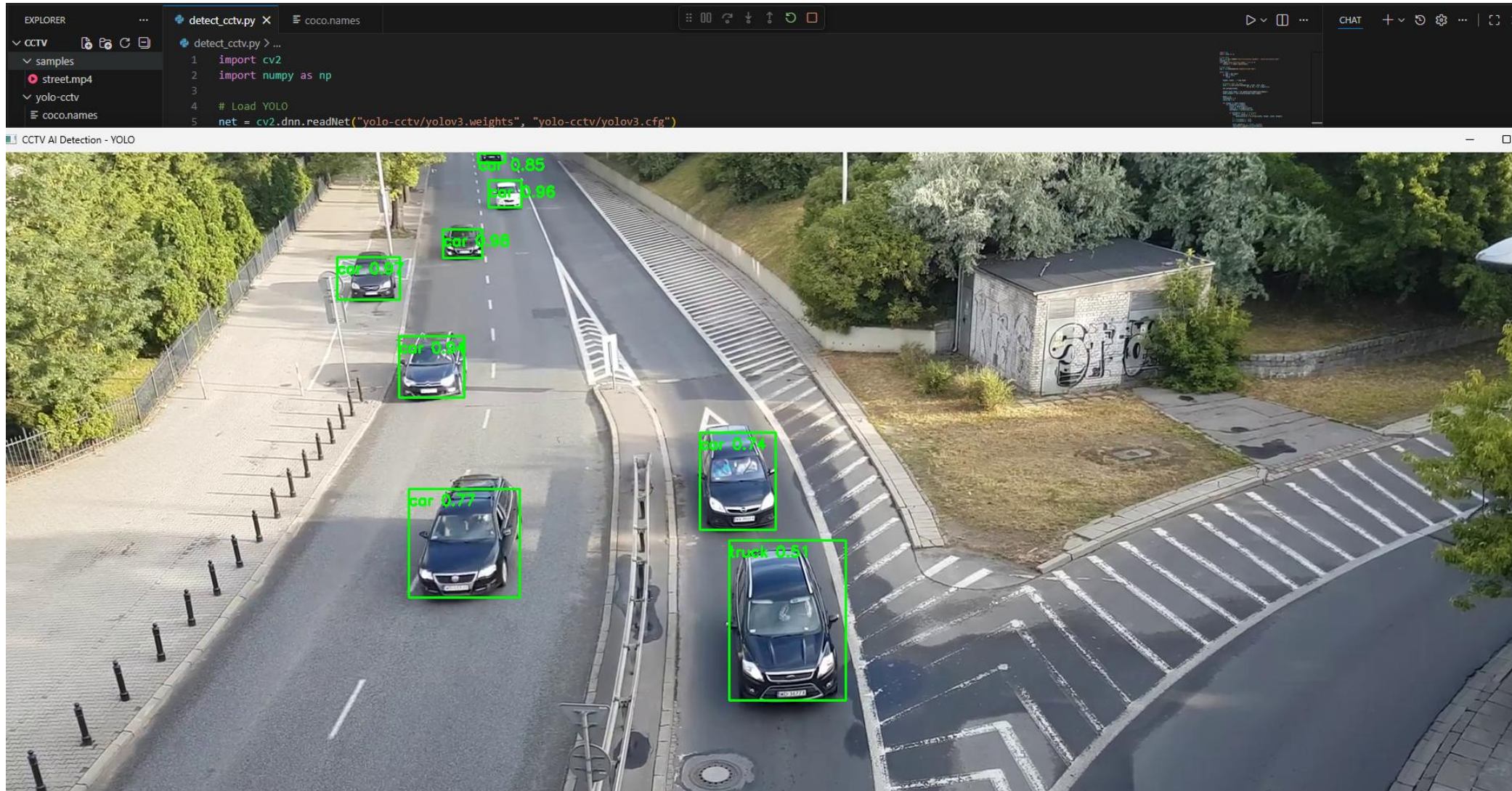


YOLO : You Only Look Once - Real Time Object Detection

Krijimi i programit për identifikim të objekteve që përdorë modelin e trajnuar me YOLO



