# Deep Learning and Convolutional Neural Network (42028)

Object Detection- 2

# Frameworks

Object Detection

Region Proposal Based

Regression/Classification Based

# Object Detection Techniques History



Images source: https://arxiv.org/pdf/1807.05511.pdf

# Object Detection Techniques Recap

## Sliding Window technique

# Object Detection Techniques Recap
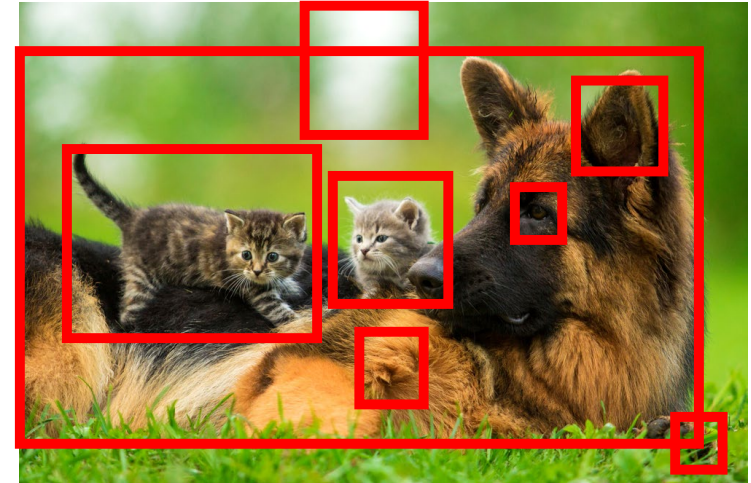
**Sliding Window technique**
- Crop images and classify using CNN
- Try different sizes of the sliding window

**Issues:**
- Slow
- Computationally very expensive
- Less accurate

# Object Detection Techniques Recap

## Region Proposals

# Predicting Bounding Boxes

Currently:
- Sliding Window
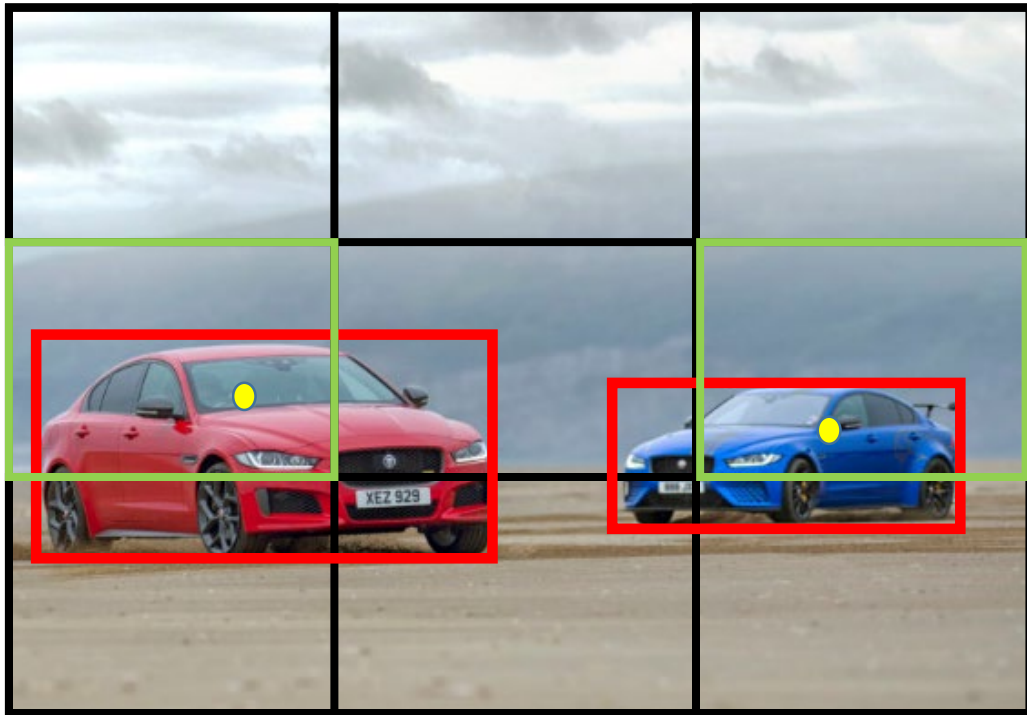- Selective Search
- Region Proposals

