# **Object Detection**

Computer Vision STUDENTS-HUB.com

Aziz M. Qaroush

Birzeit University Uploaded By: anonymous

#### Outline

#### 2

#### Introduction

- Traditional Computer Vision Techniques:
  - Sliding Window:
  - Region proposal approach (Selective Search)
- Deep Learning-Based Techniques:
  - Two-Stage Detectors:
    - R-CNN (Regions with CNN features)
  - One-Stage Detectors:
    - YOLO (You Only Look Once)
- Other Detectors
- Evaluation Measures
- Benchmarks

### **Common CV Applications**



(a) Object Classification



#### (b) Generic Object Detection (Bounding Box)



(c) Semantic Segmentation STUDENTS-HUB.com



(d) Object Instance Segmetation Uploaded By: anonymous

# **Object Detection: Task Definition**

#### 4

- Problem: Detecting and localizing generic objects from various categories, such as cars, people, etc.
  - Input: Single RGB Image
  - Output: A set of detected objects; For each object predict:

1. Category label (from fixed, known set of categories)

2. Bounding box (four numbers: x, y, width, height)

#### Challenges:

- Illumination, Viewpoint, deformations, Intra-class variability, ...
- Multiple outputs: Need to output variable numbers of objects per image
- Multiple types of output: Need to predict "what" (category label) as well as "where" (bounding box)
- Large images: Classification works at 224x224; need higher resolution for detection, often ~800x600



#### STUDENTS-HUB.com

# **Object Detection: Bounding Boxes**

- 5
- Bounding boxes (usually)
   cover only the visible portion
   of the object
- Bounding boxes are typically axis-aligned
- Oriented boxes are much less
   common



### Comparing Boxes: Intersection over Union (IoU)

- 6
- How can we compare our prediction to the groundtruth box?
- Intersection over Union

   (IoU) (Also called "Jaccard similarity" or "Jaccard index"):

Area of Intersection Area of Union Our Prediction Ground Truth

IoU > 0.5 is "decent",IoU > 0.7 is "pretty good",IoU > 0.9 is "almost perfect"

STUDENTS-HUB.com

# Detecting a single object



STUDENTS-HUB.com

7

# **Detecting Multiple Objects**









#### Need different numbers of outputs per image

CAT: (x, y, w, h)

4 numbers

DOG: (x, y, w, h) DOG: (x, y, w, h) CAT: (x, y, w, h)

....

12 numbers



pooling

Max pooling

> DUCK: (x, y, w, h) N DUCK: (x, y, w, h) r

Many numbers!

STUDENTS-HUB.com

#### Outline

#### 9

- Introduction
- Traditional Computer Vision Techniques:
  - Sliding Window
  - Region proposal approach (Selective Search)
- Deep Learning-Based Techniques:
  - Two-Stage Detectors:
    - R-CNN (Regions with CNN features)
  - One-Stage Detectors:
    - YOLO (You Only Look Once)
- Other Detectors
- Evaluation Measures
- Benchmarks

# **Sliding Window**

- 10
- In the sliding window approach, we slide a box or window over an image to select a patch and classify each image patch covered by the window using the object recognition model.
- Not only do we need to search all possible locations in the image, we have to search at different scales.
- The problem doesn't end here. Sliding window approach is good for fixed aspect ratio objects such as faces or pedestrians.
- Number of windows is large and grows quadratically with the number of pixels, and the need to search over multiple scales and aspect ratios further increases the search space.
- □ The huge search space results in high computational complexity

### **Sliding Window with HOG method**





#### sliding window

STUDENTS-HUB.com

#### **Recap – HOG features**

12



#### Find a HOG template and use as filter

STUDENTS-HUB.com

### **Sliding window + HOG features**



 Slide through the image and check if there is an object at every location

#### No person here

STUDENTS-HUB.com

**14** 



 Slide through the image and check if there is an object at every location

#### **YES!!** Person match found

STUDENTS-HUB.com



#### But what if we were looking for buses?

#### No bus found

STUDENTS-HUB.com



#### But what if we were looking for buses?

#### No bus found

STUDENTS-HUB.com

17



We will never find the object we don't choose our window size wisely!

