

# Object Detection

Tomasz Stańczyk  
[tomasz.stanczyk@inria.fr](mailto:tomasz.stanczyk@inria.fr)

*With slides from Andrew Ng and other sources (referenced)*



# Today's Agenda

Object detection fundamentals - based on DeepLearningAI materials by Andrew Ng

+ references for more information/self-study if desired

YOLOX object detection algorithm

+ references for more information/self-study if desired

Later today: YOLOX-based practical assignment



DeepLearning.AI



LearnOpenCV.com

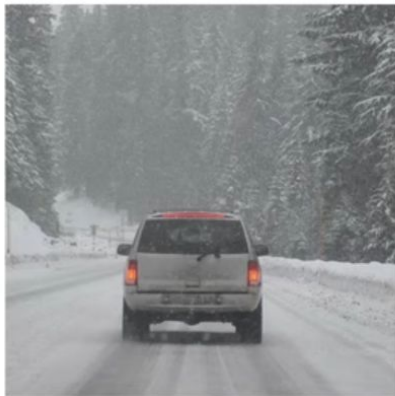
EXPERT AI BLOG

# Object Detection Fundamentals

Selected slides by Andrew Ng

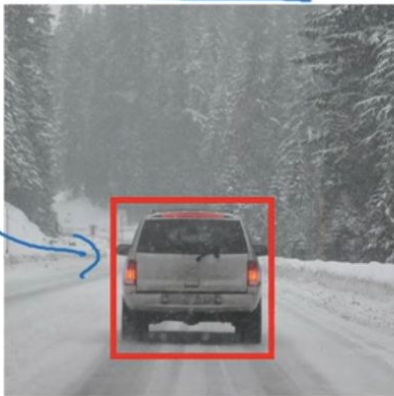
# What are localization and detection?

Image classification



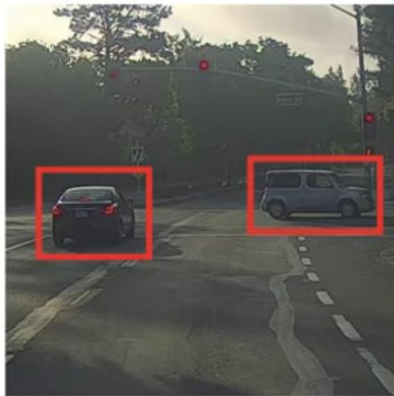
"Car"

Classification with localization

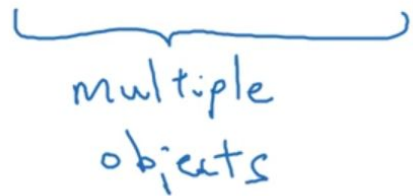


"Car"

Detection



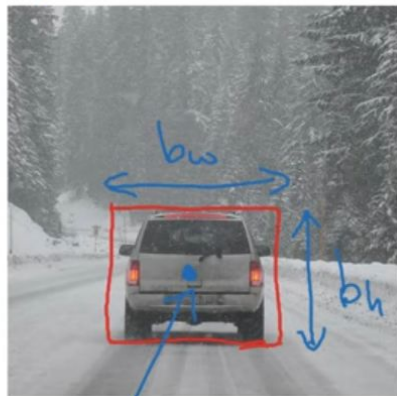
multiple objects



# Classification with localization

Object Localization

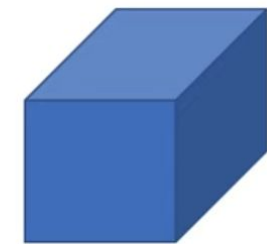
$(0,0)$



$(1,1)$

$b_x, b_y$

- 1 - pedestrian ←
- 2 - car ←
- 3 - motorcycle ←
- 4 - background



ConvNet

...



softmax (4)

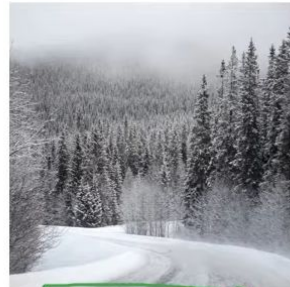
$b_x, b_y, b_h, b_w$   
bounding box

$b_x = 0.5$   
 $b_y = 0.7$   
 $b_h = 0.3$   
 $b_w = 0.4$

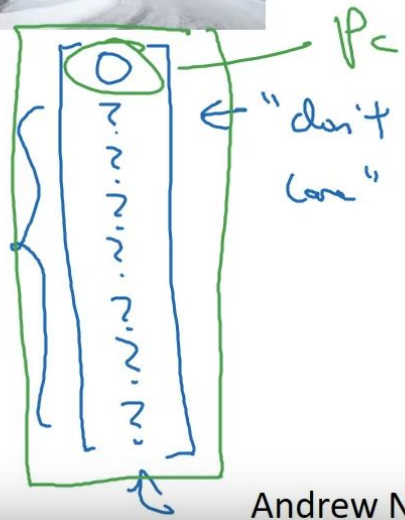
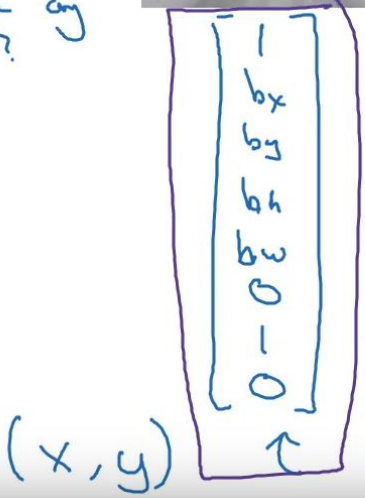
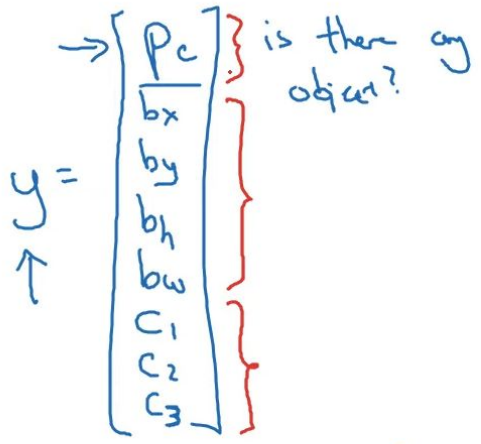
# Defining the target label $y$

- 1 - pedestrian
- 2 - car ←
- 3 - motorcycle
- 4 - background ←

Need to output  $b_x, b_y, b_h, b_w$ , class label (1-4)



$$L(\hat{y}, y) = \begin{cases} (\hat{y}_1 - y_1)^2 + (\hat{y}_2 - y_2)^2 + \dots + (\hat{y}_8 - y_8)^2 & \text{if } y_i = 1 \\ (\hat{y}_1 - y_1)^2 & \text{if } y_i = 0 \end{cases}$$



# Car detection example

Training set:

X

y



1



1



1



0



0



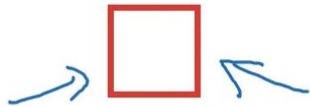
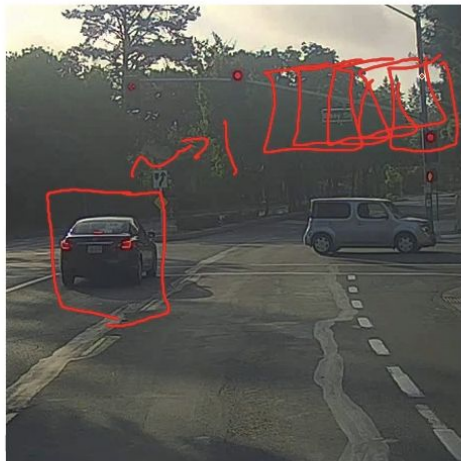
→ ConvNet → y

