

YOU LOOK ONLY ONCE (YOLO)

A Hands-on Introduction to the State of the Art
Realtime Object Detection Deep Learning Model

Bradley Erickson

MD PhD

Pouria Rouzrokh

MD MPH MHPE

Radiology Informatics Laboratory
Mayo Clinic

What to Expect

Our goal for today:

To learn about an object detection deep learning model, called “YOLO”!

What we cannot do in 60 minutes:

To discuss every single line of notebook code / basics of DL programming!

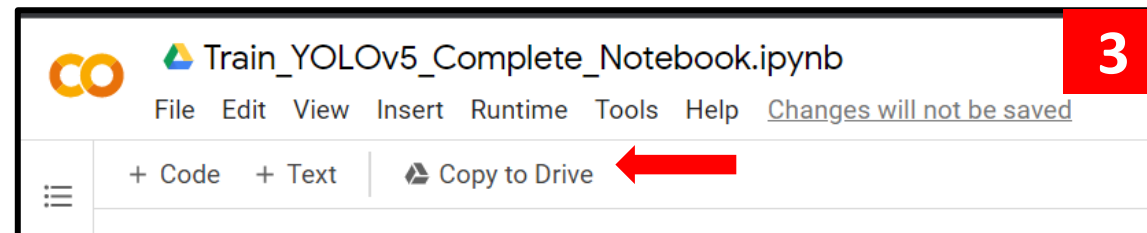
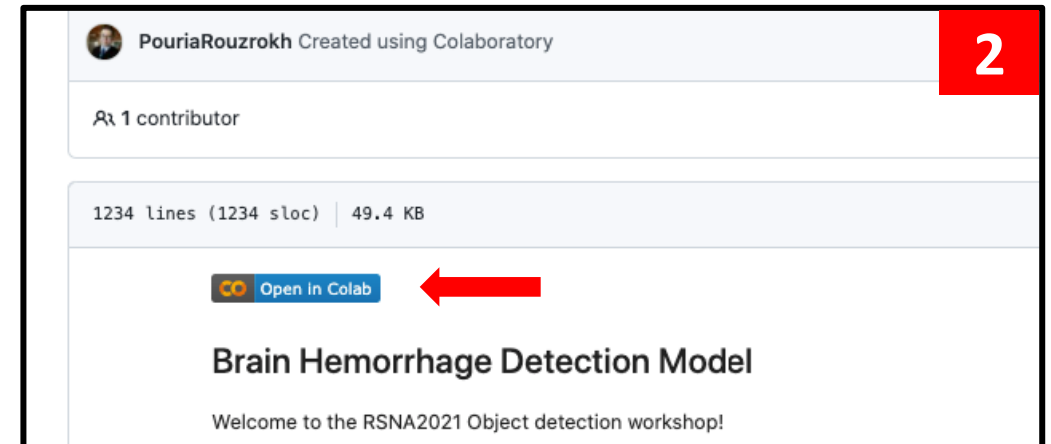
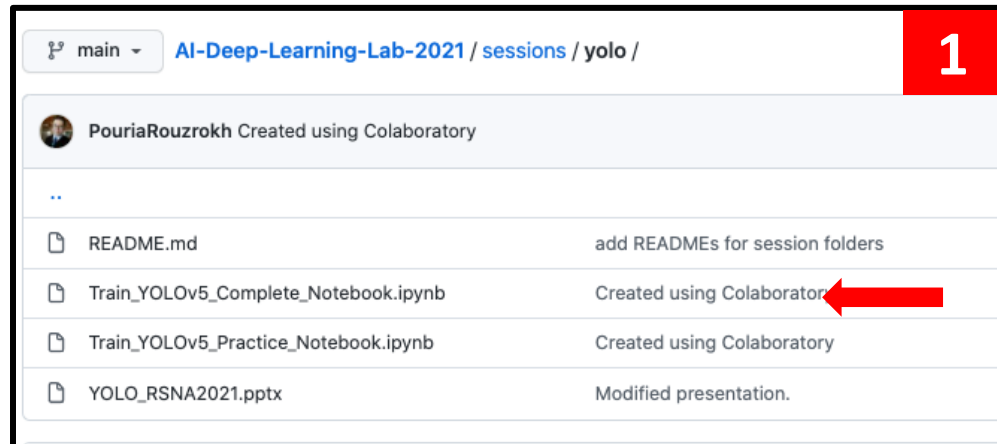
- The notebook will be accessible to you with full explanation in addition to the code.
- Ultralytics documentations include all the code details for implementing YOLO.

