

# 深度学习-目标检测篇

bilibili: 霹雳吧啦

作者：神秘的wz

# YOLO系列

## 规划

- YOLO v1 (理论简单介绍)
- YOLO v2 (理论简单介绍)
- YOLO v3 (理论详细介绍)
- YOLO v3 SPP (trick扩充+代码讲解)



代码参考: <https://github.com/ultralytics/yolov3>

# YOLO v1

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

Joseph Redmon\*, Santosh Divvala\*<sup>†</sup>, Ross Girshick<sup>¶</sup>, Ali Farhadi\*<sup>†</sup>

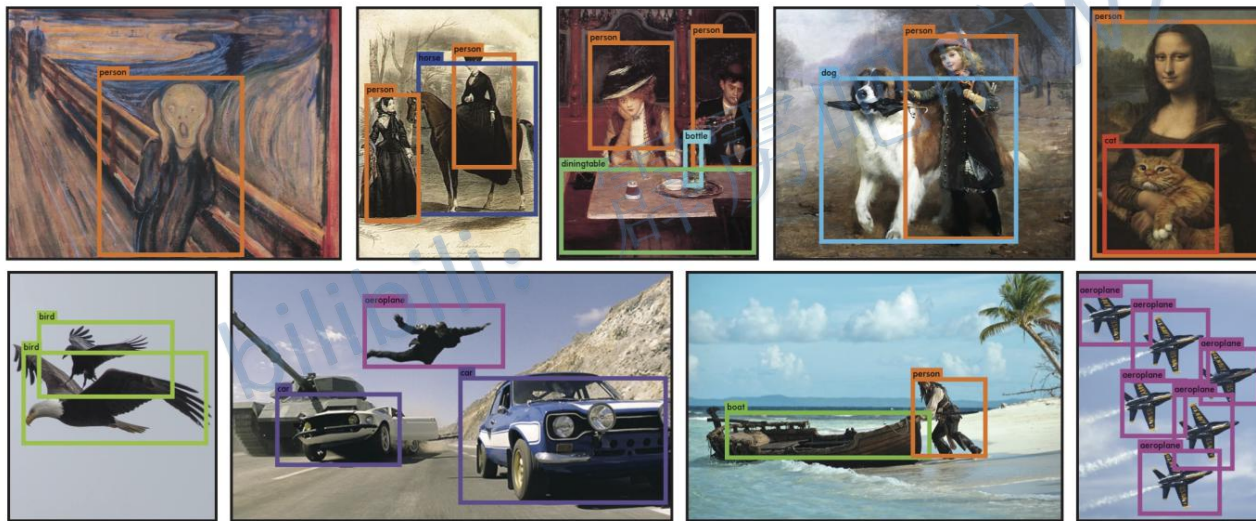
University of Washington\*, Allen Institute for AI<sup>†</sup>, Facebook AI Research<sup>¶</sup>

