

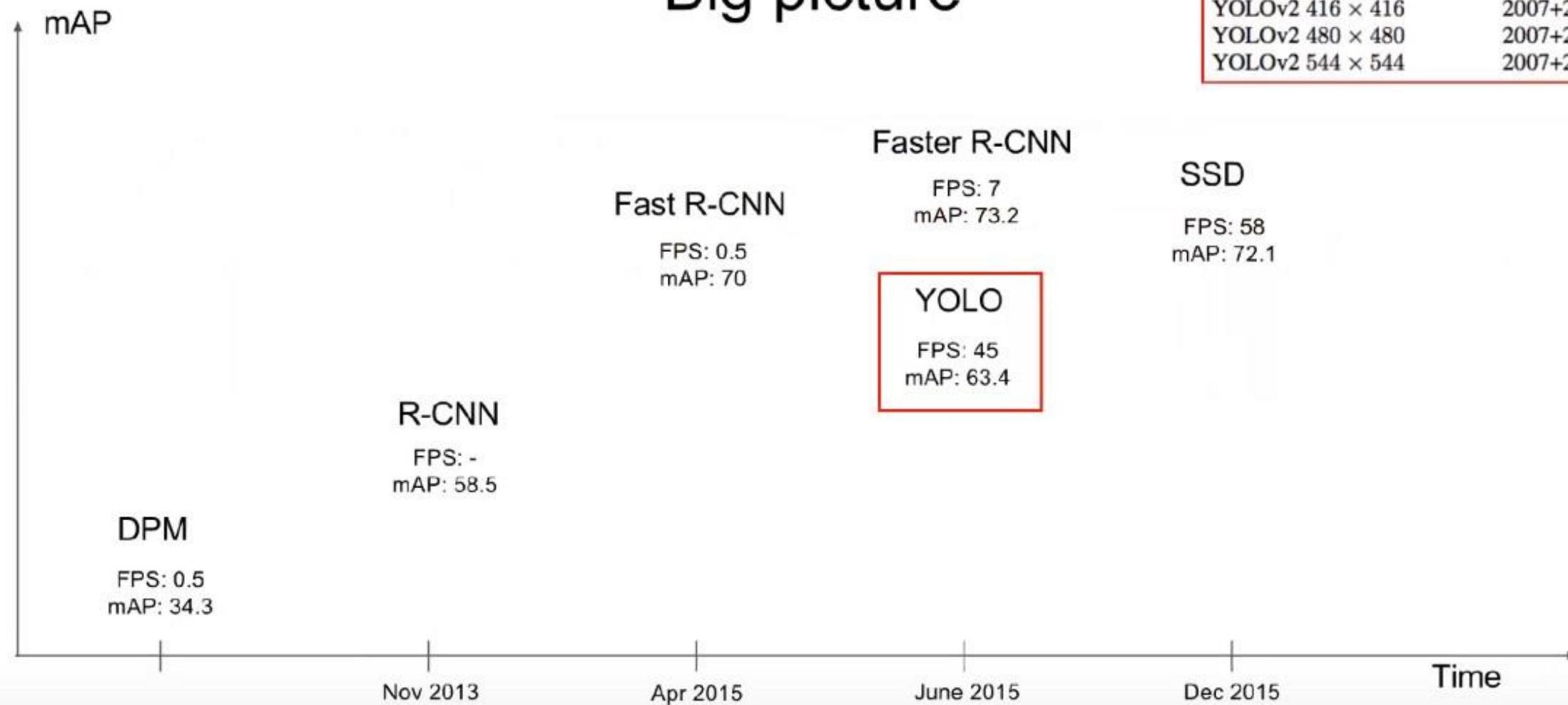
Object Detection

2019.01.21

김용규

Evaluation on VOC2007

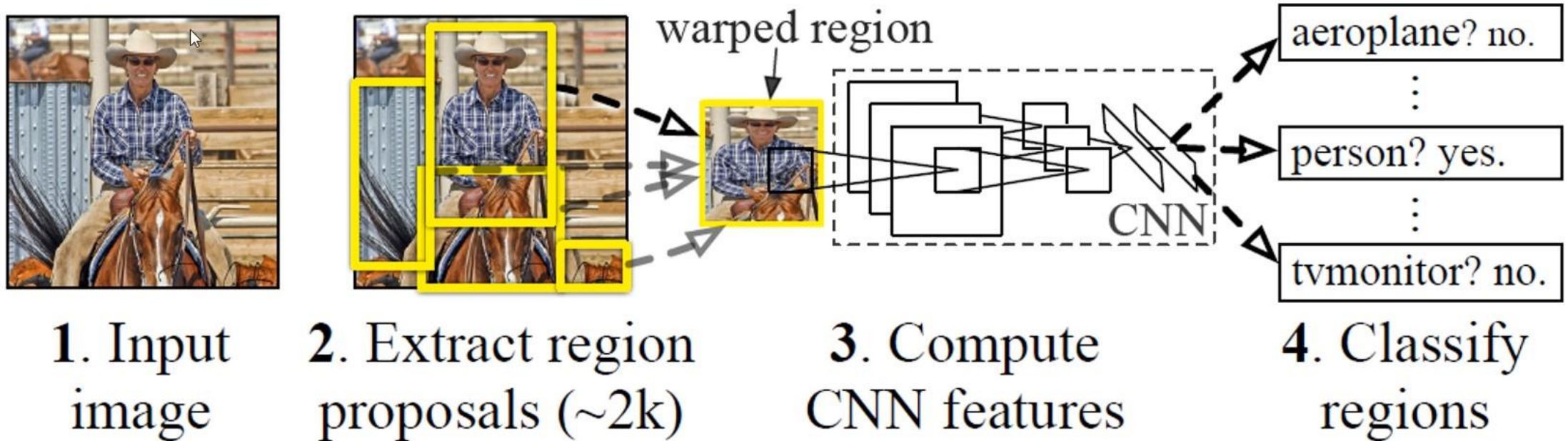
Big picture



YOLOv2 416 × 416	2007+2012	76.8	67
YOLOv2 480 × 480	2007+2012	77.8	59
YOLOv2 544 × 544	2007+2012	78.6	40

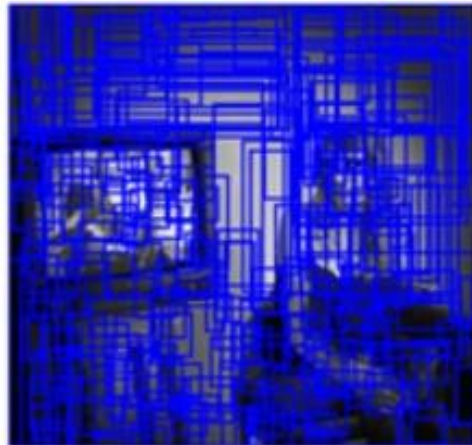
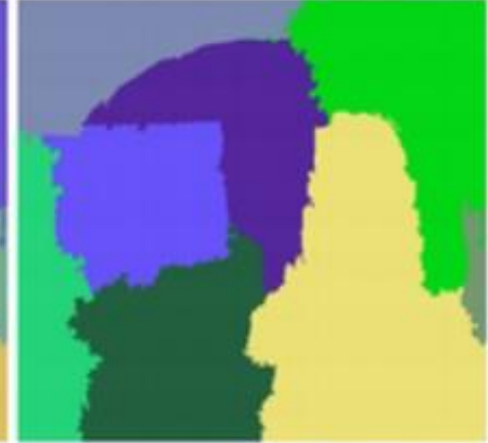
R-CNN Architecture

