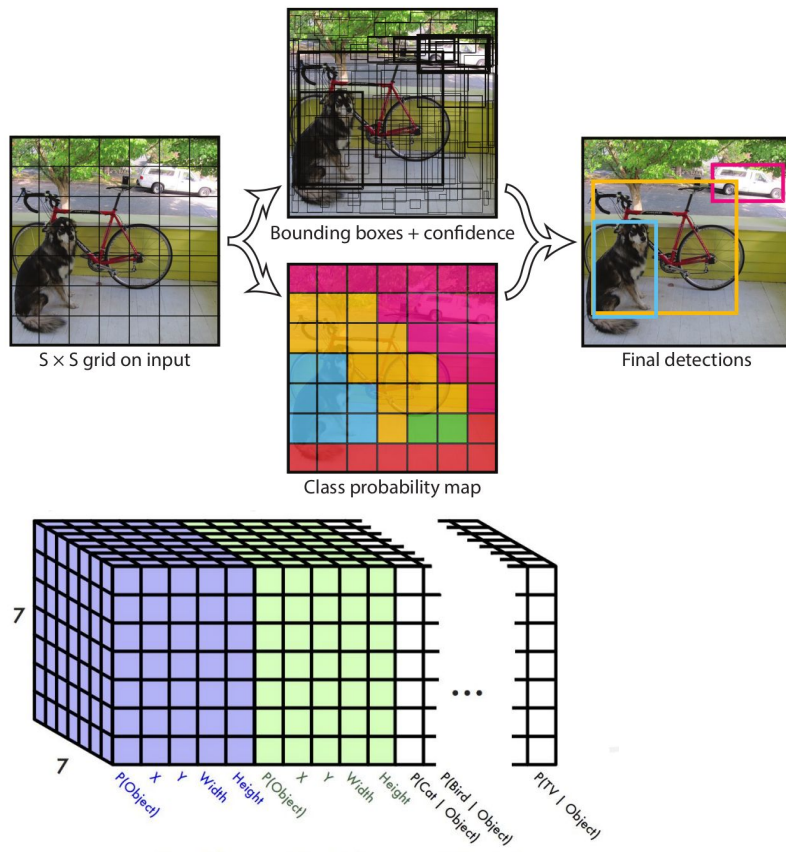P(Object)  X  Y  Width  Height  P(Object)  X  Y  Width  Height  P(Cat | Object)  P(Bird | Object)  P(TV | Object)