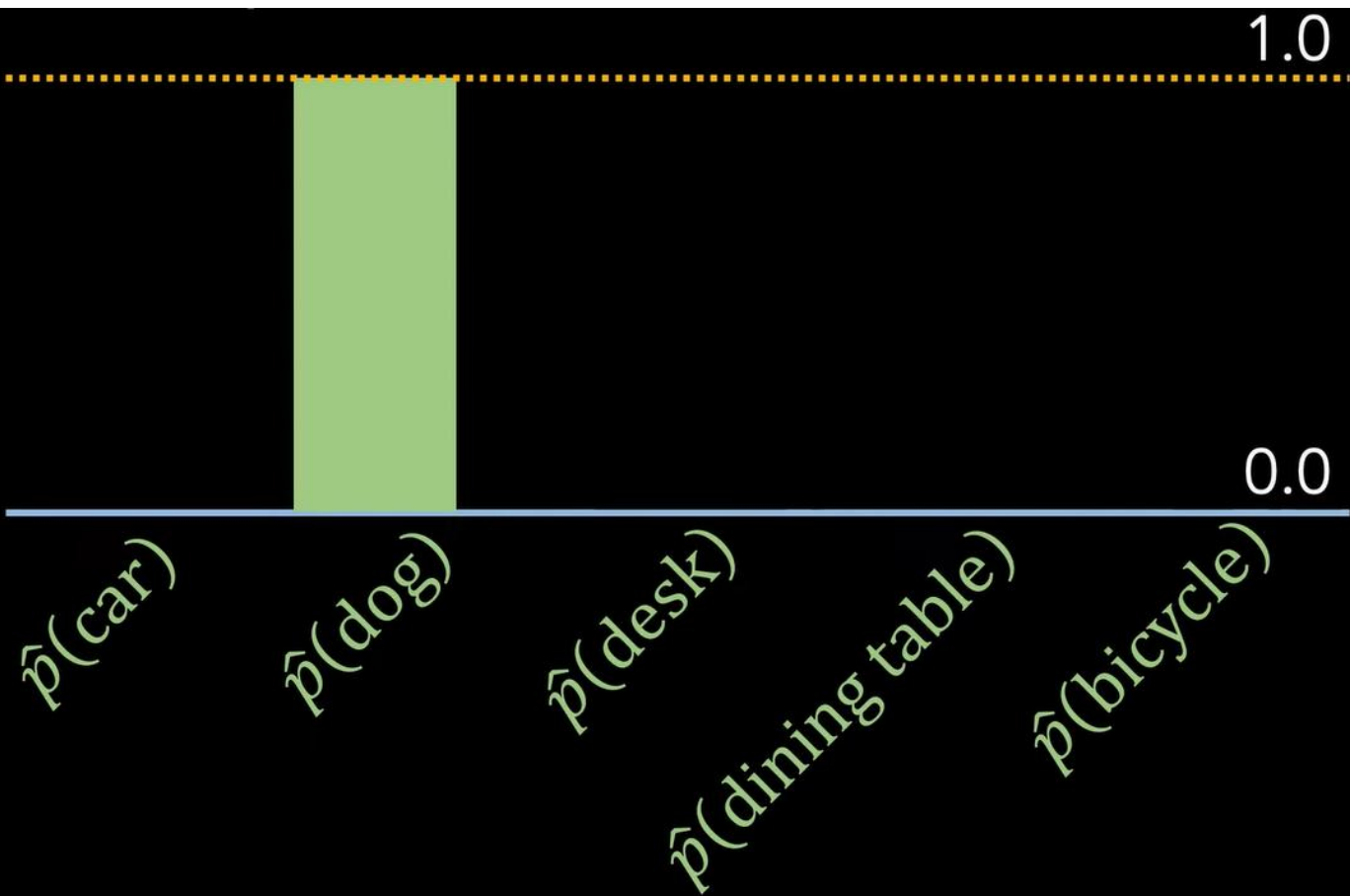
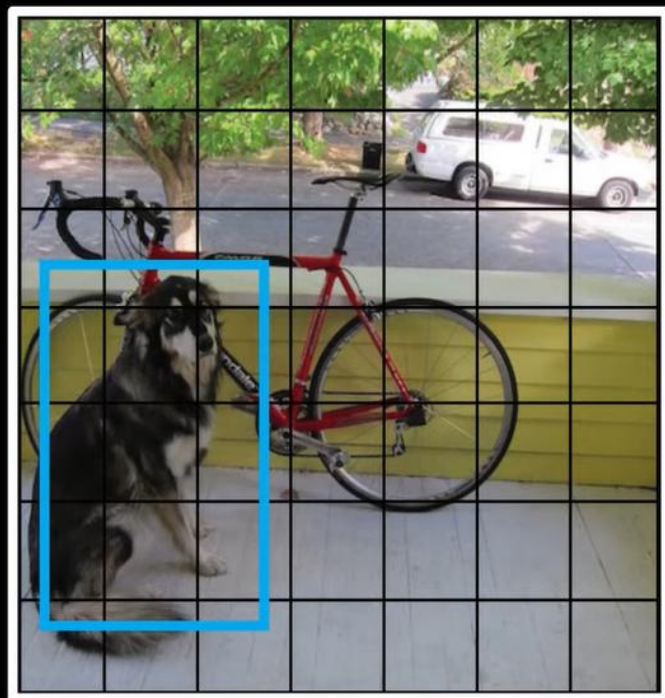
Komponentët kryesorë të programit



```
detect_cctv.py > ...
1  import cv2
2  import numpy as np
3
4  # Load YOLO
5  net = cv2.dnn.readNet("yolo-cctv/yolov3.weights", "yolo-cctv/yolov3.cfg")
6  classes = []
7  with open("yolo-cctv/coco.names", "r") as f:
8      classes = f.read().splitlines()
9
10 # Input Video
11 cap = cv2.VideoCapture("samples/street.mp4")
12
13 while True:
14     _, img = cap.read()
15     if img is None:
16         break
17
18     height, width, _ = img.shape
19
20     # Prepare input for YOLO
21     blob = cv2.dnn.blobFromImage(img, 1/255, (416, 416),
22                                   (0, 0, 0), True, crop=False)
23     net.setInput(blob)
24
25     output_layer_names = net.getUnconnectedOutLayersNames()
26     layer_outputs = net.forward(output_layer_names)
27
28     boxes = []
29     confidences = []
30     class_ids = []
31
```