What we can do in 60 minutes:

- To talk about the **big picture** of YOLO and how to implement YOLOv5.
- To talk about **handy tips** for obtaining a better performance from YOLOv5!

Let's Open Our Notebook!

<https://tinyurl.com/28bp6nm8>



Defining Object Detection

Semantic Segmentation



GRASS, CAT,
TREE, SKY

No objects, just pixels

Classification



CAT

Single Object

Object Detection



DOG, DOG, CAT

Multiple Object

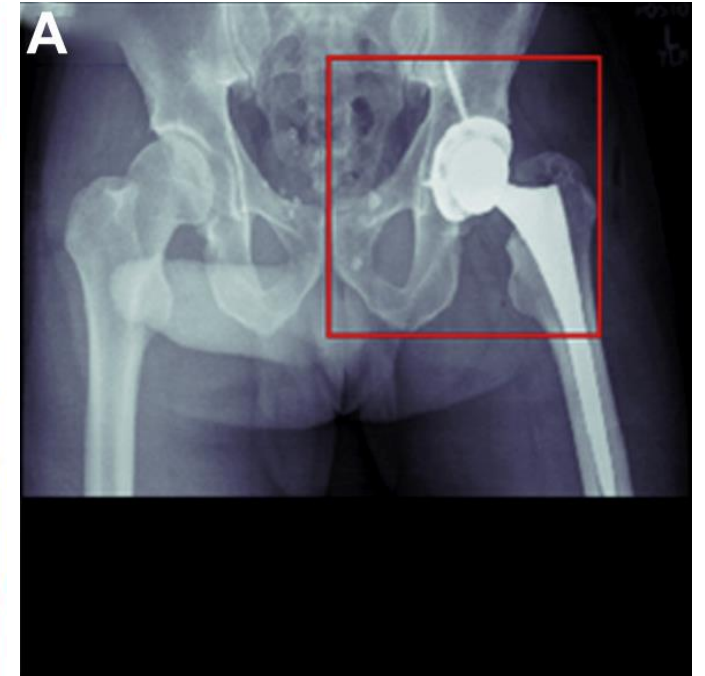
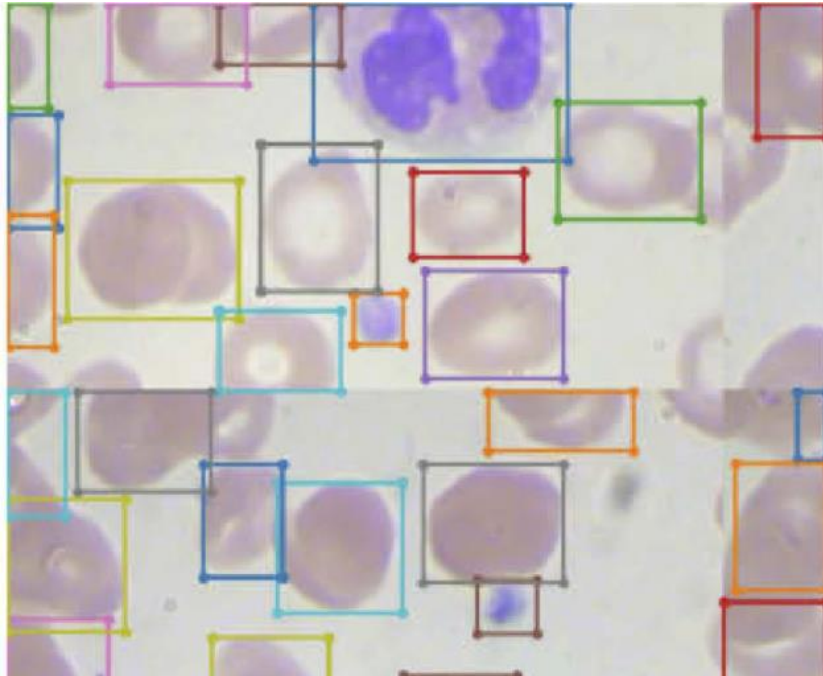
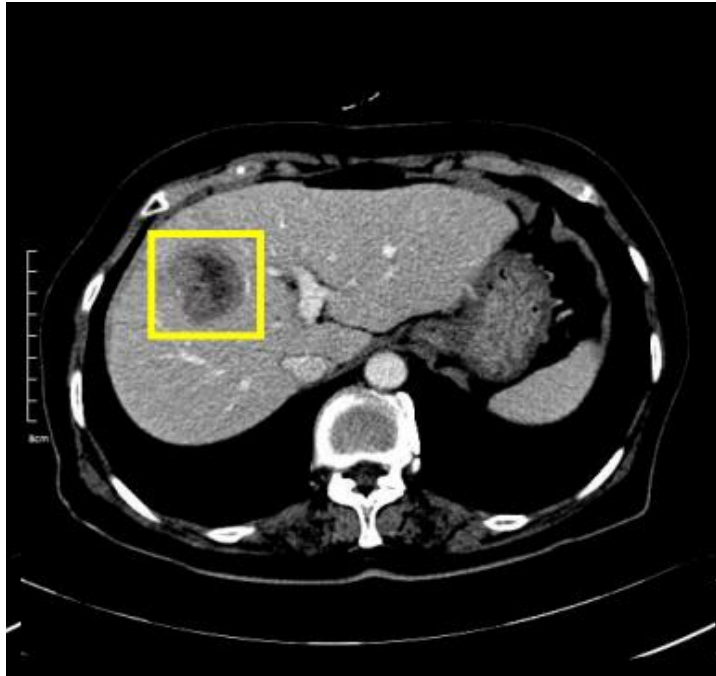
Instance Segmentation



