



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

Redmon et al.

2015

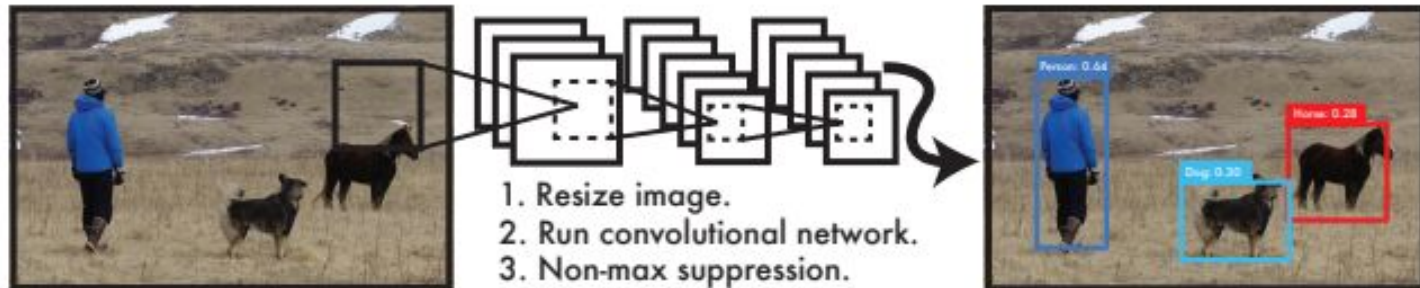


Presented by Eli Lifland, 3/22/2020

# Prior Work

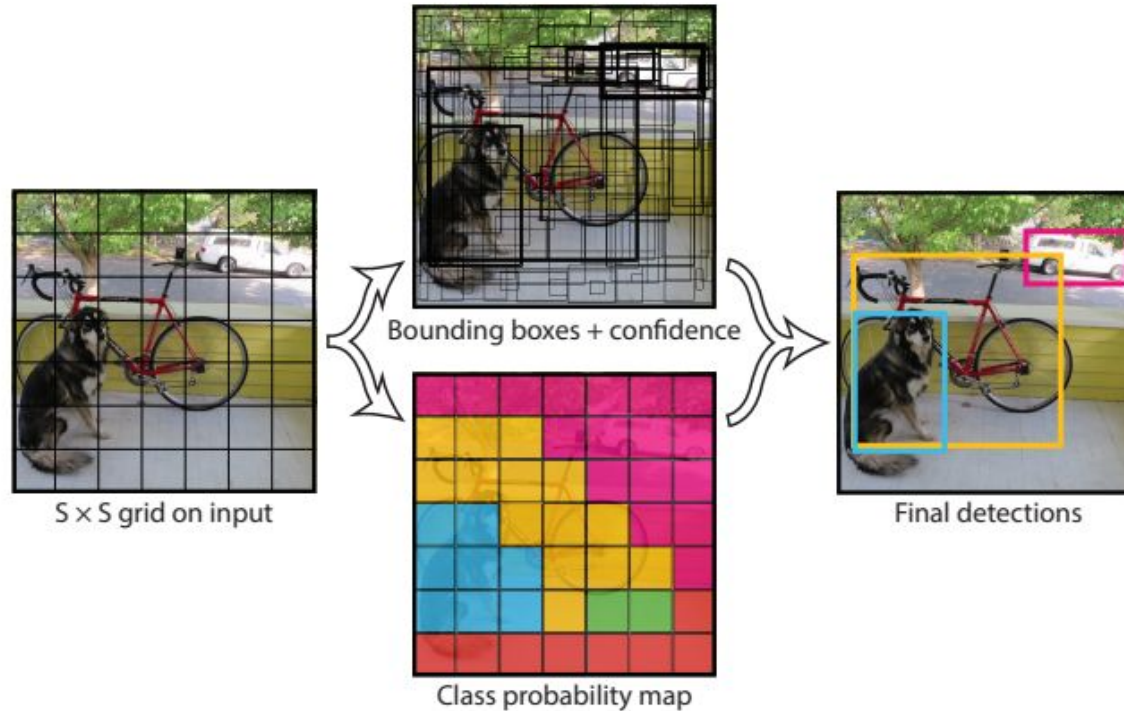
- Prior work repurposed classifiers to perform detection, running classifier on various regions of image
- Deformable Parts Model (DPM)
  - Sliding window approach, pipeline to:
    - Extract features
    - Classify regions
    - Predict bounding boxes
- R-CNN:
  - Region proposal instead of sliding windows
  - Fast/Faster R-CNN use NNs to propose regions