- yolov3.cfg = si ndërtohet rrjeti
- yolov3.weights = rezultati i trajnimit
- coco.names = emrat e objekteve

Trajnimi i modelit : Ground truth(class probabilities)



Ground Truth është vlera e vërtetë dhe e saktë e objektit në imazh.

Bounding box-i i Ground Truth konsiderohet 100% i saktë dhe shërben si standard për të matur sa afër janë parashikimet e modelit gjatë trajnimit

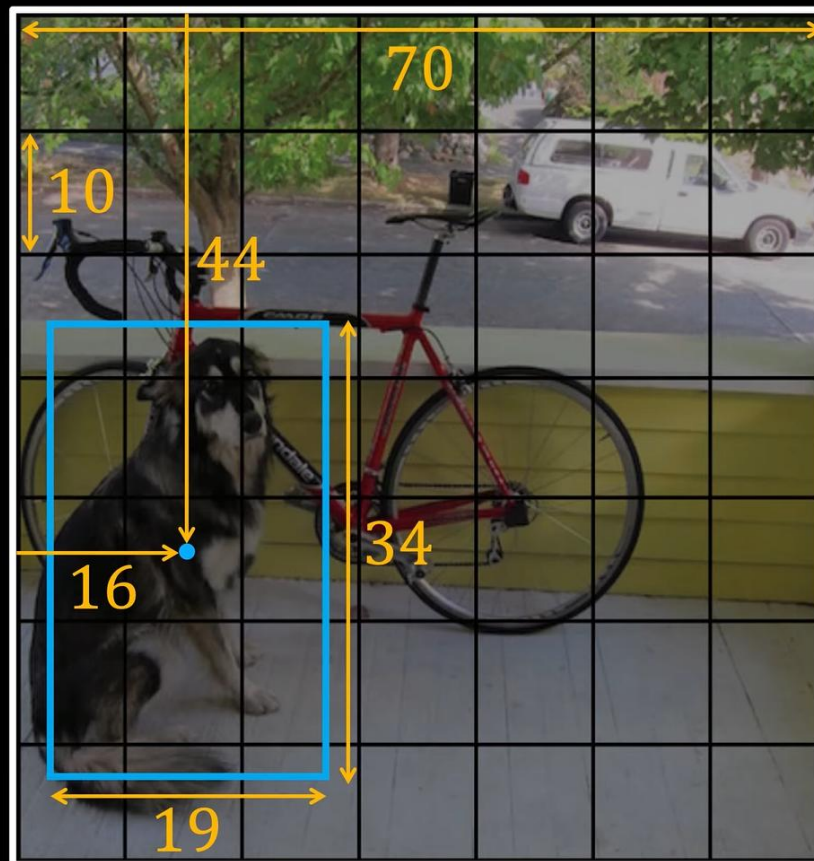
Trajnimi i modelit : Ground truth(box coordinates)