DOG, DOG, CAT

[This image is CC0 public domain](#)

Medical Applications of Object Detection

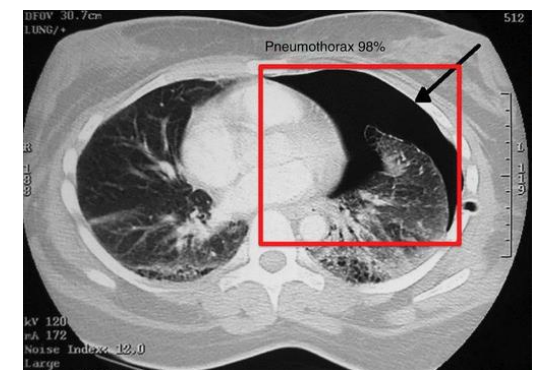
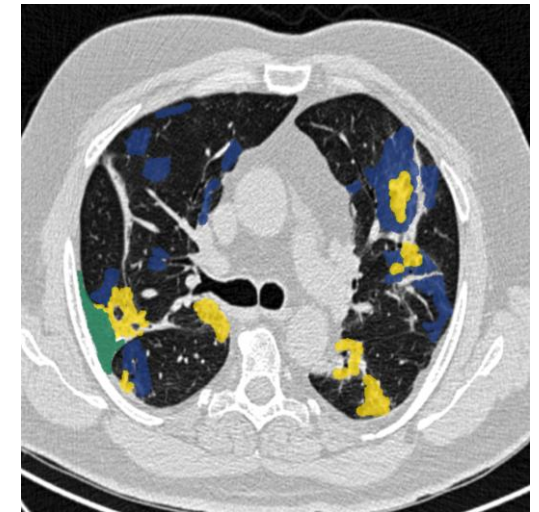


Sources:

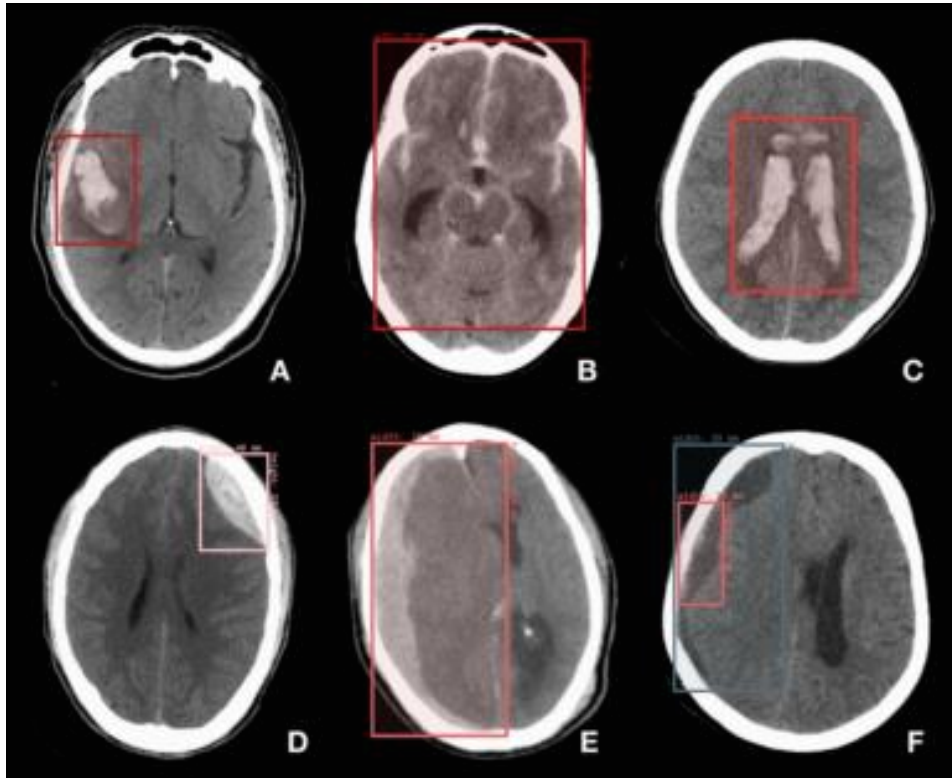
- [Liver Lesion Detection from Weakly-labeled Multi-phase CT Volumes with a Grouped Single Shot MultiBox Detector](#)
- [Improved detection performance in blood cell count by an attention-guided deep learning method](#)
- [Deep Learning Artificial Intelligence Model for Assessment of Hip Dislocation Risk Following Primary Total Hip Arthroplasty From Postoperative Radiographs](#)