# Sliding windows detection

Object  
Detection

ConvNet  $\rightarrow$  0

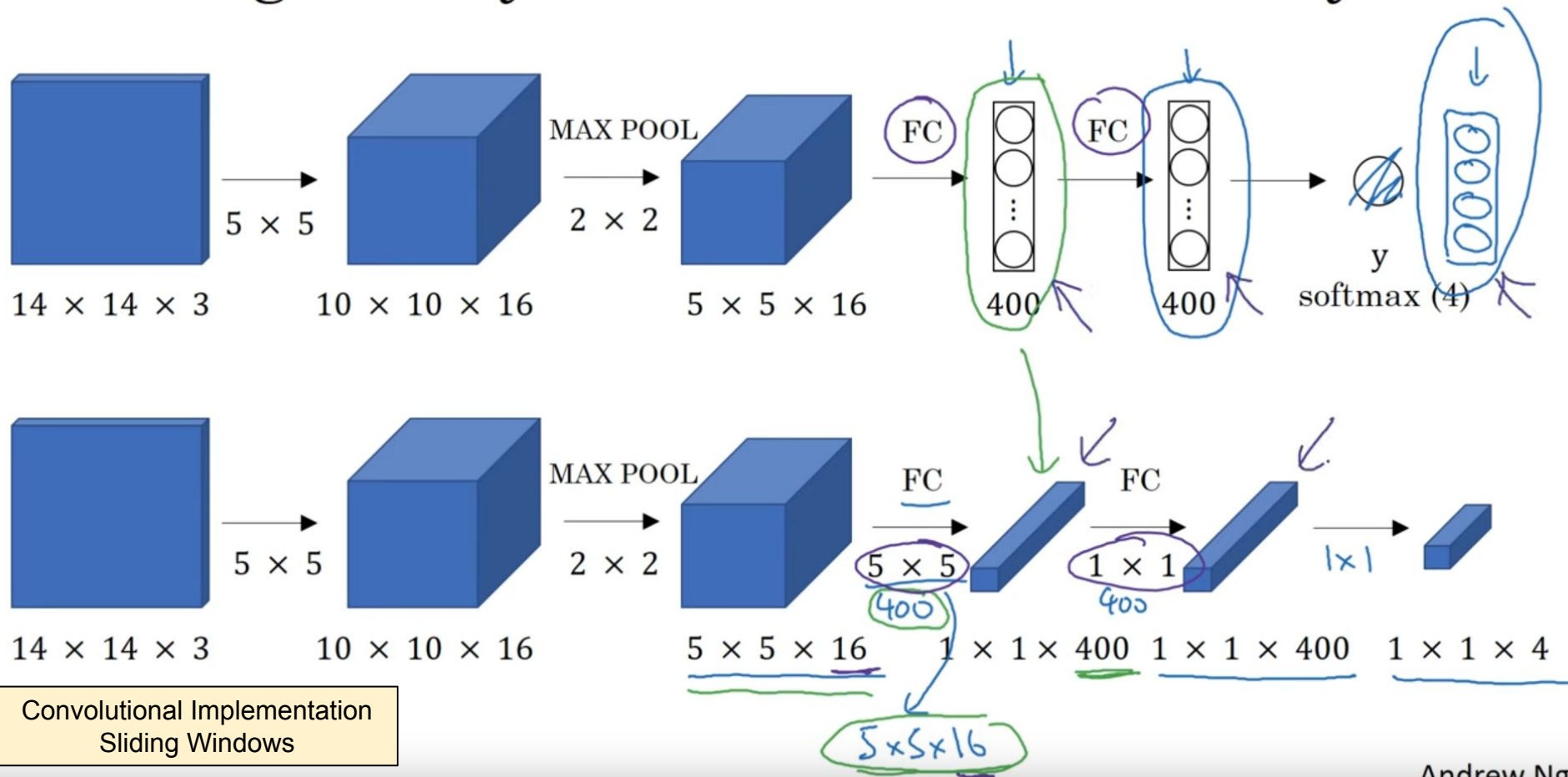
ConvNet



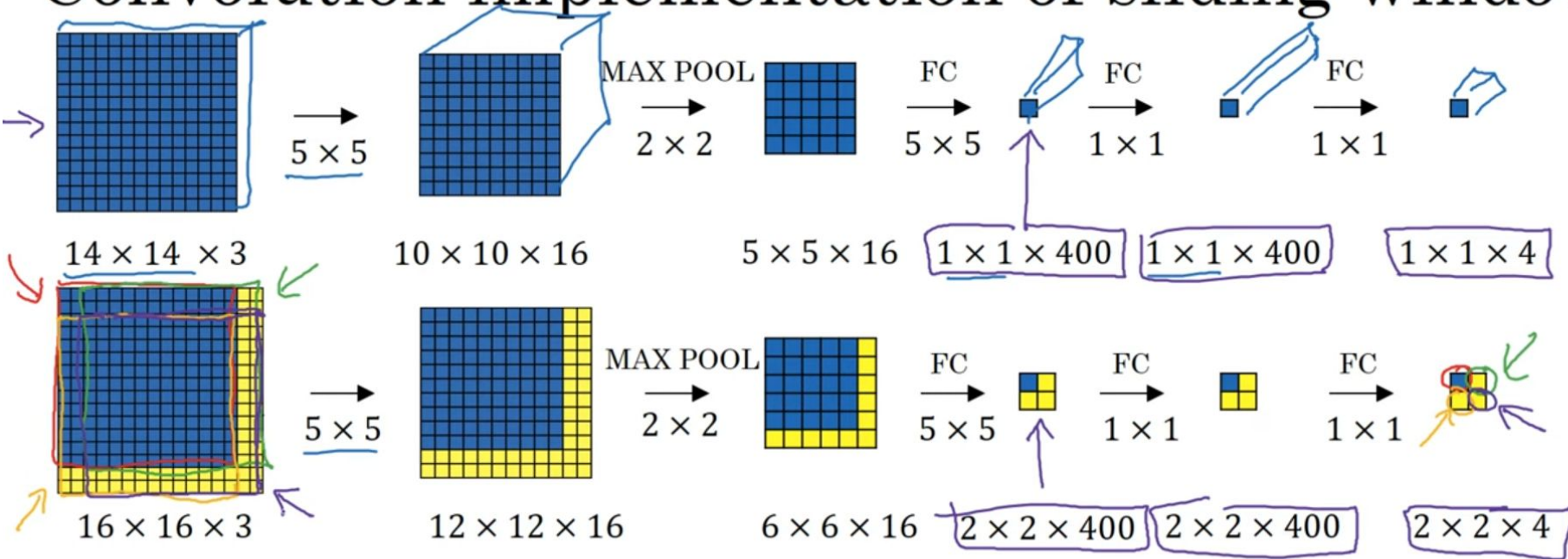
Computation cost



# Turning FC layer into convolutional layers

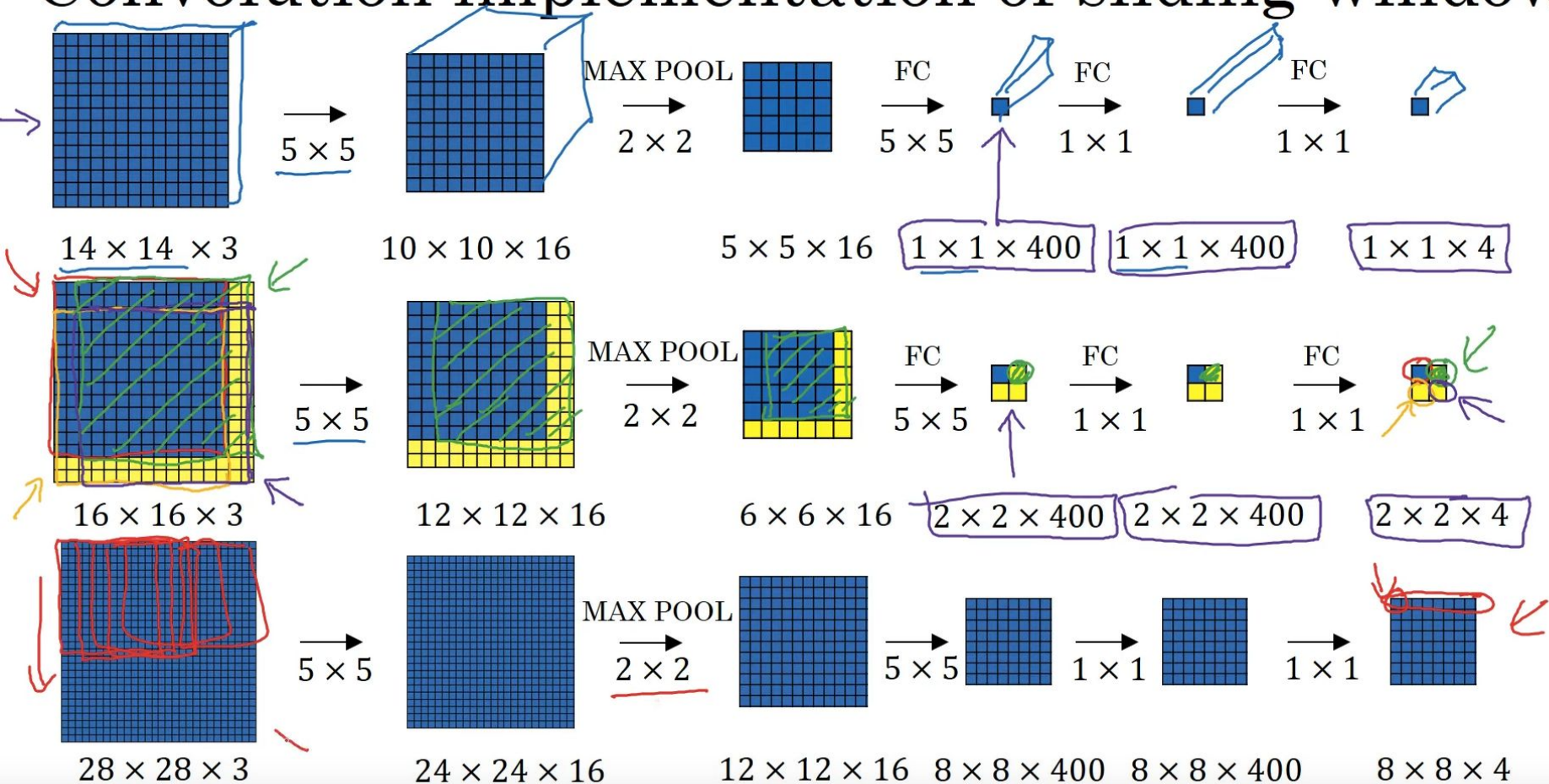


# Convolution implementation of sliding windows

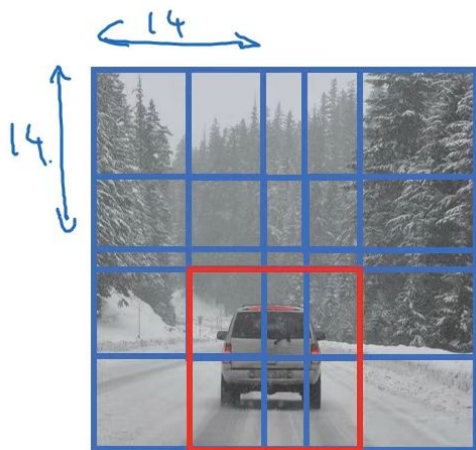
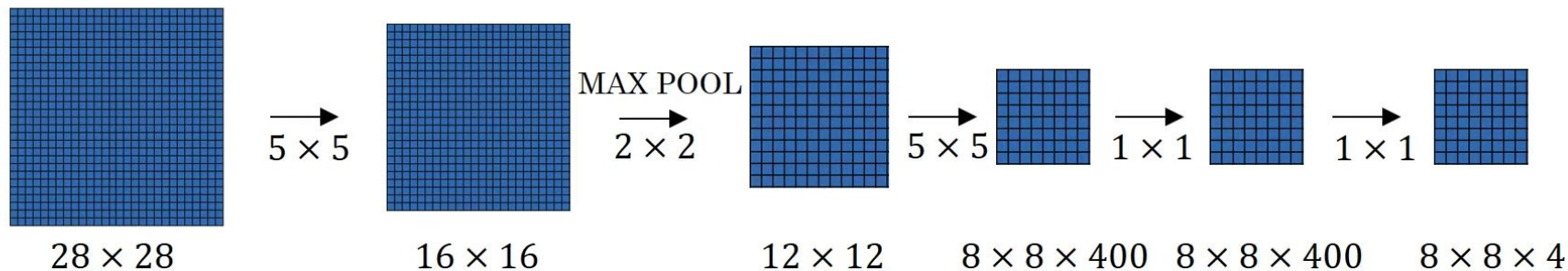