$$x = \frac{16 \% 10}{10} = 0.60$$

$$y = \frac{44 \% 10}{10} = 0.40$$

$$w = \frac{19}{70} = 0.27$$

$$h = \frac{34}{70} = 0.49$$




Box Coordinates janë vlerat që përshkruajnë vendndodhjen dhe madhësinë e një objekti në imazh. Ato përdoren nga YOLO për të vizatuar bounding boxes dhe për të llogaritur saktësinë e parashikimeve kundrejt Ground Truth.

Pozita dhe madhësia e objektit i lejojnë YOLO-s jo vetëm të identifikojë çfarë objekti është në imazh, por gjithashtu edhe ku ndodhet dhe sa i madh është. Kjo e bën YOLO-n shumë të vlefshëm për aplikime reale si CCTV, monitorim trafiku dhe robotikë.

Trajnimi i modelit : Ground truth(confidence)

- Does an object appear in this grid cell or not?

 = 0

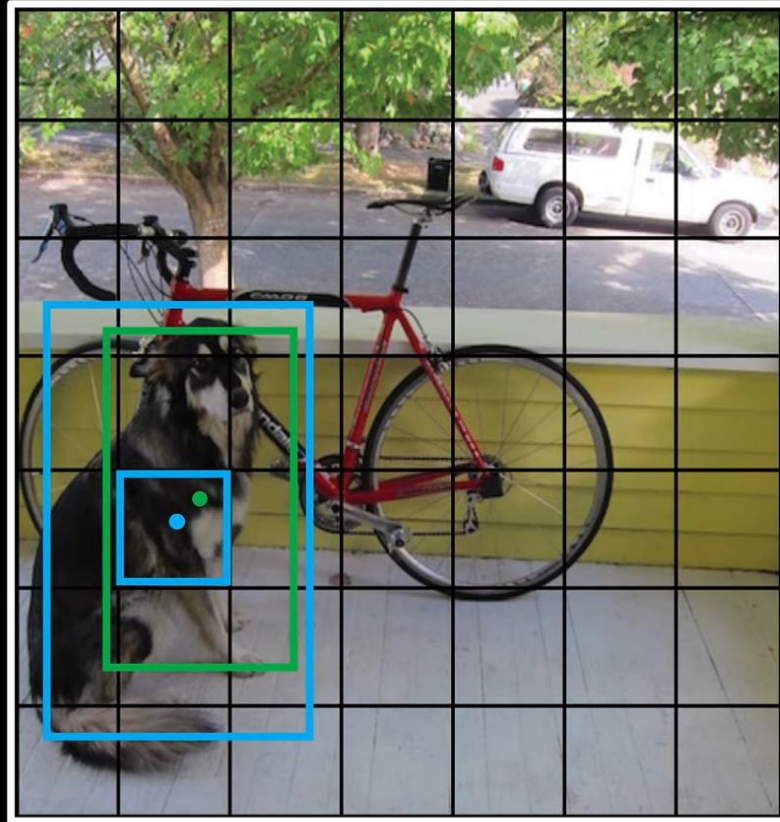
 = $\text{IoU}(\text{pred}, \text{true})$



Ground Truth Confidence është vlera binare nga dataset-i që tregon nëse një objekt ekziston në një qelizë specifike të grid-it. Përdoret për të llogaritur gabimin e konfidencës gjatë trajnimit — nëse modeli thotë se ka objekt aty ku nuk ka, ai penalizohet.

Trajnimi i modelit : Ground truth(confidence)

$$\hat{C} = \text{IoU}(\text{pred}, \text{true})$$



Gjatë trajnimit, IoU mat sa përputhen kutitë e parashikuara nga modeli me Ground Truth. Sa më i lartë IoU, aq më i saktë modeli. IoU përdoret për të llogaritur humbjen dhe për të udhëhequr përmirësimin e modelit gjatë mësimin.