# YOLO Detection System



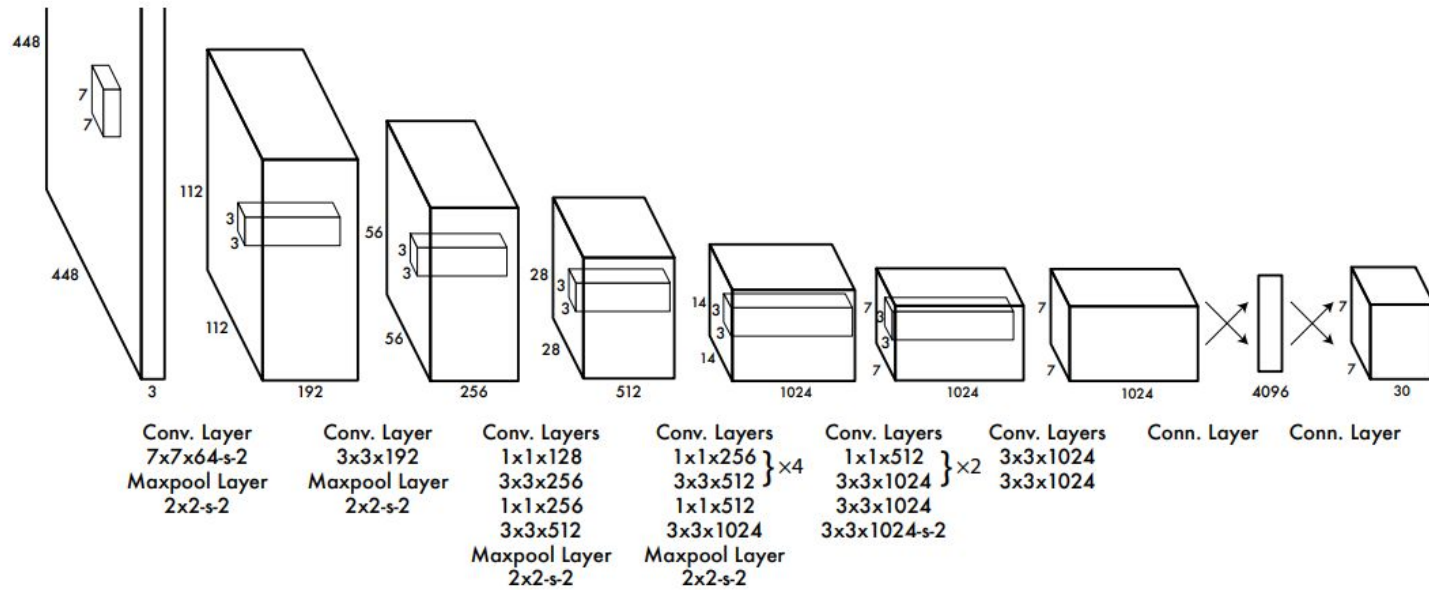
**Figure 1: The YOLO Detection System.** Processing images with YOLO is simple and straightforward. Our system (1) resizes the input image to  $448 \times 448$ , (2) runs a single convolutional network on the image, and (3) thresholds the resulting detections by the model's confidence.

# YOLO Model



**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.

# Network Design



**Figure 3: The Architecture.** Our detection network has 24 convolutional layers followed by 2 fully connected layers. Alternating  $1 \times 1$  convolutional layers reduce the features space from preceding layers. We pretrain the convolutional layers on the ImageNet classification task at half the resolution ( $224 \times 224$  input image) and then double the resolution for detection.

# Training

- Pretrain first 20 layers on ImageNet, then convert to detection by adding last 4 conv, 2 FC layers
- Use sum-squared error because easy to optimize
- To avoid model instability due to gradient from cells w/o objects, weight loss from bounding box predictions higher and weight confidence predictions for boxes without objects lower
- Change in width/height of bounding box matters more for smaller objects than larger
  - Take square root to reflect this