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

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

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

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

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

# Loss Function of YOLO v1



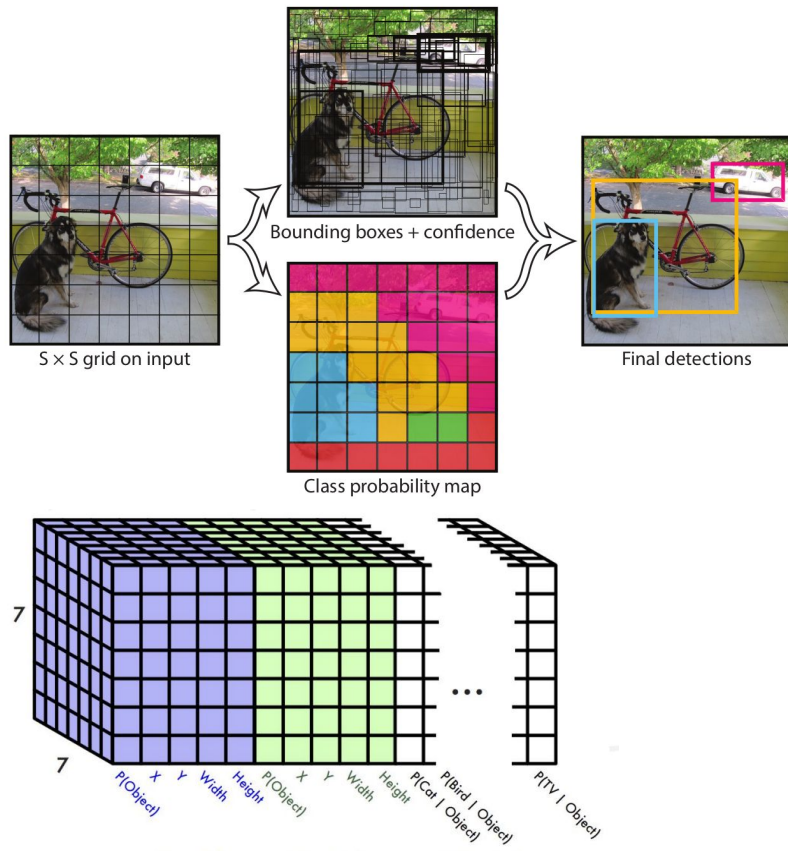S × S grid on input

Bounding boxes + confidence

Class probability map

Final detections

P(Object)  X  Y  Width  Height  P(Object)  X  Y  Width  Height  P(Bird | Object)  P(Cat | Object)  P(TV | Object)

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

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

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

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

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

# Loss Function of YOLO v1



Bounding boxes + confidence
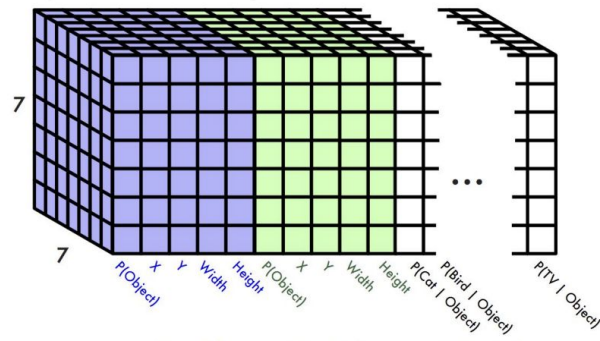
S × S grid on input

Class probability map

Final detections

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

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

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

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

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

# Loss Function of YOLO v1



S × S grid on input
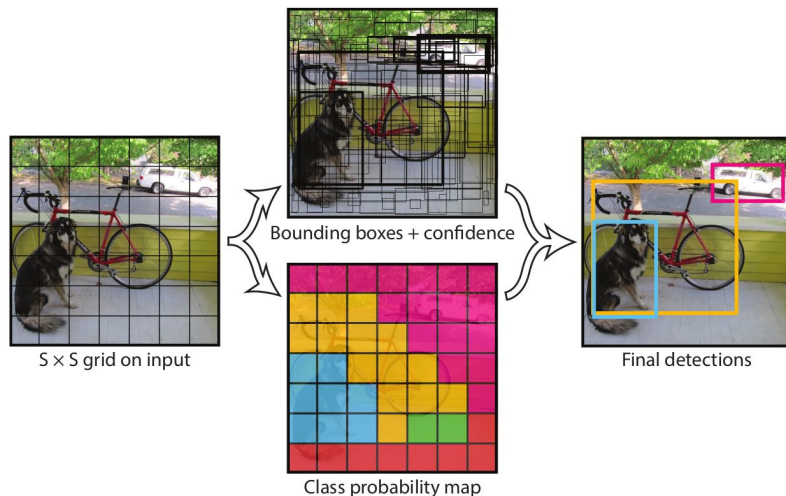
Bounding boxes + confidence