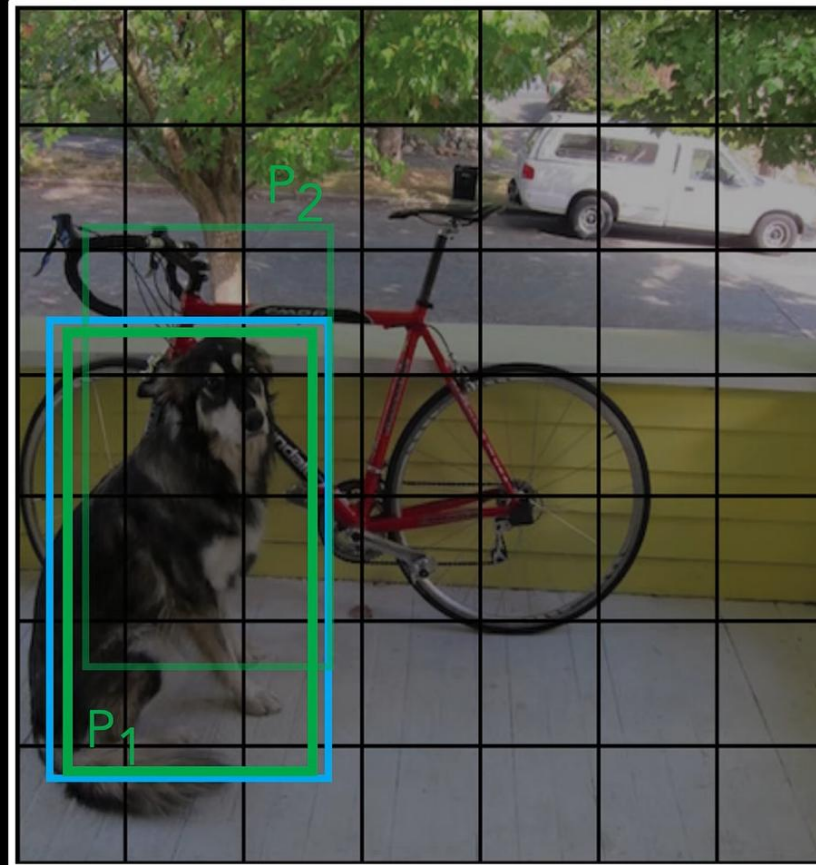
Gjatë trajnimit, modeli parashikon kutitë fillestare. Ne e krahasojmë secilën me Ground Truth duke përdorur IoU. Sa më shumë të përputhen kutitë, humbja është më e ulët. Modeli vazhdimisht përditëson parametrat derisa të arrijë konvergjencë — pra të mësojë të detektojë objektet me saktësi të lartë.

Trajnimi i modelit : Box selection(inference)

If $\text{IoU}(P_1, P_2) > \text{Threshold}$:

$$P = \text{argmax}(C(P_1), C(P_2))$$

(Non max suppression)



Non-Max Suppression është një teknikë që heq kutitë e tepërta të detektuara për të njëjtin objekt, duke zgjedhur vetëm kutinë me probabilitetin më të lartë dhe përputhjen më të mirë