#### No bus found

STUDENTS-HUB.com

18



#### We need to do multi scale sliding window

STUDENTS-HUB.com

# Sliding window + HoG with feature pyramid

- 19
- □ Sliding window technique is applied as usual over all the pyramid levels.
- The window that produces the highest similarity score out of the resizing's is used as the location of the detected object



Filter F

Score of *F* at position *p* is  $F \cdot \phi(p, H)$ 

 $\phi(p, H)$  = concatenation of HOG features from subwindow specified by p

A feature pyramid of different image resizing STUDENTS-HUB.com

#### Detecting Multiple Objects: Sliding Window

- 20
- Question: How many possible boxes are there in an image of size H x W?
  - Consider a box of size h x w:
  - Possible x positions: W w + 1
  - Possible y positions: H h + 1
  - Possible positions: (W w + 1) \* (H h + 1)
  - Total possible boxes =

$$\sum_{h=1}^{H} \sum_{w=1}^{W} (W - w + 1)(H - h + 1)$$

$$=\frac{H(H+1)}{2}\frac{W(W+1)}{2}$$

800 x 600 image has ~58M boxes! No way we can evaluate them all Uploaded By: anonymous

STUDENTS-HUB.com

#### Outline

#### 21

- Introduction
- Traditional Computer Vision Techniques:
  - Sliding Window:
  - **Region proposal approach** (Selective Search)
- Deep Learning-Based Techniques:
  - Two-Stage Detectors:
    - R-CNN (Regions with CNN features)
  - One-Stage Detectors:
    - YOLO (You Only Look Once)
- Other Detectors
- Evaluation Measures
- Benchmarks

# **Region proposal approach**

- 22
- Takes an image as the input and output bounding boxes corresponding to all patches in an image that are most likely to be objects.
- These region proposals can be noisy, overlapping and may not contain the object perfectly but amongst these region proposals, there will be a proposal which will be very close to the actual object in the image.
- We can then classify these proposals using the object recognition model. The region proposals with the high probability scores are locations of the object.
- An important property of a region proposal method is to have a very **high recall**. This is just a fancy way of saying that the regions that contain the objects we are looking have to be in our list of region proposals.

STUDENTS-HUB.com

# **Region Proposals via Selective Search**

- 23
- It aims to generate a diverse set of potential regions in an image that may contain objects.
- It create region proposals from smaller segments to larger segments in a bottomup approach.
- Algorithm steps:
  - Step 1: Starts by over-segmenting the image based on intensity of the pixels using Felzenszwalb and Huttenlocher graph-based segmentation method. The goal of this step is to over-segment the image into a large number of small regions.



### **Region Proposals via Selective Search**

- 24
- Step 2: From set of regions, choose two adjacent segments that are most similar based on their similarity.
  - **The similarity measure often considers color, texture, size, and shape features.**
- **Step 3**: Combine them into a single, larger region.
  - Use a hierarchical grouping algorithm (e.g., agglomerative clustering) based on the similarity measure.
  - The goal of this step is to merge similar segments into larger regions.
- **Step 4:** Repeat step 2. Continue grouping until a stooping criteria satisfied.
  - **c** Common stopping criteria used in the Selective Search algorithm:
    - Number of Region Proposals: Set a predetermined fixed number of proposals to generate.
    - **Hierarchy Size:** Limit the size of the hierarchy or the number of levels considered.
    - **Segment Size:** Smaller segments are likely to capture fine details but may not be as meaningful.
- Step 5: Generate Bounding Boxes: For each group of similar regions, create a bounding box that tightly encapsulates the regions in the group. These bounding boxes represent the final region proposals.
- Step 6: Object Detection: Pass the selected region proposals to an object detection model (e.g., a classifier or a deep neural network) to identify and classify objects within those regions.