Class probability map

Final detections

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

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

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

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

$$+ \sum_{i=0}^{S^2} \mathbb{1}_{i}^{\mathrm{obj}} \sum_{c \in \mathrm{classes}} \left( p_i(c) - \hat{p}_i(c) \right)^2$$
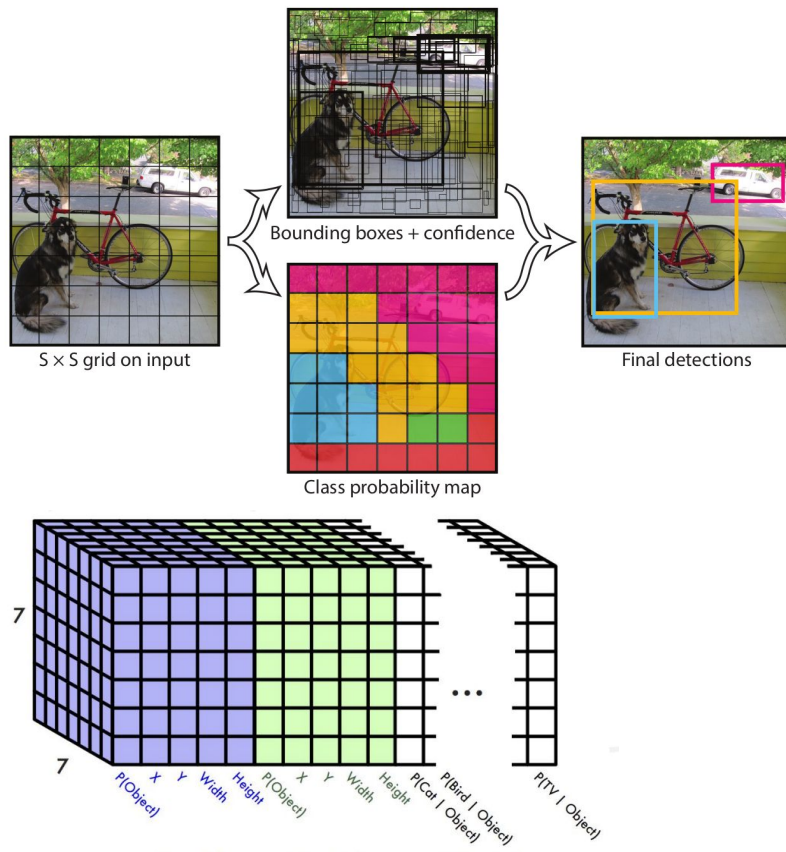
# Loss Function of YOLO v1



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

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

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

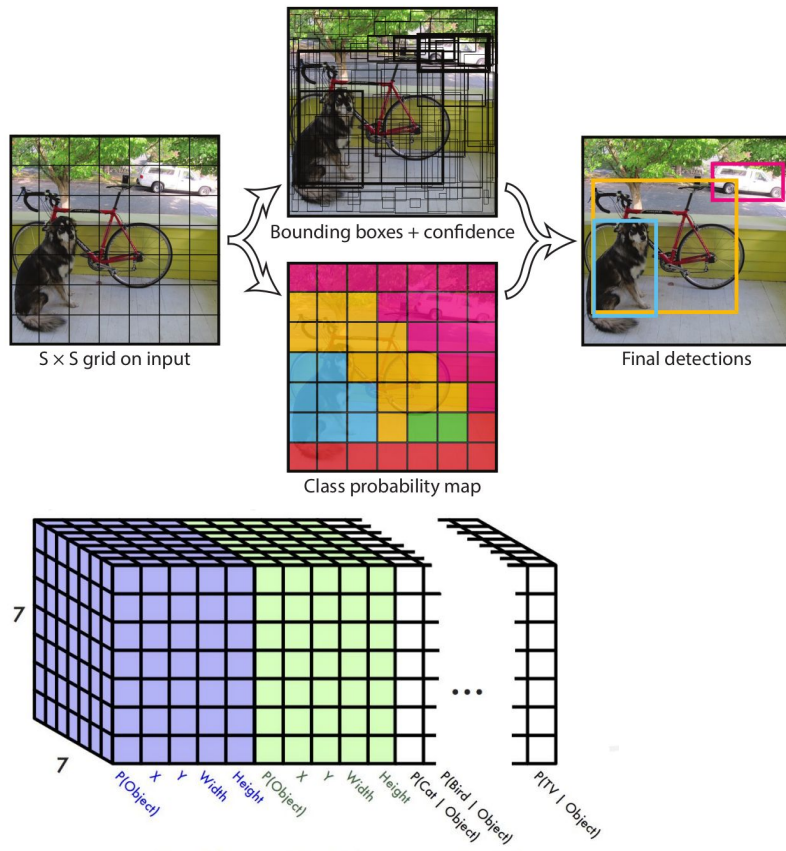$$+ \lambda_{\text{noobj}} \sum_{i=0}^{S^2} \sum_{j=0}^{B} \mathbb{1}_{ij}^{\text{noobj}} \left( C_i - \hat{C}_i \right)^2$$

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

# Loss Function of YOLO v1



S × S grid on input

Bounding boxes + confidence

Class probability map

Final detections

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

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

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

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

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

# End of Training the model

# Let's look at the inference !

# Non Maximum Suppression

- **This is applied during inference**
- **All the bounding box predictions are sorted according to their objectness scores in descending order**