Task:
- Predict Bounding boxes from CNN

# Predicting Bounding Boxes



- Place a grid over the image

- Apply image classification and localization to each of the grid cells

# Predicting Bounding Boxes



**Class : {car, bike}**
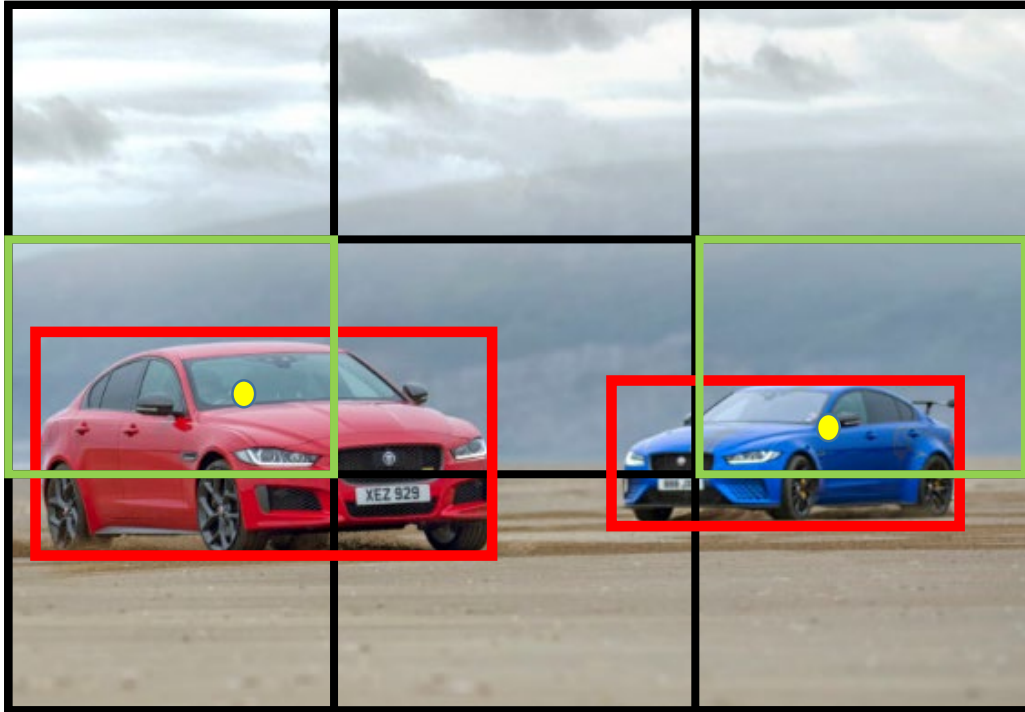
**Training Strategy:**

Input:
-   Image: (ht x wd x 3)

Target:
-   Bounding box information for each object
-   Class for each object

# Predicting Bounding Boxes



**Class : {car, bike}**

**Idea:** Take the mid-point of the object and
Assign it to a grid cell based on its location

**Training Strategy:**

Target:

$Y = \{p_o, x, y, h, w, c_1, c_2\}$ for each cell

e.g:

$Cell(1,1) = \{0, ?, ?, ?, ?, ?, ?\}$

:

$Cell(2,1) = \{1, x, y, h, t, 1, 0\}$

$Cell(2,2) = \{0, ?, ?, ?, ?, ?, ?\}$

$Cell(2,3) = \{1, x, y, h, t, 1, 0\}$

:

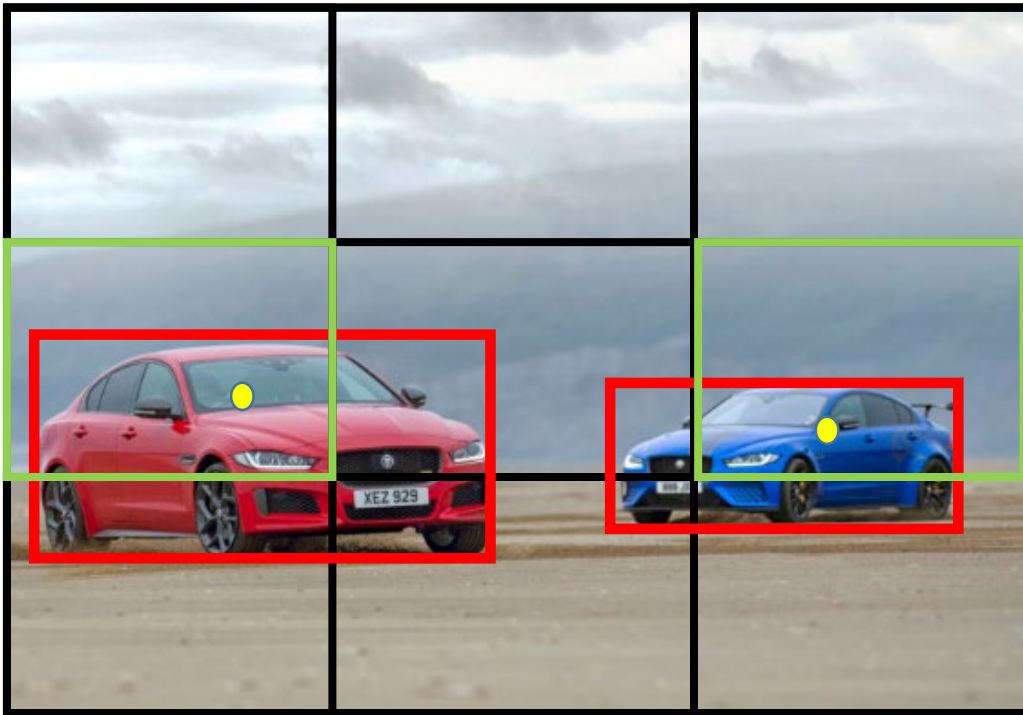$Cell(3,3) = \{0, ?, ?, ?, ?, ?, ?\}$

# Predicting Bounding Boxes



**Class : {car, bike}**

**Training Strategy:**

Target output vector:

3 X 3 X 7

3 X 3: Grid size

7: (5 + Number-of-Classes)