Object Detection vs. Object Segmentation

Object Detection	Semantic Segmentation
Does patch (bounding box)-level predictions.	Does pixel-level prediction.
Does not identify the exact borders or volume of an object.	Can identify the exact borders (and volume) of an object.
Less tedious manual labeling: Simply draw a bounding box!	More tedious manual labeling: Manually segment each pixel of interest!



What We Do Today!



To train a deep learning model to detect brain hemorrhage lesion on Head CT scans!

For our data, we use the publicly available [CQ500 Head CT-scan dataset](#) of patients with brain hemorrhage.

For our labels, we use the bounding box annotation on CQ500 dataset by [Physionet](#).

Our Notebook Workflow

- **Part 0: Setting up the working directory**
- **Part 1: Extracting CT images from DICOM files**
- **Part 2: Collecting and visualizing the annotated bounding box labels**
- **Part 3: Data splitting and setting up the training and validation sets**
- **Part 4: Downloading the YOLO model and configuring it**
- **Part 5: View and modify the hyperparameters (optional)**
- **Part 6: Training**
- **Part 7: Inference**



80% of code!



20% of code!

You Look Only Once (YOLO)



- The original YOLO model was the first object detection network to combine the problem of drawing bounding boxes and identifying class labels **in a single pass of training** (different than two-stage models like Faster RCNN)
- YOLO models are often **very fast and very small** – often making them faster to train and easier to deploy, especially to compute-limited edge devices.

You Look Only Once (YOLO)

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

Joseph Redmon*, Santosh Divvala
University of Washington*, Allen Institute of AI
<http://pjreddie.com/yolo/>

Abstract

We present YOLO, a new approach to object detection. Prior work on object detection repurposes classifiers to perform detection. Instead, we frame object detection as a regression problem to spatially separated bounding boxes and associated class probabilities. A single neural network predicts bounding boxes and class probabilities directly from full images in one evaluation. Since the whole detection pipeline is a single network, it can be optimized end-to-end directly on detection performance.

Our unified architecture is extremely fast. Our baseline YOLO model processes images in real-time at 45 frames per second. A smaller version of the network, Fast YOLO, processes an astounding 155 frames per second while still achieving double the mAP of other real-time detectors. Compared to state-of-the-art detection systems, YOLO makes more localization errors but is less likely to predict false positives on background. Finally, YOLO learns very general representations of objects. It outperforms other detection methods, including DPM and R-CNN, when generalizing from natural images to other domains like artwork.

YOLO9000:

Better, Faster

Joseph Redmon
University of Washington
<http://pjreddie.com/yolo/>

Abstract

We introduce YOLO9000, a state-of-the-art, real-time object detection system that can detect over 9000 object categories. First we propose various improvements to YOLO detection method, both novel and drawn from prior work. The improved model, YOLOv2, is state-of-the-art on standard detection tasks like PASCAL VOC and COCO. Then, we propose a novel, multi-scale training method the same YOLO model can run at varying sizes, offering an easy trade-off between speed and accuracy. At 67 FPS, YOLOv2 achieves 76.8 mAP on VOC 2007. At 40 FPS, YOLOv2 gets 71.2 mAP, outperforming state-of-the-art methods like Faster R-CNN with ResNet and SSD while still running significantly faster. Finally we propose a method to jointly train on object detection and classification. Using this method we train YOLO9000 simultaneously on the COCO detection data and the ImageNet classification dataset. Our joint training allows YOLO9000 to predict detections for object classes that don't have labelled detection data. We validate our approach on the ImageNet detection task. YOLO9000 achieves 19.7 mAP on the ImageNet detection validation set despite only having detection data for 44 of the 200 classes. On 156 classes not in COCO, YOLO9000 gets 16.0 mAP. YOLO can detect more than just 200 classes; it predicts