<http://pjreddie.com/yolo/>

2016 CVPR

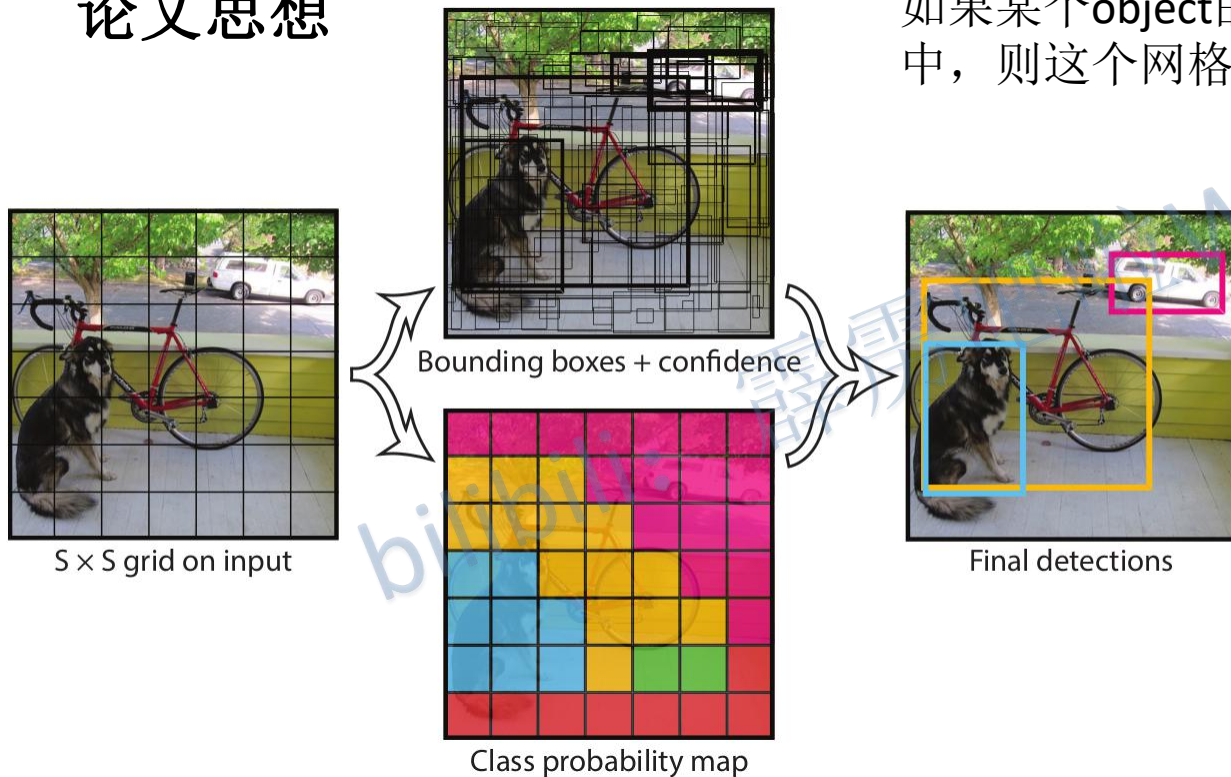
45FPS 448x448

63.4mAP



# YOLO v1

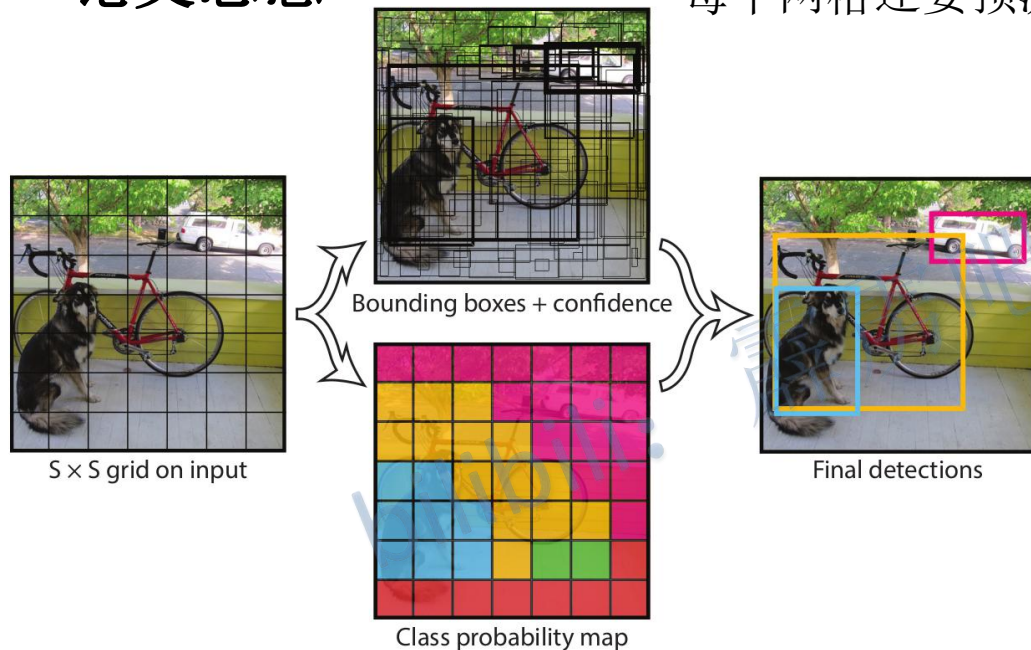
## 论文思想



1) 将一幅图像分成 $S \times S$ 个网格(grid cell), 如果某个object的中心落在这个网格中, 则这个网格就负责预测这个object。

# YOLO v1

## 论文思想



2)每个网格要预测B个bounding box，每个bounding box除了要预测位置之外，还要附带预测一个confidence值。每个网格还要预测C个类别的分数。

For evaluating YOLO on PASCAL VOC, we use  $S = 7$ ,  $B = 2$ . PASCAL VOC has 20 labelled classes so  $C = 20$ . Our final prediction is a  $7 \times 7 \times 30$  tensor.