R-CNN: *Regions with CNN features*

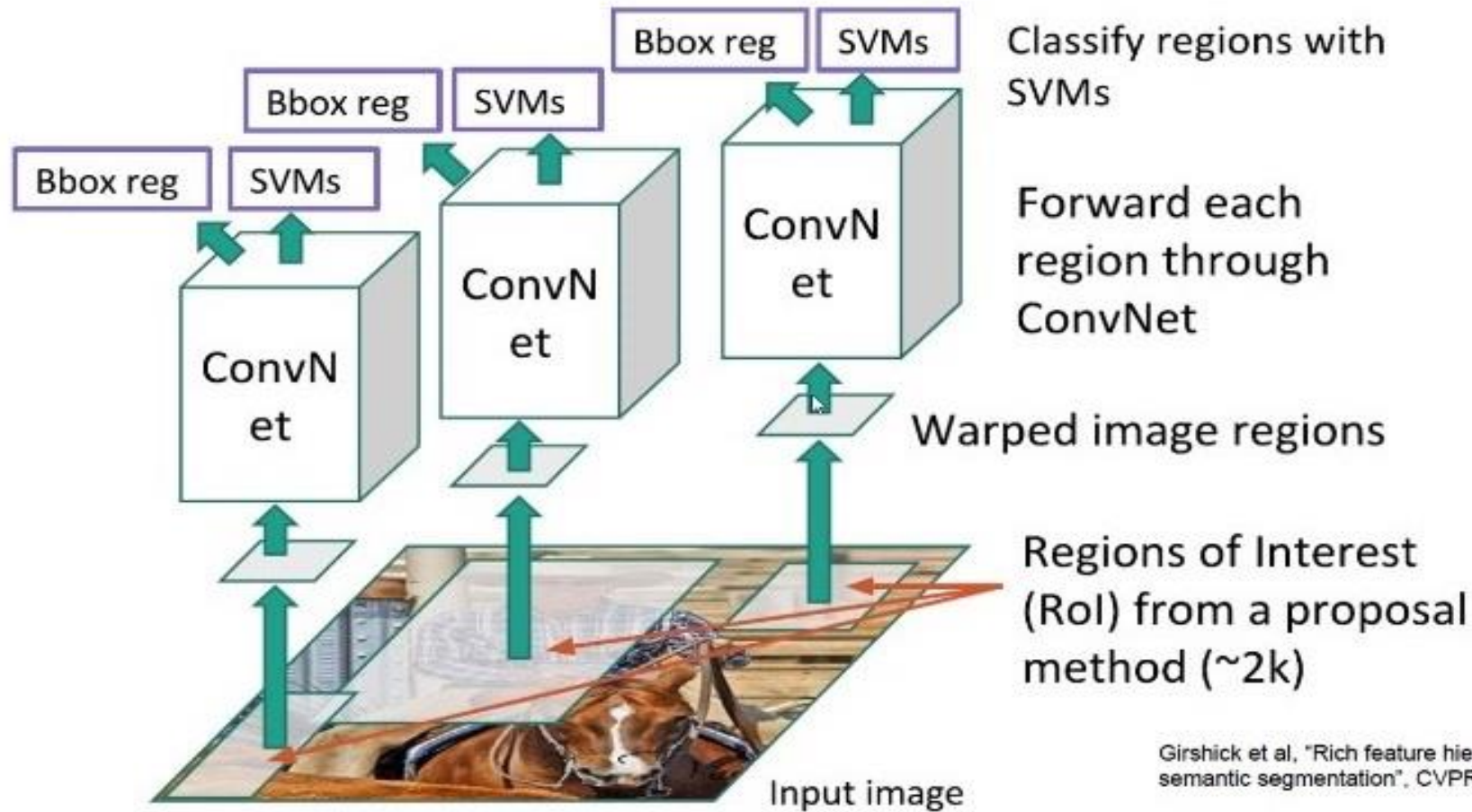




Input Image

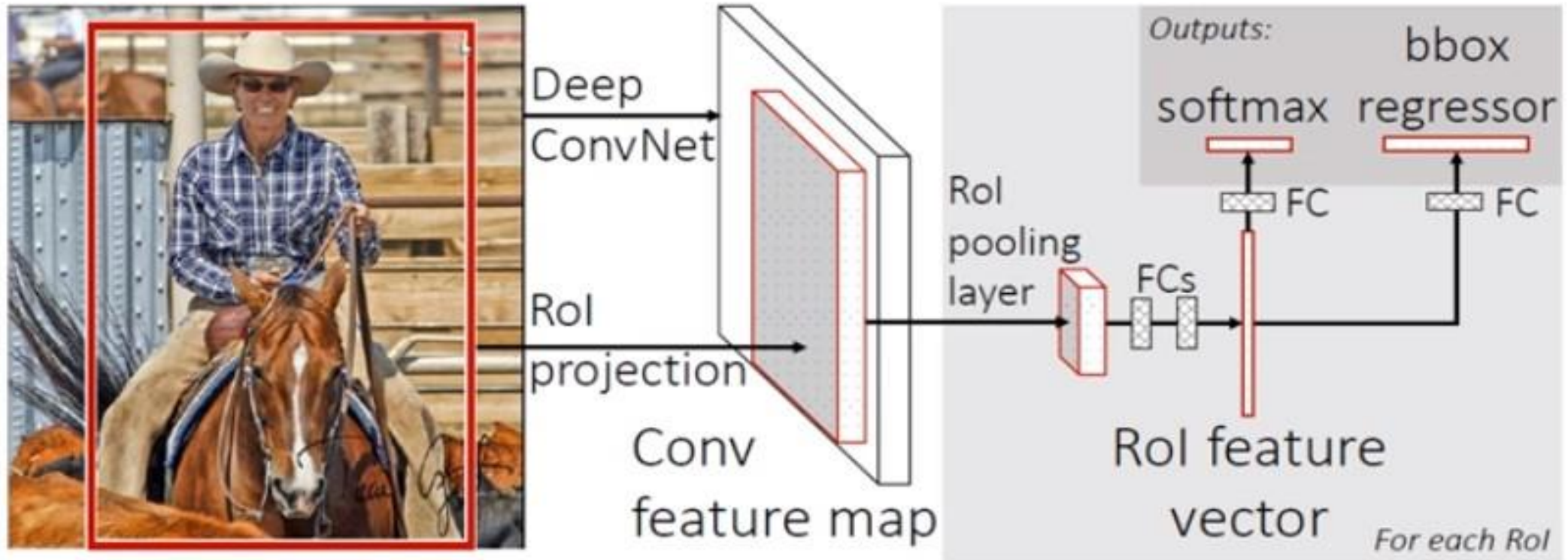


R-CNN

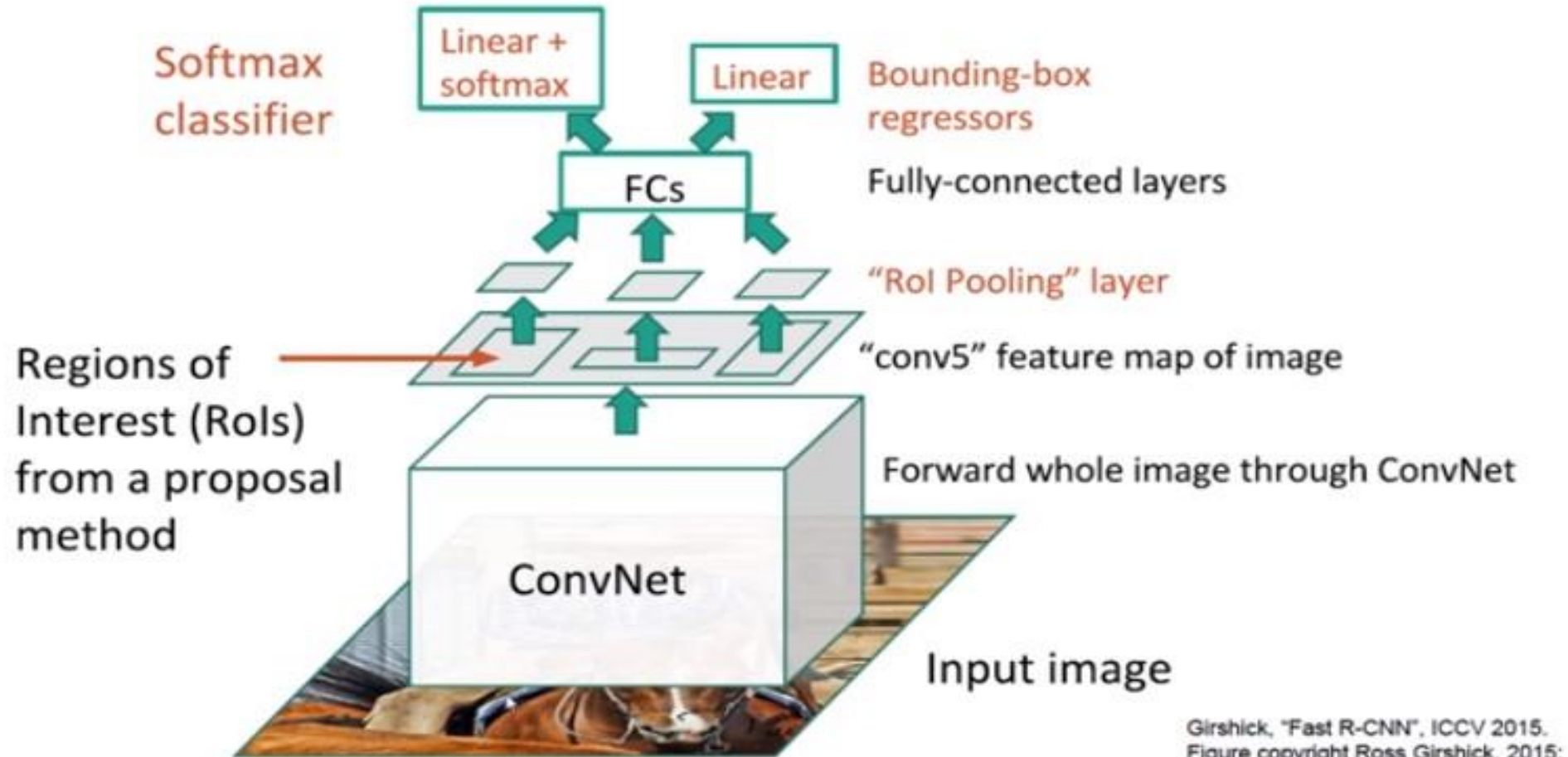


Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.

Fast R-CNN Architecture

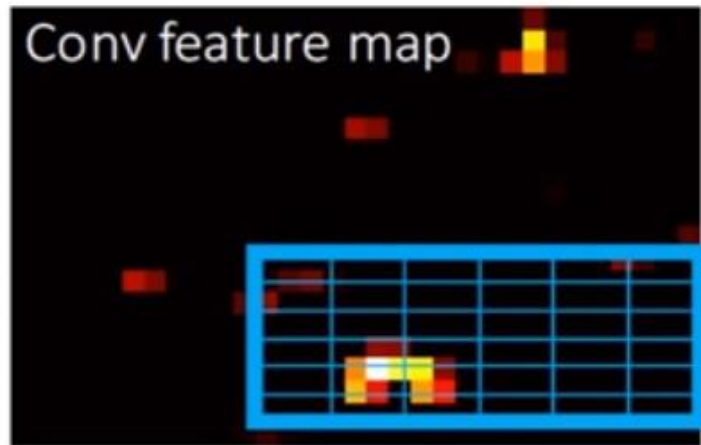