YOLOv3: An Incremental Improvement

Abstract

We present some updates to YOLO. It's a lot of little design changes to make it better. It's the last time but more accurate. It's still a little worry. At 320×320 YOLOv3 runs in 100 ms as accurate as SSD but three times faster. At the old .5 IOU mAP detection method, it achieves 57.9 AP₅₀ in 51 ms compared to 57.5 AP₅₀ in 198 ms by RetinaNet. It's 3.8x faster. As always, all

1. Introduction

Sometimes you just kinda phone it in. I didn't do a whole lot of research. I spent a lot of time on Twitter. Played around

YOLOv4: Optimal Speed and Accuracy of Object Detection

Alexey Bochkovskiy*
alexeyab84@gmail.com

Chien-Yao Wang*
Institute of Information Science
Academia Sinica, Taiwan
kinyiu@iis.sinica.edu.tw

Hong-Yuan Mark Liao
Institute of Information Science
Academia Sinica, Taiwan
liao@iis.sinica.edu.tw

Abstract

There are a huge number of features which are said to improve Convolutional Neural Network (CNN) accuracy. Practical testing of combinations of such features on large datasets, and theoretical justification of the result, is required. Some features operate on certain models exclusively and for certain problems exclusively, or only for small-scale datasets; while some features, such as batch-normalization and residual-connections, are applicable to the majority of models, tasks, and datasets. We assume that such universal features include Weighted-Residual-Connections (WRC), Cross-Stage-Partial-connections (CSP), Cross mini-Batch Normalization (CmBN), Self-adversarial-training (SAT) and Mish-activation. We use new features: WRC, CSP, CmBN, SAT, Mish activation, Mosaic data augmentation, CmBN, DropBlock regularization, and CIoU loss, and combine some of them to achieve state-of-the-art results: 43.5% AP (65.7% AP₅₀) for the MS COCO dataset at a real-time speed of ~65 FPS on Tesla V100. Source code is at <https://github.com/AlexeyAB/darknet>.

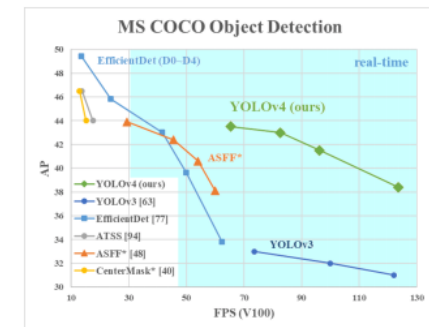
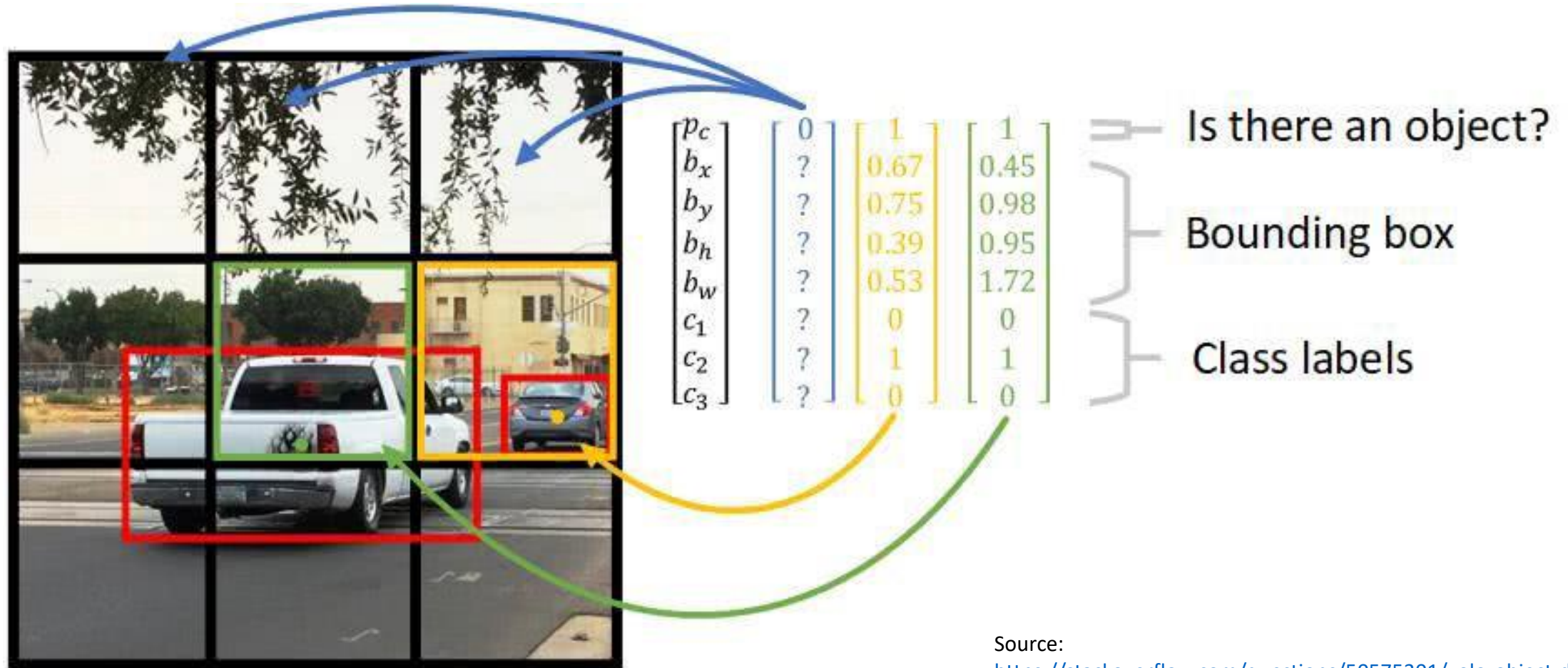


