

Computer Vision: Fall 2022 — Lecture 15

Dr. Karthik Mohan

Univ. of Washington, Seattle

November 19, 2022

References

Generic ML/DL

- ① [Good Book for Machine Learning Concepts](#)
- ② [Deep Learning Reference](#)

CNN

- ① [Convolutional Neural Networks for Visual Recognition](#)
- ② [Convolutional Neural Net Tutorial](#)
- ③ [CNN Transfer Learning](#)
- ④ [PyTorch Transfer Learning Tutorial](#)

CNN Publication References

CNN surveys

- ① Convolutional Neural Networks: A comprehensive survey, 2019
- ② A survey of Convolutional Neural Networks: Analysis, Applications, and Prospects, 2021

CNN Archs

- ① GoogLeNet
- ② Top models on ImageNet
- ③ ResNet ILSVRC paper

Object Detection and Image Segmentation References

Object Detection

- ① A survey of modern deep learning based object detection methods
- ② YOLO Survey ✓
- ③ YOLO Original Paper ✓ }

↳ product of UW → ALI FARHADI

Learnings from the Mini-Project - Breakout Session!

Breakout and Discuss - Peer Learning (5 mins)

Breakout and discuss in your zoom room - What were your key learnings from the mini-project? What strategies worked and what didn't? How much did hyper-param tuning play a role in the result? Did you get to build your intuition with the models you tested?

Last Lecture

- 1 Introduction to Object Detection and Instance Segmentation
- 2 R-CNN model and metrics for Object Detection
& data sets

Today

- ① Object Detection Recap ✓
- ② R-CNN variants - Fast and Faster R-CNN
- ③ YOLO - Single Stage Object Detection
- ④ Results and Benchmarking on the data sets

Today

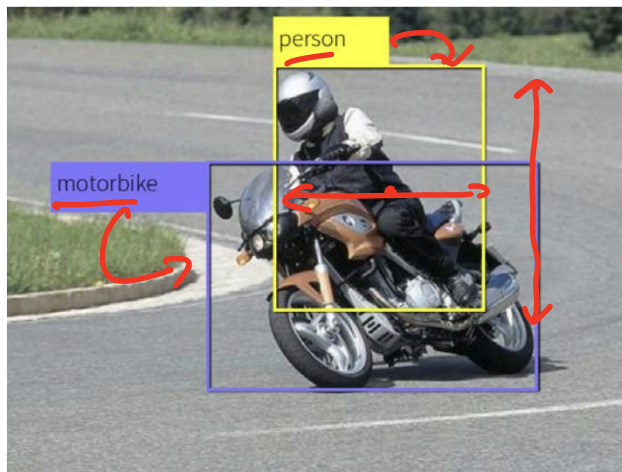
- ① **Object Detection Recap**
- ② R-CNN variant - Fast R-CNN
- ③ YOLO - Single Stage Object Detection
- ④ Results and Benchmarking on the data sets

Object Detection vs Image Segmentation

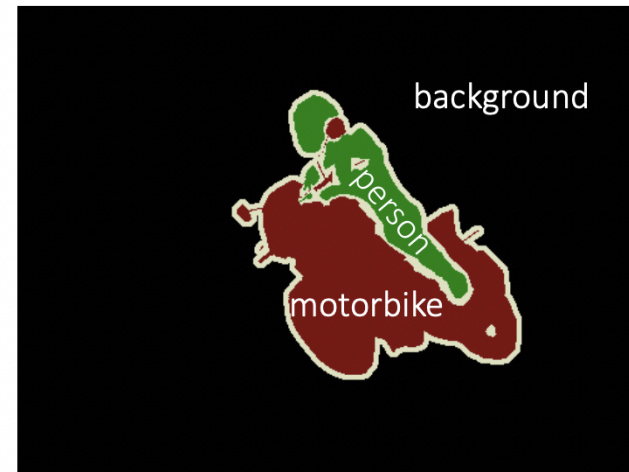
input image



object detection

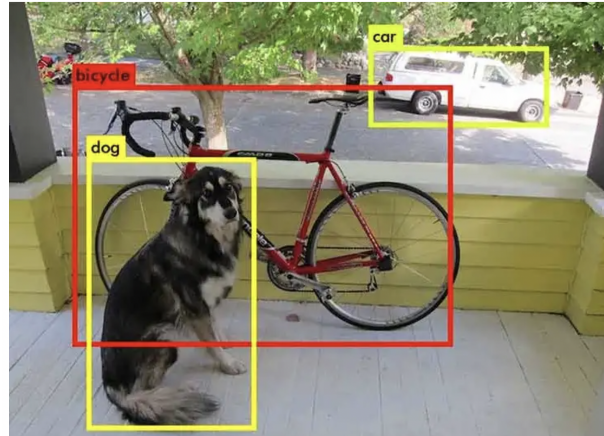


segmentation



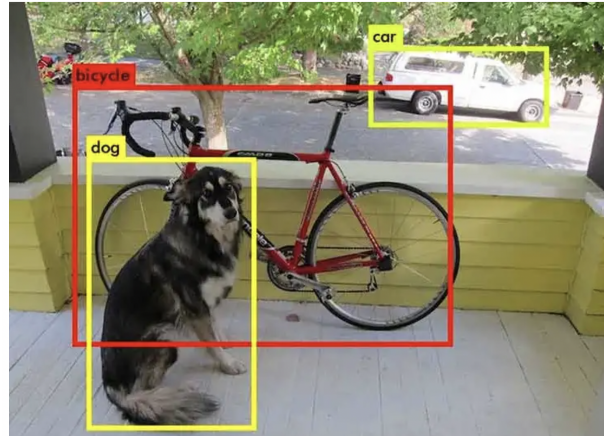
$(c_x, c_y, w, h), class$

Object Detection History



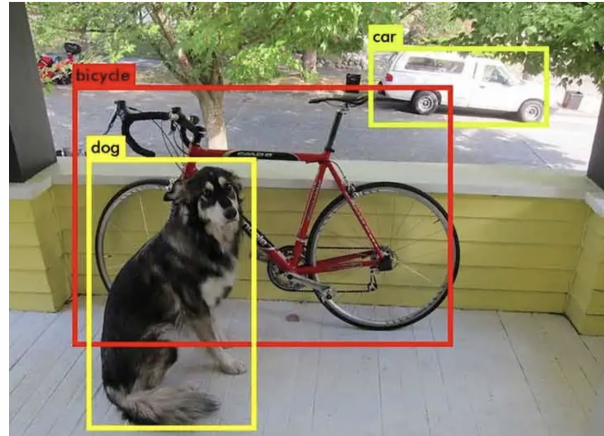
- 1 Has been an uphill task until 2012

Object Detection History



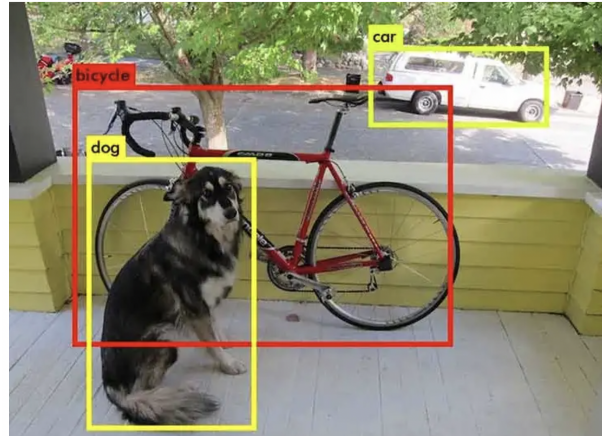
- 1 Has been an uphill task until 2012
- 2 Early detectors for objects - Ensemble of hand-crafted ones

Object Detection History



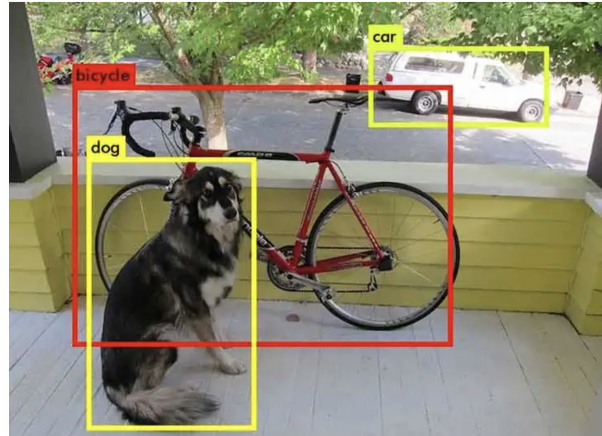
- 1 Has been an uphill task until 2012
- 2 Early detectors for objects - Ensemble of hand-crafted ones
- 3 Early detectors: Low accuracy and cumbersome/time-consuming

Object Detection History



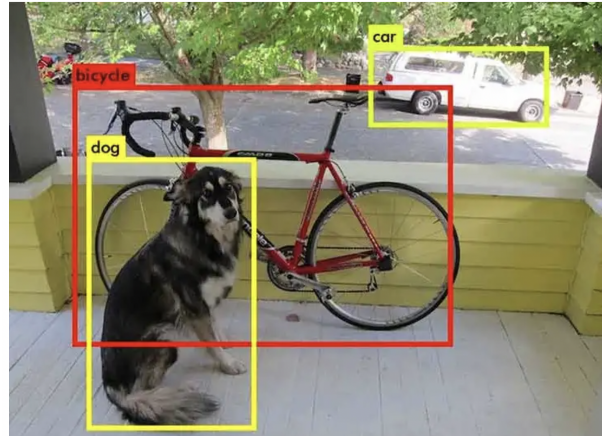
- ① Has been an uphill task until 2012
- ② Early detectors for objects - Ensemble of hand-crafted ones
- ③ Early detectors: Low accuracy and cumbersome/time-consuming
- ④ CNN changed the landscape - Better Accuracy, faster train, generalizability