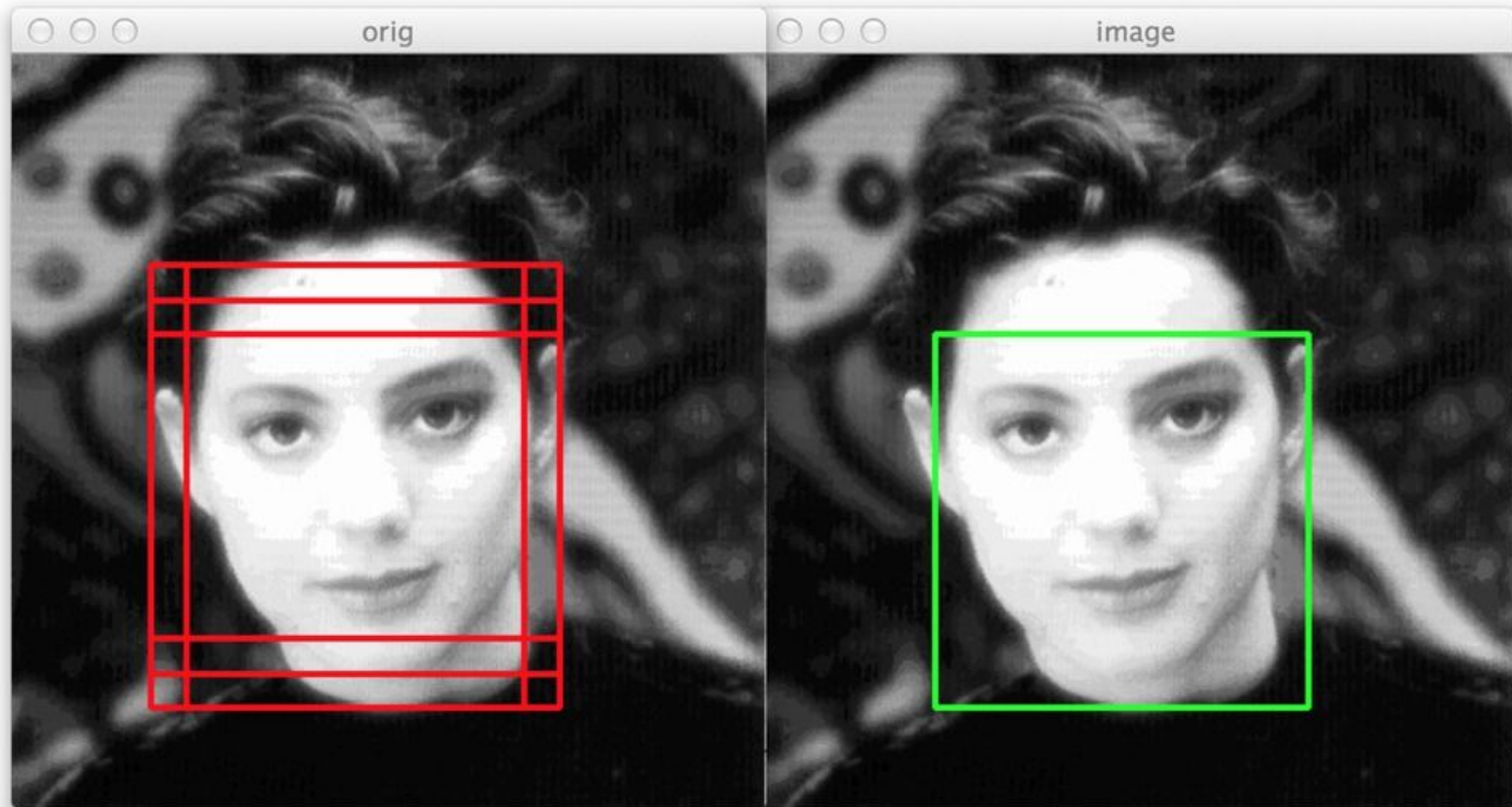
# Training Loss

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



# Non-Maximal Suppression

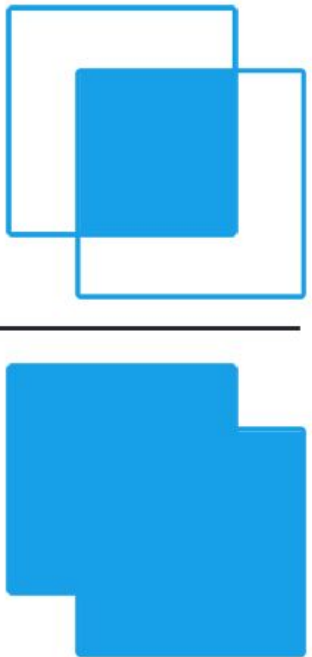
- Use during inference to avoid overlapping predictions