Raw predictions from YOLO

# Non Maximum Suppression

- This is applied during inference
- All the bounding box predictions are sorted according to their objectness scores in descending order
- **All boxes below a certain objectness score are removed**
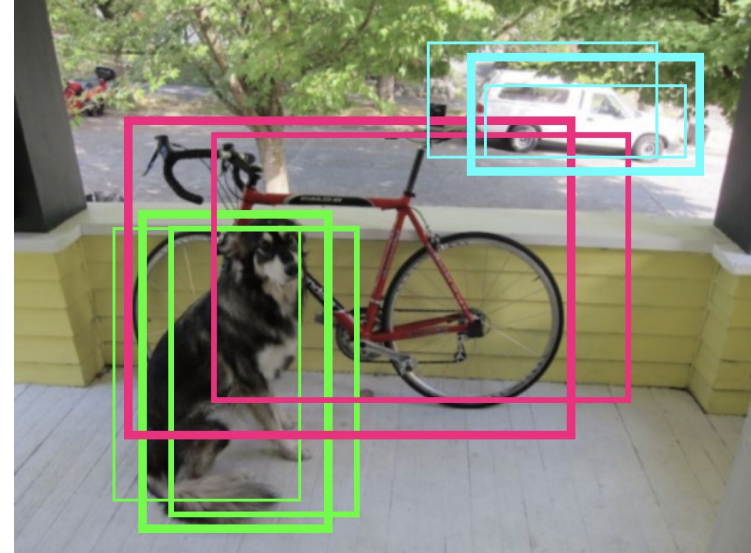


Not so confident predictions removed

# Non Maximum Suppression

- This is applied during inference
- All the bounding box predictions are sorted according to their objectness scores in descending order
- All boxes below a certain objectness score are removed
- **Given a BBox prediction, all BBox with a lower objectness score and belonging to the same class having an IOU above a user specified threshold are removed from the prediction set.**
- **This is applied to all remaining predictions iteratively**



Intra-class overlapping predictions removed

# End of YOLO v-1

# YOLO V2

## Improvements over V1

- 

$$b_x = \sigma(t_x) + c_x$$
$$b_y = \sigma(t_y) + c_y$$
$$b_w = p_w e^{t_w}$$
$$b_h = p_h e^{t_h}$$

$$Pr(\text{object}) * IOU(b, \text{object}) = \sigma(t_o)$$

The network predicts 5 bounding boxes at each cell in the output feature map. The network predicts 5 coordinates for each bounding box, $t_x$, $t_y$, $t_w$, $t_h$, and $t_o$. If the cell is offset from the top left corner of the image by $(c_x, c_y)$ and the bounding box prior has width and height $p_w$, $p_h$, then the predictions correspond to:
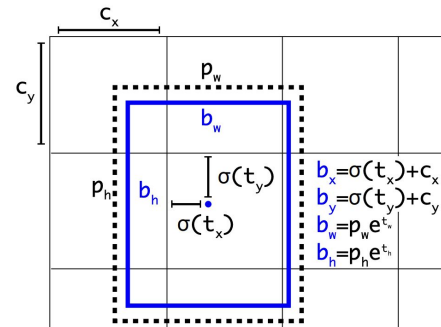


**Figure 3: Bounding boxes with dimension priors and location prediction.** We predict the width and height of the box as offsets from cluster centroids. We predict the center coordinates of the box relative to the location of filter application using a sigmoid function.
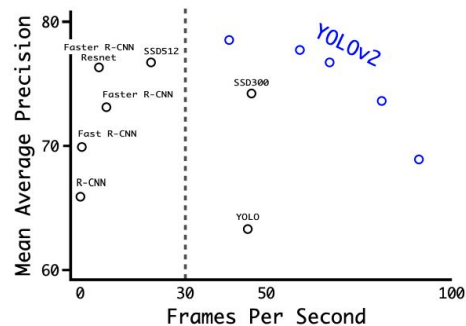


**Figure 4: Accuracy and speed on VOC 2007.**

# YOLO V2

- The model has a new backbone darknet 19
- The model is purely convolutional even till the last layers
- Training is done at multiple scales of input all at multiples of 32 ranging from 320X320 to 608X608
- The biggest gain is achieved by direct location prediction (discussed in previous slide)
- Pre training was done at a higher resolution of 448X448