Object Detection History



- 1 Has been an uphill task until 2012
- 2 Early detectors for objects - Ensemble of hand-crafted ones
- 3 Early detectors: Low accuracy and cumbersome/time-consuming
- 4 CNN changed the landscape - Better Accuracy, faster train, generalizability
- 5 AlexNet (2012) - First CNN archs to be applied to Obj. Detection

Object Detection History

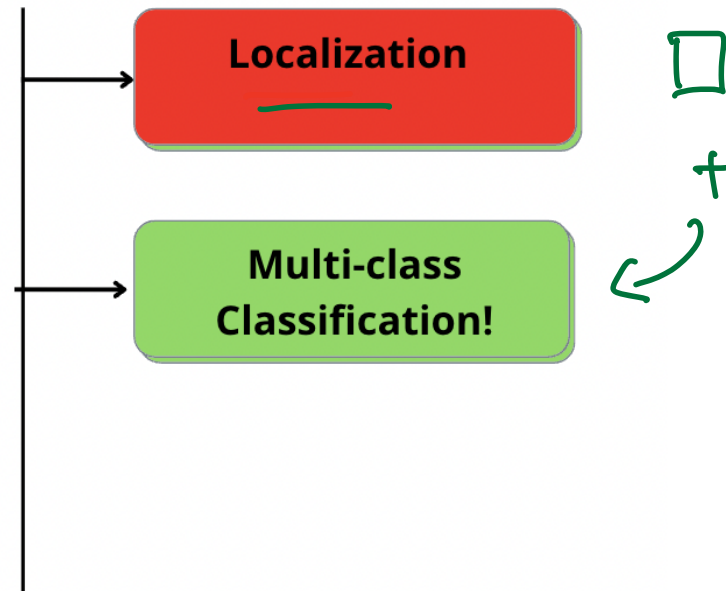


- 1 Has been an uphill task until 2012
- 2 Early detectors for objects - Ensemble of hand-crafted ones
- 3 Early detectors: Low accuracy and cumbersome/time-consuming
- 4 CNN changed the landscape - Better Accuracy, faster train, generalizability
- 5 AlexNet (2012) - First CNN archs to be applied to Obj. Detection
- 6 **Real world application:** Self-driving cars

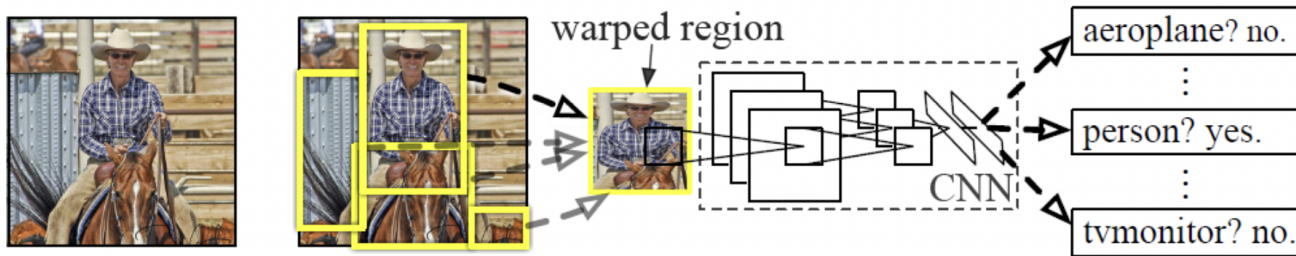
Multi-class Classification vs Object Detection