Convolutional Implementation  
Sliding Windows

# Convolution implementation of sliding windows



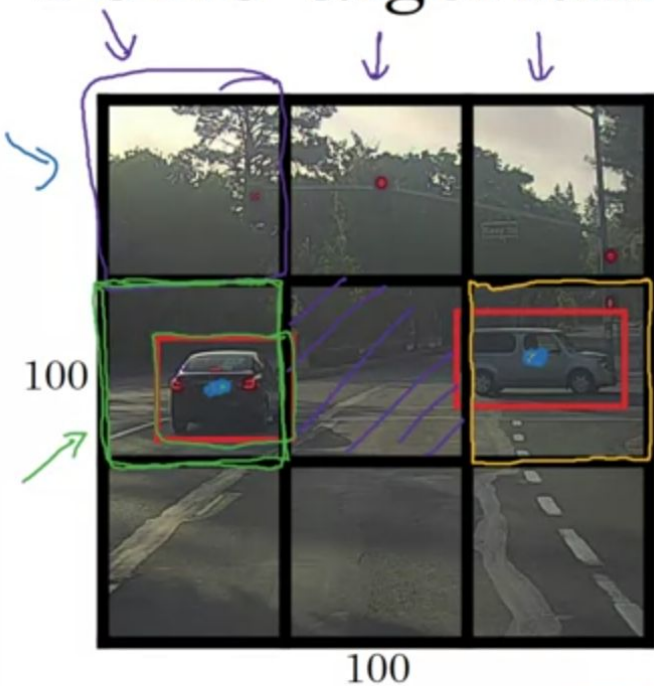
# Convolution implementation of sliding windows



Convolutional Implementation  
Sliding Windows

# YOLO algorithm

Bounding Box Predictions

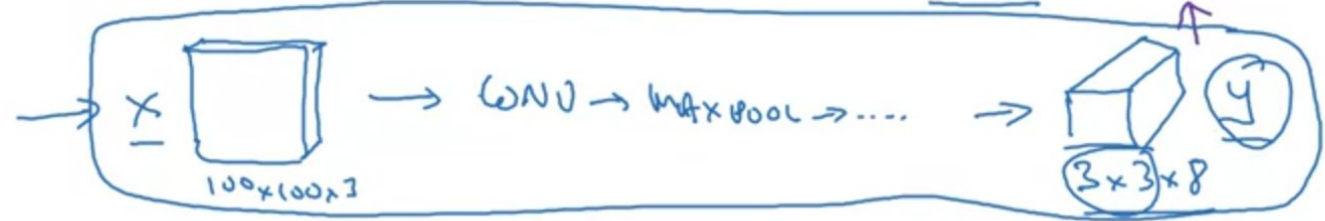
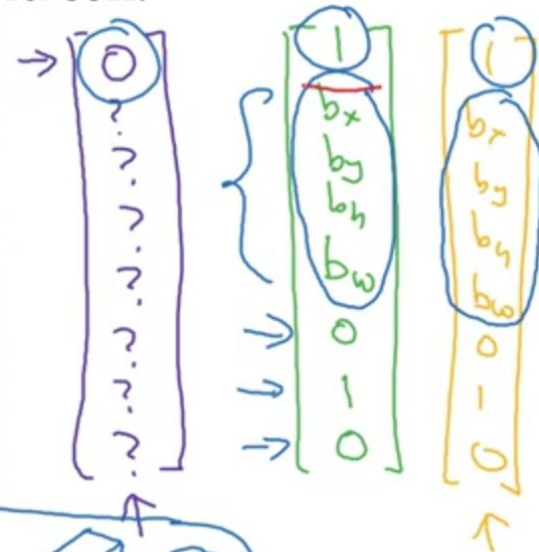


Labels for training  
For each grid cell:

Target output:  
 $3 \times 3 \times 8$

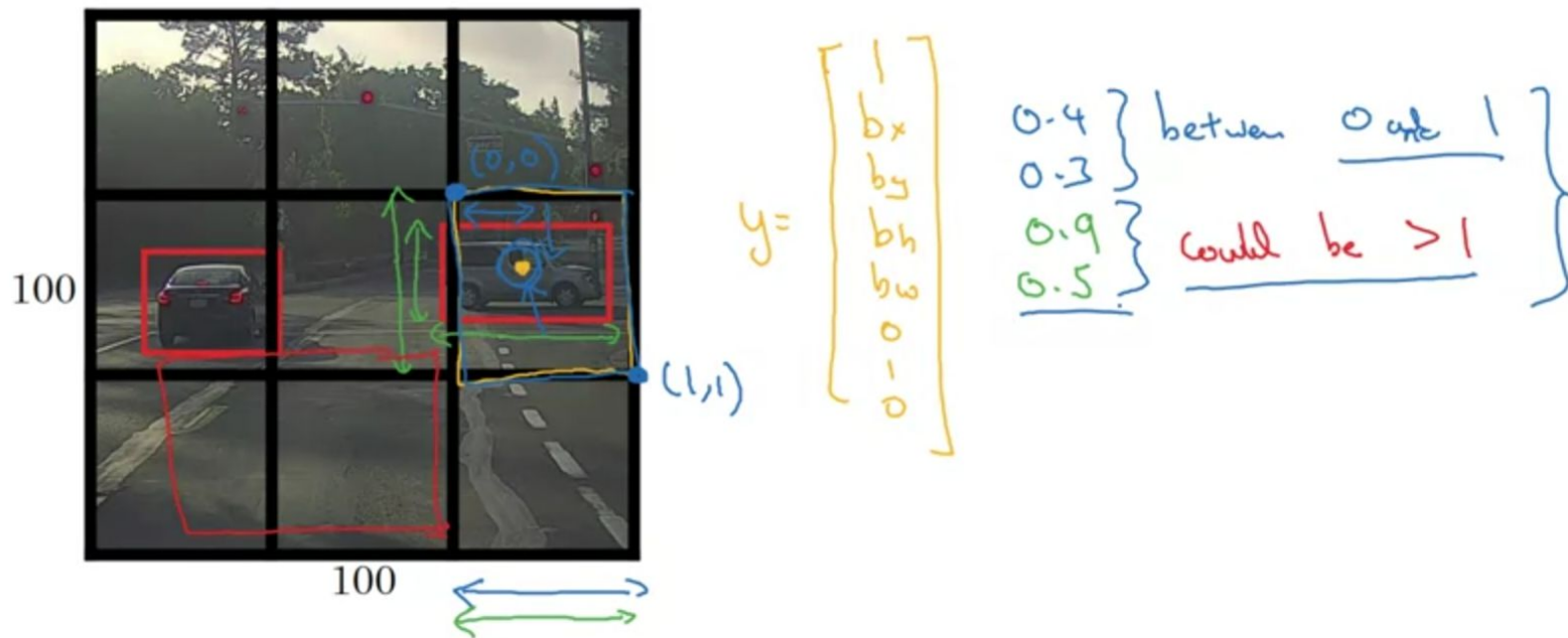


$y = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix}$



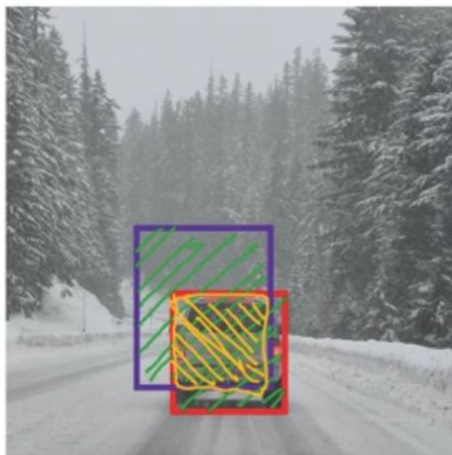
# Specify the bounding boxes

Bounding Box Predictions



# Evaluating object localization

Intersection  
Over Union



Intersection over Union (IoU)

$$= \frac{\text{Size of } \text{[yellow box]}}{\text{Size of } \text{[green box]}}$$

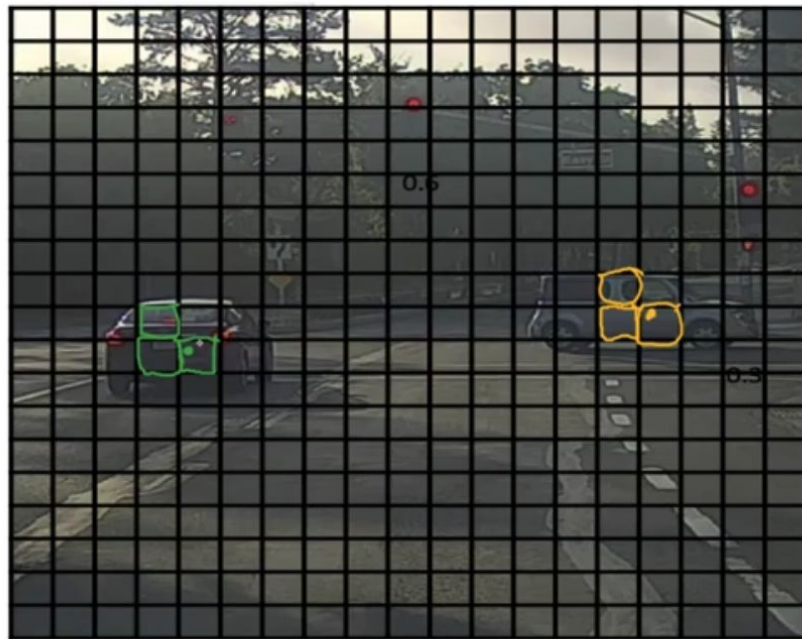
“Correct” if IoU  $\geq$  0.5 ←

0.6 ←

More generally, IoU is a measure of the overlap between two bounding boxes.

# Non-max suppression example

Non-max  
Suppression

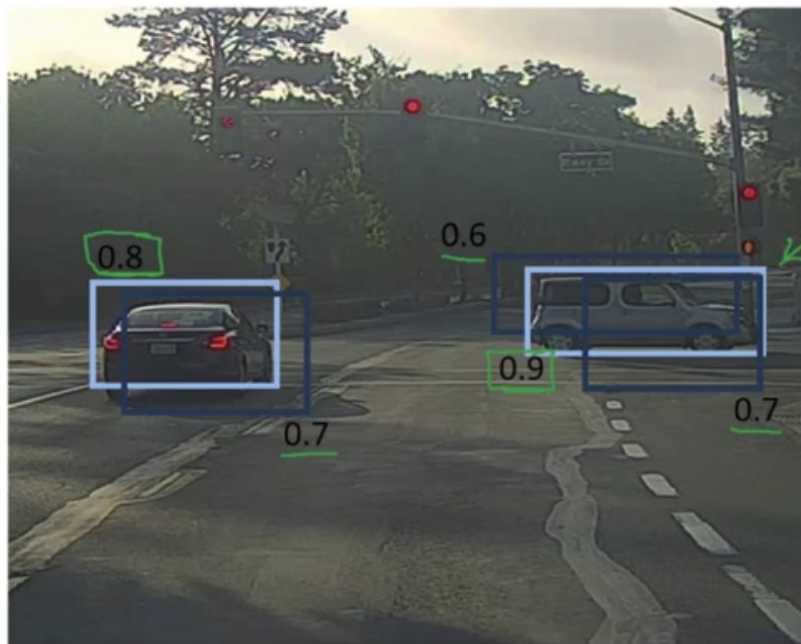


19x19



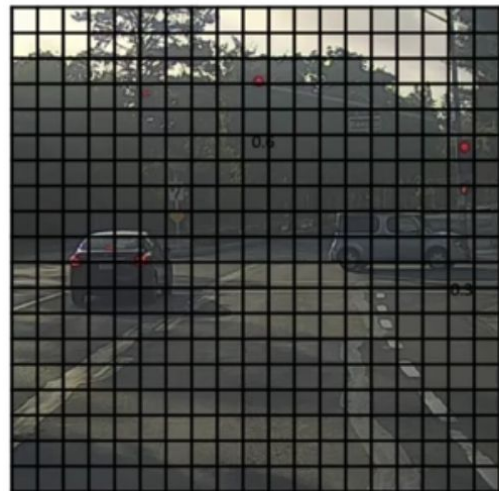
# Non-max suppression example

Non-max  
Suppression



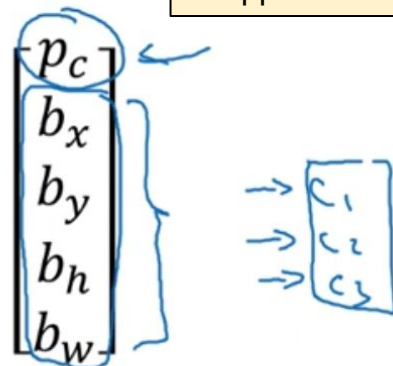
# Non-max suppression algorithm

Non-max  
Suppression



19x 19

Each output prediction is:

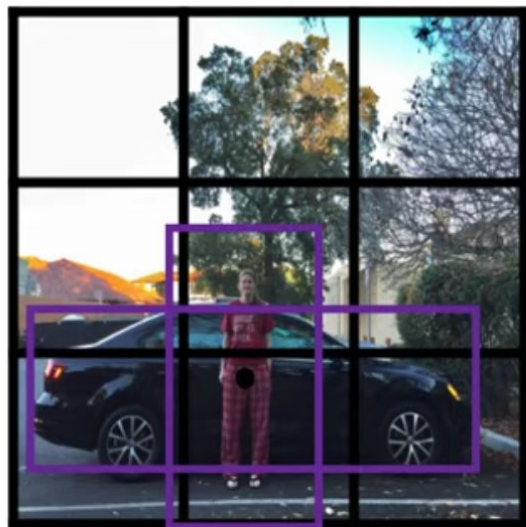


Discard all boxes with  $p_c \leq 0.6$

→ While there are any remaining boxes:

- Pick the box with the largest  $p_c$   
Output that as a prediction.
- Discard any remaining box with  $\text{IoU} \geq 0.5$  with the box output in the previous step

# Overlapping objects:

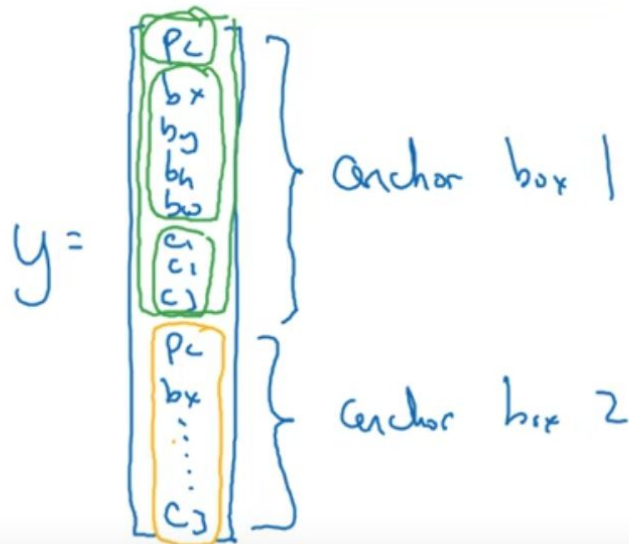


$$y = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix}$$

Anchor box 1:



Anchor box 2:

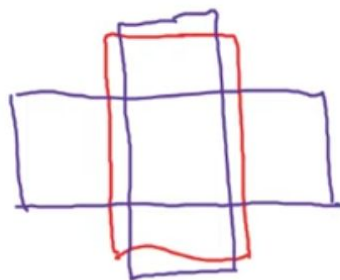


# Anchor box algorithm

Previously:

Each object in training image is assigned to grid cell that contains that object's midpoint.

Output  $y$ :  
 $3 \times 3 \times 8$



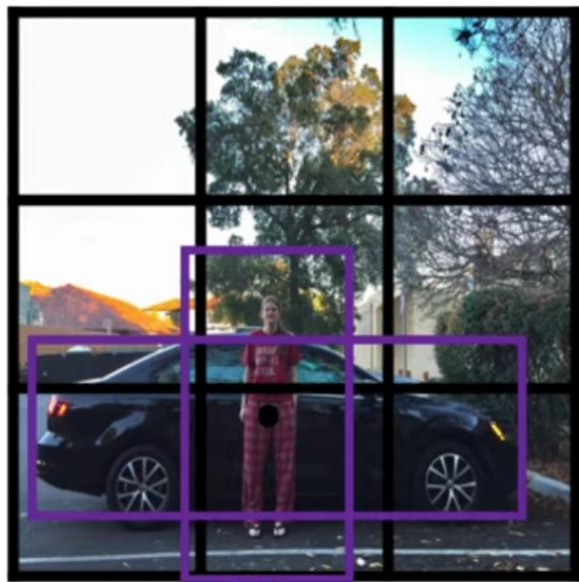
With two anchor boxes:

Each object in training image is assigned to grid cell that contains object's midpoint and anchor box for the grid cell with highest IoU.

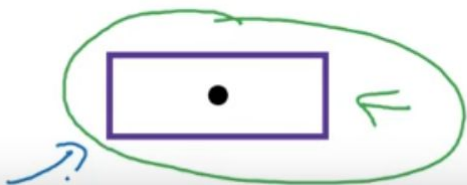
(grid cell, anchor box)

Output  $y$ :  
 $3 \times 3 \times \underline{16}$   
 $3 \times 3 \times \underline{2} \times \underline{8}$

# Anchor box example



Anchor box 1:      Anchor box 2:



$y =$

$$\begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \\ p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix}$$

Handwritten annotations for the first vector  $y$  (orange and green):

$$\begin{bmatrix} 1 \\ b_x \\ b_y \\ b_h \\ b_w \\ 0 \\ 0 \\ 0 \\ 1 \\ b_x \\ b_y \\ b_h \\ b_w \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

Handwritten annotations for the second vector  $y$  (green):

$c_1$  only?

$$\begin{bmatrix} 0 \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ 1 \\ b_x \\ b_y \\ b_h \\ b_w \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

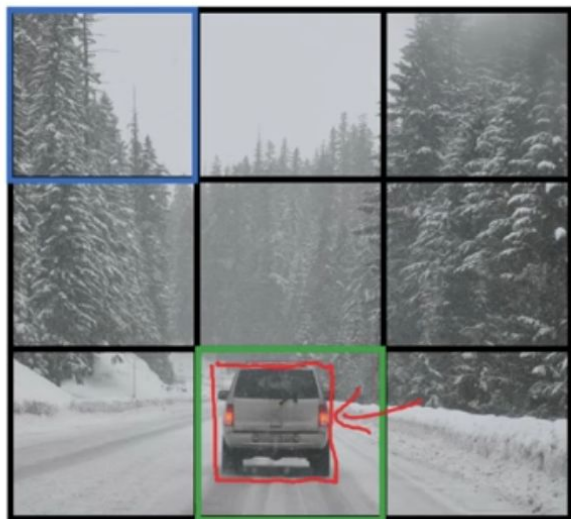
anchor box 1

anchor box 2

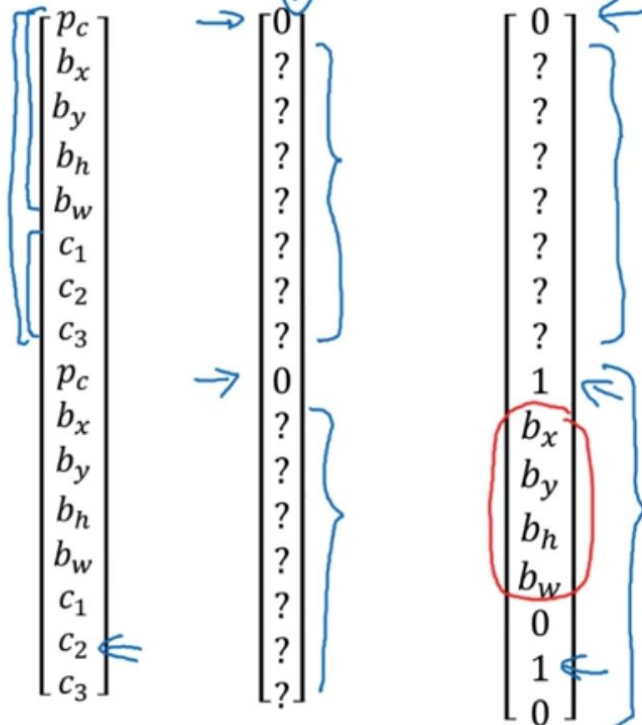
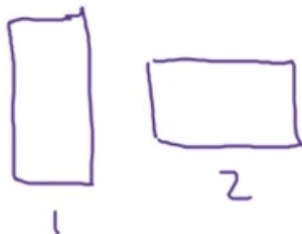
# Training