Fast R-CNN

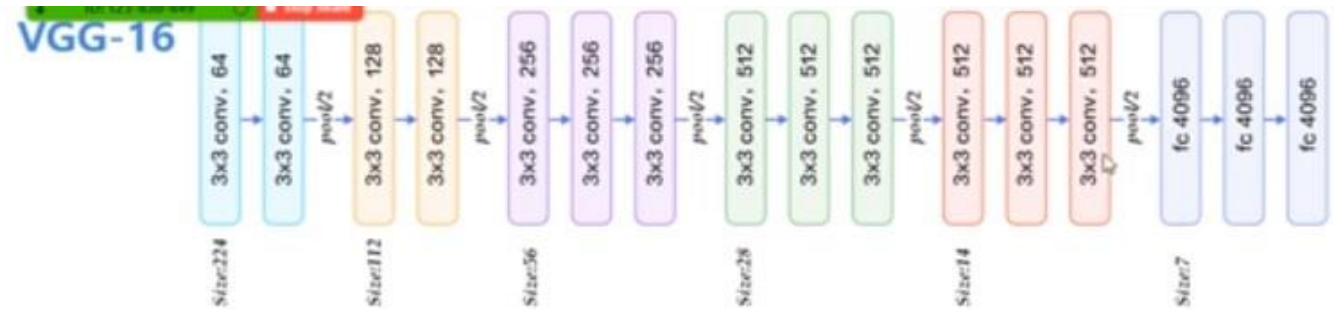


Girshick, "Fast R-CNN", ICCV 2015.
Figure copyright Ross Girshick, 2015;

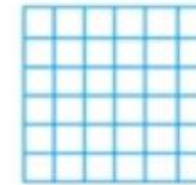
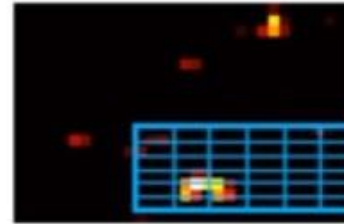
Rol Pooling



Region of Interest (RoI)



Rol
pooling
layer



fc layers ...