**Multi-Class
Classification**

Object Detection



Multi-label Classification vs Object Detection



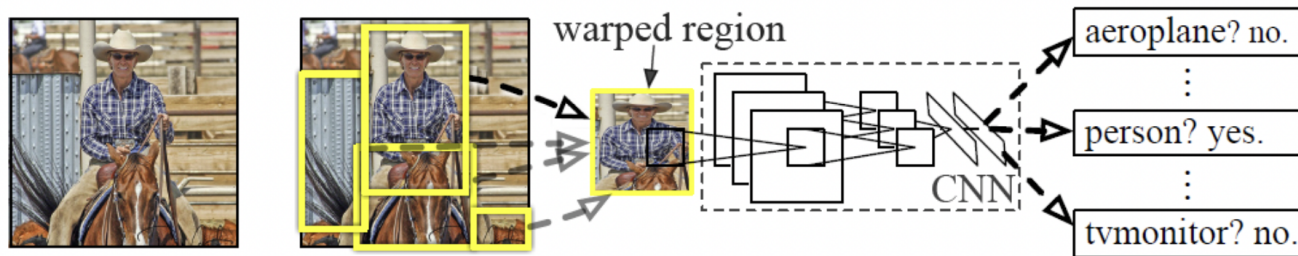
**Multi-Class
Classification**

Object Detection

Localization

**Multi-class
Classification!**

Multi-class Classification vs Object Detection



**Multi-Class
Classification**

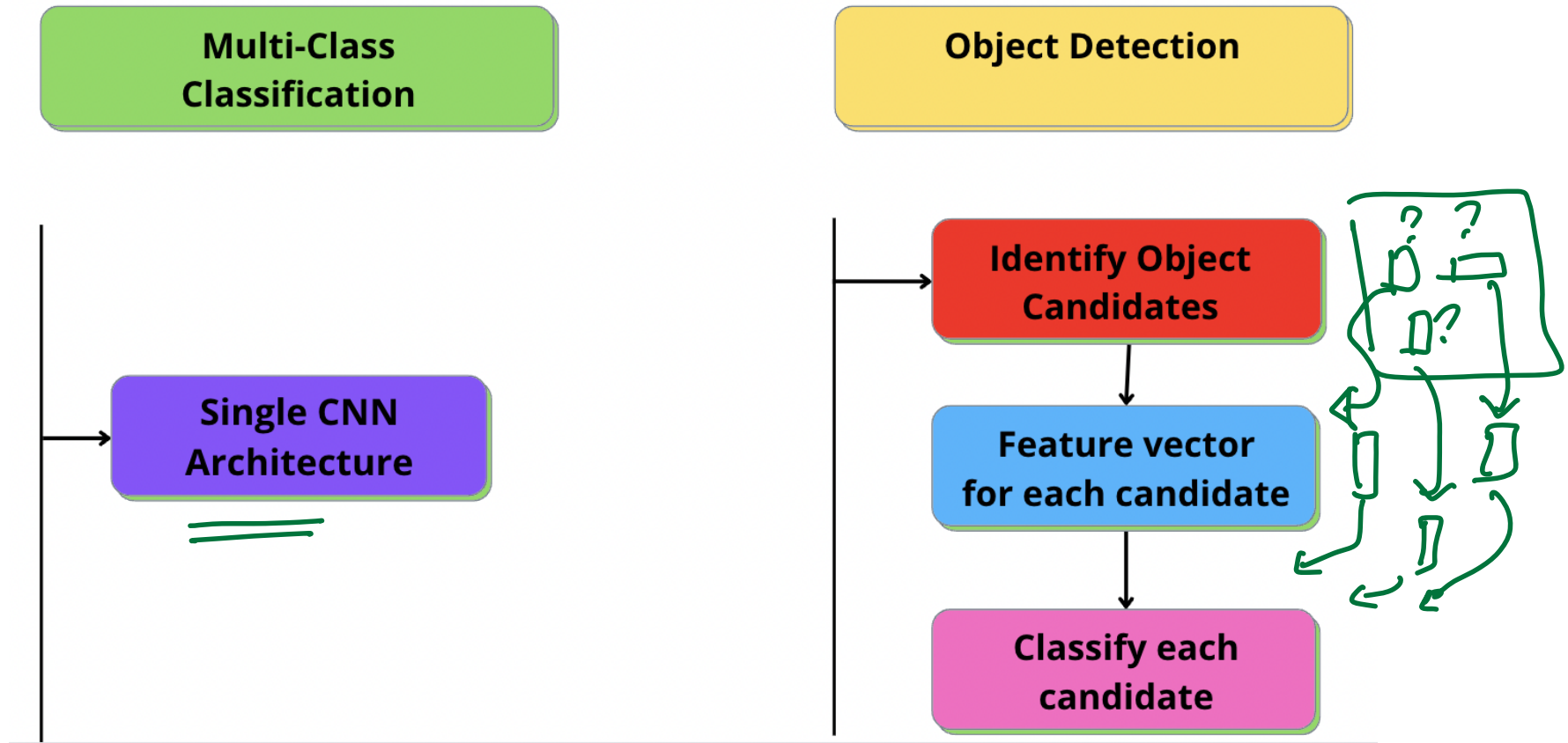
Object Detection

**Multi-label
Classificaiton**

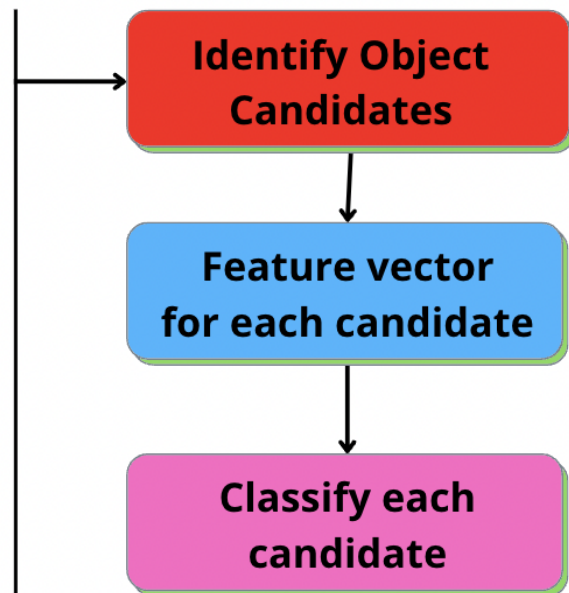
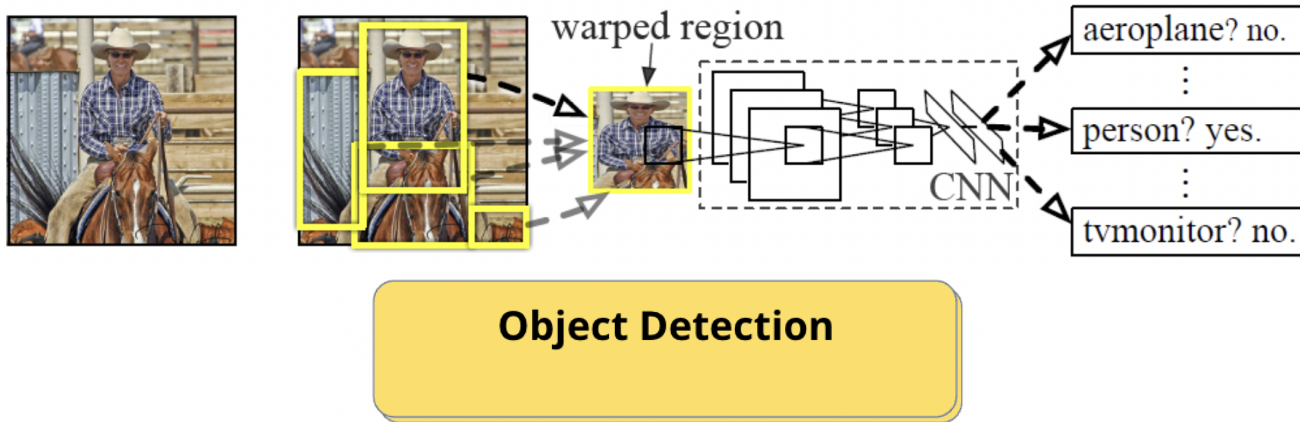
Localization

**Multi-class
Classification!**

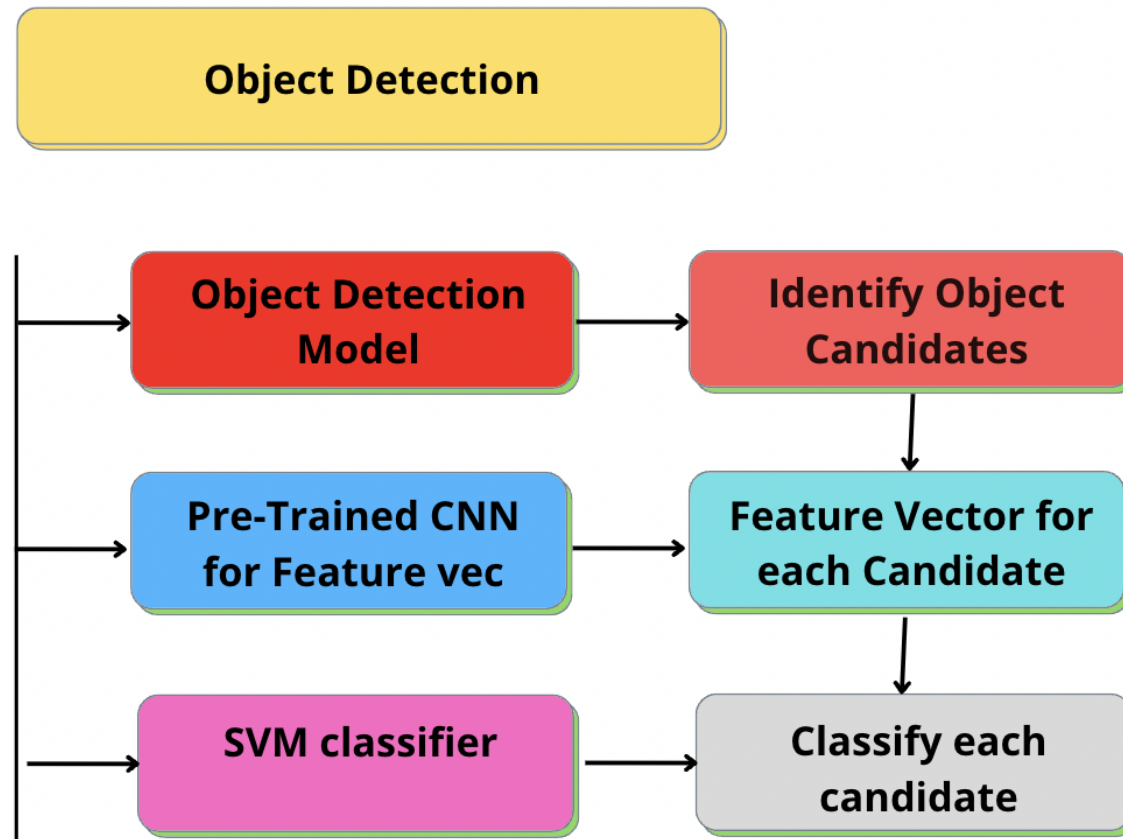
Multi-class Classification vs Object Detection



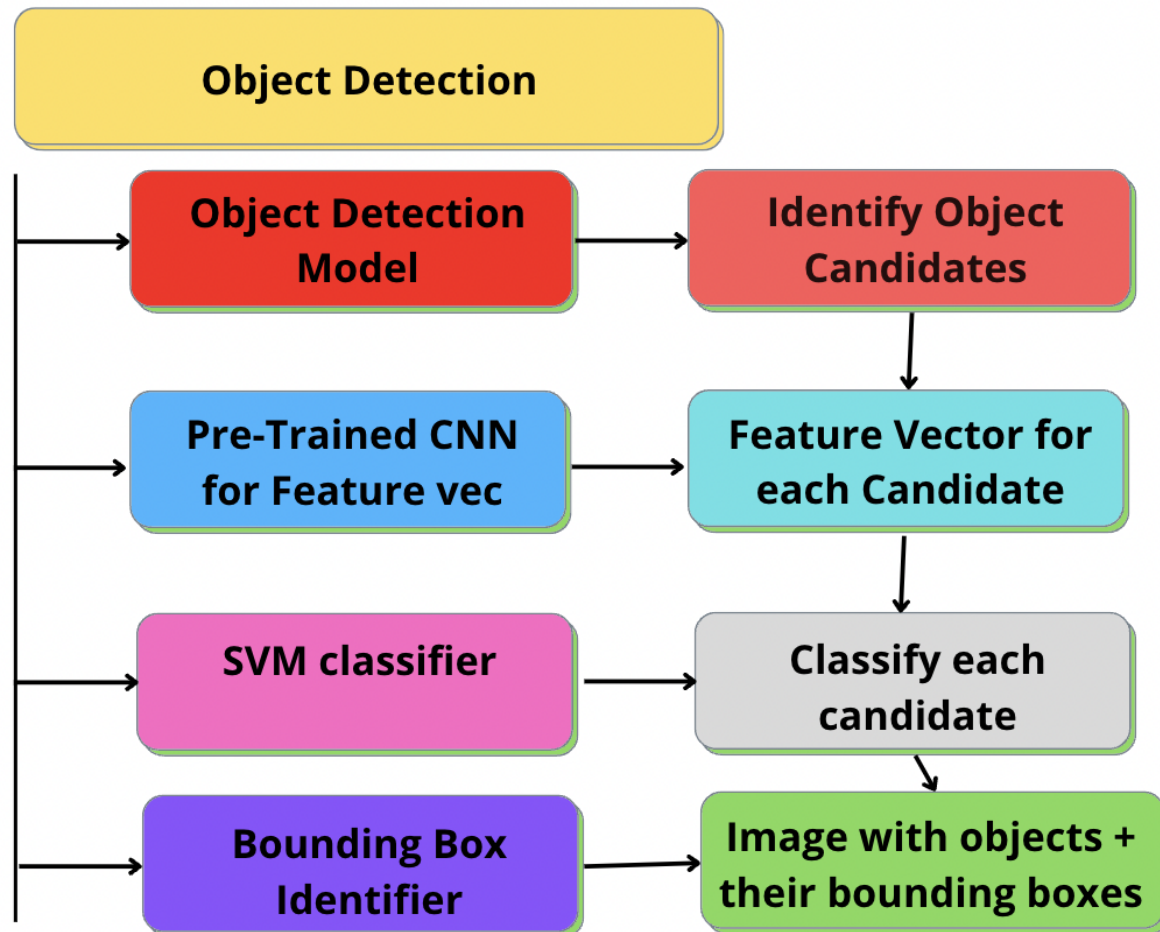
Object Detection Intuition



Object Detection Model Framework



Object Detection Model Framework

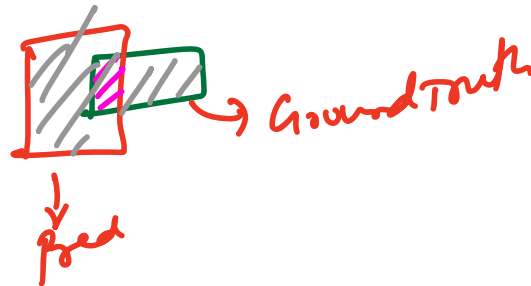


*R-CNN
↳ region*

Metrics for Classification

IOU

Intersection over Union: Is the ratio of the area of the *intersection* between the predicted bounding box and the ground truth bounding box **over** the union of the area between the predicted bounding box and the ground truth bounding box



Metrics for Classification

IOU

Intersection over Union: Is the ratio of the area of the *intersection* between the predicted bounding box and the ground truth bounding box **over** the union of the area between the predicted bounding box and the ground truth bounding box