3 X 3 X 7

# Predicting Bounding Boxes

## Training Strategy:
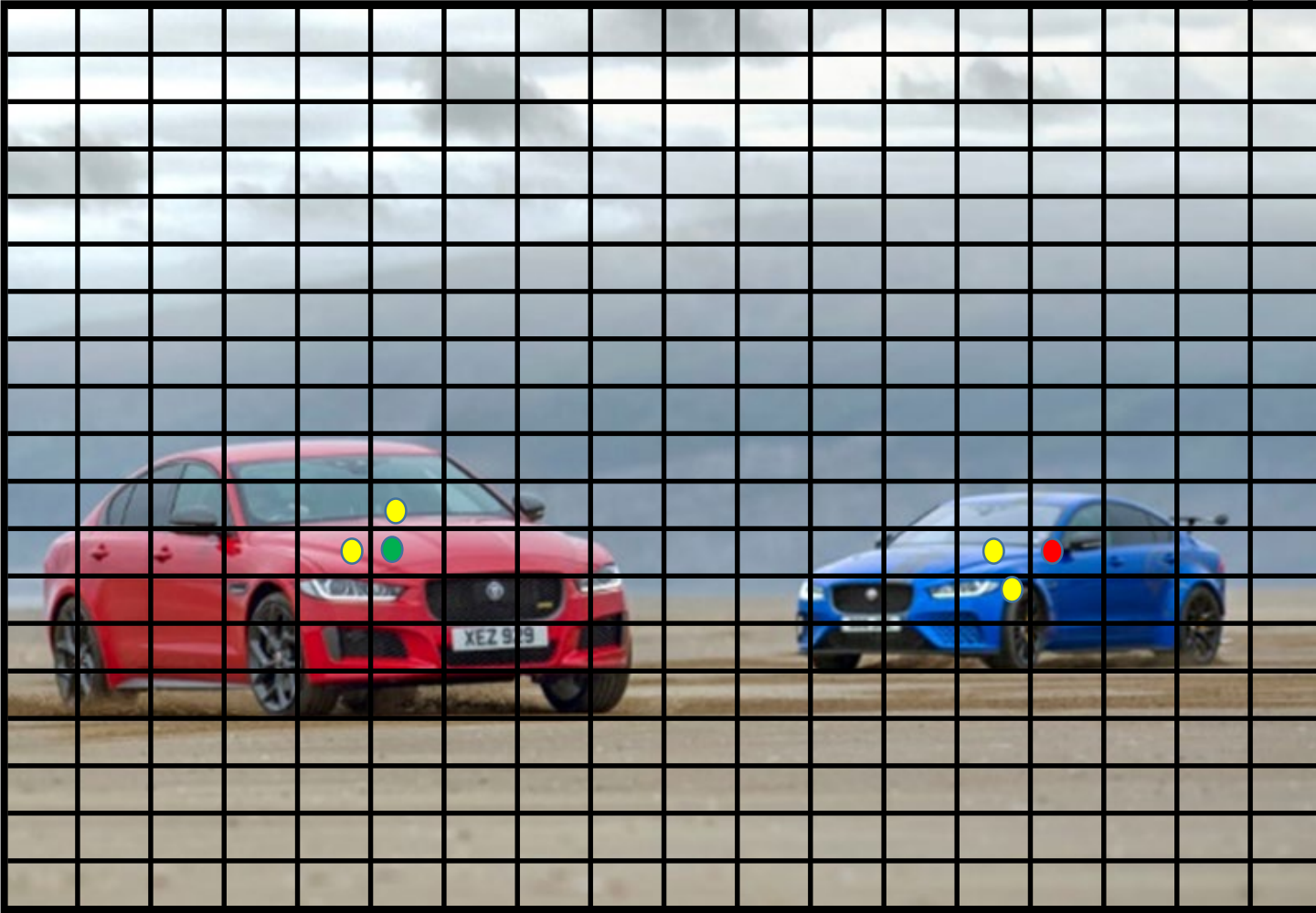
**Input: X**



**Class : {car, bike}**

**Target: Y**

**CNN**

3 X 3 X 7

**In practice:** The grid is finer, 19 X 19 instead of 3 X 3

So, Target will be of size: 19 X 19 X 7