STUDENTS-HUB.com

# **Region Proposals via Selective Search**

25



STUDENTS-HUB.com

#### Outline

#### 26

- Introduction
- Traditional Computer Vision Techniques:
  - Sliding Window:
  - Region proposal approach (Selective Search)

#### Deep Learning-Based Techniques:

- Two-Stage Detectors:
  - R-CNN (Regions with CNN features)
- One-Stage Detectors:
  - YOLO (You Only Look Once)
- Other Detectors
- Evaluation Measures
- Benchmarks



STUDENTS-HUB.com

27

#### 28

- Anchor-based object detection relies on predefined anchors, which are small rectangular boxes of various sizes and aspect ratios, to represent potential object locations.
  - The object detection model predicts the bounding box of an object relative to the anchor that best fits it.
  - Anchor boxes serve as reference templates at different scales and aspect ratios, helping the model efficiently handle objects of various sizes and shapes.
  - This approach is widely used and has proven to be effective in various object detection tasks.
  - Advantages of Anchor-Based:
    - Effective for handling objects of various scales and aspect ratios.
    - Provides a systematic way to propose candidate regions for object detection.
- Anchor-free object detection, on the other hand, does not use predefined anchors. Instead, it directly predicts the bounding box coordinates of objects, typically represented by center points and dimensions.
  - Advantages of Anchor-free:
    - Simplifies the model architecture and training process.
    - Avoids the need for anchor design, making the model more flexible.

#### STUDENTS-HUB.com

29



The left side is the anchor-based method which uses the fixed different ratio aspects anchors to locate the location of an object, and the right side is the anchor-free method that directly estimate the bounding box.

STUDENTS-HUB.com

- 30
- **Two-Stage Object Detection**: consist of two main stages: region proposal and classification.
  - **Region Proposal Stage:** The first stage generates a set of region proposals, which 1. are potential locations of objects in the image. This stage typically uses a specialized algorithm, such as Selective Search or Region Proposal Network (RPN), to identify regions that are likely to contain objects.
  - 2. Classification Stage: The second stage classifies each region proposal as either an object or background. This stage typically uses a convolutional neural network (CNN) classifier to extract features from each region proposal and classify it based on those features.
- **One-Stage Object Detection:** perform both region proposal and classification in a single stage. One-stage object detection algorithms typically use a CNN to predict bounding boxes and class labels directly from the input image. The CNN is trained on a dataset of images that have been labeled with bounding boxes and class labels. STUDENTS-HUB.com

- Both two-stage and one-stage detectors have their strengths and trade-offs, making them suitable for different applications based on requirements such as speed, accuracy, and computational resources.
  - **Speed:** One-stage detectors are generally faster.
  - Accuracy: Two-stage detectors often provide higher accuracy due to the separate region proposal and classification stages.
  - Complexity: One-stage detectors tend to be simpler and computationally more efficient.
- If accuracy is the primary concern, then a two-stage algorithm is typically the best choice.
- One-stage object detection algorithms are often used in applications where speed is critical, such as real-time object detection in video streams.

STUDENTS-HUB.com

31

#### **Object Detection - Past vs. Present**



#### Outline

#### 33

- Introduction
- Traditional Computer Vision Techniques:
  - Sliding Window:
  - Region proposal approach (Selective Search)
- Deep Learning-Based Techniques:
  - Two-Stage Detectors:
    - R-CNN (Regions with CNN features)
  - One-Stage Detectors:
    - YOLO (You Only Look Once)
- Other Detectors
- Evaluation Measures
- Benchmarks

# R-CNN (Regions with CNN features)

- 34
- R-CNN (Region-based Convolutional Neural Network) is a two-stage object detection framework that combines region proposals with convolutional neural networks.
- □ R-CNN consists of four main steps:
  - Region Proposal: R-CNN uses a region proposal algorithm (such as selective search) to generate candidate regions in an image that are likely to contain objects.
  - Feature Extraction: Each proposed region is cropped and warped, and the features within these regions are extracted using a pre-trained Convolutional Neural Network (CNN).
  - Object Classification: The CNN-extracted features are used to train Support Vector Machines (SVMs) for object classification, determining the presence of objects and their respective classes.
- Bounding Box Regression: A separate regression model is trained to refine the coordinates of the proposed bounding boxes, improving the localization accuracy of the detected objects. This ensures that the bounding boxes tightly enclose the detected objects.
   STUDENTS-HUB.com

# **R-CNN (Regions with CNN features)**



### **R-CNN: Bounding Box Regression**

- 36
- Bounding Box regression is a metric for measuring how well predicted bounding boxes captures objects.
- It aims to refine the predicted bounding boxes generated by the region proposal stage to accurately enclose the detected objects (ground truth bounding boxes).
- Bounding Box Regression utilizes machine learning regression technique to learn the offsets between the predicted bounding boxes and the true object locations.
- The target labels for regression are the corrections needed for the coordinates of the proposed bounding boxes to match the coordinates of the ground truth bounding boxes.
Bounding box regression: Predict "transform" to correct the Bbox Rol: 4 numbers (t<sub>x</sub>, t<sub>y</sub>, t<sub>h</sub>, t<sub>w</sub>)