Figure adapted
from Kaiming He

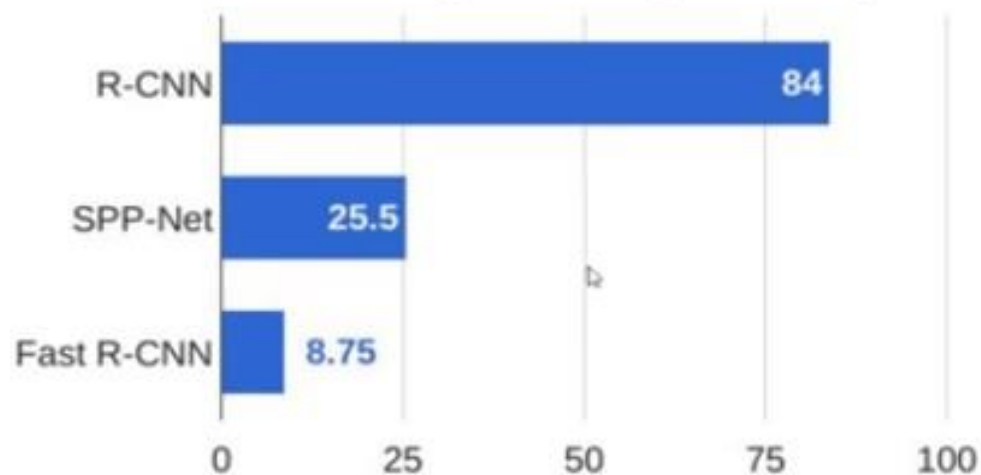
Just a special case of the SPP layer with one pyramid level

Rol in Conv feature map : $21 \times 14 \rightarrow 3 \times 2$ max pooling with stride(3, 2) \rightarrow output : 7×7

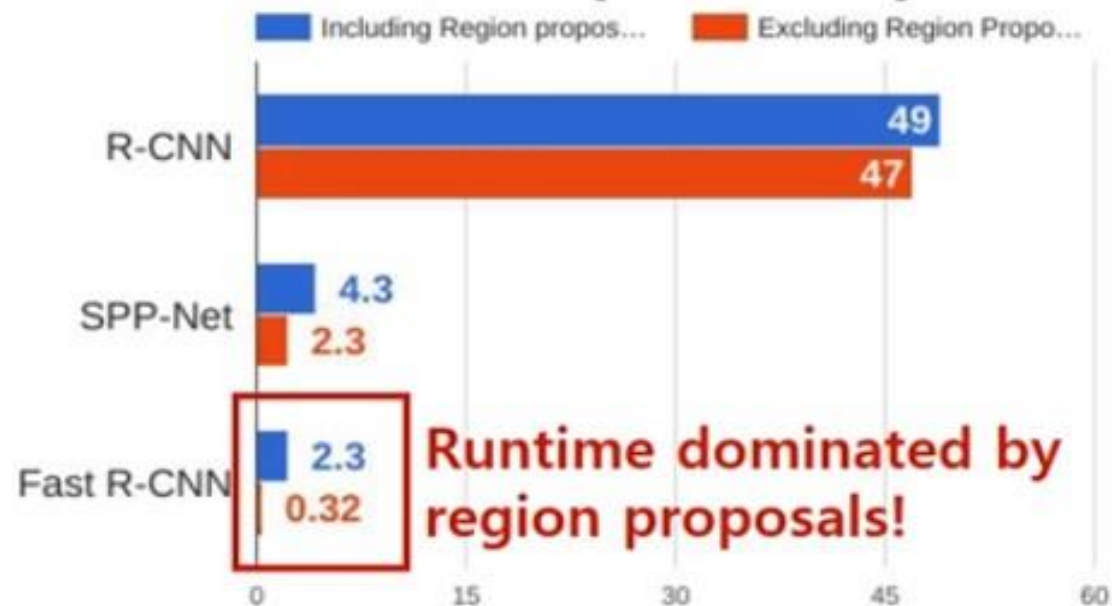
Rol in Conv feature map : $35 \times 42 \rightarrow 5 \times 6$ max pooling with stride(5, 6) \rightarrow output : 7×7

R-CNN vs SPP-net vs Fast R-CNN

Training time (Hours)

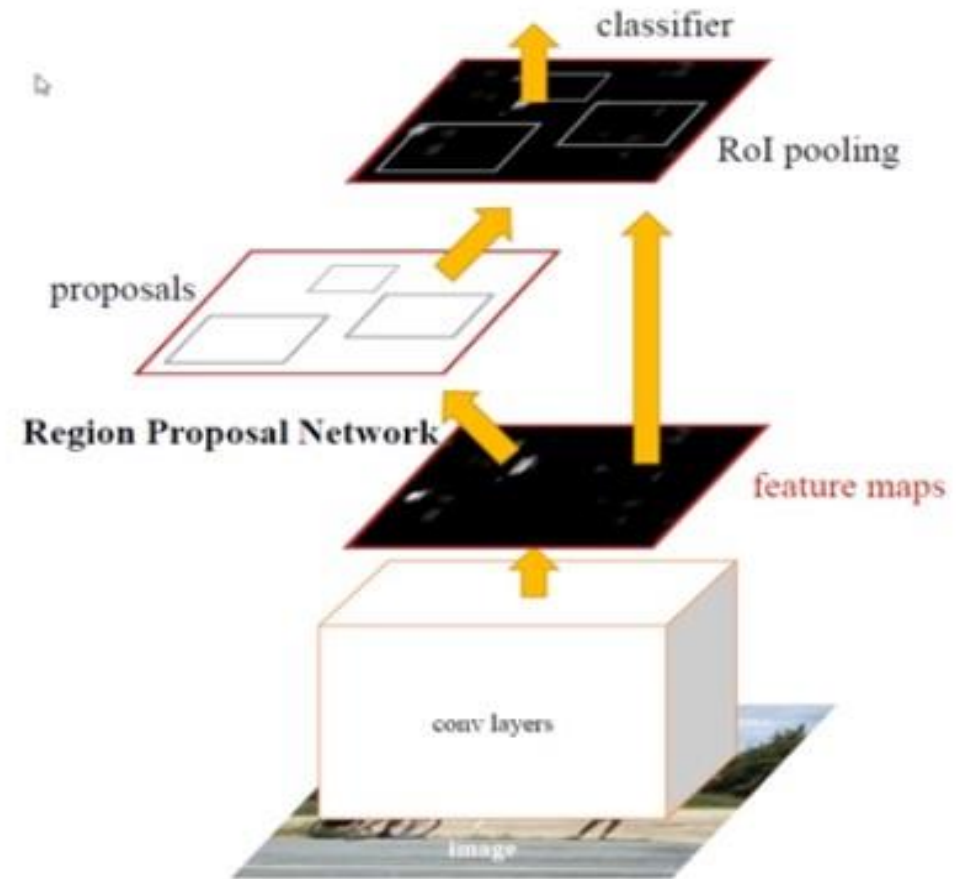


Test time (seconds)