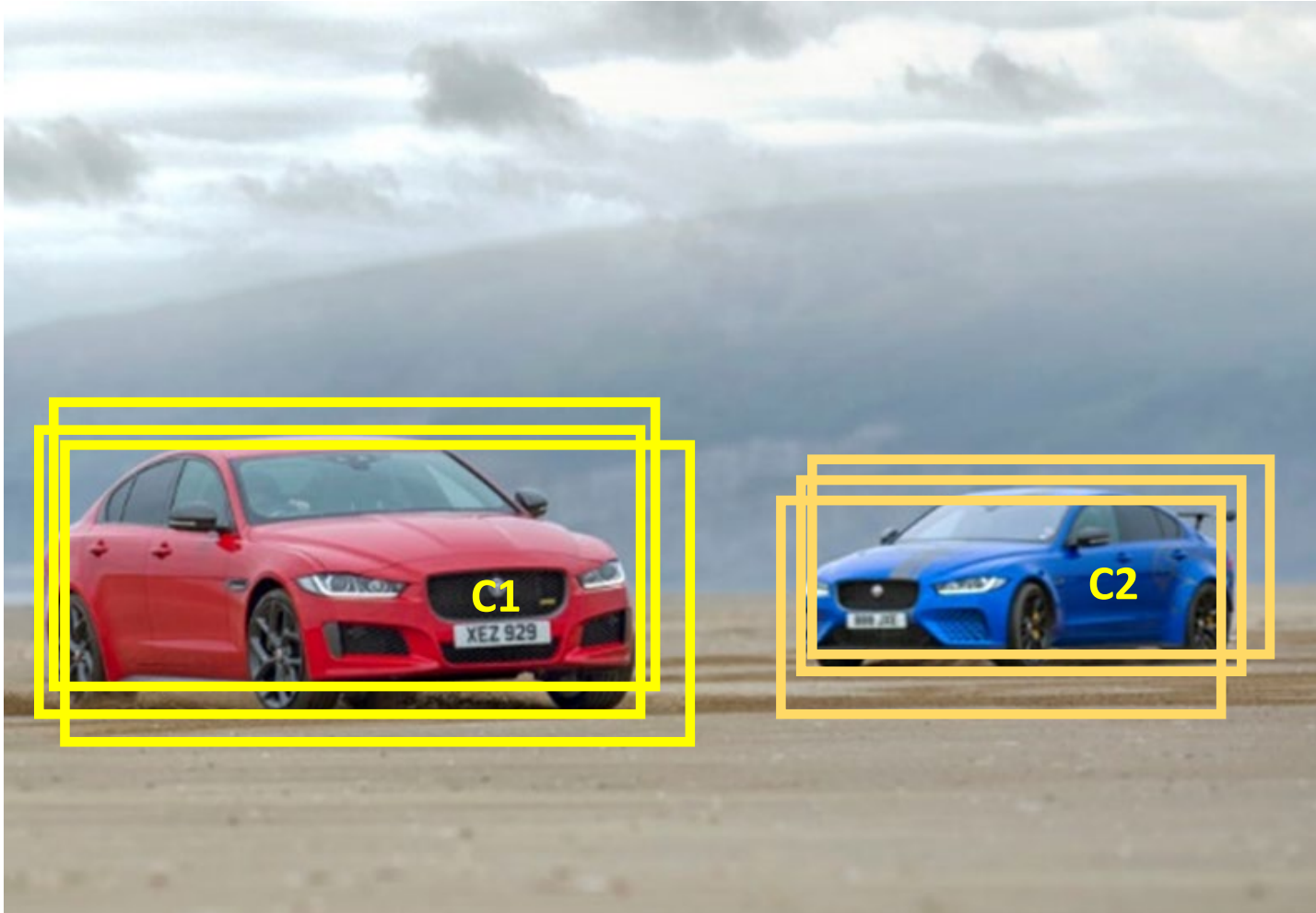Works well for non-overlapping objects

# Non Maxima Suppression (NMS)



**Issues with Object Detection:**

1. Each object has one mid-point
2. As each cells are subjected to object+localization classification
3. Hence, neighbouring cells might assume that it has the mid-point
4. Hence, Multiple detection bounding box