MAP @ 0.5 IOU

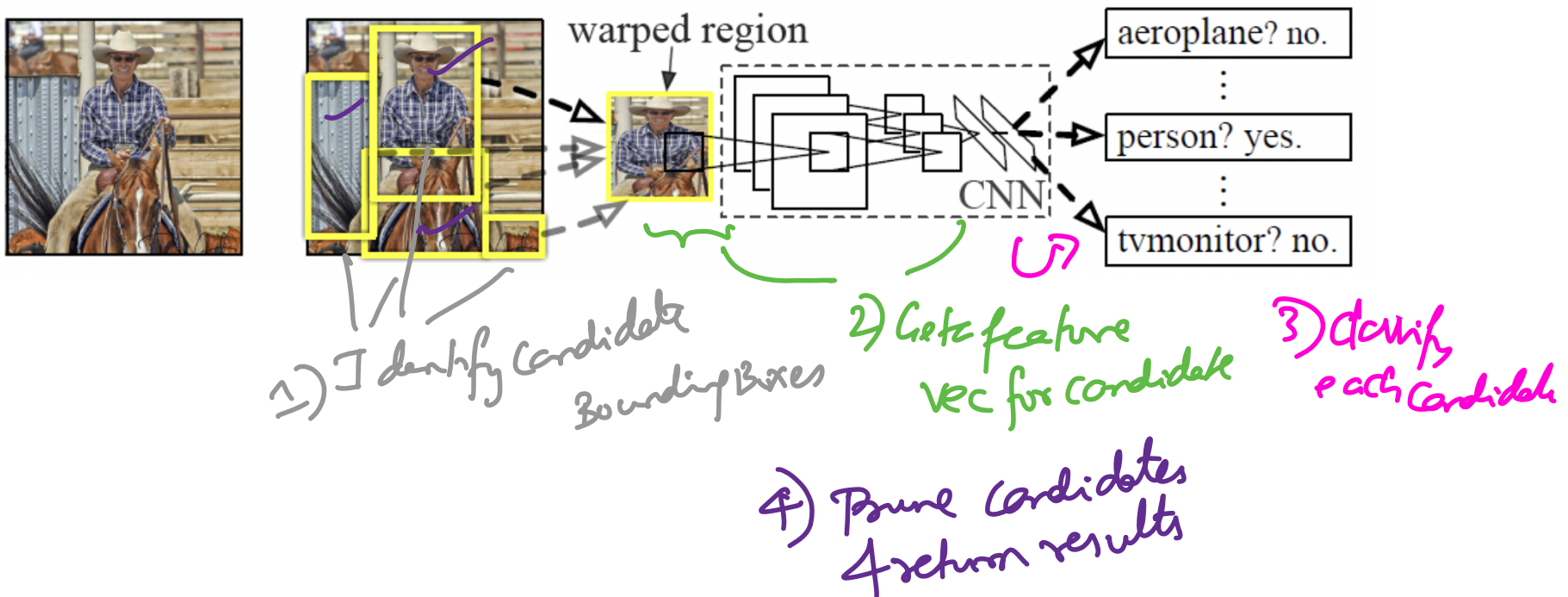
Average Precision (AP) @ 0.5 IOU: If

$$\text{IOU} > \underline{0.5} \rightarrow t \text{ (0.6, 0.7, \dots)}$$

across examples of a given class, count the precision as 1, else 0. Average Precision is the average of all the precisions in a given class.

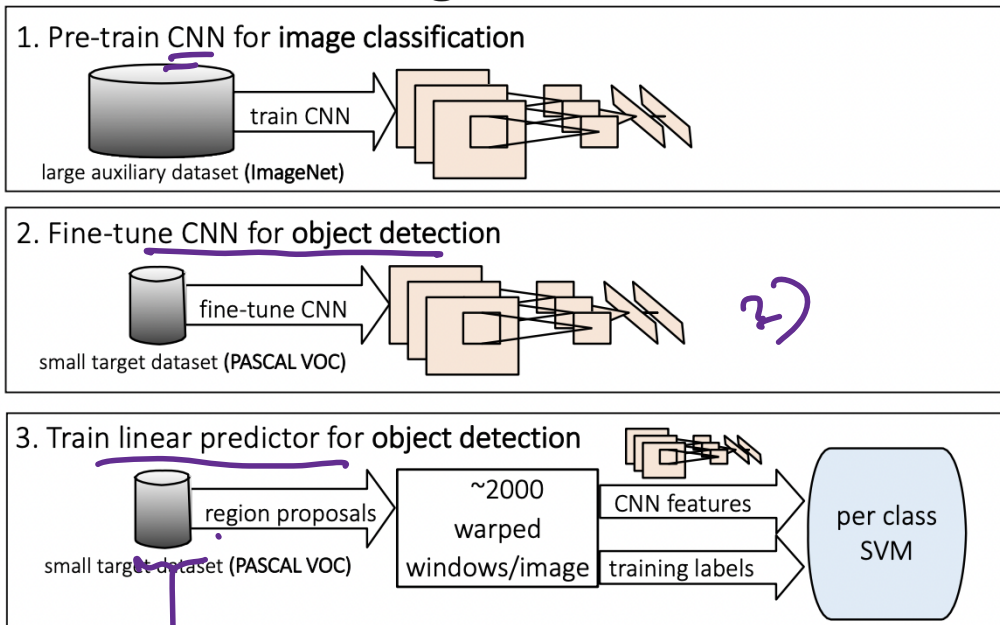
First CNN Model for Object Detection: R-CNN model

Region of Interest



First CNN Model for Object Detection: R-CNN model

R-CNN: Training

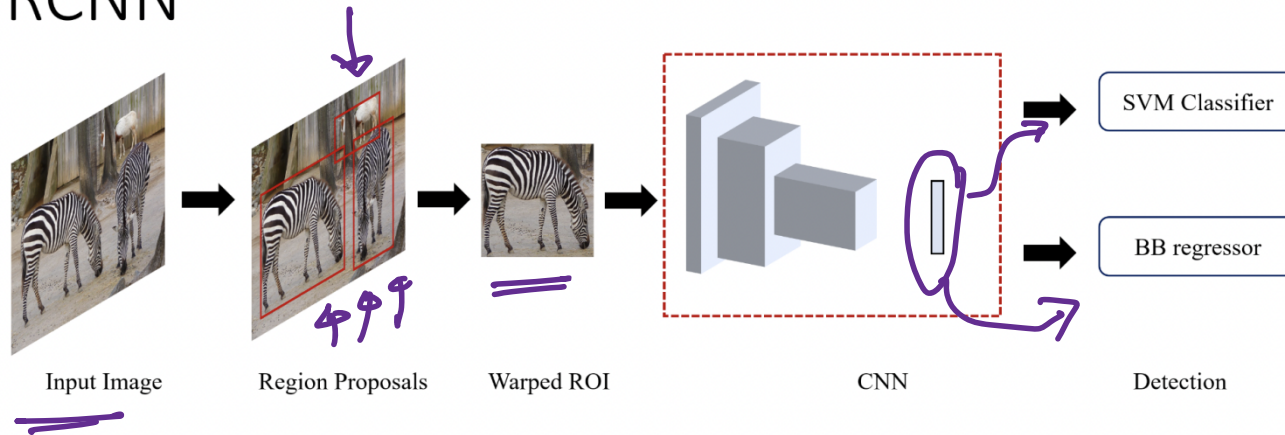


1) Region Proposal Model

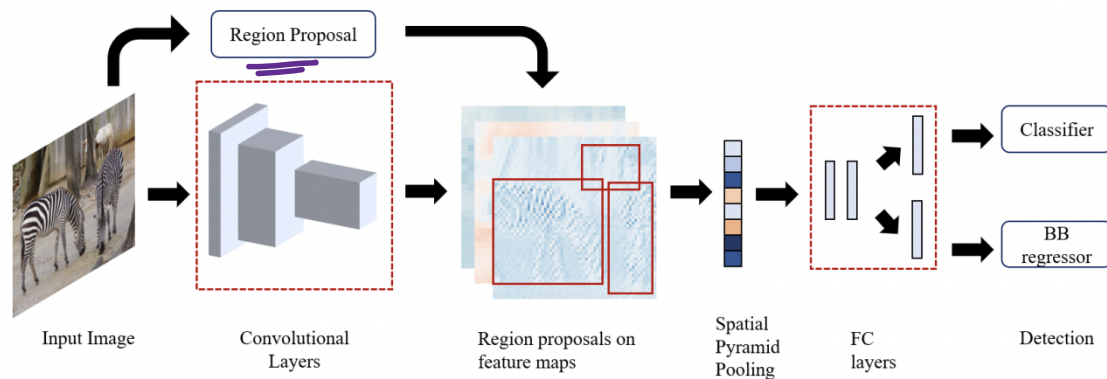
3)

R-CNN variant

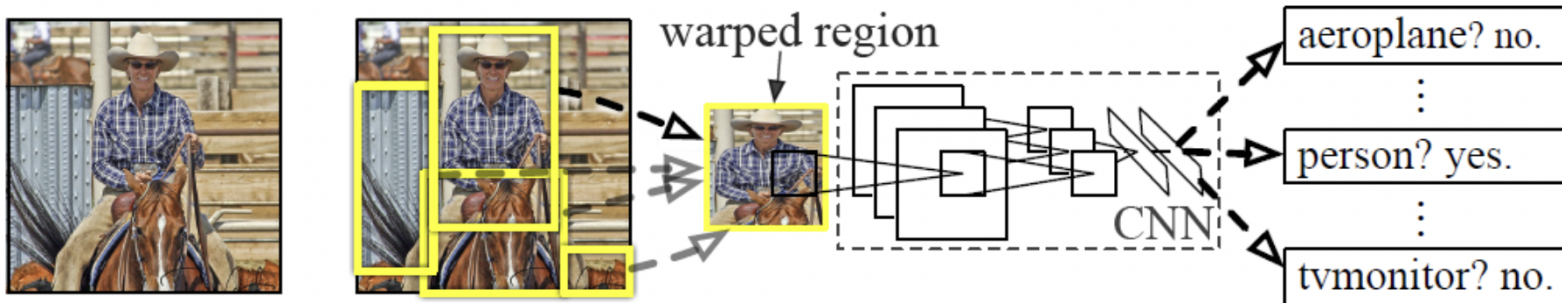
RCNN



Fast RCNN

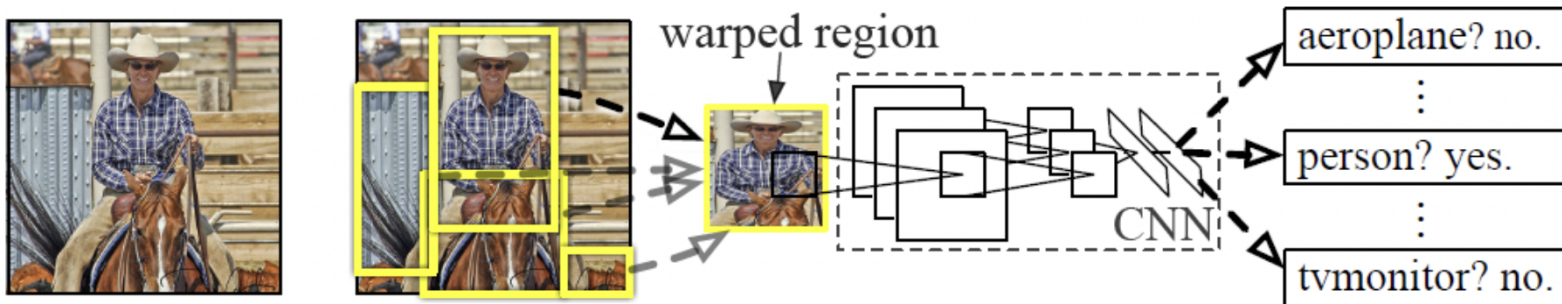


Drawbacks of Two Stage Detection



- 1 **R-CNN variants:** a) Generate Bounding Boxes b) Classify these boxes and c) Merge bounding boxes to eliminate duplicates - 3 steps, 3 models to train!