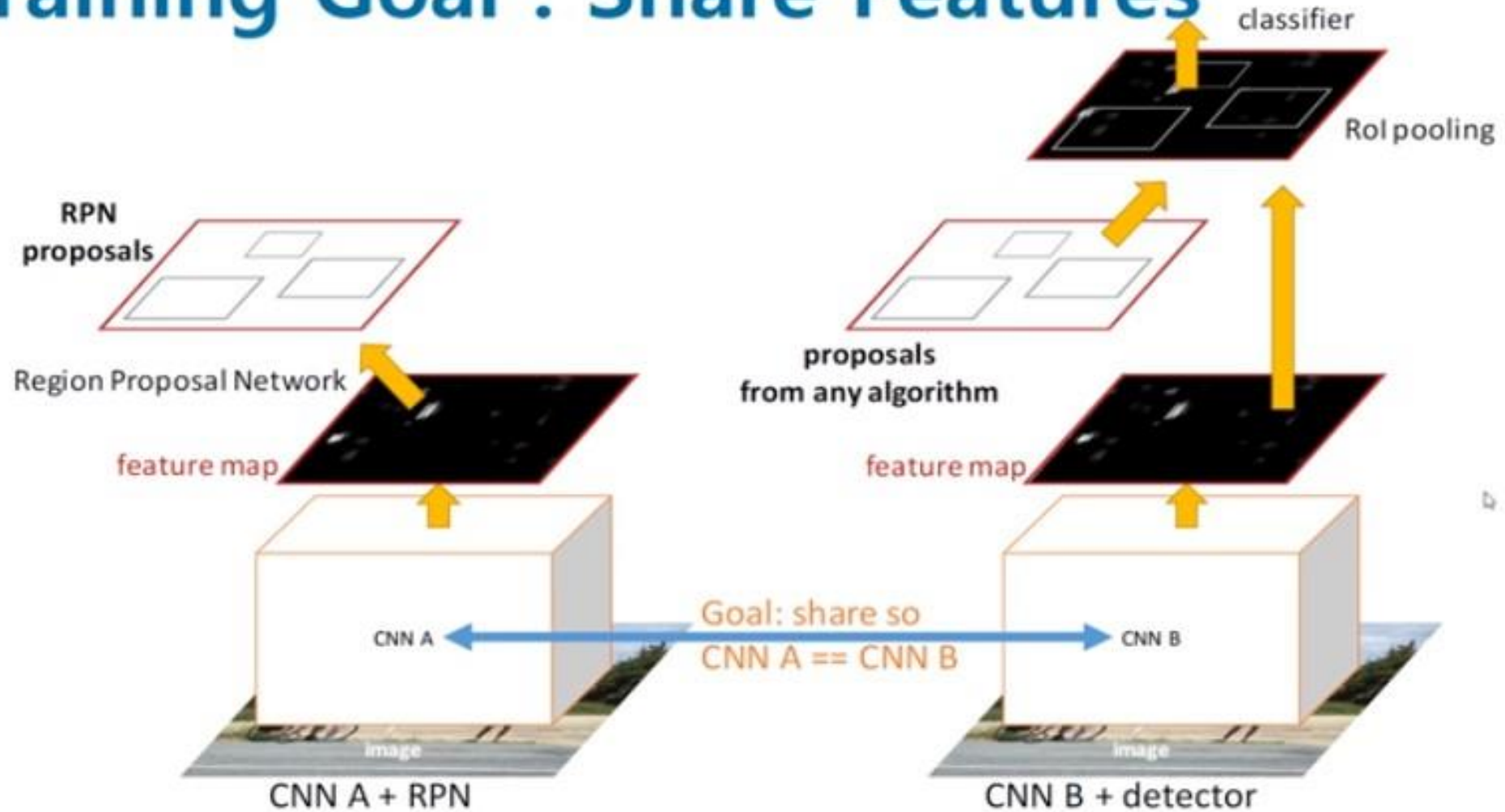


Faster R-CNN(RPN + Fast R-CNN)

- Insert a Region Proposal Network (RPN) after the last convolutional layer → using GPU!
- RPN trained to produce region proposals directly; no need for external region proposals
- After RPN, use RoI Pooling and an upstream classifier and bbox regressor just like Fast R-CNN



Training Goal : Share Features



1

YOLO

- Extremely Fast
- Global reasoning
- Generalizable representation

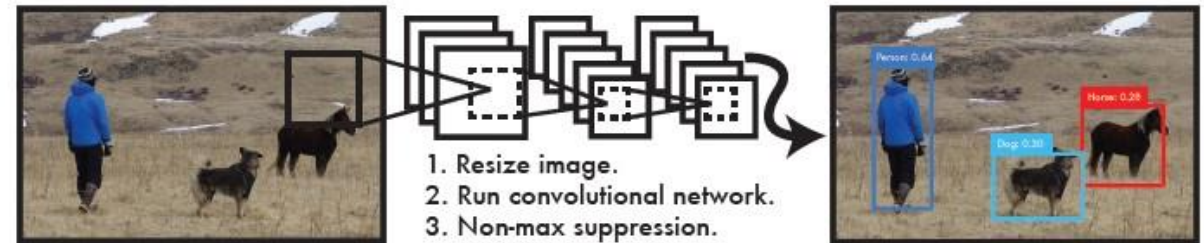


Figure 1: The YOLO Detection System. Processing images with YOLO is simple and straightforward. Our system (1) resizes the input image to 448×448 , (2) runs a single convolutional network on the image, and (3) thresholds the resulting detections by the model's confidence.

1

YOLO

- * B : Bboxes and Confidence score
(confidence score : $\text{Pr}(\text{Object}) * \text{IOU}$)
- * C : class probabilities w.r.t # classes
 $\text{Pr}(\text{Class} \mid \text{Object})$

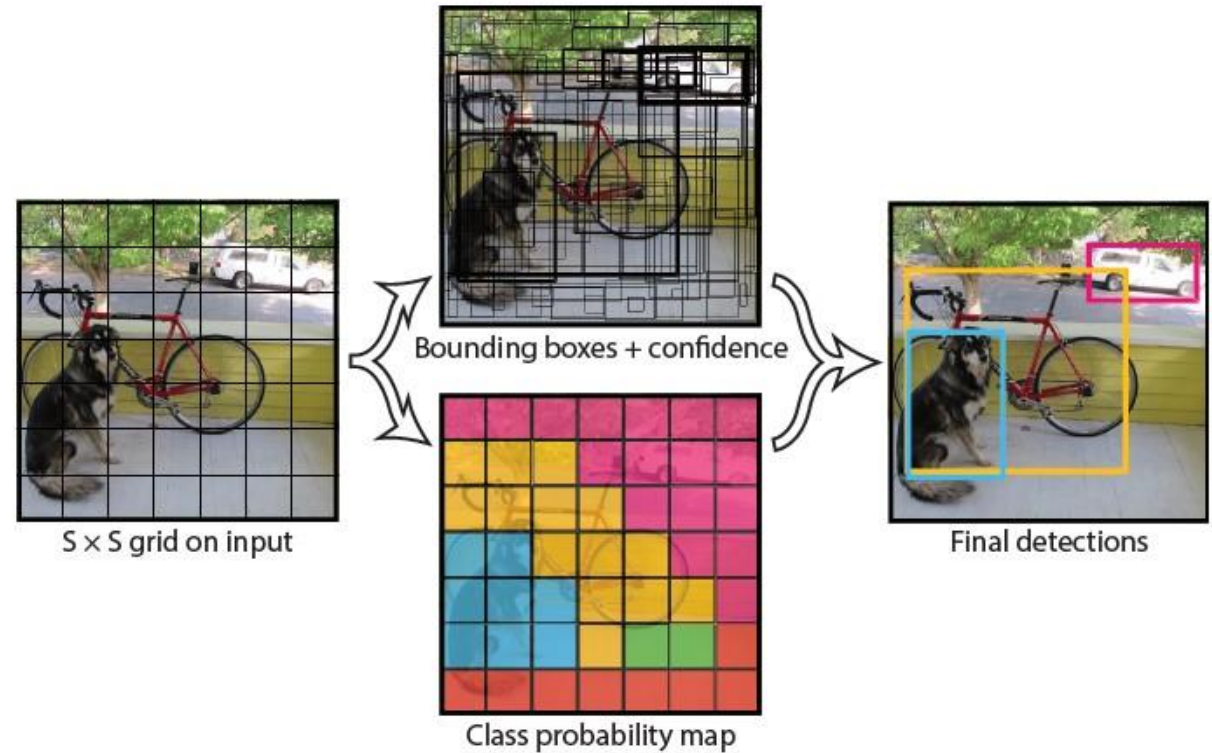


Figure 2: The Model. Our system models detection as a regression problem. It divides the image into an $S \times S$ grid and for each grid cell predicts B bounding boxes, confidence for those boxes, and C class probabilities. These predictions are encoded as an $S \times S \times (B * 5 + C)$ tensor.