Trajnimi i modelit : Loss function

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

Loss function

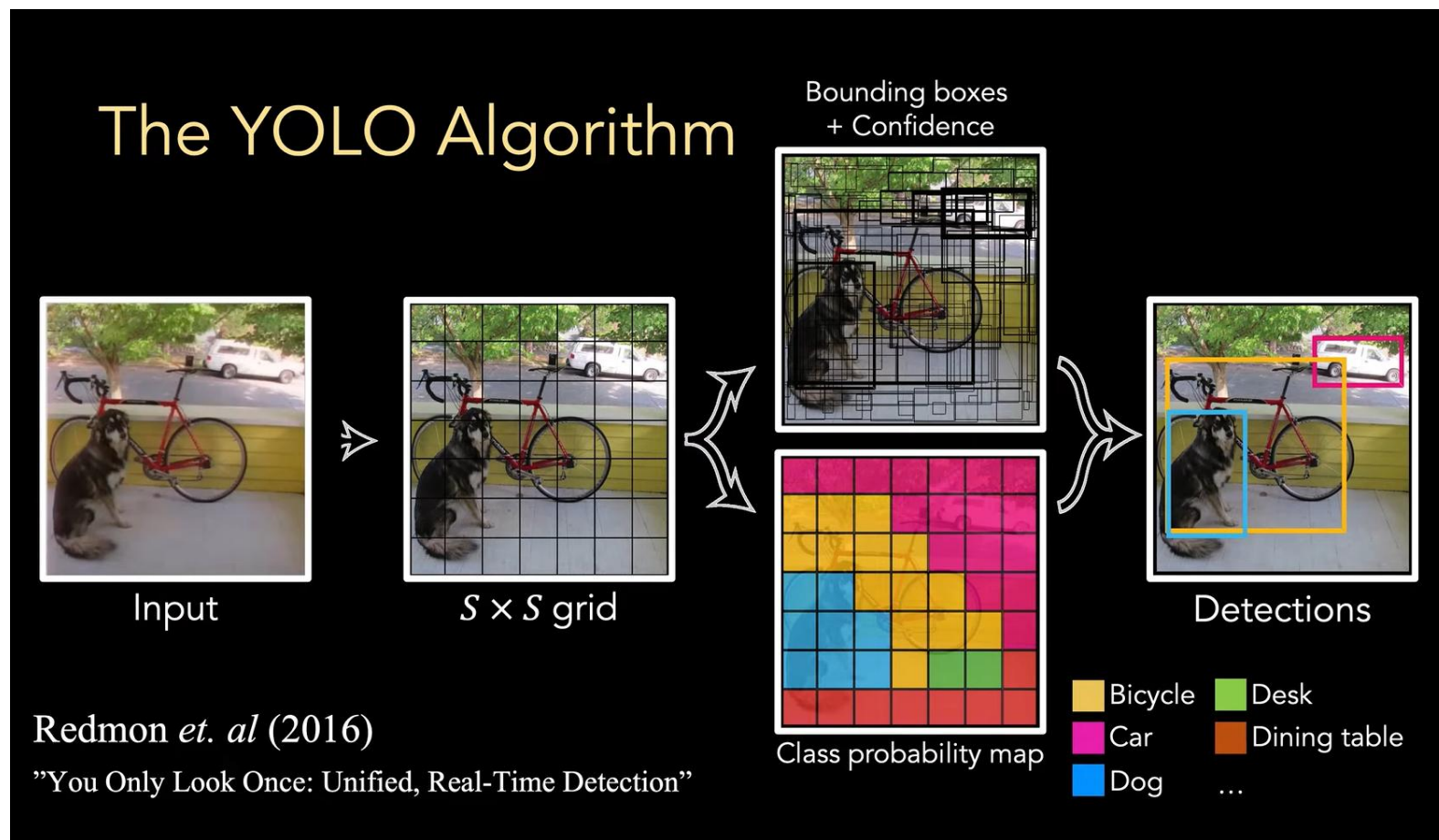
(Object detection as regression)

Loss Function në YOLO përbëhet nga tre komponentë kryesorë:

- Localization Loss për pozicionin dhe përmasat e kutive,
- Confidence Loss për vlerësimin e pranisë së objektit dhe
- Classification Loss për të identifikuar saktë klasën e objektit.

Gjatë trajnimit, modeli përditëson parametrat në mënyrë që të minimizohet vlera totale e humbjes dhe të përmirësojë saktësinë e detektimit.

Inference



Imazhi ndahet në 7×7 grid që çdo qelizë të jetë përgjegjëse për detektimin e objektit në zonën e saj. Kjo e thjeshton detektimin, e ul kompleksitetin dhe e bën YOLO-n jashtëzakonisht të shpejtë për aplikime në kohë reale si CCTV.

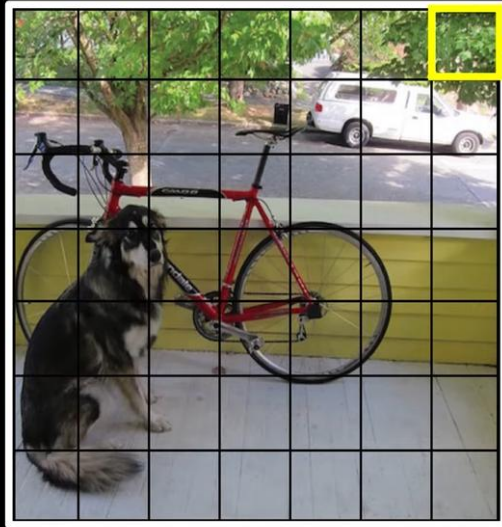
Bounding boxes tregojnë vendndodhjen e objektit, ndërsa confidence score tregon sa i sigurt është modeli në këtë parashikim.