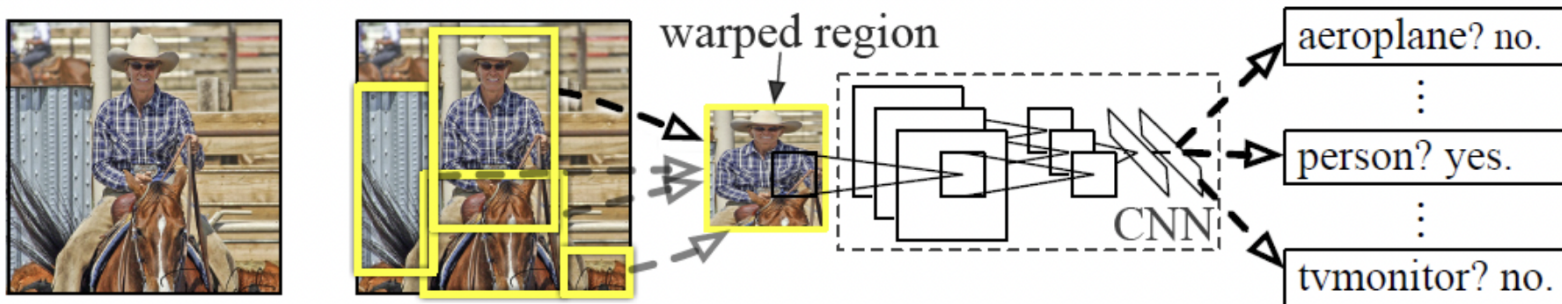
Drawbacks of Two Stage Detection



- 1 **R-CNN variants:** a) Generate Bounding Boxes b) Classify these boxes and c) Merge bounding boxes to eliminate duplicates - 3 steps, 3 models to train!
- 2 **Inference Time Taken:** FPS can be low - E.g. even Faster R-CNN has a speed of (5) FPS at inference time

↳ slow for video processing!

Drawbacks of Two Stage Detection



- 1 **R-CNN variants:** a) Generate Bounding Boxes b) Classify these boxes and c) Merge bounding boxes to eliminate duplicates - 3 steps, 3 models to train!
- 2 **Inference Time Taken:** FPS can be low - E.g. even Faster R-CNN has a speed of 5 FPS at inference time
- 3 **Train Time:** 3 separate models implies more time to train

Single Stage Detection

Simplify!

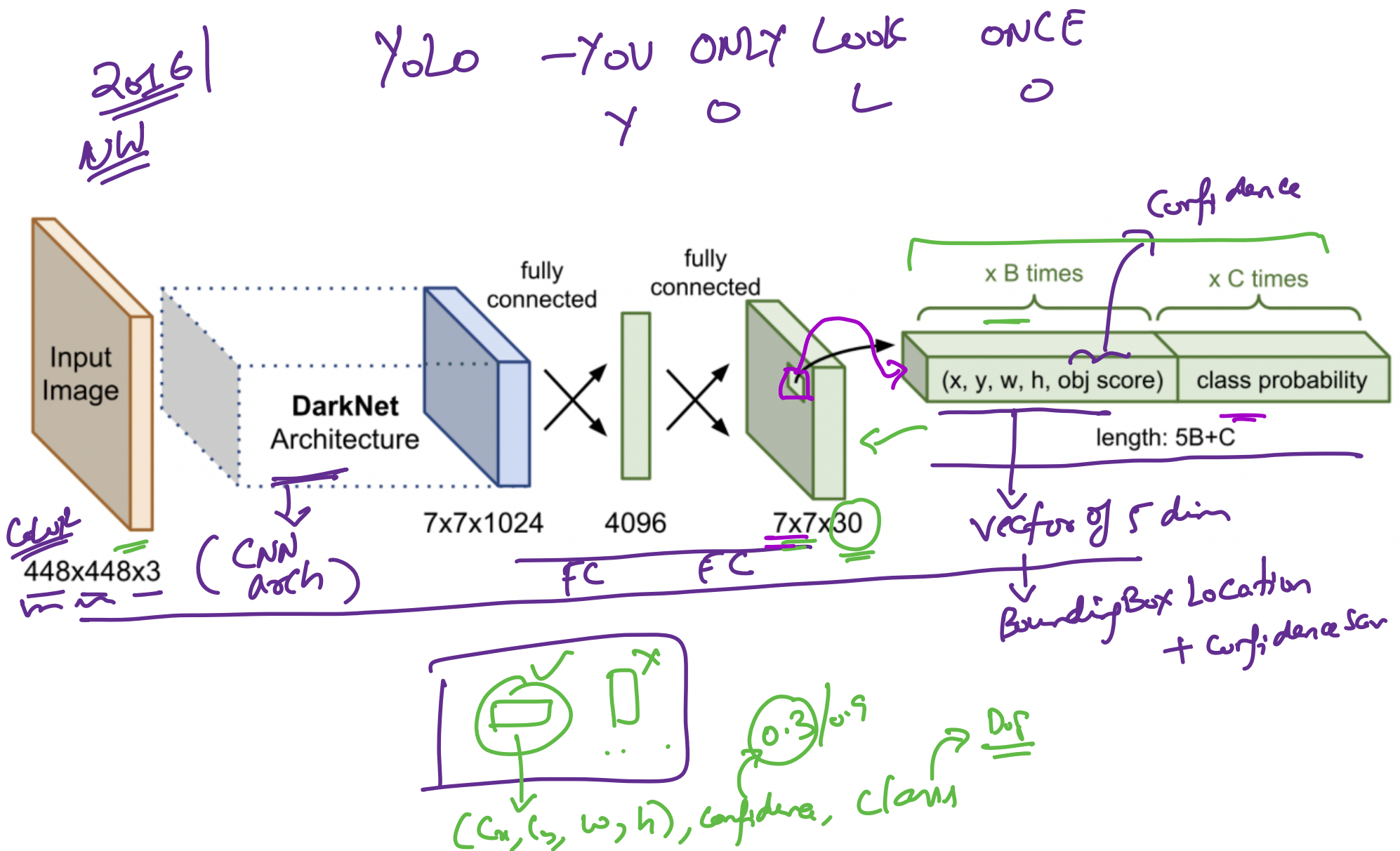
1 Model

What if we could simplify the detection process and also make training and inference more efficient in the process? Enter **YOLO**