Images source: https://arxiv.org/pdf/1807.05511.pdf

# Non Maxima Suppression (NMS)



**Sample prediction:**

**For C1:**

**Box1: 0.9 (Confidence Score)**
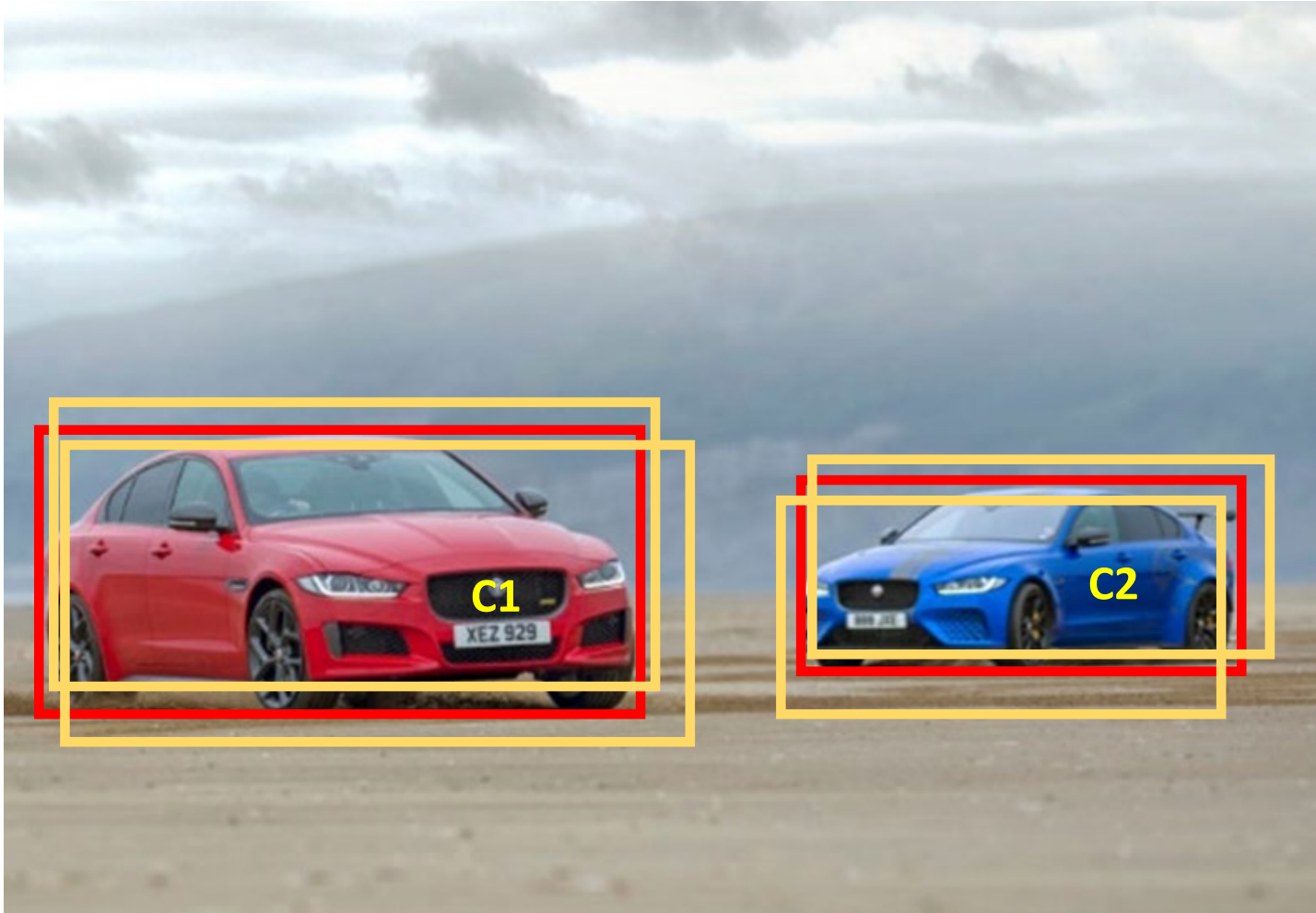
**Box2: 0.79**

**Box3: 0.82**

**For C2:**

**Box1: 0.92**

**Box2: 0.85**

**Box3: 0.7**

**NMS cleans/removes the multiple detection and only keeps the one with very high confidence**

Images source: https://arxiv.org/pdf/1807.05511.pdf

# Non Maxima Suppression (NMS)



1. Check the probabilities of each detection and keep ones with
*score > Threshold (0.7)*

2. For remaining boxes:
- Box with highest score is the detection results.
- Discard any remaining boxes with IoU > 0.5 with final detected box,
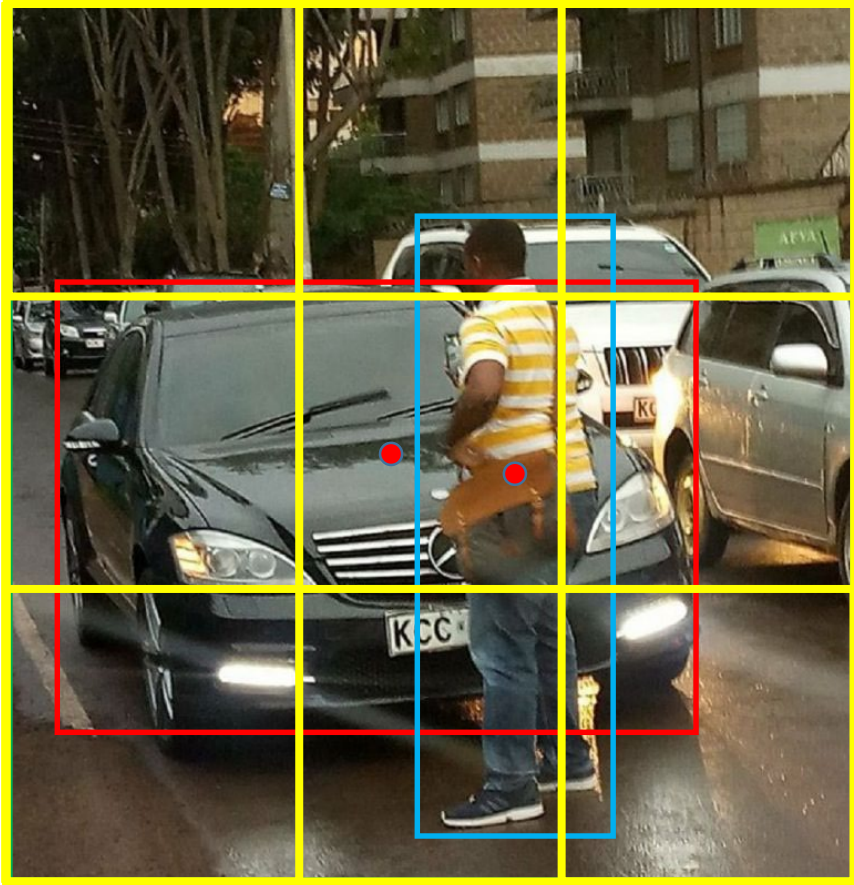i.e: overlap with the box with highest score.