Figure 1: Comparison of the proposed YOLOv4 and other state-of-the-art object detectors. YOLOv4 runs twice faster than EfficientDet with comparable performance. Improves YOLOv3's AP and FPS by 10% and 12%, respectively.

The main goal of this work is designing a fast operating

Fundamental Idea of YOLO

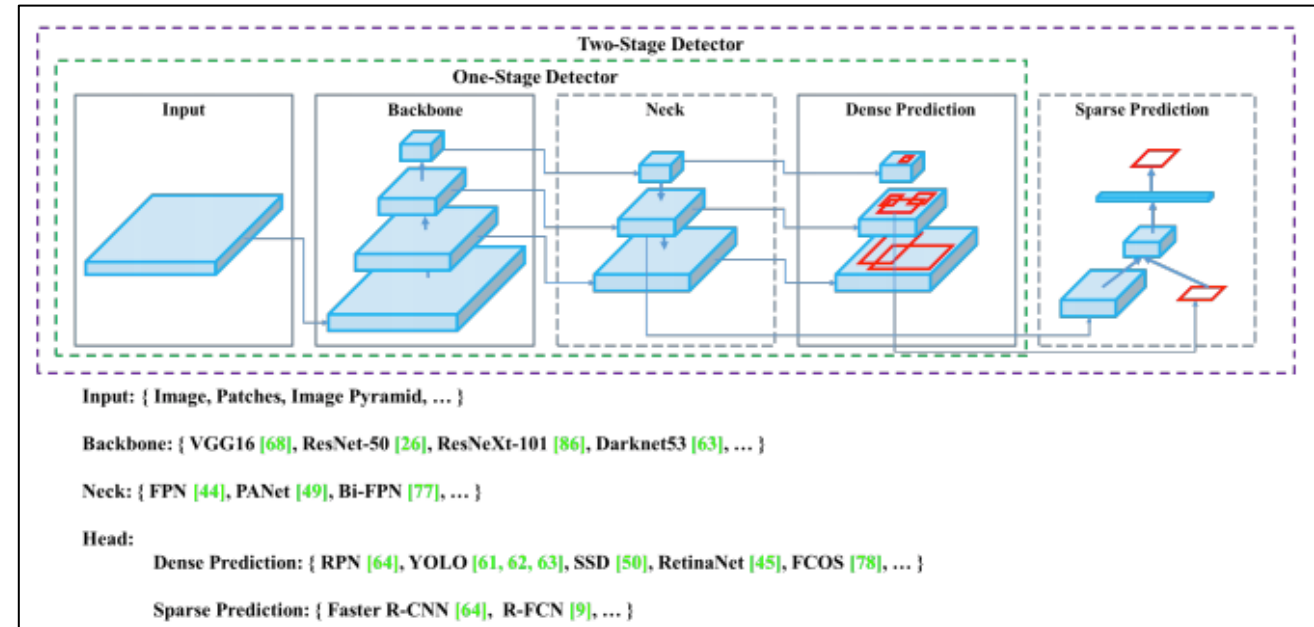


Source:

<https://stackoverflow.com/questions/50575301/yolo-object-detection-how-does-the-algorithm-predict-bounding-boxes-larger-than>

Well, YOLO has evolved a lot!