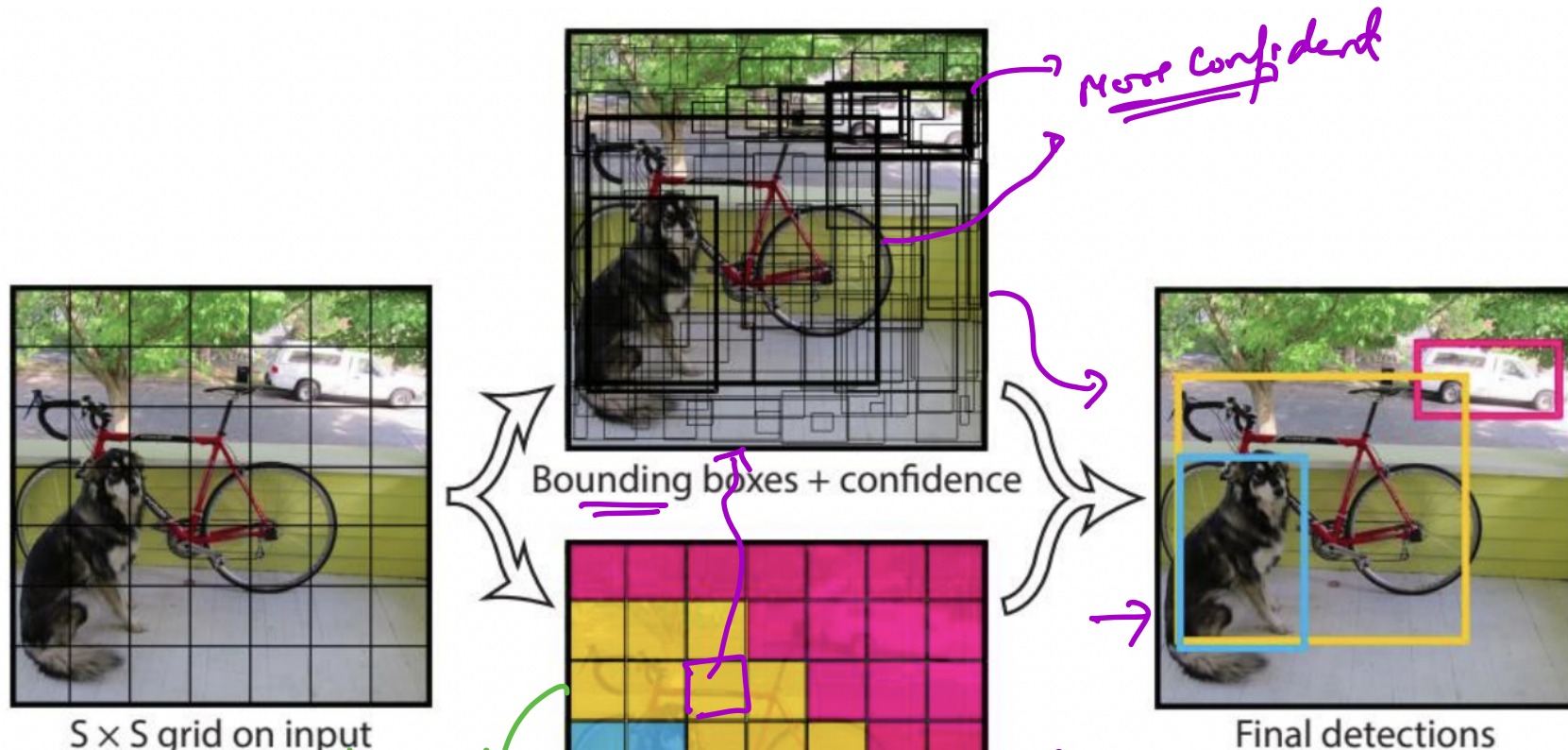
Today

- ① Object Detection Recap
- ② R-CNN variants - Fast and Faster R-CNN
- ③ **YOLO - Single Stage Object Detection** ✓
- ④ Results and Benchmarking on the data sets

YOLO - Single Stage Detection



YOLO Breakdown



S × S grid on input
S = 7 (default)

Bounding boxes + confidence

Class probability map

Final detections

More confident

car

Bicycle

Dog

7x7 = 49 grids!!!

B = 3 ⇒ Total candidate Bounding Boxes = 7x7x3 = 147

YOLO v3 Video

YOLO Real-time Detection



ICE #1


YOLO breakdown

YOLO implicitly divides the input into 7×7 regions and makes bounding box predictions for each of the 49 grids. How would this be different from 'explicitly' breaking up the input into 49 grids and passing each grid through an object detection method?

- 1 It would be the same thing and give similar results
- 2 It would be different
- 3 Surrounding context doesn't get encoded in the latter and impacts accuracy of object detection
- 4 The explicit breakup would yield higher accuracy on the IOU metric as compared to the implicit breakup

YOLO benefits

① **Single pass**

 Less chunky

YOLO benefits

- ① **Single pass**
- ② **Fast**

YOLO benefits

① Single pass

② Fast

③ Global reasoning

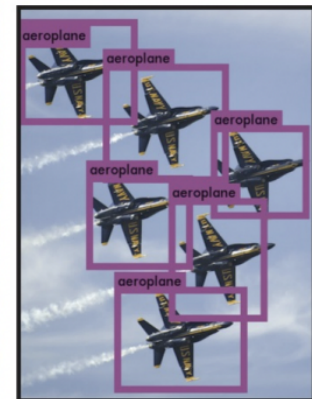
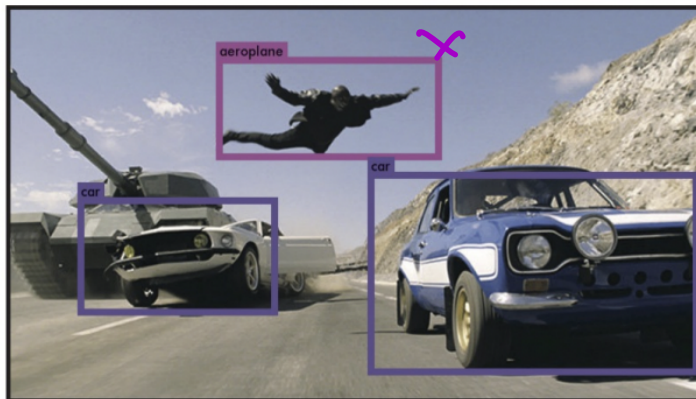
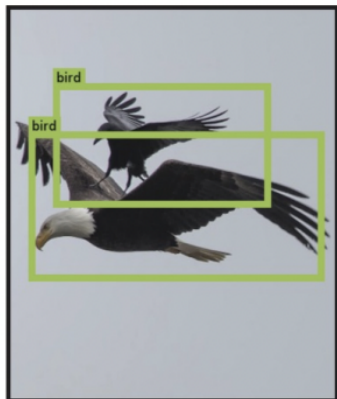
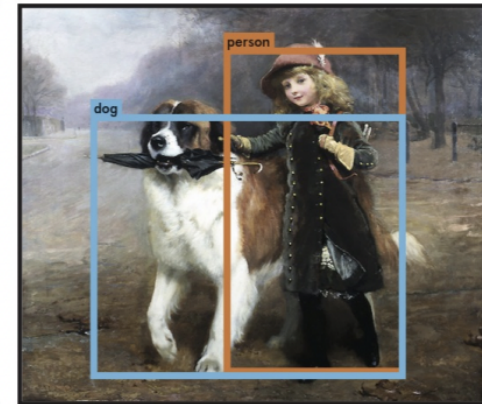
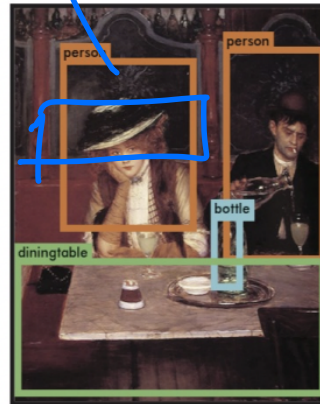
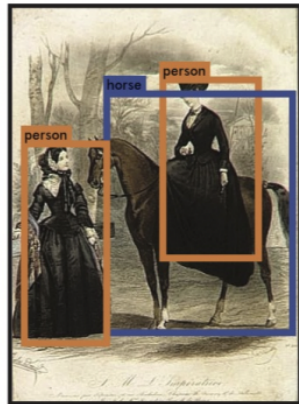
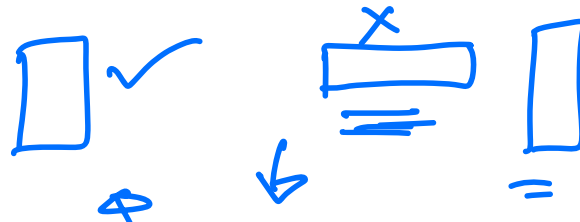
*vs (Local reasoning)
in R-CNN*

(Keeps context of surrounding grids inform object detection)

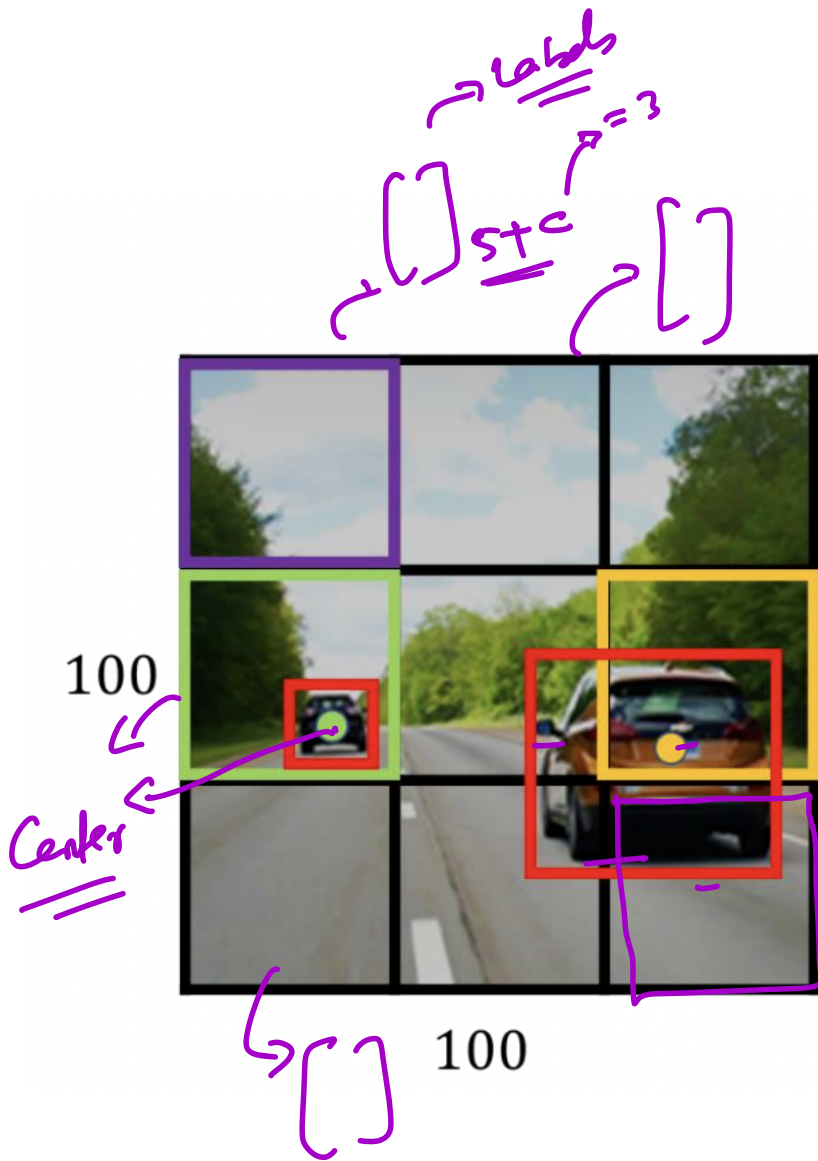
YOLO benefits

- ① **Single pass**
- ② **Fast**
- ③ **Global reasoning**
- ④ **More generalized representations**