Images source: https://arxiv.org/pdf/1807.05511.pdf

# YOLO: You Only Look Once Algorithm



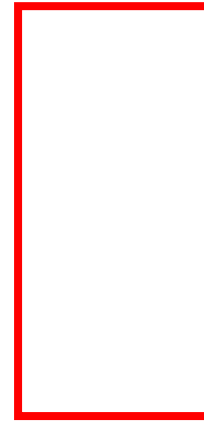**Challenges with overlapping objects**

**- Each grid cell detect only one object**
**- For multiple overlapping objects, Mid point are on the same grid cell**

# Anchor Boxes



So, Currently the Target Y = {1, x, y, h, w, C1, C2},
As the mid-points for both the objects are on the
same grid cell, only one of the objects will be associated

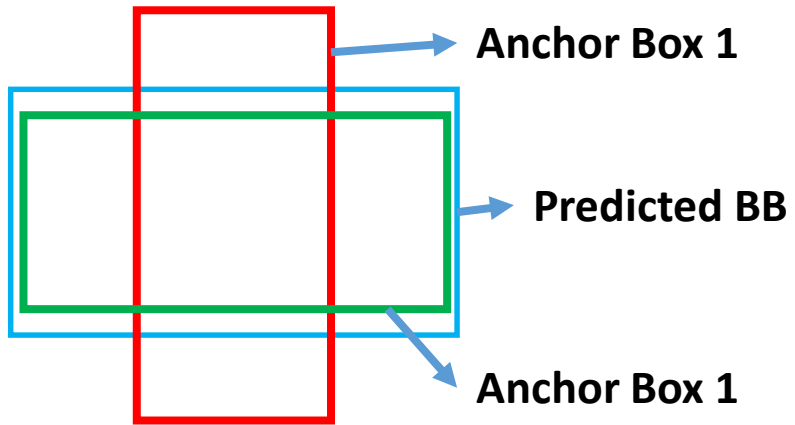**Anchor Box 1**          **Anchor Box 2**

# Anchor Boxes

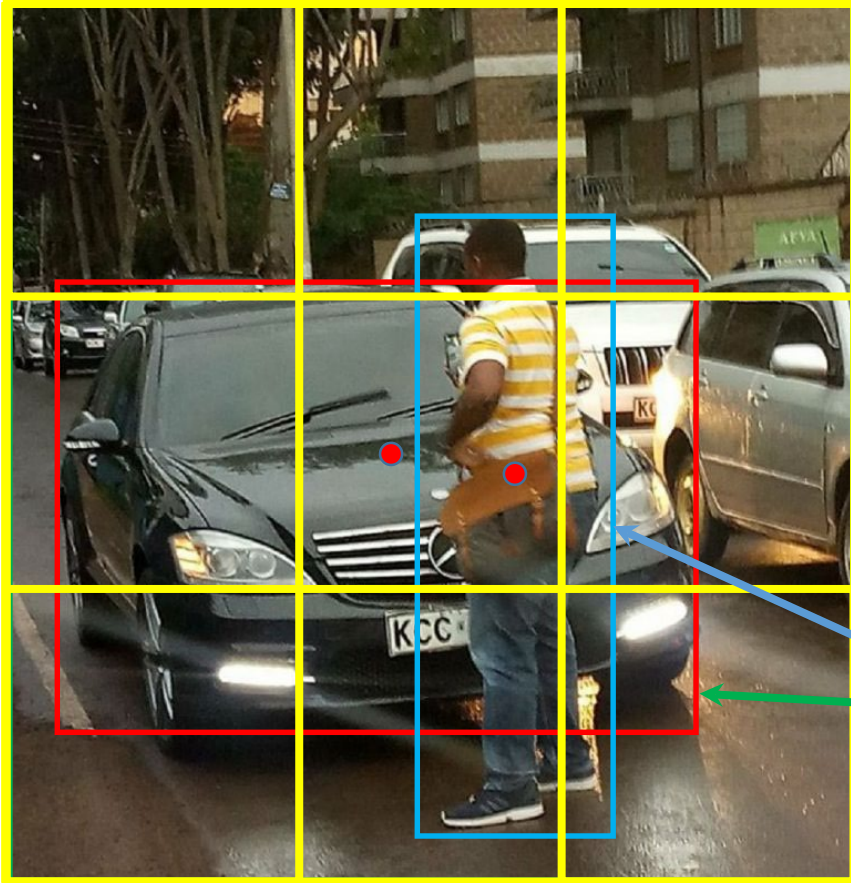Anchor Box 1

Predicted BB

Anchor Box 1

**Calculate the IoU of**
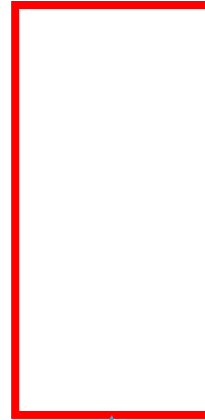Anchor boxes and predicted BB

**Associate each object to:**

1. **A cell which contains its mid-point and**
2. **Anchor box for the cell with highest IoU**