- New versions of YOLO use advanced architectures, data augmentation techniques, and training strategies! Please refer to their papers for understanding the novelties.



YOLOv5

- One of the latest versions of YOLO.
- No scientific manuscript released yet.
- Based on YOLOv4 + some technical improvements.
- Based on **Pytorch** (as opposed to Darknet), and therefore significantly smaller than YOLOv4.
- Easy implementation for custom training and helpful documentation/tutorials by **Ultralytics**.



Big Picture of Training a YOLOv5 Model

- Ultralytics has implemented YOLOv5 as a ready to train package, i.e., training and inference with it take only **1 to 2 lines of code**!
- However, most of your time will be spent on **collecting and organizing your data** as YOLOv5 expects it to be.
- When your data is ready, you should **clone the YOLOv5 model directory from the Ultralytics GitHub page**, and then configure it to know your data.
- You can select **different size variants of YOLOv5**, with or without **transfer learning** from the models trained on ImageNet.

Data Organization

▼ model_data

▼ images

▶ train

▶ valid

Image1.png
Image2.png
...

▼ labels

▶ train

▶ valid

Image1.txt
Image2.txt
...

train.txt

val.txt

Path to image1.png
Path to image2.png
...

▶ yolov5

Each file

Inside Image1.txt

- One row per object
- Each row is `class x_center y_center width height` format.
- Box coordinates must be in **normalized xywh** format (from 0 - 1). If your boxes are in pixels, divide `x_center` and `width` by image width, and `y_center` and `height` by image height.
- Class numbers are zero-indexed (start from 0).

```
0 0.21 0.45 0.18 0.23  
1 0.6 0.74 0.2 0.1
```

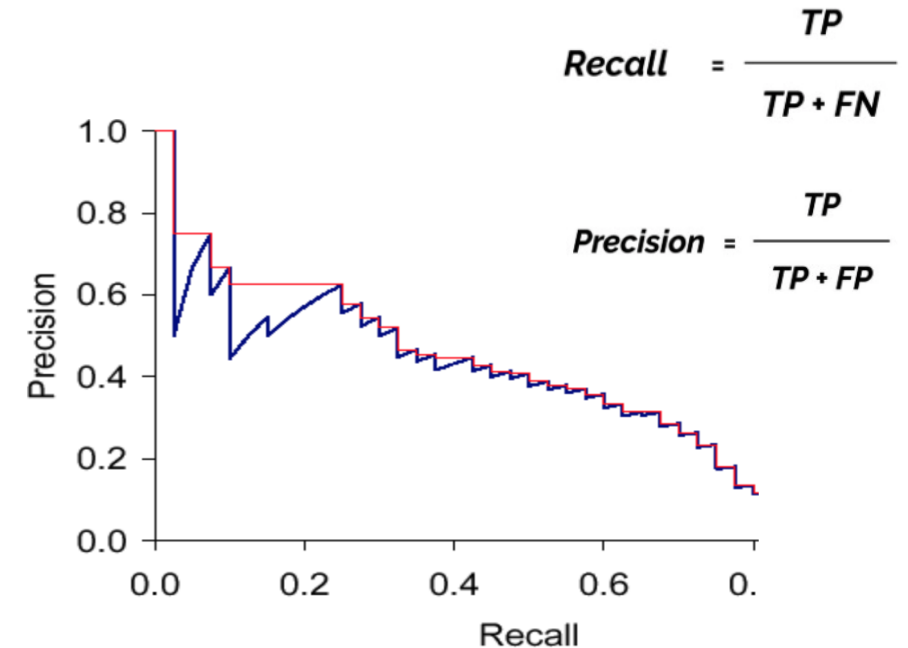
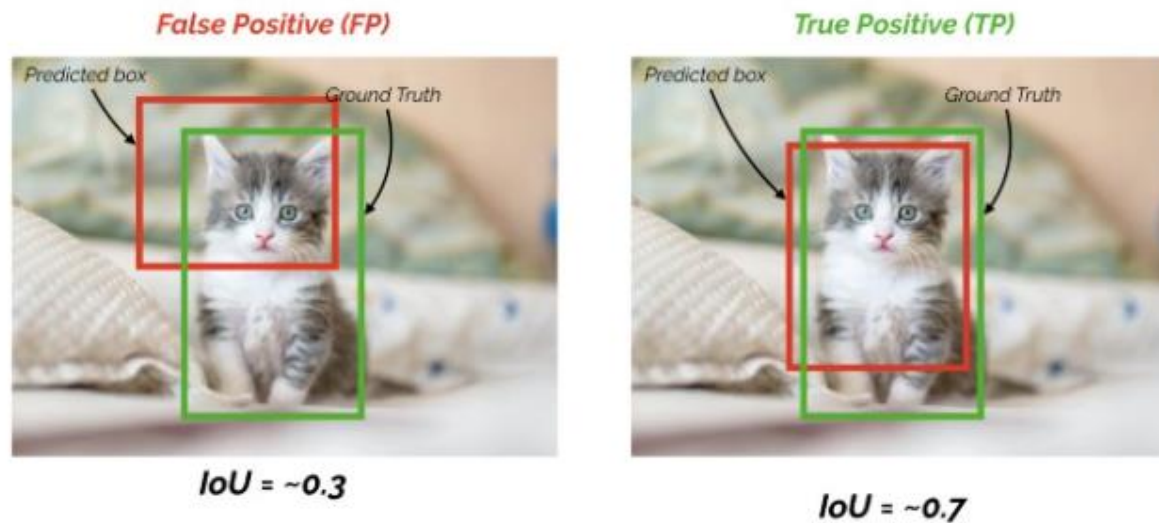
Intersection over Union (IOU)