YOLO examples



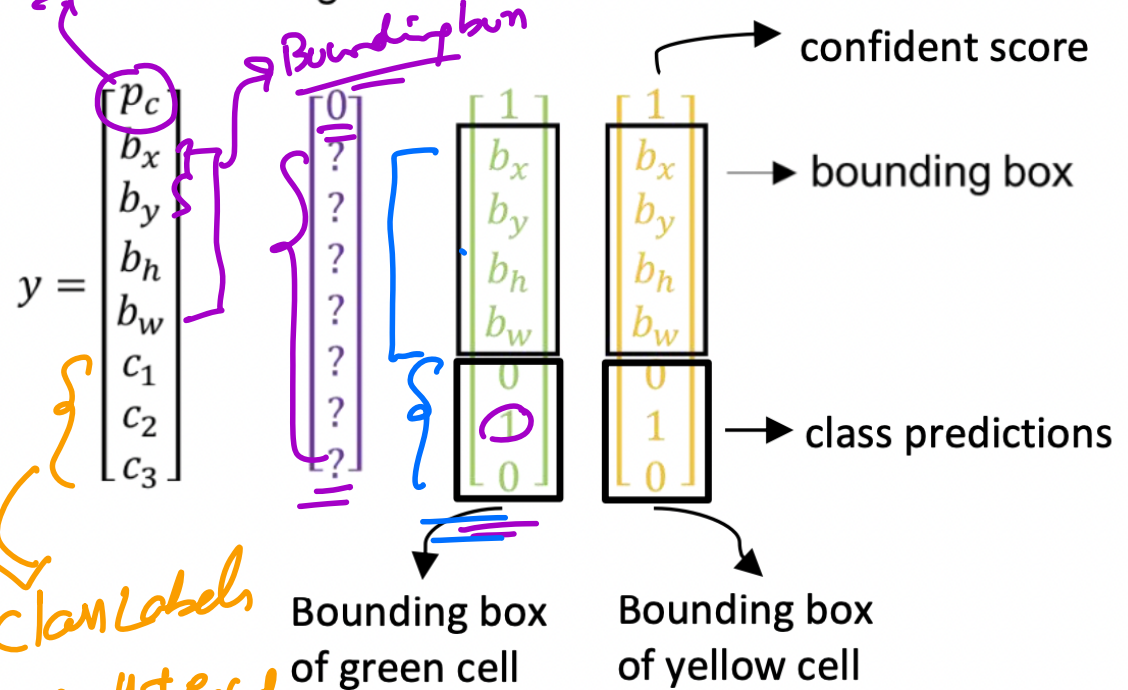
YOLO labels



$$49 \times 8 = \underline{\underline{392}}$$

Confidence

Labels for training for each grid cell:



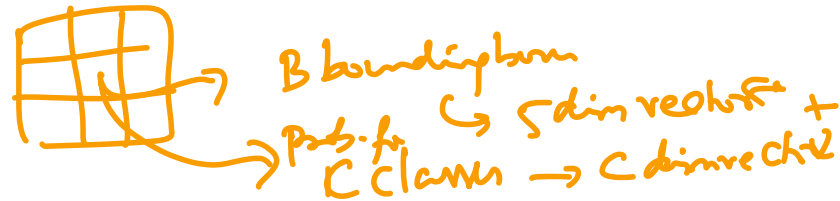
class labels

- one Hot encoding

only one component is non-zero

ICE #2

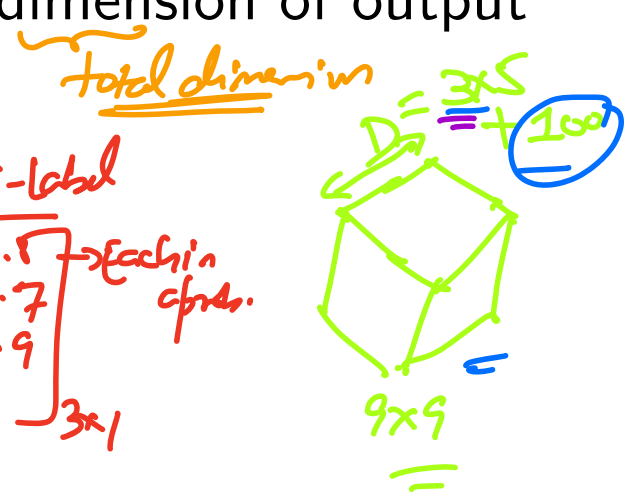
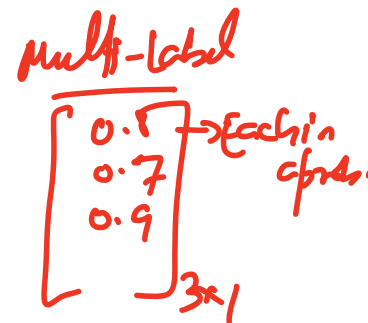
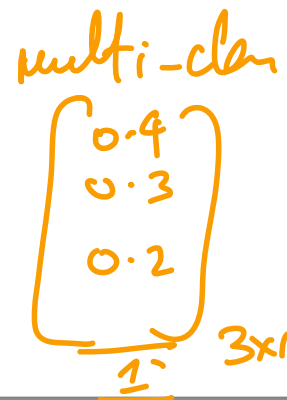
$s \rightarrow [c_1, c_2, w, h, \text{confidence}]_{5 \times 1}$



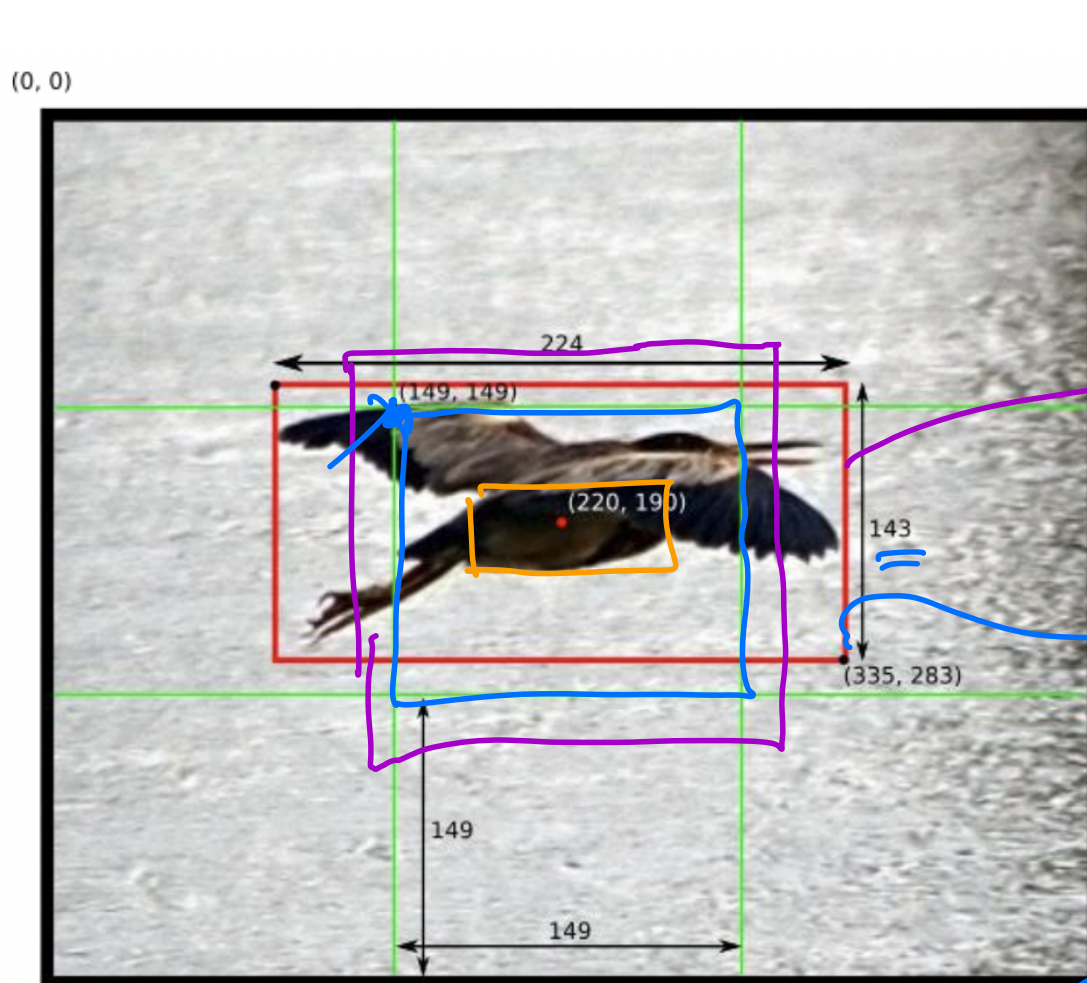
Label Dimensions (2 mins)

Consider a data set of Images, let's say a subset of the MS Coco data set. Let's look at 3 different but related ML problem formulations on this data set: Multi-Class Classification, Multi-Label Classification, and Object Detection (YOLO style). Let's say there are 100 classes to pick from. For object detection, assume that YOLO uses a 9x9 grid with each grid producing 3 different bounding boxes. What's the dimension of output vector for these 3 problem formulations?