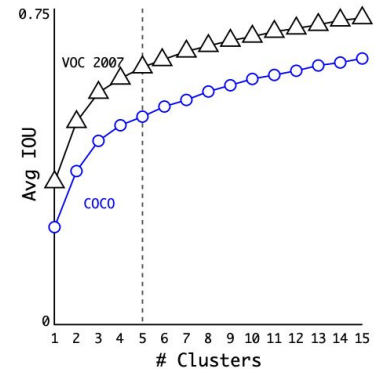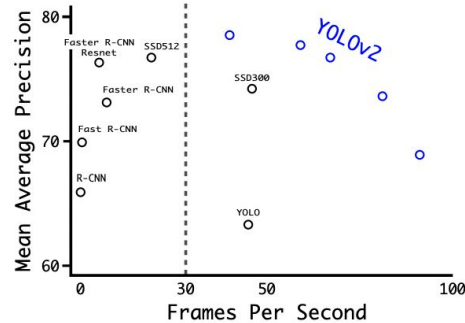|  | YOLO |  |  |  |  |  |  |  | YOLOv2 |
|---|---|---|---|---|---|---|---|---|---|
| batch norm? |  | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| hi-res classifier? |  |  | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| convolutional? |  |  |  | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| anchor boxes? |  |  |  | ✓ | ✓ |  |  |  |  |
| new network? |  |  |  |  | ✓ | ✓ | ✓ | ✓ | ✓ |
| 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 | **78.6** |



**Figure 4: Accuracy and speed on VOC 2007.**

# YOLO V3

Improvements over V2

- V3 comes with a new backbone darknet 53 over darknet 19
- The model does multi scale predictions and not just multiscale training.
  - This allows predicting larger objects in initial layers and smaller objects in final layers
  - 3 box predictions are made at each scale
  - Each box prediction has its own classification vector

| | Type | Filters | Size | Output |
|---|---|---|---|---|
| | Convolutional | 32 | 3 × 3 | 256 × 256 |
| | Convolutional | 64 | 3 × 3 / 2 | 128 × 128 |
| 1× | Convolutional | 32 | 1 × 1 | |
| | Convolutional | 64 | 3 × 3 | |
| | Residual | | | 128 × 128 |
| | Convolutional | 128 | 3 × 3 / 2 | 64 × 64 |
| 2× | Convolutional | 64 | 1 × 1 | |
| | Convolutional | 128 | 3 × 3 | |
| | Residual | | | 64 × 64 |
| | Convolutional | 256 | 3 × 3 / 2 | 32 × 32 |
| 8× | Convolutional | 128 | 1 × 1 | |
| | Convolutional | 256 | 3 × 3 | |
| | Residual | | | 32 × 32 |
| | Convolutional | 512 | 3 × 3 / 2 | 16 × 16 |
| 8× | Convolutional | 256 | 1 × 1 | |
| | Convolutional | 512 | 3 × 3 | |
| | Residual | | | 16 × 16 |
| | Convolutional | 1024 | 3 × 3 / 2 | 8 × 8 |
| 4× | Convolutional | 512 | 1 × 1 | |
| | Convolutional | 1024 | 3 × 3 | |
| | Residual | | | 8 × 8 |
| | Avgpool | | Global | |
| | Connected | | 1000 | |
| | Softmax | | | |

Table 1. **Darknet-53.**

# YOLO V3



## Improvements over V2

- **The model does multi scale predictions and not just multiscale training.**
  - This allows predicting larger objects in initial layers and smaller objects in final layers
  - 3 box predictions are made at each scale
  - Each box prediction has its own classification vector



| Method | mAP | time |
|---|---|---|
| [B] SSD321 | 28.0 | 61 |
| [C] DSSD321 | 28.0 | 85 |
| [D] R-FCN | 29.9 | 85 |
| [E] SSD513 | 31.2 | 125 |
| [F] DSSD513 | 33.2 | 156 |
| [G] FPN FRCN | 36.2 | 172 |
| RetinaNet-50-500 | 32.5 | 73 |
| RetinaNet-101-500 | 34.4 | 90 |
| RetinaNet-101-800 | **37.8** | 198 |
| **YOLOv3-320** | 28.2 | **22** |
| **YOLOv3-416** | 31.0 | 29 |
| **YOLOv3-608** | 33.0 | 51 |