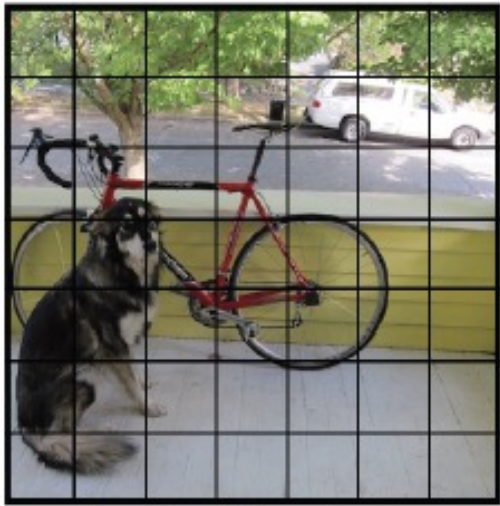


# Object detection: Where and What together

D.A. Forsyth

# Yolo

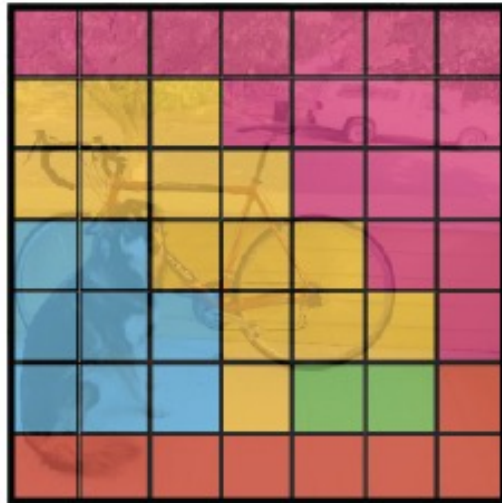
- Predict box objectness score, label scores indep.
- Process:
  - split image into grid
  - each cell predicts box location, confidence
  - each cell predicts a class probability
  - multiply box and class probabilities
  - cleanup
    - threshold, NMS