Consider a region proposal with center  $(p_x, p_y)$ , width  $p_w$ , height  $p_h$ 

Model predicts a transform  $(t_x, t_y, t_w, t_h)$ to correct the region proposal

The output box is defined by:

 $b_{x} = p_{x} + p_{w}t_{x}$   $b_{y} = p_{y} + p_{h}t_{y}$   $b_{w} = p_{w}\exp(t_{w})$   $b_{h} = p_{h}\exp(t_{h})$ STUDENTS-HUB.com

37

Shift center by amount relative to proposal size

Scale proposal; exp ensures that scaling factor is > 0



Consider a region proposal with center  $(p_x, p_y)$ , width  $p_w$ , height  $p_h$ 

Model predicts a transform  $(t_x, t_y, t_w, t_h)$ to correct the region proposal

The output box is:  $b_x = p_x + p_w t_x$   $b_y = p_y + p_h t_y$   $b_w = p_w \exp(t_w)$  $b_h = p_h \exp(t_h)$ 

38

When transform is 0, output = proposal

L2 regularization encourages leaving proposal unchanged



39

Consider a region proposal with center  $(p_x, p_y)$ , width  $p_w$ , height  $p_h$ 

Model predicts a <u>transform</u>  $(t_x, t_y, t_w, t_h)$ to correct the region proposal

The output box is:  $b_x = p_x + p_w t_x$   $b_y = p_y + p_h t_y$   $b_w = p_w \exp(t_w)$   $b_h = p_h \exp(t_h)$  Scale / Translation invariance: Transform encodes *relative* difference between proposal and output; important since CNN doesn't see absolute size or position after cropping



Consider a region proposal with center  $(p_x, p_y)$ , width  $p_w$ , height  $p_h$ 

Model predicts a transform  $(t_x, t_y, t_w, t_h)$ to correct the region proposal

The output box is:  $b_x = p_x + p_w t_x$   $b_y = p_y + p_h t_y$   $b_w = p_w \exp(t_w)$  $b_h = p_h \exp(t_h)$ 

40

Given proposal and target output, we can solve for the transform the network should output:

 $t_x = (b_x - p_x)/p_w$   $t_y = (b_y - p_y)/p_h$   $t_w = \log(b_w/p_w)$  $t_h = \log(b_h/p_h)$ 



### R-CNN Training

41

#### Input Image



Categorize each region proposal as positive, negative, or neutral based on overlap with ground-truth boxes:

#### Positive: > 0.5 IoU with a GT box Negative: < 0.3 IoU with all GT boxes Neutral: between 0.3 and 0.5 IoU with GT boxes

STUDENTS-HUB.com



STUDENTS-HUB.com

42

### **R-CNN** Test-Time

### Input Image



- 1. Run proposal method
- 2. Run CNN on each proposal to get class scores, transforms
- 3. Threshold class scores to get a set of detections

2 problems:

- CNN often outputs overlapping boxes
- How to set thresholds?

STUDENTS-HUB.com

# **Overlapping Boxes: Non-Max Suppression (NMS)**

44

**Problem**: Object detectors often output many overlapping detections:

# Solution: Post-process raw detections using Non-Max Suppression (NMS)

- 1. Select next highest-scoring box
- Eliminate lower-scoring boxes with IoU > threshold (e.g. 0.7)
- 3. If any boxes remain, GOTO 1



# **Overlapping Boxes: Non-Max Suppression (NMS)**

45

**Problem**: Object detectors often output many overlapping detections:

Solution: Post-process raw detections using Non-Max Suppression (NMS)

- 1. Select next highest-scoring box
- Eliminate lower-scoring boxes with IoU > threshold (e.g. 0.7)
- 3. If any boxes remain, GOTO 1

loU(∎, ∎) = 0.74



# **Overlapping Boxes: Non-Max Suppression (NMS)**

**Problem**: Object detectors often output many overlapping detections:

Solution: Post-process raw detections using Non-Max Suppression (NMS)

46

- 1. Select next highest-scoring box
- Eliminate lower-scoring boxes with IoU > threshold (e.g. 0.7)
- 3. If any boxes remain, GOTO 1



### Fast R-CNN

- Instead of extracting CNN feature vectors independently for each region proposal, this model aggregates them into one CNN forward pass over the entire image and the region proposals share this feature matrix.
- Then the same feature matrix is branched out to be used for learning the object classifier and the bounding-box regressor. In conclusion, computation sharing speeds up R-CNN.