Each bounding box consists of 5 predictions:  $x, y, w, h$ , and confidence. The  $(x, y)$  coordinates represent the center of the box relative to the bounds of the grid cell. The width and height are predicted relative to the whole image. Finally the confidence prediction represents the IOU between the predicted box and any ground truth box.

also how accurate it thinks the box is that it predicts. Formally we define confidence as  $\Pr(\text{Object}) * \text{IOU}_{\text{pred}}^{\text{truth}}$ . If no

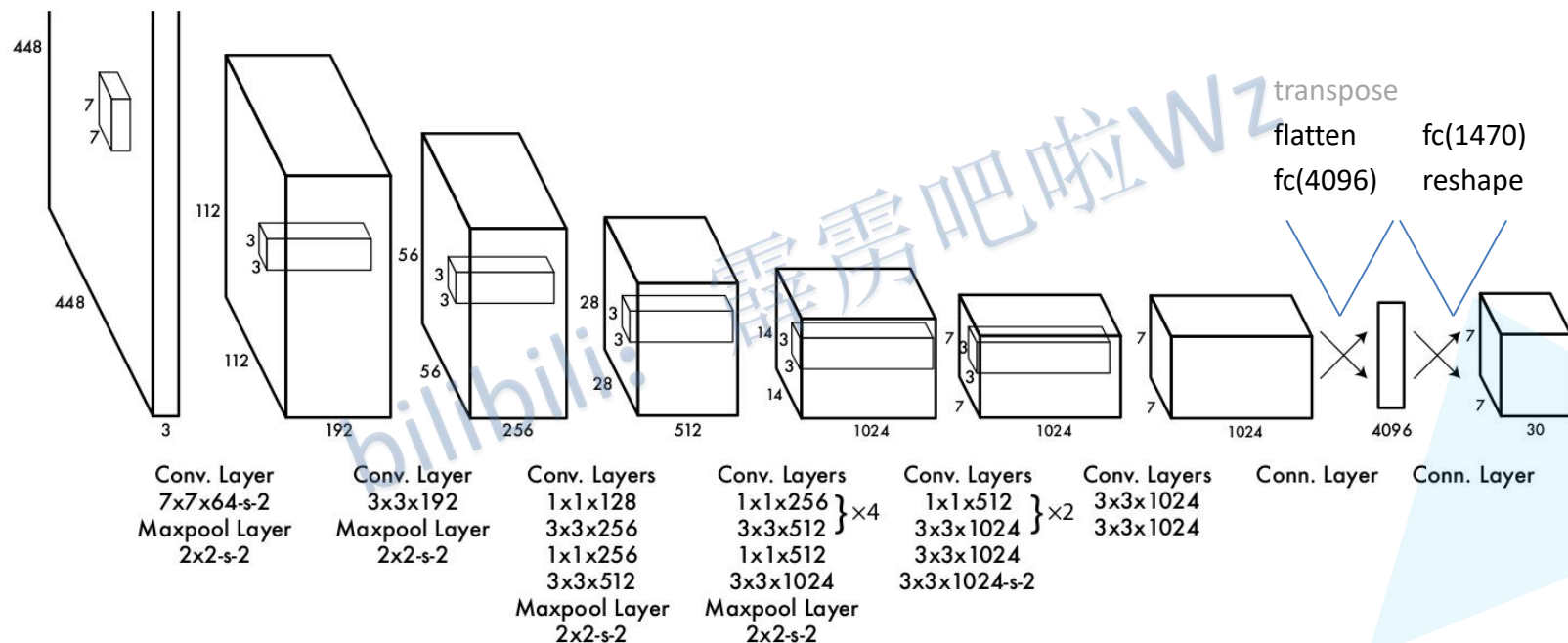
At test time we multiply the conditional class probabilities and the individual box confidence predictions,

$$\Pr(\text{Class}_i | \text{Object}) * \Pr(\text{Object}) * \text{IOU}_{\text{pred}}^{\text{truth}} = \Pr(\text{Class}_i) * \text{IOU}_{\text{pred}}^{\text{truth}} \quad (1)$$