# Anchor Boxes



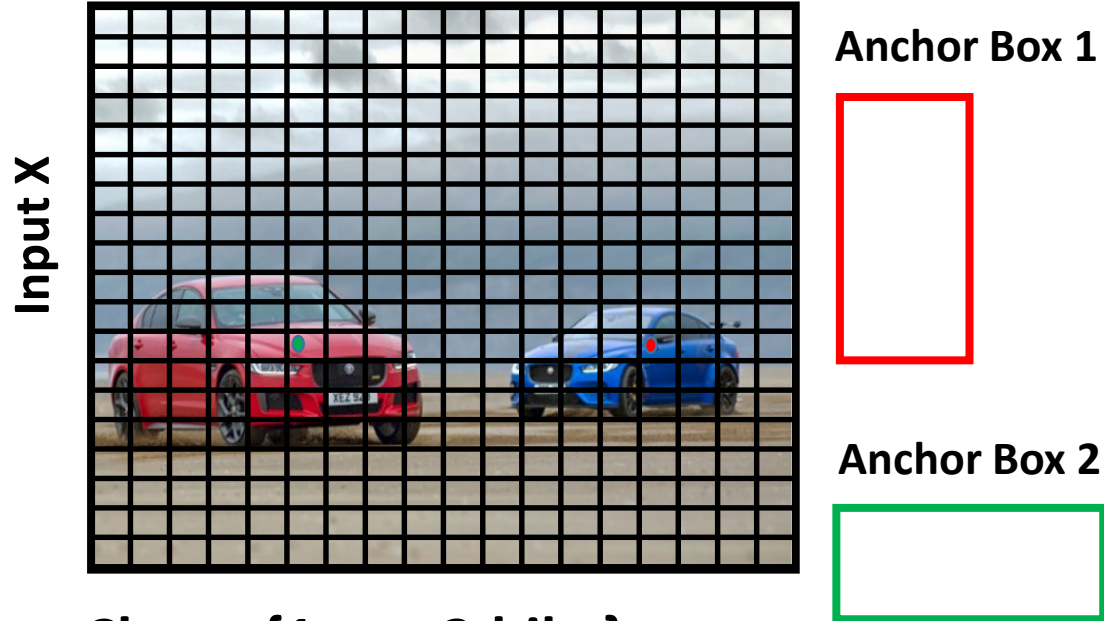Anchor Box 1

Anchor Box 2

Similar Shape

So, with Anchor boxes:

Target Y = {$P_o$, x, y, h, w, C1, C2, $P_o$, x, y, h, w, C1, C2},

Anchor Box 1          Anchor Box 2

Images source: https://arxiv.org/pdf/1807.05511.pdf

# YOLO: You Only Look Once Algorithm

**Input X**

**Anchor Box 1**

**Anchor Box 2**

Class : {1:car, 2:bike}

Y size : ( 19 X 19 X 2 X 7 )

**Grid Size**

**#Anchor Box**

$5(P_o, x,y,h,w) + \#Classes(2)$

**Training Set**

$Y = \{P_o, x, y, h, w, C1, C2, P_o, x, y, h, w, C1, C2\}$

$Cell(1,1) = \{0, ?, ?, ?, ?, ?, ?, 0, ?, ?, ?, ?, ?, ?\}$
⋮
$Cell(12,6) = \{0, ?, ?, ?, ?, ?, ?, 1, x, y, h, w, 1, 0\}$
⋮
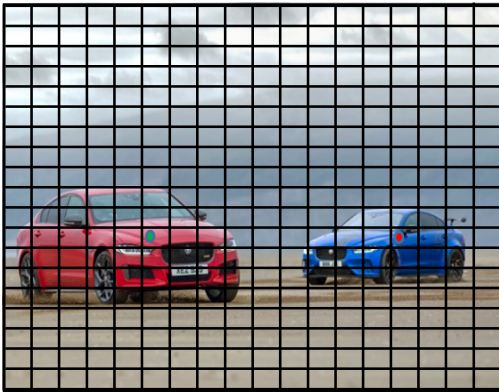$Cell(12,15) = \{0, ?, ?, ?, ?, ?, ?, 1, x, y, h, w, 1, 0\}$
⋮
$Cell(19,19) = \{0, ?, ?, ?, ?, ?, ?, 0, ?, ?, ?, ?, ?, ?\}$

# YOLO: You Only Look Once Algorithm

**Training:**

**Input: X**

**Target: Y**



CNN

19 X 19 X 2 X 7

**Class : {car, bike}**

# YOLO: You Only Look Once Algorithm

**Testing:**

**Input: X**



**Class : {car, bike}**

CNN

19 X 19 X 2 X 7

Y = {P_o, x, y, h, w, C1, C2, P_o, x, y, h, w, C1, C2}

{0, ?, ?, ?, ?, ?, ?, 0, ?, ?, ?, ?, ?, ?}
:
{0, ?, ?, ?, ?, ?, ?, 1, x, y, h, w, 1, 0}
:
{0, ?, ?, ?, ?, ?, ?, 1, x, y, h, w, 1, 0}
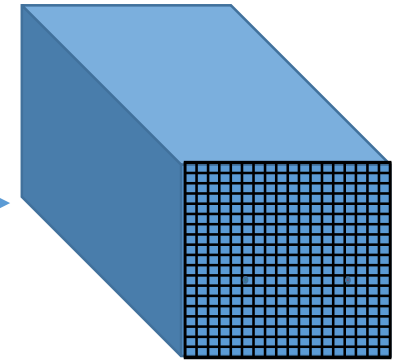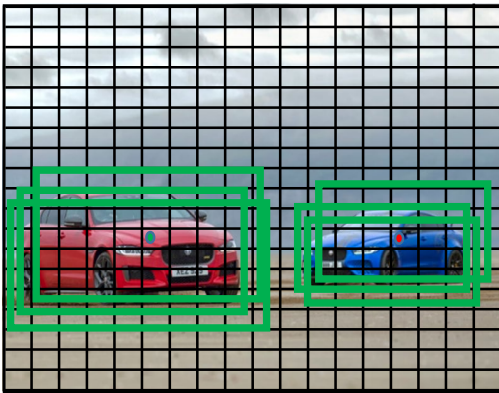:
{0, ?, ?, ?, ?, ?, ?, 0, ?, ?, ?, ?, ?, ?}