Fast R-CNN is much faster in both training and testing time. However, the improvement is not dramatic because the region proposals are generated separately by another model and that is very expensive. Uploaded By: anonymous

# Fast R-CNN vs "Slow" R-CNN

# **Fast R-CNN**: Apply differentiable cropping to shared image features



**"Slow" R-CNN**: Apply differentiable cropping to shared image features



STUDENTS-HUB.com

**48** 

# **Faster R-CNN**

An intuitive speedup solution is to integrate the region proposal algorithm into the CNN model. Faster R-CNN is doing exactly this: construct a single, unified model composed of RPN (region proposal network) and fast R-CNN with shared convolutional feature layers.



STUDENTS-HUB.com

### **Comparison**





STUDENTS-HUB.com

### Outline

#### 51

- Introduction
- Traditional Computer Vision Techniques:
  - Sliding Window:
  - Region proposal approach (Selective Search)
- Deep Learning-Based Techniques:
  - Two-Stage Detectors:
    - R-CNN (Regions with CNN features)
  - One-Stage Detectors:
    - YOLO (You Only Look Once)
- Other Detectors
- Evaluation Measures
- Benchmarks

# YOLO (You Only Look Once)

- 52
- The R-CNN family of techniques primarily use regions to localize the objects within the image. The network does not look at the entire image, only at the parts of the images which have a higher chance of containing an object.
- The YOLO framework (You Only Look Once) on the other hand, deals with object detection in a different way. It takes the entire image in a single instance and predicts the bounding box coordinates and class probabilities for these boxes.
- You Only Look Once (YOLO) is a state-of-the-art, real-time object detection algorithm introduced in 2015
- The authors frame the object detection problem as a regression problem instead of a classification task by spatially separating bounding boxes and associating probabilities to each of the detected images using a single convolutional neural network (CNN).

# How does the YOLO Framework Function?

#### 53

- First, image is split into a SxS grid
- For each grid square, generate B bounding boxes
- For each bounding box, there are
   5 predictions
- Image classification and localization are applied on each cell grid.
- Each cell in the grid is responsible for localizing and predicting the class of the object that it covers, along with the probability/confidence value.



STUDENTS-HUB.com

# How does the YOLO Framework Function?

#### 54

 YOLO predicts both bounding box parameters and class probabilities for each grid cell in a single forward pass through the network.

#### Bounding Box Predictions:

 For each grid cell, YOLO predicts bounding box parameters. These parameters include the coordinates (x, y) of the box's center, its width (w), and its height (h).

#### Class Predictions:

 YOLO predicts class probabilities for each bounding box within a grid cell. The model assigns a probability score for each possible class that the object might belong to (e.g., person, car, dog).

#### Confidence Score:

Each bounding box is associated with a confidence score. This score reflects the model's confidence that the predicted box contains an object.

# How does the YOLO Framework Function?



STUDENTS-HUB.com

### **YOLO Architecture**

56

YOLO architecture is similar to <u>GoogleNet</u>. As illustrated below, it has overall 24 convolutional layers, four max-pooling layers, and two fully connected layers.



Uploaded By: anonymous

STUDENTS-HUB.com

## **YOLO Architecture**

#### 57

- The architecture works as follows:
  - Resizes the input image into 448x448 before going through the convolutional network.
  - A 1x1 convolution is first applied to reduce the number of channels, which is then followed by a 3x3 convolution to generate a cuboidal output.
  - The activation function under the hood is ReLU, except for the final layer, which uses a linear activation function.
  - Some additional techniques, such as batch normalization and dropout, respectively regularize the model and prevent it from overfitting.