YOLO Algorithm

- 1 - pedestrian
- 2 - car
- 3 - motorcycle



$$y =$$



$3 \times 3 \times 16$

y is  $3 \times 3 \times 2 \times 8$

$19 \times 19 \times 16$   
 $19 \times 19 \times 40$

#anchors

$5 + \#classes$

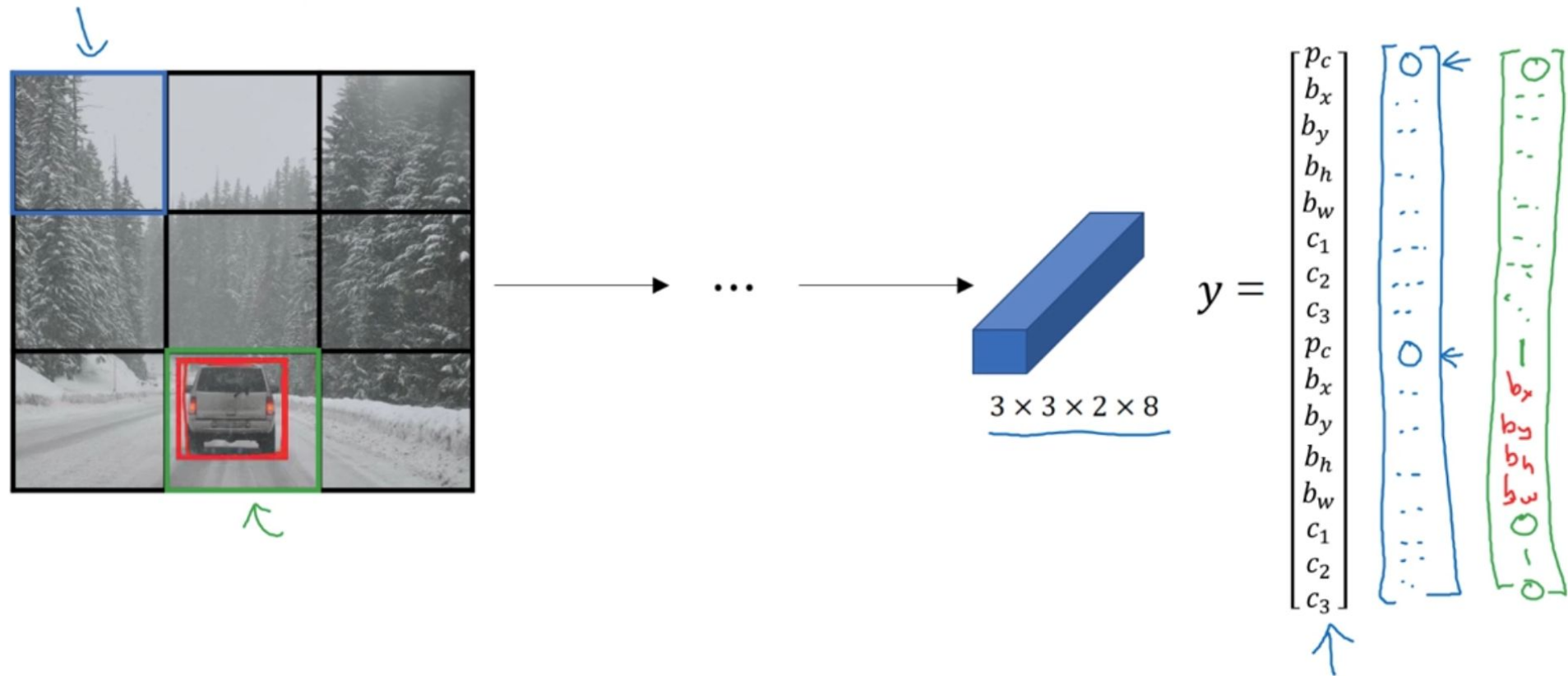


ConvNet



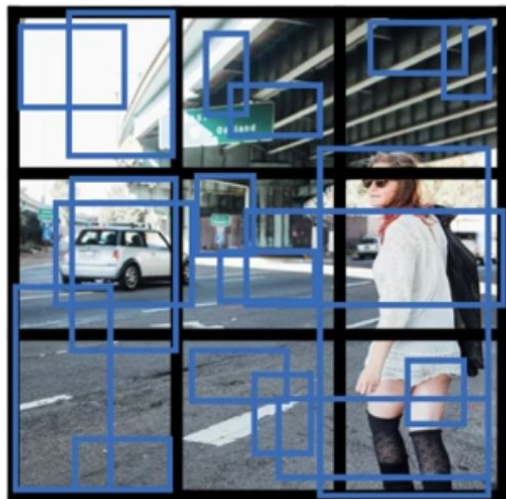
# Making predictions

YOLO  
Algorithm



# Outputting the non-max suppressed outputs

YOLO  
Algorithm



- For each grid cell, get 2 predicted bounding boxes.



# Outputting the non-max suppressed outputs

YOLO  
Algorithm



- For each grid cell, get 2 predicted bounding boxes.
- Get rid of low probability predictions.
- For each class (pedestrian, car, motorcycle) use non-max suppression to generate final predictions.

# For more details...

Check the Andrew Ng's videos on object detection

Available on YouTube

See the following playlist:

<https://www.youtube.com/playlist?list=PLkDaE6sCZn6GI29AoE31iwdVwSG-KnDzF>

Videos: C4W3L01, C4W3L03, C4W3L04, C4W3L06, C4W3L07, C4W3L08,  
C4W3L09

# YOLOX

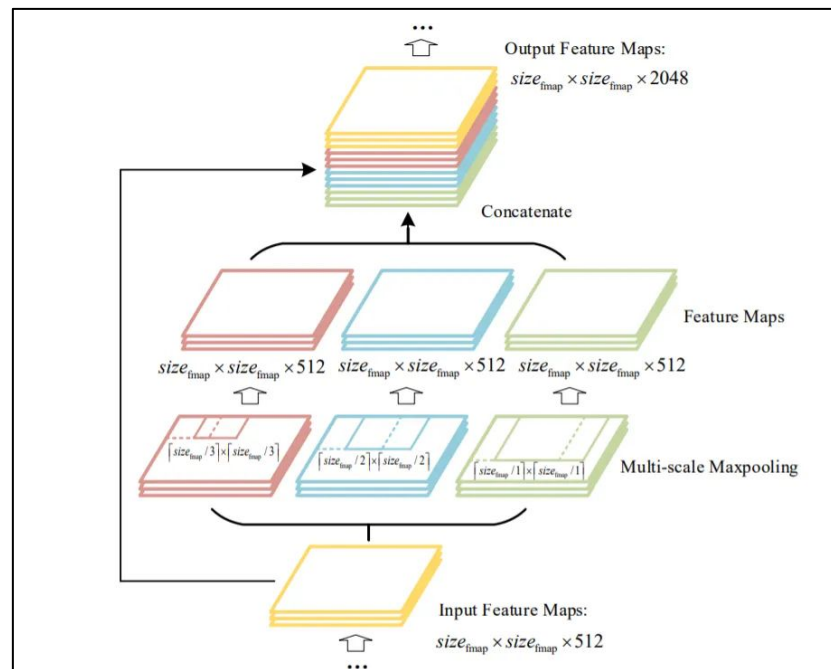
## Object Detection Algorithm

Some information and figures  
Courtesy of LearnOpenCV

# YOLOX

Built on top of YOLOv3 with Darknet-53 backbone and SPP layer (YOLOv3-SPP)

Spatial Pyramid Pooling (SPP) layer →



# YOLOX

YOLOX' distinctive features:

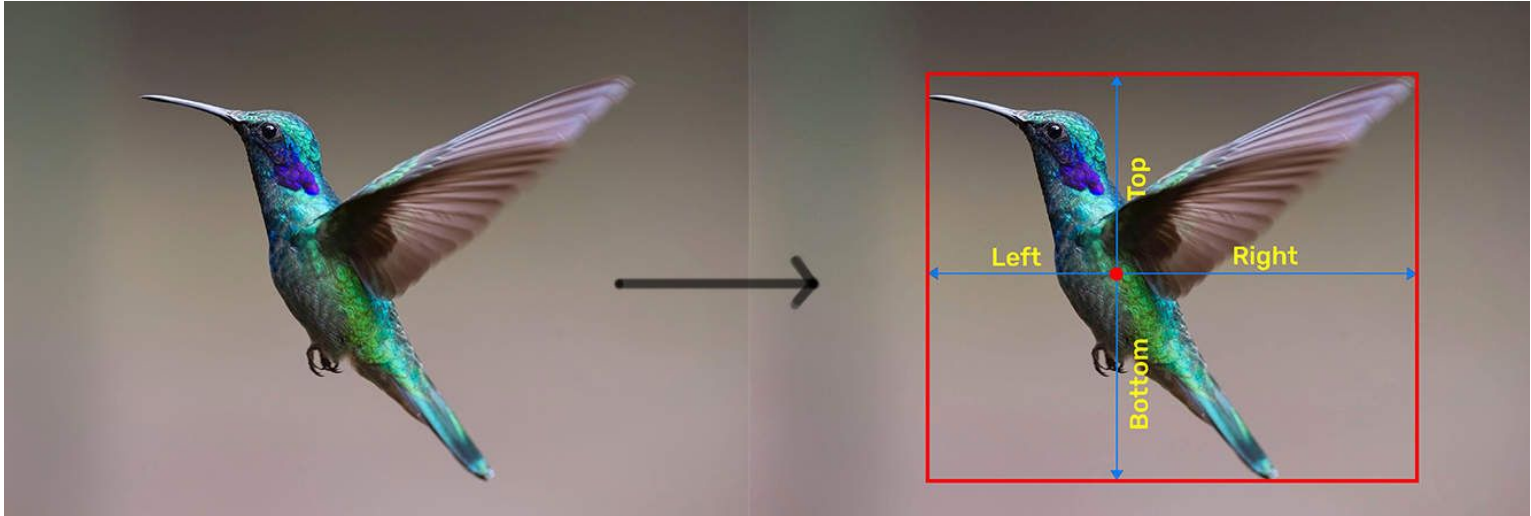
- Anchor free design
- Decoupled head
- simOTA label assignment strategy
- Advanced Augmentations: Mixup and Mosaic