# YOLO: You Only Look Once Algorithm

**Testing:**

**Input: X**



**Class : {car, bike}**

CNN

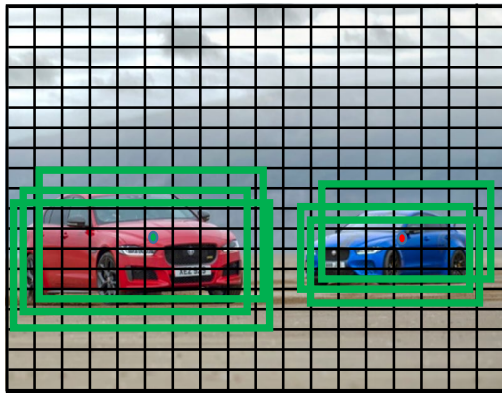19 X 19 X 2 X 7

Y = {$P_O$, x, y, h, w, C1, C2, $P_O$, x, y, h, w, C1, C2}

{0, ?, ?, ?, ?, ?, ?, 0, ?, ?, ?, ?, ?, ?}
:
{0, ?, ?, ?, ?, ?, ?, 1, x, y, h, w, 1, 0}
:
{0, ?, ?, ?, ?, ?, ?, 1, x, y, h, w, 1, 0}
:
{0, ?, ?, ?, ?, ?, ?, 0, ?, ?, ?, ?, ?, ?}

**Apply NMS**
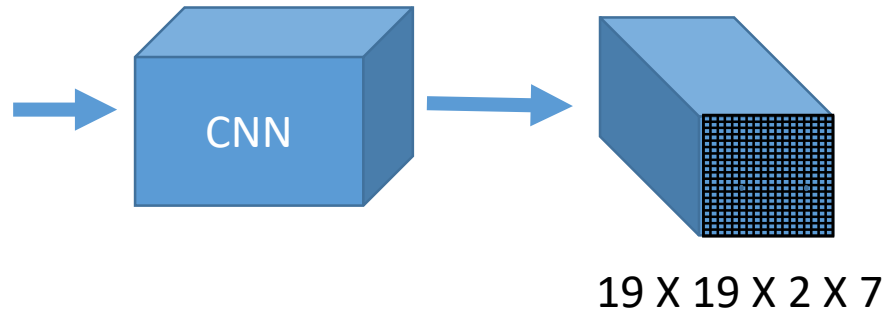
# YOLO: You Only Look Once Algorithm

- Real-time performance with 45 frames per sec, 0.02 sec per image

- Not suitable for small objects

- Issues with new or multiple aspect ratios and unable to generalize

# Single Shot Detector(SSD)

- Similar to YOLO

- VGG16 base CONV layers

- Take advantage of Anchor boxes with different aspect ratios

- Large number of anchors boxes are chosen

- Not suitable for small objects

- 3 times faster than Faster-RCNN

- With ResNet101 base SSD may be help in detecting small objects with better features from the CONV base

# Single Shot Detector(SSD)

SSD300 architecture:

# Object Detection State-of-the-Art

**Dataset: PASCAL VOC 2007 and 2017**

**Test Dataset : PASCAL VOC 2007**

| Method | Train Dataset | mAP | Time in sec/image | Time Frame /sec |
|---|---|---|---|---|
| RCNN (VGG16) | Pascal VOC 2007 | 66.0 | 50 | - |
| Fast RCNN | VOC 2007+2012 | 70.0 | 2 | - |
| Faster RCNN (VGG16) | VOC 2007+2012 | 73.2 | 0.11 | 9 |
| Faster RCNN (ResNet101) | VOC 2007+2012 | **83.8** | 2.24 | 0.4 |
| Yolo | VOC 2007+2012 | 63.4 | 0.02 | 45 |
| SSD300 | VOC 2007+2012 | 74.3 | 0.02 | 45 |
| SSD512 | VOC 2007+2012 | **76.8** | 0.05 | 19 |

# Object Detection Summary

**Base Networks:**
- VGG16
- REsNet101
- Inception V2
- Inception V3
- ResNet
- MobileNet
- Alexnet
- ZFNet

Etc.

**Object Detection FrameWorks:**
- RCNN Family (RCNN, Fast/Faster RCNN)
- Yolo
- SSD
- F-RCN

**Summary:**
- Faster-RCNN is more accurate but slower
- Yolo/SSD are faster/real-time but not much accurate

Source: http://cs231n.stanford.edu/slides/2017/cs231n_2017_lecture11.pdf