# Limitations

- Each grid cell limited to  $B=2$  boxes and 1 class
  - Struggles with many small objects such as flocks of birds
- Loss function treats errors in small, large boxes the same even though errors in small boxes have bigger effect on IOU
  - Localization isn't great

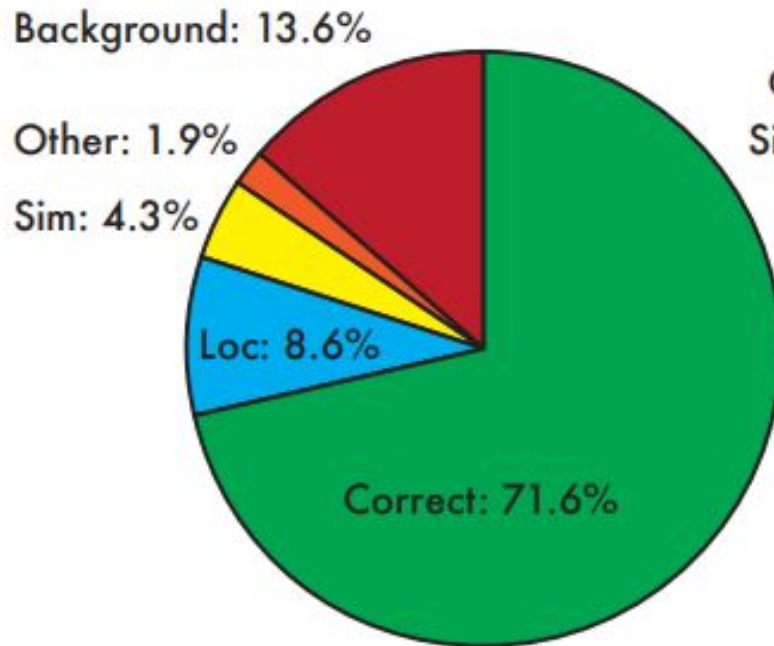
$$\text{IoU} = \frac{\text{Area of Overlap}}{\text{Area of Union}}$$


# Results: PASCAL VOC 2007

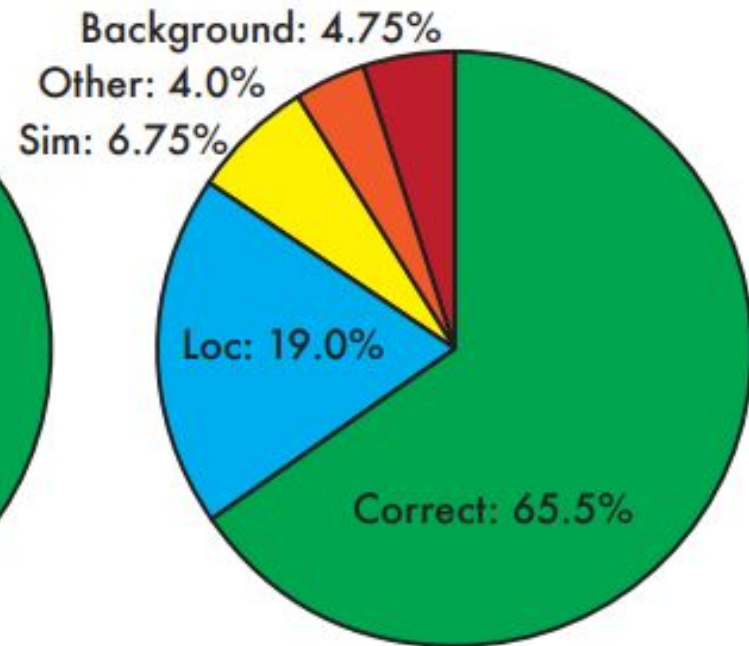
Real-Time Detectors	Train	mAP	FPS
100Hz DPM [31]	2007	16.0	100
30Hz DPM [31]	2007	26.1	30
Fast YOLO	2007+2012	52.7	<b>155</b>
YOLO	2007+2012	<b>63.4</b>	45
<hr/>			
Less Than Real-Time			
Fastest DPM [38]	2007	30.4	15
R-CNN Minus R [20]	2007	53.5	6
Fast R-CNN [14]	2007+2012	70.0	0.5
Faster R-CNN VGG-16[28]	2007+2012	73.2	7
Faster R-CNN ZF [28]	2007+2012	62.1	18
YOLO VGG-16	2007+2012	66.4	21