# Center based detector

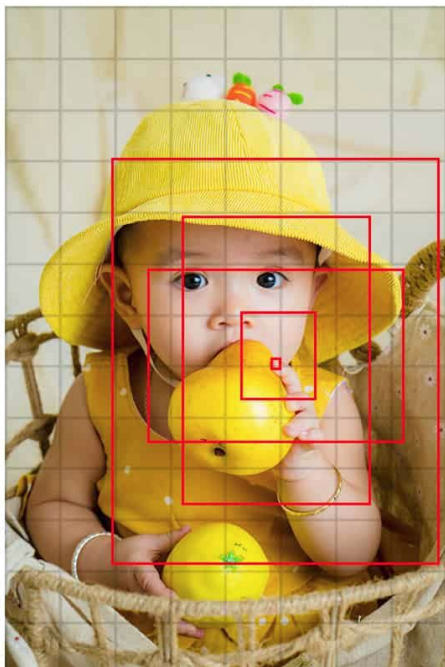
Find positive point in the center

Predict four distances from the positive to the boundary

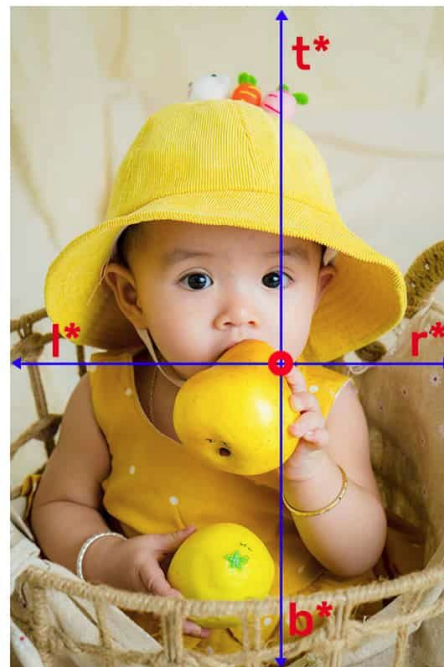


# Anchor Free YOLOX

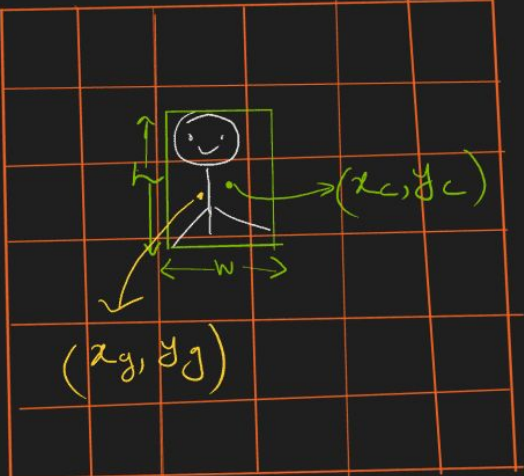
Anchor based



Anchor free



# What is Anchor Free Object Detection



The diagram shows a 5x5 grid with a stick figure in the center. A green bounding box is drawn around the figure. The center of the bounding box is marked with a dot and labeled  $(x_c, y_c)$ . The grid cell center is marked with a dot and labeled  $(x_g, y_g)$ . The width of the bounding box is labeled  $w$  and the height is labeled  $h$ . Arrows indicate the dimensions and the positions of the center points.

Anchor-free prediction

- ground truth bbox & its center
- grid cell center

A per-cell anchor-free model should predict

$$\delta x_g = x_g - x_c$$
$$\delta y_g = y_g - y_c$$

height =  $h$ , width =  $w$   
class = ground-truth class

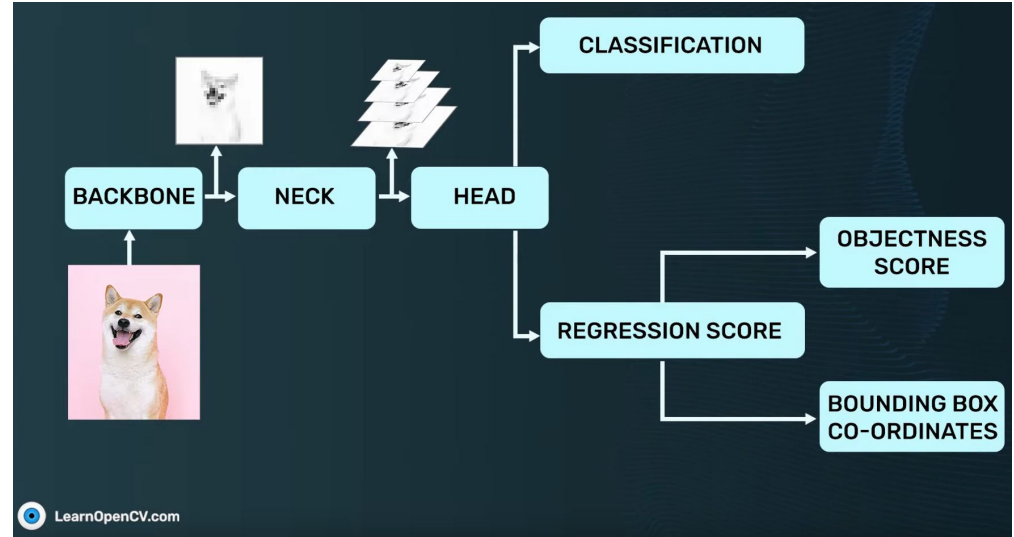


# YOLO architecture

Backbone: extracts features of an image

Neck: producing feature maps with multiple scales

Head: outputs localization and classification scores



# YOLOX - decoupled head

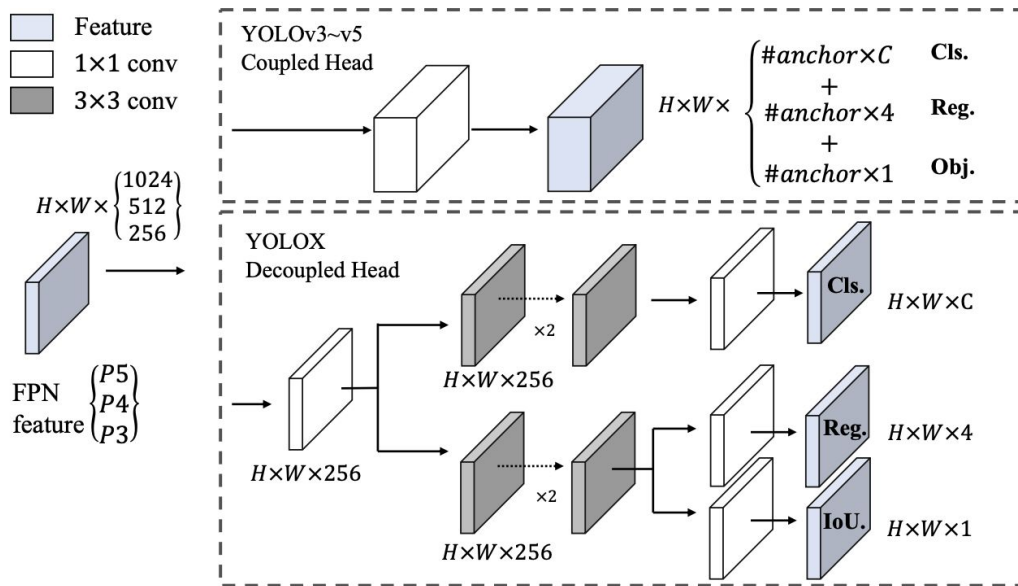
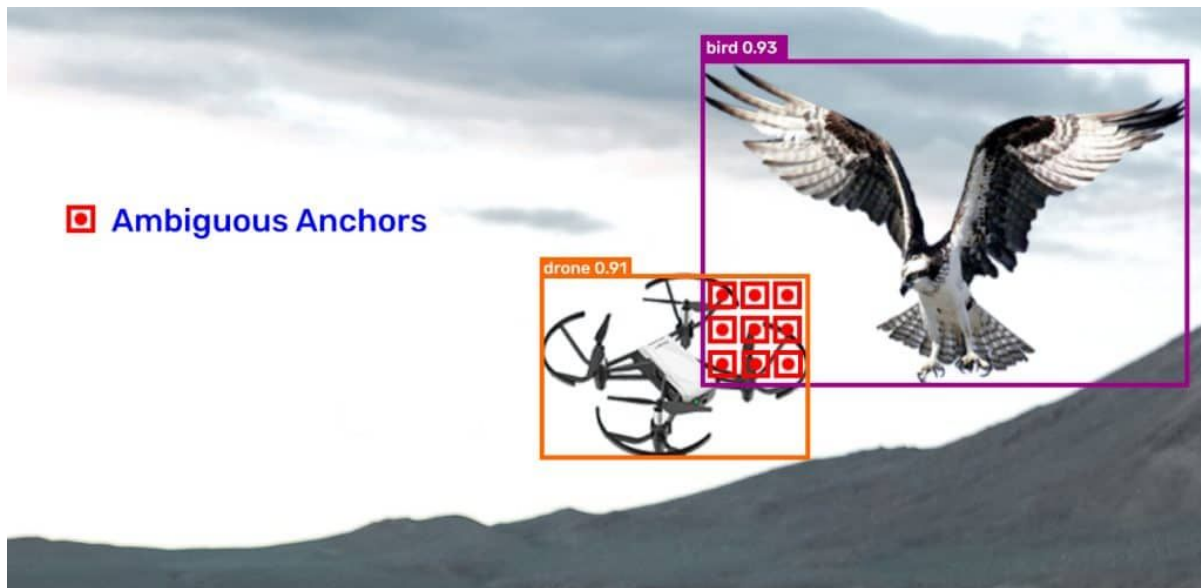