Intersection over Union (IoU) = $\frac{\text{Area of Overlap}}{\text{Area of Union}}$

Source: <https://towardsdatascience.com/map-mean-average-precision-might-confuse-you-5956f1bfa9e2>

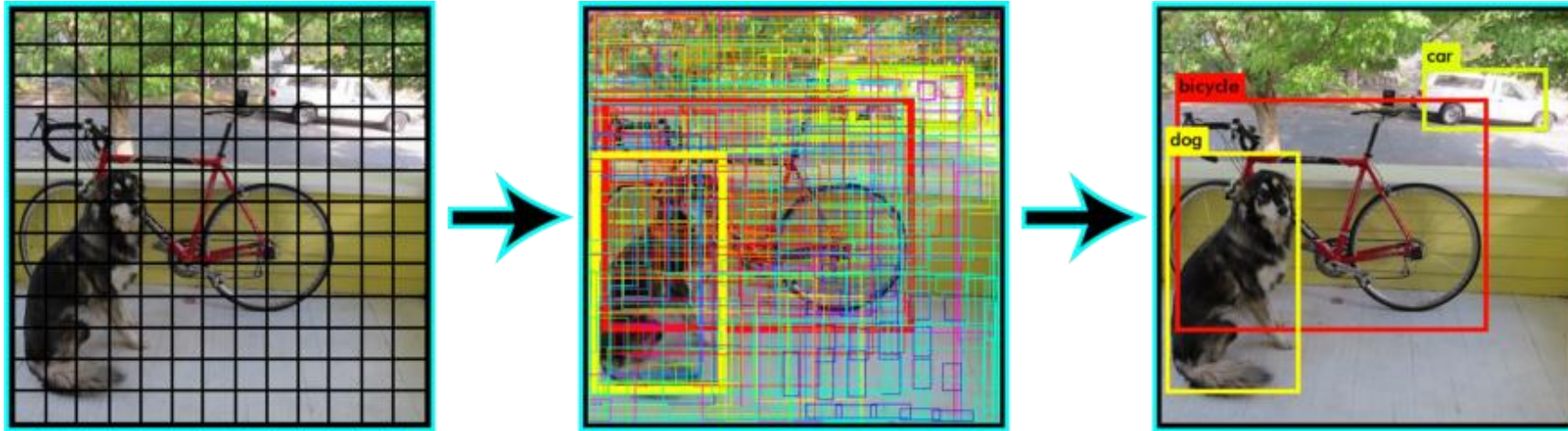
Mean Average Precision (mAP)

If IoU threshold = 0.5



The area under the Precision-Recall curve for an object detection model is called the “Average Precision” of that model. If you take the mean of these “Average Precision” scores across all classes which the model predicts, you will achieve a metric for evaluating your model: mean Average Precision (mAP)!

Non-max Suppression



Non-Max Suppression:

To select the best bounding box, from the multiple predicted bounding boxes, these object detection algorithms use non-max suppression.

In this technique, if two bounding boxes have an IoU over a certain threshold (e.g., 60%), the bounding box with the lower confidence score will be discarded.

Source: <https://www.pyimagesearch.com/2018/11/12/yolo-object-detection-with-opencv/>

Thank you for your attention!