- 1 100, 300, and 8100
- 2 100, 100, and 9315
- 3 100, 8100, and 9315
- 4 100, 100 and 8100



YOLO Bounding Box



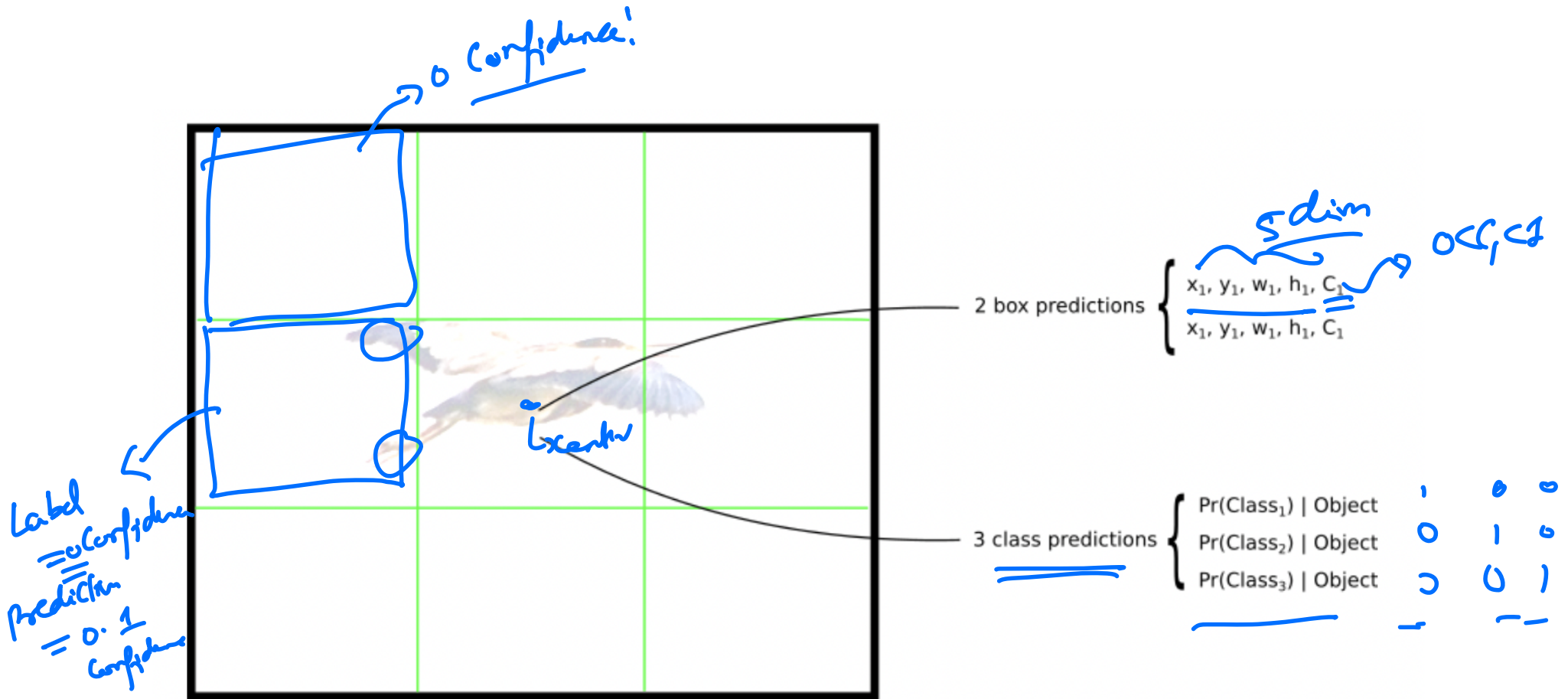
Blue Grid \Rightarrow Location + 5dim vector
 \Rightarrow 3x 5dim vector \Rightarrow Bounding Box

$$x = (220 - 149) / 149 = 0.48$$
$$y = (190 - 149) / 149 = 0.28$$
$$w = 224 / 448 = 0.50$$
$$h = 143 / 448 = 0.32$$

$0 < x < 1$
 $0 < y < 1$
 $0 < w < 1$
 $0 < h < 1$

Key pre-processing step: - normalize x, y, w, h (Labels)

YOLO Bounding Box 2



YOLO and Regression

YOLO Loss Function - Regression!

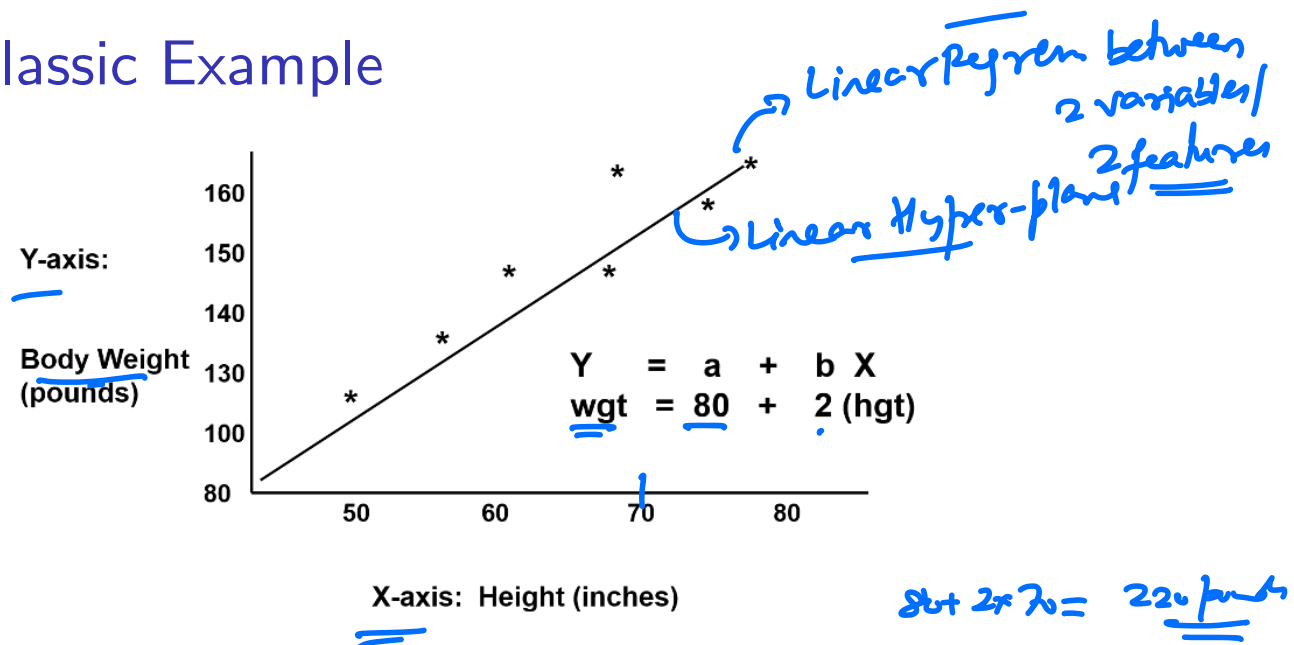
YOLO loss function turns out to be just like a Regression Loss! Why Regression?

YOLO and Regression

YOLO Loss Function - Regression!

YOLO loss function turns out to be just like a Regression Loss! Why Regression?

Linear Regression Classic Example



ICE #3

Regression (2 mins)

You want to predict the ‘sharpness’ of an image when the input is an image. Sharpness for this exercise is defined on a continuous scale between 0 and 1. The training data looks like $\{Image, Sharpness\}$ where Image is the input and Sharpness (on a continuous scale) is the output. You devise an ingenious loss function as follows: Take the prediction \hat{y}_i of the sharpness, subtracts it from the ground truth sharpness y_i , and obtain the error, e_i . Define the loss, $L = \sum_i e_i$. You then minimize the loss as you hope a good model for sharpness would give zero errors and hence a close to zero loss. Optimizing the loss function:

- 1 Will help you train a good model for sharpness
- 2 Is a good idea but may have to watch out for overfitting
- 3 Would not be a good idea
- 4 Could result in a model with overall zero error but poor individual predictions

YOLO Loss Function

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

Today

- ① Object Detection Recap
- ② R-CNN variants - Fast and Faster R-CNN
- ③ YOLO - Single Stage Object Detection
- ④ **Results and Benchmarking on the data sets**

Data Sets

Dataset	Classes	Train			Validation			Test
		Images	Objects	Objects/Image	Images	Objects	Objects/Image	
PASCAL VOC 12	20	5,717	13,609	2.38	5,823	13,841	2.37	10,991
MS-COCO	80	118,287	860,001	7.27	5,000	36,781	7.35	40,670
ILSVRC	200	456,567	478,807	1.05	20,121	55,501	2.76	40,152
OpenImage	600	1,743,042	14,610,229	8.38	41,620	204,621	4.92	125,436

Results for Object Detection

Model	Year	Backbone	Size	AP _[0.5:0.95]	AP _{0.5}	FPS
R-CNN*	2014	AlexNet	224	-	58.50%	~0.02
SPP-Net*	2015	ZF-5	Variable	-	59.20%	~0.23
Fast R-CNN*	2015	VGG-16	Variable	-	65.70%	~0.43
Faster R-CNN*	2016	VGG-16	600	-	67.00%	5
R-FCN	2016	ResNet-101	600	31.50%	53.20%	~3
FPN	2017	ResNet-101	800	36.20%	59.10%	5
Mask R-CNN	2018	ResNeXt-101-FPN	800	39.80%	62.30%	5
DetectoRS	2020	ResNeXt-101	1333	53.30%	71.60%	~4
YOLO*	2015	(Modified) GoogLeNet	448	-	57.90%	45
SSD	2016	VGG-16	300	23.20%	41.20%	48
YOLOv2	2016	DarkNet-19	352	21.60%	44.00%	81
RetinaNet	2018	ResNet-101-FPN	400	31.90%	49.50%	12
YOLOv3	2018	DarkNet-53	320	28.20%	51.50%	45
CenterNet	2019	Hourglass-104	512	42.10%	61.10%	7.8
EfficientDet-D2	2020	Efficient-B2	768	43.00%	62.30%	41.7
YOLOv4	2020	CSPDarkNet-53	512	43.00%	64.90%	31
Swin-L	2021	HTC++	-	57.70%	-	-

^aModels marked with * are compared on PASCAL VOC 2012, while others on MS COCO. Rows colored gray are real-time detectors (>30 FPS).

vs! → 2020

Breakout for Takeaways!

Discuss Takeaways (5 mins)

From today's lecture in your zoom group

Next Lecture

- ① Newer Variants of YOLO
- ② Object Detection vs Instance Segmentation
- ③ Image Captioning Models