Class Probability Map është grupi i probabiliteteve për secilën klasë objekti në çdo qelizë të grid-it. Ai i tregon modelit se çfarë objekti është

Inference - vektORIZIMI

The YOLO Algorithm $(B \times 5 + n)$

Output vector length



$S \times S$ grid

$[p(c_1), p(c_2), \dots, p(c_n)]$

Class probabilities

$[x_1, y_1, \sqrt{w_1}, \sqrt{h_1}, C_1]$

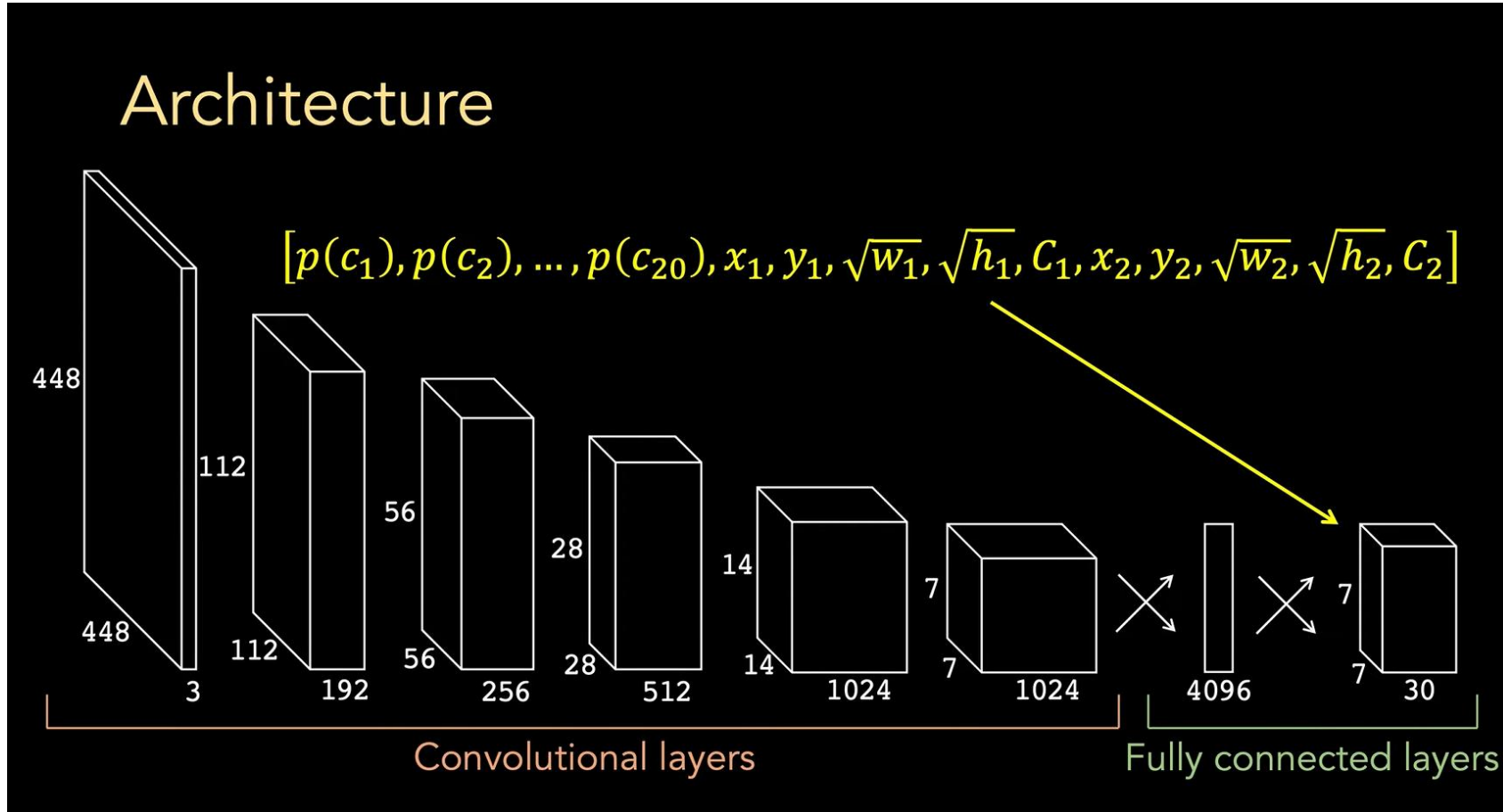
$[x_2, y_2, \sqrt{w_2}, \sqrt{h_2}, C_2]$

...