Figure 2: Illustration of the difference between YOLOv3 head and the proposed decoupled head. For each level of FPN feature, we first adopt a  $1 \times 1$  conv layer to reduce the feature channel to 256 and then add two parallel branches with two  $3 \times 3$  conv layers each for classification and regression tasks respectively. IoU branch is added on the regression branch.

# SimOTA Advanced Label Assignment Strategy

OTA: Optimal Transport Assignment for Object Detection



# SimOTA Advanced Label Assignment Strategy

Briefly explained by LearnOpenCV

What is simOTA in YOLOX?

Simplified OTA or simOTA is the redesigned Optimal Transport Assignment strategy. The training cost does not increase but average precision(AP) is definitely improved. It is shown with empirical evidence in the paper.

In simOTA, iteration is not performed for every positive label. A strategy called **Dynamic Top K** is used to estimate the approximate number of positive anchors for each ground truth. Here, only the top **K** number of positive labels are selected. This reduces the number of iterations by many folds.

The number of positive labels per ground truth (GT) varies due to the following factors.

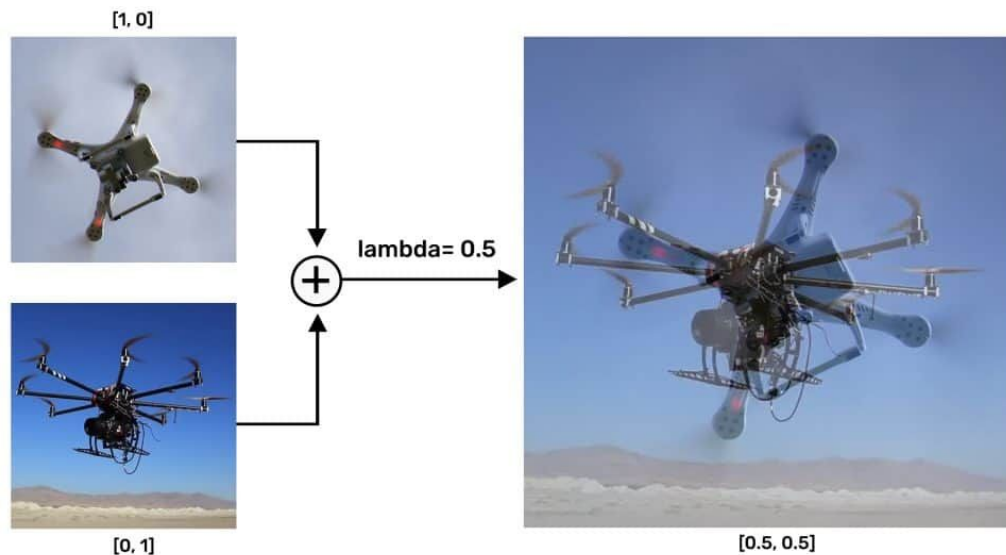
- Size
- Scale
- Occlusion conditions etc.

However, it is difficult to model a mapping function from these factors to the positive anchor number **K**. Hence it is done on the basis of IoU value. The [IoU values](#) of the **anchors** to the ground truth(GT) are summed up to represent the GT's estimated number of positive anchors.

**The intuition is such that the number of positive anchors for a certain GT should be positively correlated with the number of anchors that have well regressed.**

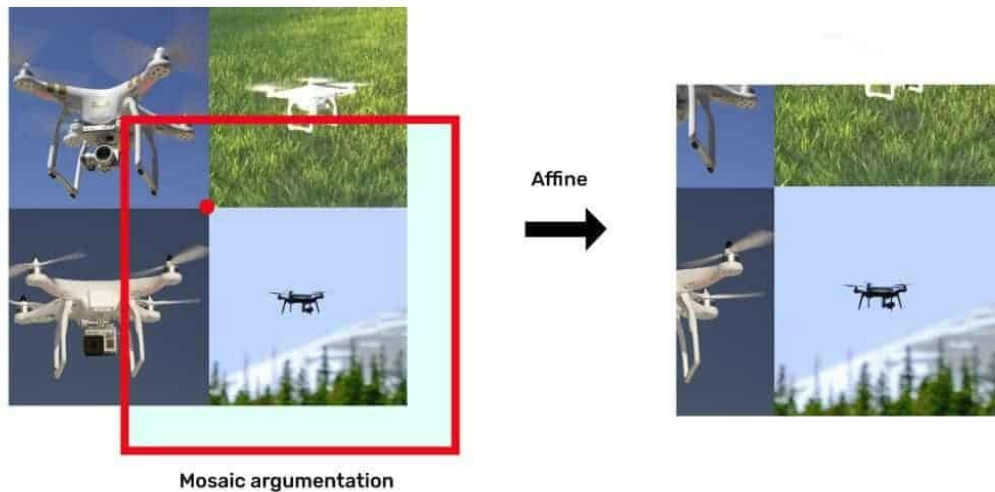
# Strong Data Augmentation in YOLOX

## Mixup Augmentation



# Strong Data Augmentation in YOLOX

## Mosaic Augmentation



# Performance gain - step by step

Methods	AP (%)	Parameters	GFLOPs	Latency	FPS
YOLOv3-ultralytics <sup>2</sup>	44.3	63.00 M	157.3	10.5 ms	95.2
YOLOv3 baseline	38.5	63.00 M	157.3	10.5 ms	95.2
+decoupled head	39.6 (+1.1)	63.86 M	186.0	11.6 ms	86.2
+strong augmentation	42.0 (+2.4)	63.86 M	186.0	11.6 ms	86.2
+anchor-free	42.9 (+0.9)	63.72 M	185.3	11.1 ms	90.1
+multi positives	45.0 (+2.1)	63.72 M	185.3	11.1 ms	90.1
+SimOTA	<b>47.3 (+2.3)</b>	63.72 M	185.3	11.1 ms	90.1
+NMS free (optional)	46.5 (-0.8)	67.27 M	205.1	13.5 ms	74.1

Table 2: Roadmap of YOLOX-Darknet53 in terms of AP (%) on COCO *val*. All the models are tested at  $640 \times 640$  resolution, with FP16-precision and batch=1 on a Tesla V100. The latency and FPS in this table are measured without post-processing.

# For more details...

Check the following:

YOLOX Object Detector Paper Explanation and Custom Training

<https://learnopencv.com/yolox-object-detector-paper-explanation-and-custom-training/>

CenterNet: Objects as Points – Anchor Free Object Detection Explained

<https://learnopencv.com/centernet-anchor-free-object-detection-explained/>

Paper Review: “YOLOX: Exceeding YOLO Series in 2021”

<https://medium.com/mlearning-ai/paper-review-yolox-exceeding-yolo-series-in-2021-ffc1bd94a1f3>

The YOLOX paper

<https://arxiv.org/pdf/2107.08430.pdf>

...more interesting blog posts of your choice! Plenty of stuff available online!



# Questions and Answers