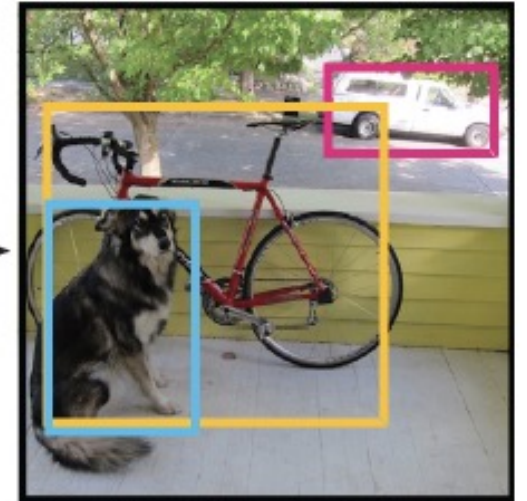
$S \times S$  grid on input



Bounding boxes + confidence



Class probability map

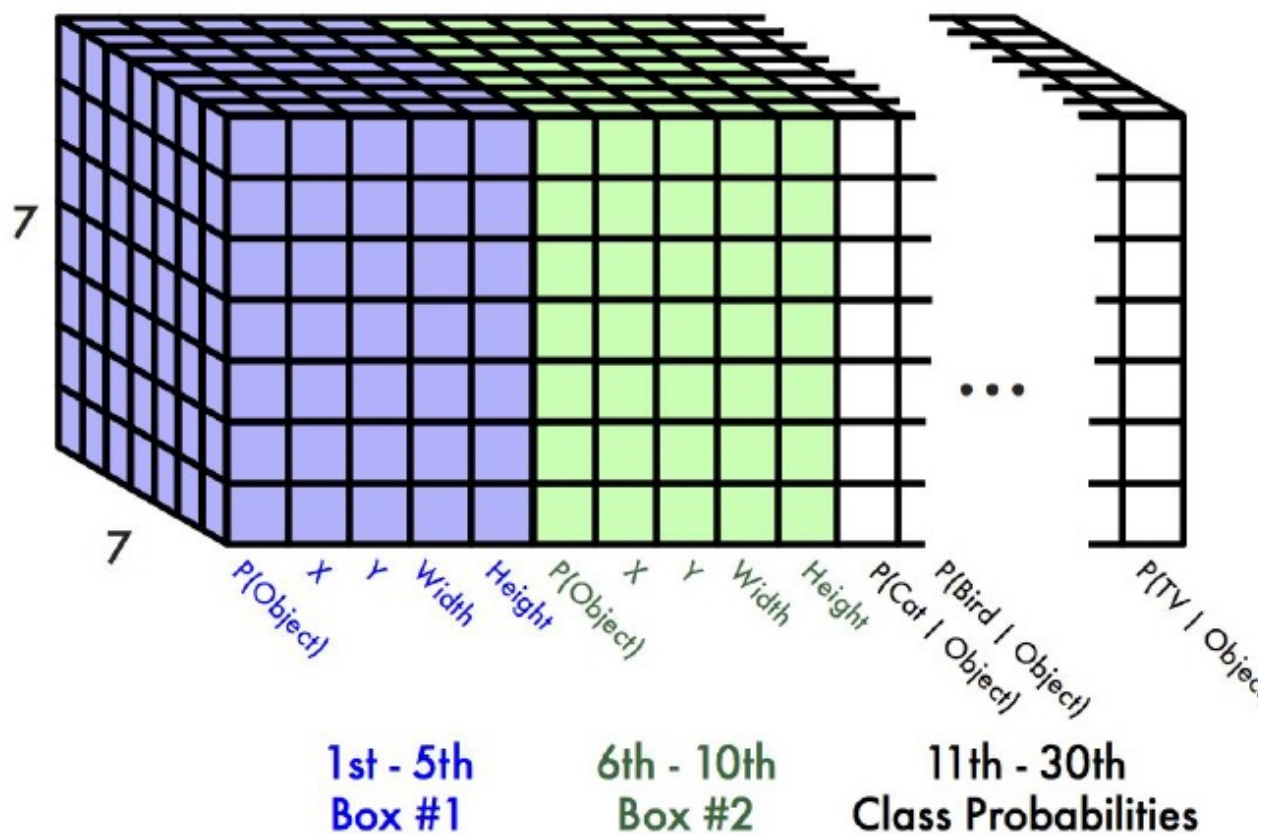


Final detections

# This parameterization fixes the output size

Each cell predicts:

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



For Pascal VOC:

- 7x7 grid
- 2 bounding boxes / cell
- 20 classes

$$7 \times 7 \times (2 \times 5 + 20) = 7 \times 7 \times 30 \text{ tensor} = \mathbf{1470 \text{ outputs}}$$



- Original caption: *“YOLO running on sample artwork and natural images from the internet. It is mostly accurate although it does think one person is an airplane”*

Credit: Figure 5 of You Only Look Once: Unified, Real-Time Object Detection  
Joseph Redmon, Santosh Divvala, Ross Girshick, Ali Farhadi

# Speed/accuracy payoffs

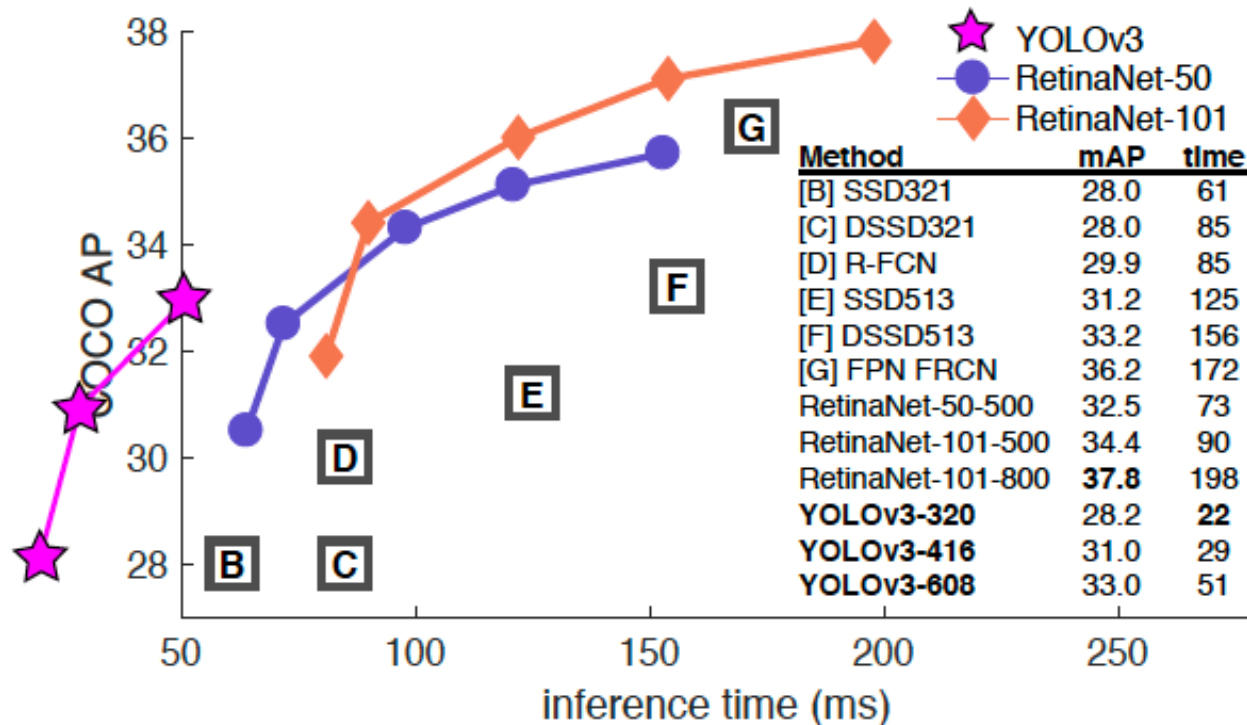


Figure 1. We adapt this figure from the Focal Loss paper [9]. YOLOv3 runs significantly faster than other detection methods with comparable performance. Times from either an M40 or Titan X, they are basically the same GPU.

Figure from: “YOLOv3: An Incremental Improvement”

# Yolo variants are maintained

- At V.11 now
- Still very widely used
- Training is relatively straightforward