$[x_B, y_B, \sqrt{w_B}, \sqrt{h_B}, C_B]$

Bounding box predictions

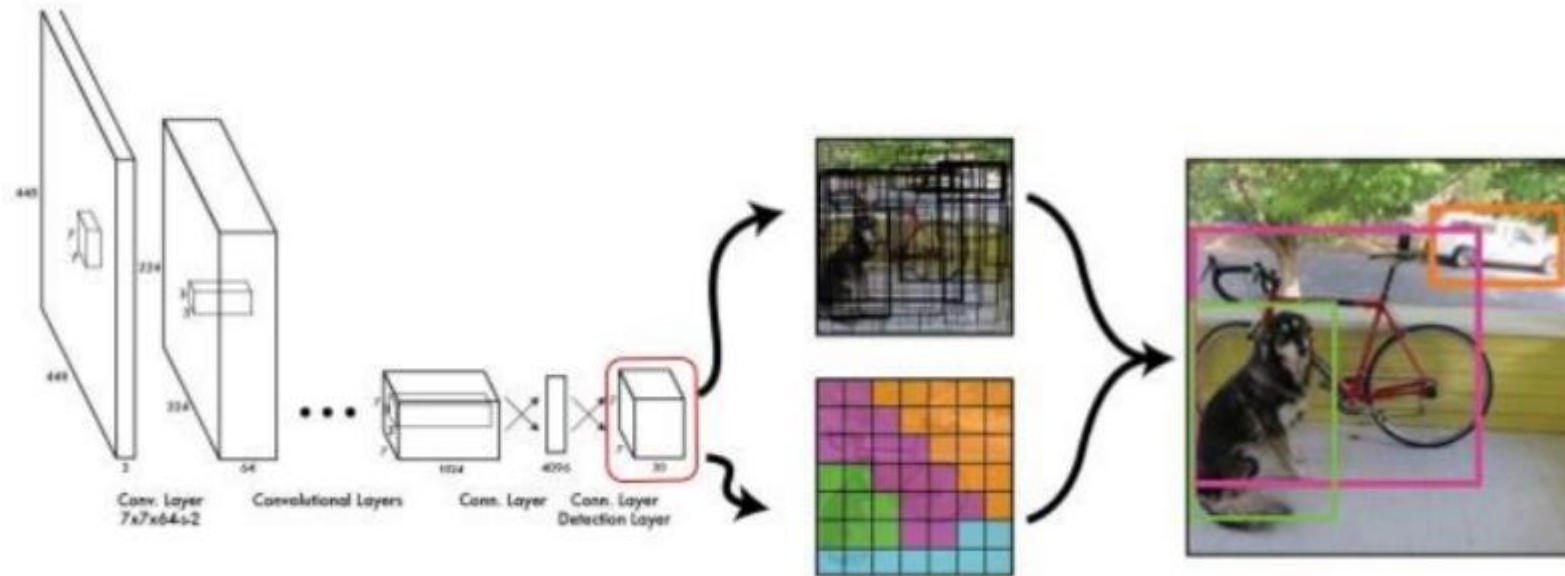
Inference - CNN



YOLO transformon imazhin fillestar përmes shtresave konvolicionale për të nxjerrë tipare, dhe më pas përmes fully connected layers prodhon vektorët e parashikimit të objekteve për secilën nga 49 qelizat e grid-it (7×7), duke përfshirë klasa, pozicion dhe madhësi të bounding box-eve.

Infernce - CNN

YOLO: You Only Look Once



Imazhi përpunohet nga CNN për të nxjerrë tipare vizuale.

Për secilën qelizë të grid-it, modeli parashikon coordinates e bounding box-it, confidence score dhe probabilitetet e klasave.

Pas filtrimit të kutive të tepërta, YOLO identifikon objektet finale dhe pozicionin e tyre në imazh në kohë reale.