YOLO

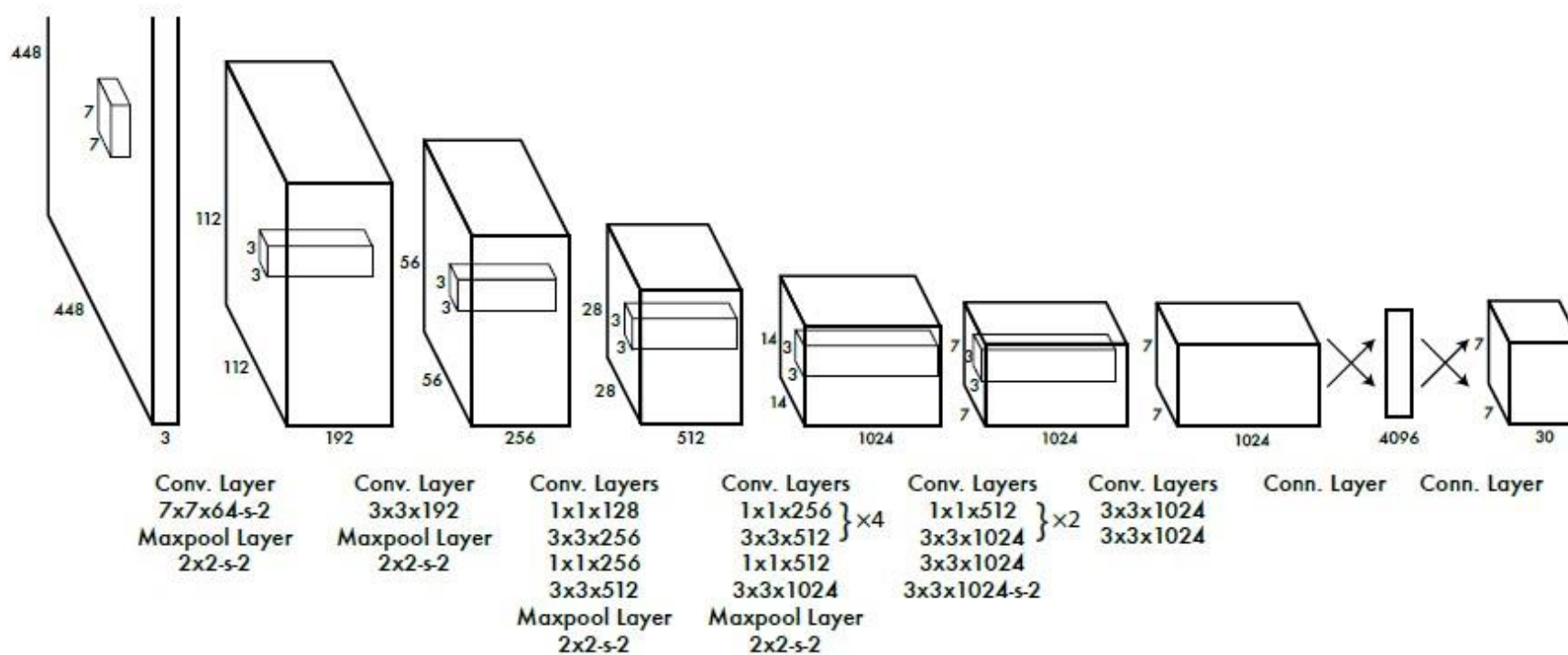
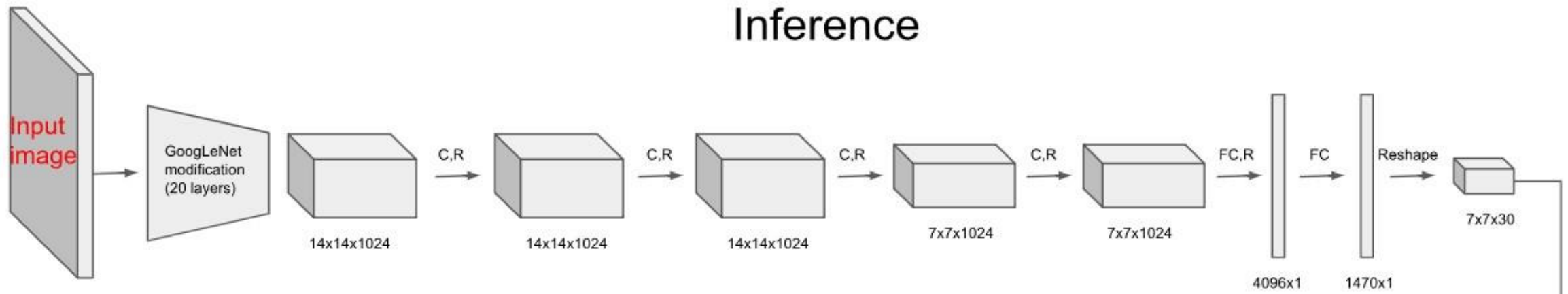
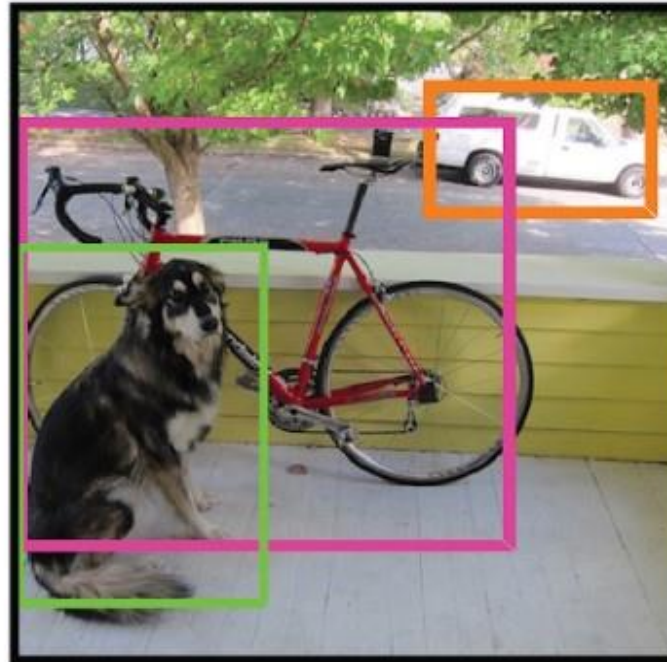


Figure 3: The Architecture. Our detection network has 24 convolutional layers followed by 2 fully connected layers. Alternating 1×1 convolutional layers reduce the features space from preceding layers. We pretrain the convolutional layers on the ImageNet classification task at half the resolution (224×224 input image) and then double the resolution for detection.