# **YOLO Training**

- **58**
- We need to pass the labelled data to the model in order to train it.
- Suppose we have divided the image into a grid of size 3 X 3 and there are a total of 2 classes which we want the objects to be classified into.
- YOLO determines the attributes of the bounding boxes using a single regression module in the following format, where Y is the final vector representation for each bounding box.

Y = [pc, bx, by, bh, bw, c1, c2]

Here,

- pc corresponds to the probability score of the grid containing an object.
- bx, by are the x and y coordinates of the center of the bounding box with respect to the enveloping grid cell.
- bh, bw correspond to the height and the width of the bounding box with respect to the enveloping grid cell.

■ c1 and c2 correspond to the two classes. STUDENTS-HUB.com

# **YOLO Training**

#### 59

- □ YOLO is a regression algorithm. What is X? What is Y?
  - X is simple, just an image width (in pixels) \* height (in pixels) \* RGB values
  - Y is a tensor of size S \* S \* (B \* 5 + C)
  - B\*5 + C term represents the predictions + class predicted distribution for a grid block



STUDENTS-HUB.com

### How to Encode Bounding Boxes?





Bounding box centers



## **Intersection Over Unions or IOU**

- **61**
- Most of the time, a single object in an image can have multiple grid box candidates for prediction, even though not all of them are relevant.
- The goal of the IOU (a value between 0 and 1) is to discard such grid boxes to only keep those that are relevant.
- □ Here is the logic behind it:
  - The user defines its IOU selection threshold, which can be, for instance, 0.5.
  - Then YOLO computes the IOU of each grid cell which is the Intersection area divided by the Union Area.
  - Finally, it ignores the prediction of the grid cells having an IOU ≤ threshold and considers those with an IOU > threshold.
- Intuitively, the more you increase the threshold, the better the predictions become.

STUDENTS-HUB.com

### **Intersection Over Unions or IOU**

**62** 



We can observe that the object originally had two grid candidates, then only "Grid 2" was selected at the end.

STUDENTS-HUB.com

# **Non-Max Suppression or NMS**

- 63
- Setting a threshold for the IOU is not always enough because an object can have multiple boxes with IOU beyond the threshold, and leaving all those boxes might include noise.
- Non-Max suppression algorithm:
  - 1. Discard all the boxes having probabilities less than or equal to a pre-defined threshold (say, 0.5)
  - 2. For the remaining boxes:
    - 1. Pick the box with the highest probability and take that as the output prediction
    - 2. Discard any other box which has IoU greater than the threshold with the output box from the above step
  - 3. Repeat step 2 until all the boxes are either taken as the output prediction or discarded

STUDENTS-HUB.com

### **Non-Max Suppression or NMS**





STUDENTS-HUB.com

64

# **YOLO Objective Function**

For YOLO, we need to minimize the following loss
 Sum squared error is used

$$\lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^{B} \mathbb{1}_{ij}^{\text{obj}} \left[ (x_i - \hat{x}_i)^2 + (y_i - \hat{y}_i)^2 \right] + \lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^{B} \mathbb{1}_{ij}^{\text{obj}} \left[ \left( \sqrt{w_i} - \sqrt{\hat{w}_i} \right)^2 + \left( \sqrt{h_i} - \sqrt{\hat{h}_i} \right)^2 \right] \right]$$

$$+ \lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^{B} \mathbb{1}_{ij}^{\text{obj}} \left[ \left( \sqrt{w_i} - \sqrt{\hat{w}_i} \right)^2 + \left( \sqrt{h_i} - \sqrt{\hat{h}_i} \right)^2 \right]$$

$$+ \sum_{i=0}^{S^2} \sum_{j=0}^{B} \mathbb{1}_{ij}^{\text{obj}} \left( C_i - \hat{C}_i \right)^2$$

$$+ \lambda_{\text{noobj}} \sum_{i=0}^{S^2} \sum_{j=0}^{B} \mathbb{1}_{ij}^{\text{noobj}} \left( C_i - \hat{C}_i \right)^2$$

$$+ \sum_{i=0}^{S^2} \mathbb{1}_{ij}^{B} \mathbb{1}_{ij}^{\text{noobj}} \left( C_i - \hat{C}_i \right)^2$$