# Results: Error Analysis

## Fast R-CNN



## YOLO



# Results: Combining Fast R-CNN and YOLO

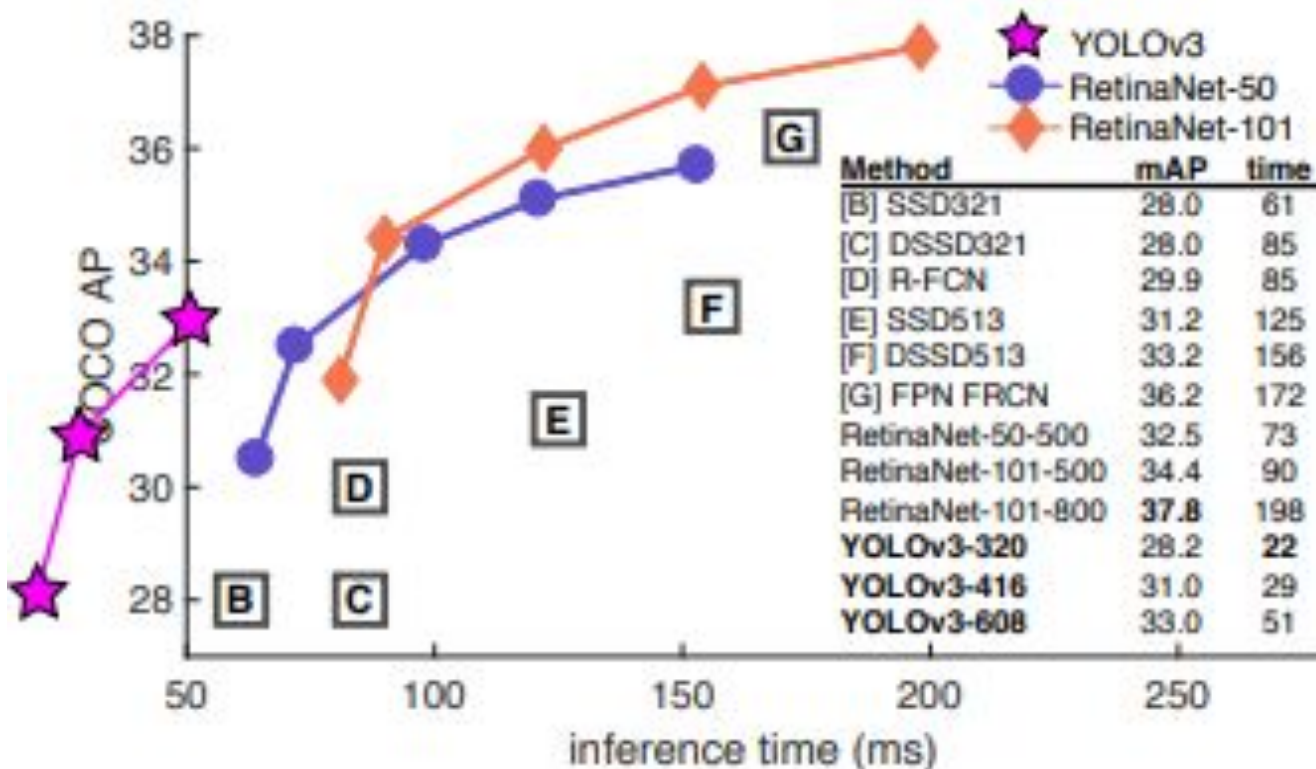
	mAP	Combined	Gain
Fast R-CNN	71.8	-	-
Fast R-CNN (2007 data)	<b>66.9</b>	72.4	.6
Fast R-CNN (VGG-M)	59.2	72.4	.6
Fast R-CNN (CaffeNet)	57.1	72.1	.3
YOLO	63.4	<b>75.0</b>	<b>3.2</b>

# Results: Generalization to Art

	VOC 2007	Picasso		People-Art
	AP	AP	Best $F_1$	AP
<b>YOLO</b>	<b>59.2</b>	<b>53.3</b>	<b>0.590</b>	<b>45</b>
R-CNN	54.2	10.4	0.226	26
DPM	43.2	37.8	0.458	32
Poselets [2]	36.5	17.8	0.271	
D&T [4]	-	1.9	0.051	



# Prologue: YOLOv3: An Incremental Improvement





# Prologue: Joe Redmon



**Joe Redmon**

@pjreddie

I stopped doing CV research because I saw the impact my work was having. I loved the work but the military applications and privacy concerns eventually became impossible to ignore.