Inference

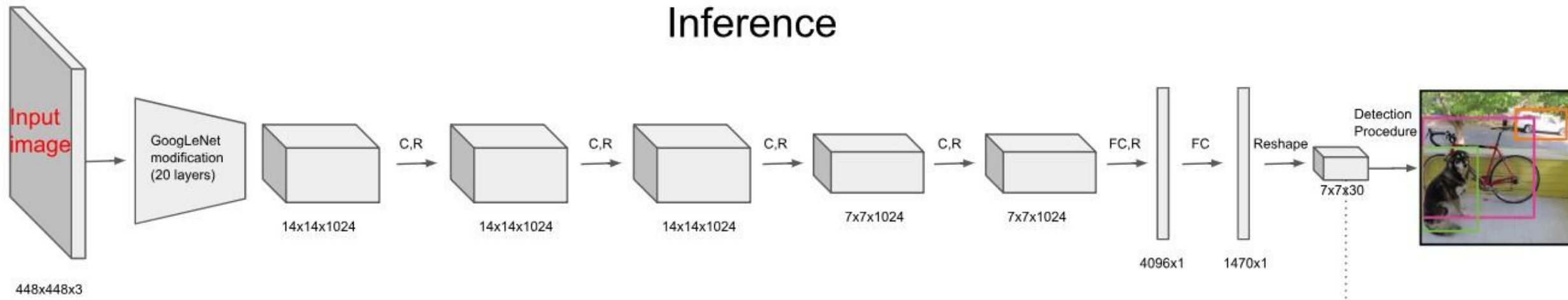


448x448x3

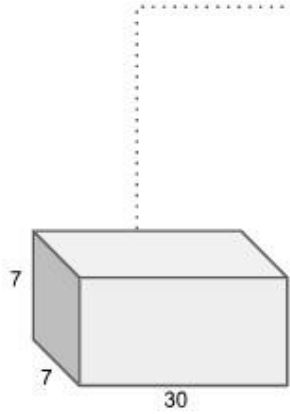


Detection Procedure

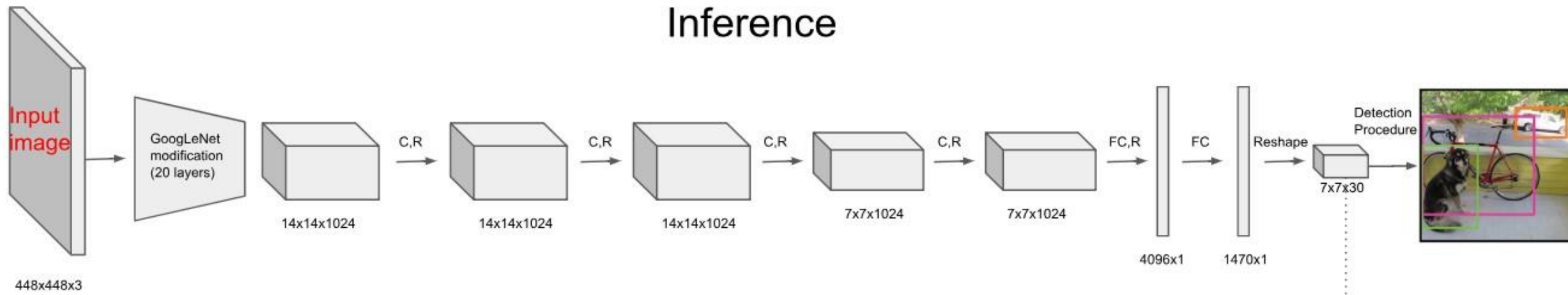
Inference



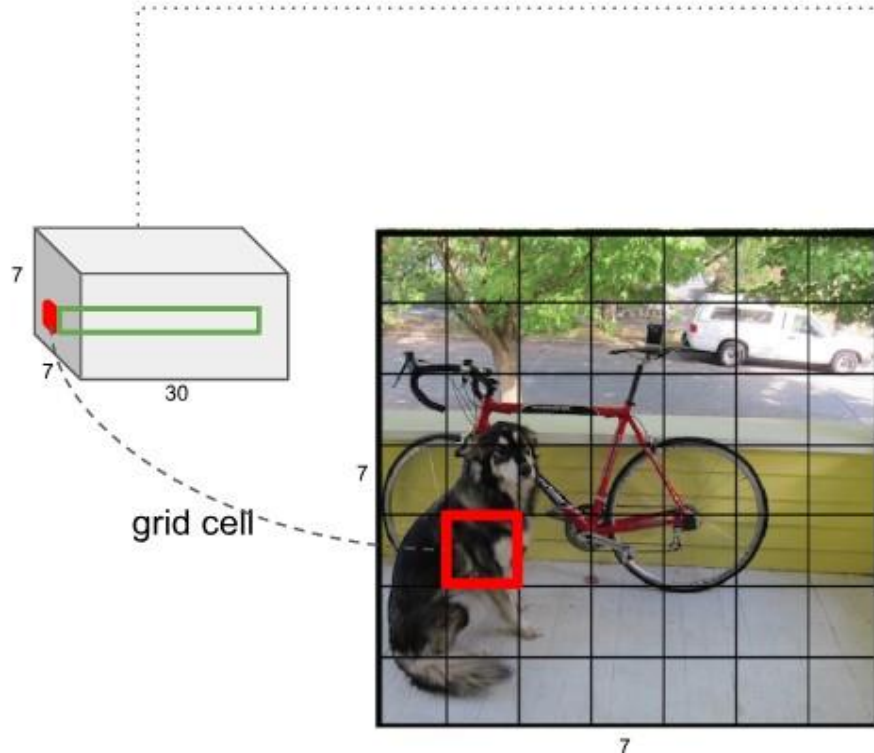
Tensor values interpretation