$$+ \sum_{i=0}^{S^2} \mathbb{1}_{ij}^{B} \mathbb{1}_{ij}^{\text{noobj}} \left( C_i - \hat{C}_i \right)^2$$

$$+ \sum_{i=0}^{S^2} \mathbb{1}_{ij}^{B} \mathbb{1}_{ij}^{\text{noobj}} \left( C_i - \hat{C}_i \right)^2$$

$$+ \sum_{i=0}^{S^2} \mathbb{1}_{ij}^{B} \mathbb{1}_{ij}^{O(i)} \left( C_i - \hat{C}_i \right)^2$$

$$+ \sum_{i=0}^{S^2} \mathbb{1}_{ij}^{B} \mathbb{1}_{ij}^{O(i)} \left( C_i - \hat{C}_i \right)^2$$

$$+ \sum_{i=0}^{S^2} \mathbb{1}_{ij}^{B} \sum_{c \in \text{classes}} \left( p_i(c) - \hat{p}_i(c) \right)^2$$

$$Class loss, minimize loss between true class of object in grid box$$

#### STUDENTS-HUB.com

65

### **YOLO Generations**





#### STUDENTS-HUB.com

### **Some Comparisons**

67



### **Some Comparisons**

**68** 



Latency A100 TensorRT FP16 (ms/img)

STUDENTS-HUB.com

### Outline

#### 69

- Introduction
- Traditional Computer Vision Techniques:
  - Sliding Window:
  - Region proposal approach (Selective Search)
- Deep Learning-Based Techniques:
  - Two-Stage Detectors:
    - R-CNN (Regions with CNN features)
  - One-Stage Detectors:
    - YOLO (You Only Look Once)

### Other Detectors

- Evaluation Measures
- Benchmarks

### **Other Detectors**

#### 70

#### Mask R-CNN

- An extension of Faster R-CNN, Mask R-CNN,
- Includes an additional branch for instance segmentation.
- It simultaneously predicts object masks along with bounding boxes and class scores.

#### SSD (Single Shot Multibox Detector)

 Similar to YOLO, SSD is a single-shot object detector that predicts bounding boxes and class scores at multiple scales.

#### RetinaNet

- Addressing the issue of class imbalance in object detection
- It is known for its strong performance on datasets with a large number of background samples.

#### EfficientDet

EfficientDet focuses on balancing model efficiency and accuracy by using a compound scaling method to optimize the model's depth, width, and resolution.

#### Cascade R-CNN

- An extension of Faster R-CNN that aims to improve the accuracy of object detection by using a cascade structure.
- It refines the bounding box predictions in a cascaded manner.

#### CenterNet

- Directly predicts object centers.
- □ It achieves high accuracy with a single-stage architecture.

#### STUDENTS-HUB.com

### **Detectors Performance**





### **Detectors Performance**



STUDENTS-HUB.com
### **Detectors Performance**



### Outline

- Introduction
- Traditional Computer Vision Techniques:
  - Sliding Window:
  - Region proposal approach (Selective Search)
- Deep Learning-Based Techniques:
  - Two-Stage Detectors:
    - R-CNN (Regions with CNN features)
  - One-Stage Detectors:
    - YOLO (You Only Look Once)
- Other Detectors
- Evaluation Measures
- Benchmarks

#### 75

- Mean Average Precision (mAP) is a commonly used metric for evaluating the performance of object detectors in computer vision tasks, particularly in the context of object detection.
- mAP is calculated based on:

#### Precision and Recall

- Precision: Precision is the ratio of true positive predictions to the total number of positive predictions (true positives + false positives). It measures the accuracy of positive predictions.
- Recall: Recall is the ratio of true positive predictions to the total number of actual positive instances (true positives + false negatives). It measures the ability of the model to capture all positive instances.

#### Intersection over Union (IoU)

 IoU is a measure of the overlap between the predicted bounding box and the ground truth bounding box.

#### Precision-Recall Curve

The precision-recall curve is generated by varying the confidence threshold of the model and calculating precision and recall at each threshold.

#### Average Precision (AP)

• AP is calculated by computing the area under the precision-recall curve. It represents the average precision across all recall levels.

#### Mean Average Precision (mAP)

mAP is the mean of the average precision values calculated for each class. It provides an overall assessment of the detector's performance across different object categories.