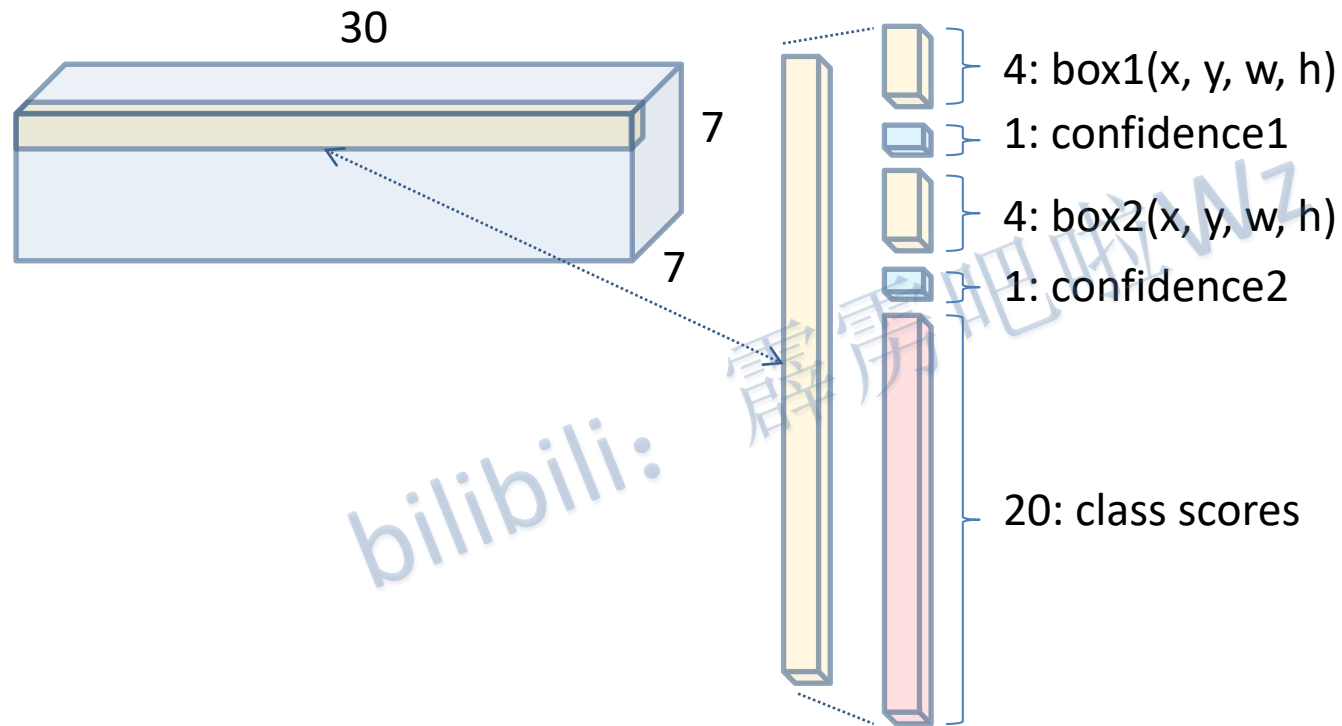
# YOLO v1

## 网络结构



# YOLO v1

## 网络结构



# YOLO v1

## 损失函数

sum-squared error

误差平方和

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

bounding box损失

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

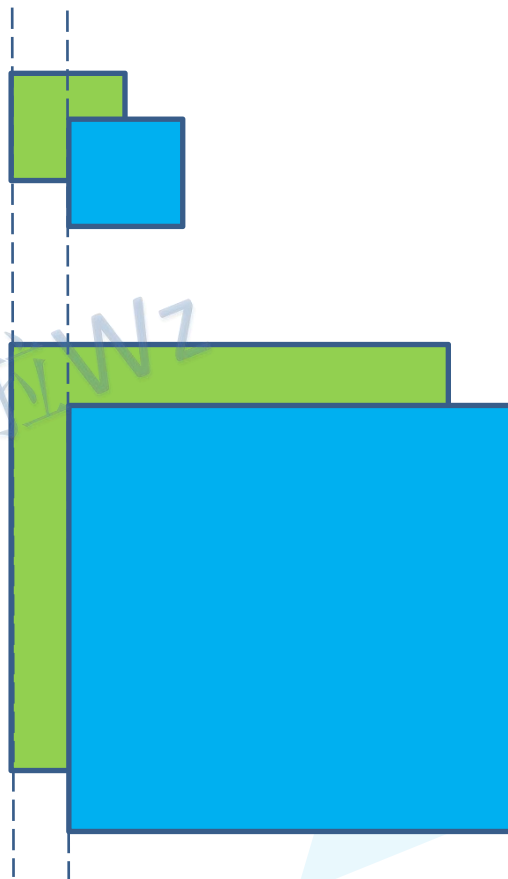
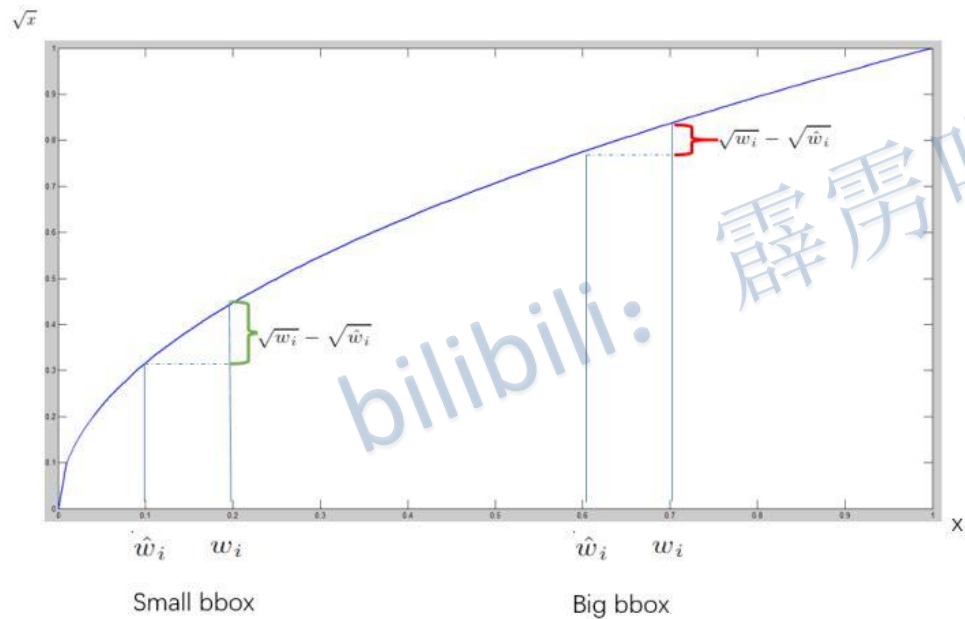
confidence损失

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

(3) classes损失



# YOLO v1



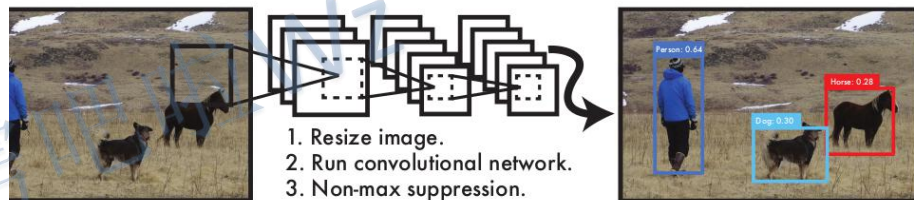
# YOLO v1

## 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.



# 沟通方式

## 1.github

<https://github.com/WZMIAOMIAO/deep-learning-for-image-processing>

## 2.CSDN

[https://blog.csdn.net/qq\\_37541097/article/details/103482003](https://blog.csdn.net/qq_37541097/article/details/103482003)

## 3.bilibili

<https://space.bilibili.com/18161609/channel/index>

尽可能每周更新