#### STUDENTS-HUB.com

# mAP: Step-by-step

- 76
- 1. Run object detector on all test images (with NMS)
- For each category, compute Average Precision (AP) = area under Precision vs Recall Curve
  - 1. For each detection (highest score to lowest score)
    - 1. If it matches some GT box with IoU > 0.5, mark it as positive and eliminate the GT
    - 2. Otherwise mark it as negative
    - 3. Plot a point on PR Curve
  - 2. Average Precision (AP) = area under PR curve
- 3. Mean Average Precision (mAP) = average of AP for each category
- For "COCO mAP": Compute mAP@thresh for each IoU threshold (0.5, 0.55, 0.6, ..., 0.95) and take average



### All dog detections sorted by score



#### All ground-truth dog boxes

STUDENTS-HUB.com





STUDENTS-HUB.com





STUDENTS-HUB.com





STUDENTS-HUB.com





STUDENTS-HUB.com





STUDENTS-HUB.com







STUDENTS-HUB.com

STUDENTS-HUB.com

### Evaluating Object Detectors: Mean Average Precision (mAP)

- 84
- Run object detector on all test images (with NMS)
- For each category, compute Average Precision
   (AP) = area under Precision vs Recall Curve
  - 1. For each detection (highest score to lowest score)
    - If it matches some GT box with IoU > 0.5, mark it as positive and eliminate the GT
    - 2. Otherwise mark it as negative
    - 3. Plot a point on PR Curve
  - 2. Average Precision (AP) = area under PR curve
- Mean Average Precision (mAP) = average of AP for each category

```
Car AP = 0.65
Cat AP = 0.80
Dog AP = 0.86
mAP@0.5 = 0.77
```

STUDENTS-HUB.com

### Evaluating Object Detectors: Mean Average Precision (mAP)

- Run object detector on all test images (with NMS)
- For each category, compute Average Precision (AP) = area under Precision vs Recall Curve
  - 1. For each detection (highest score to lowest score)
    - If it matches some GT box with IoU > 0.5, mark it as positive and eliminate the GT
    - 2. Otherwise mark it as negative
    - 3. Plot a point on PR Curve
  - 2. Average Precision (AP) = area under PR curve
- 3. Mean Average Precision (mAP) = average of AP for each category
- 4. For "COCO mAP": Compute mAP@thresh for each IoU threshold (0.5, 0.55, 0.6, ..., 0.95) and take average

```
mAP@0.5 = 0.77
mAP@0.55 = 0.71
mAP@0.60 = 0.65
...
mAP@0.95 = 0.2
```

# **Object Detection Benchmarks**

- MS COCO (Microsoft Common Objects in Context)
  - One of the largest and most widely used object detection benchmarks.
  - Consists of 328,000 images with annotations for 80 object categories.
  - Provides evaluation metrics such as Average Precision (AP) and mean Average Precision (mAP).



# **Object Detection Benchmarks**

## PASCAL VOC Challenge

- This benchmark is older than COCO, but it is still widely used.
- Consists of about 20,000 images with annotations for 20 object categories.
- Evaluation metrics include mAP.



# **Object Detection Benchmarks**

#### 88

#### KITTI Vision Benchmark Suite

- This benchmark focuses on object detection in autonomous driving scenarios.
- Consists of 80,000 images with annotations for 3 object categories (cars, pedestrians, and cyclists).
- Provides metrics like precision, recall, and F1 score.



# Acknowledgement

#### 89

- □ The material in these slides are based on:
  - Digital Image Processing: Rafael C. Gonzalez, and Richard
  - Forsythe and Ponce: Computer Vision: A Modern Approach
  - Rick Szeliski's book: Computer Vision: Algorithms and Applications
  - cs131@ Stanford University
  - cs131n@ Stanford University
  - CS198-126@ University of California, Berkely
  - CAP5415@ University of Central Florida
  - CSW182 @ University of California, Berkely
  - Deep Learning Lecture Series @UCL
  - EECS 498.008 @ University of Michigan
  - CSE576 @ Washington University
  - 11-785@ Carnegie Mellon University
  - CSCI1430@ Brown University
  - Computer Vision@ Bonn University
  - ICS 505@ KFUPM
- Digital Image Processing@ University of Jordan STUDENTS-HUB.com