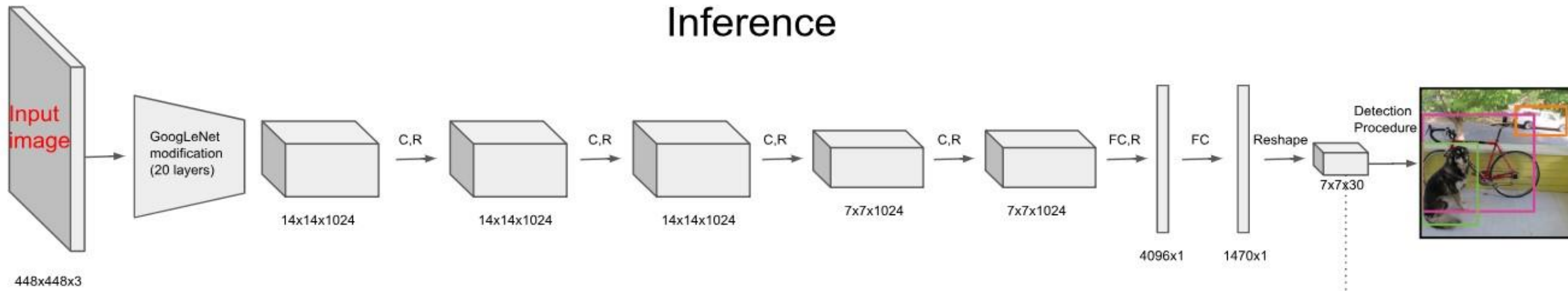
Inference



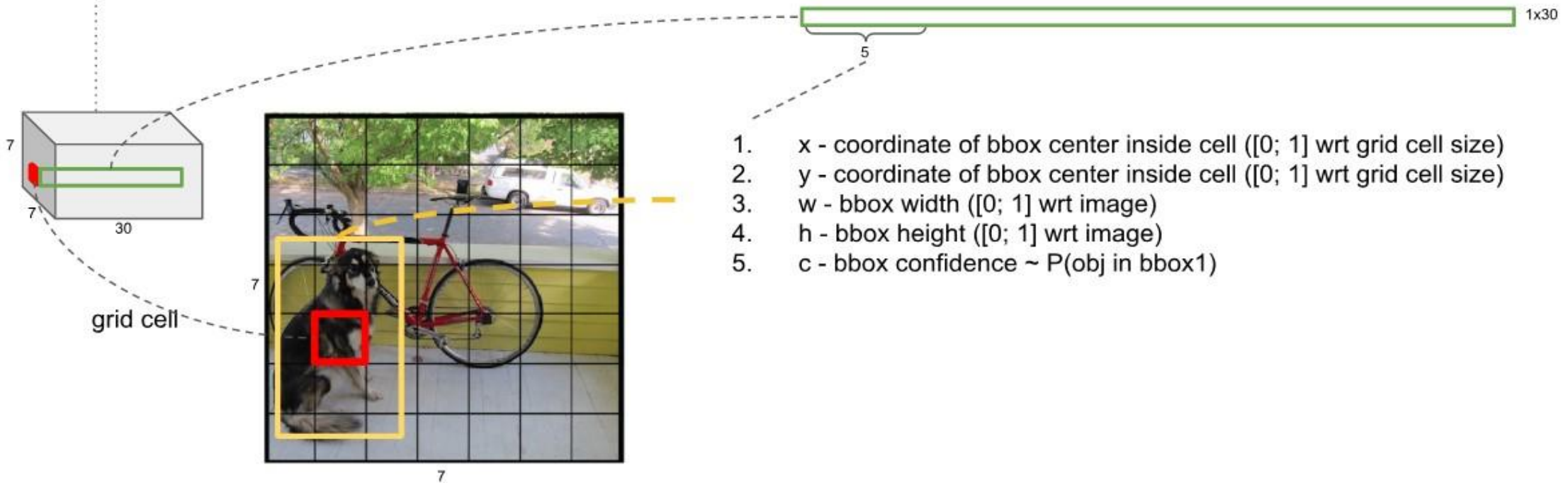
Tensor values interpretation



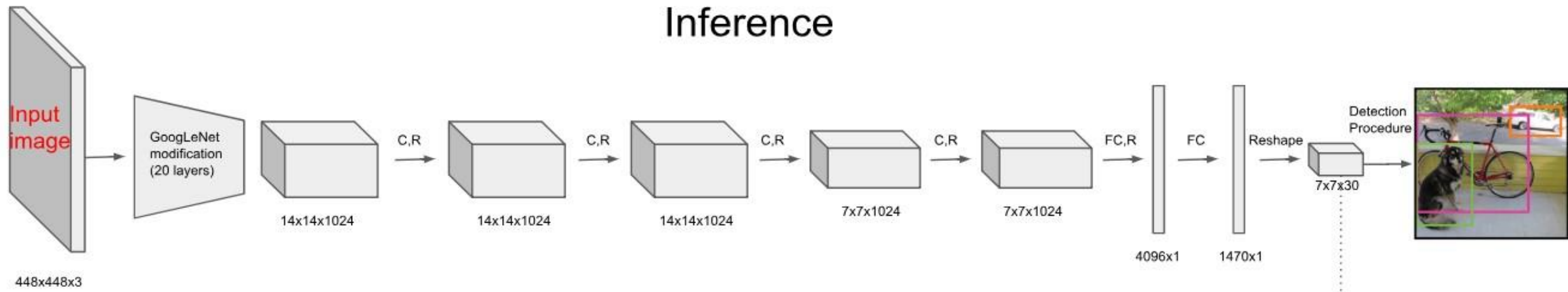
Inference



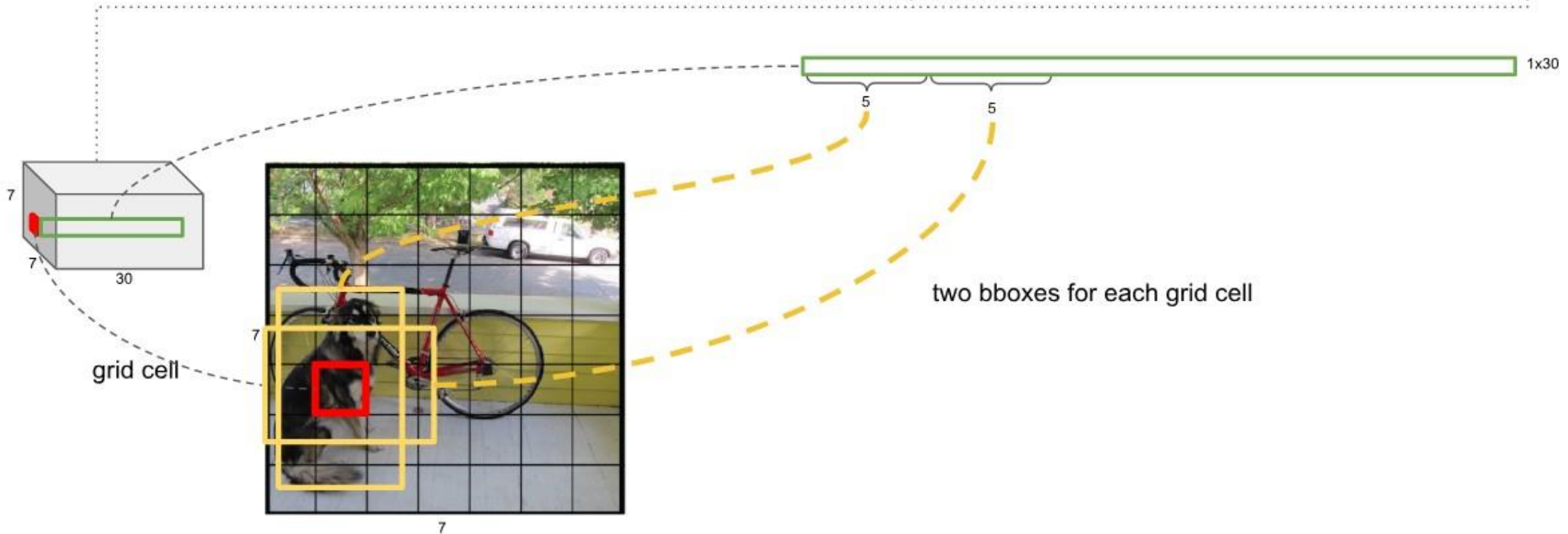
Tensor values interpretation



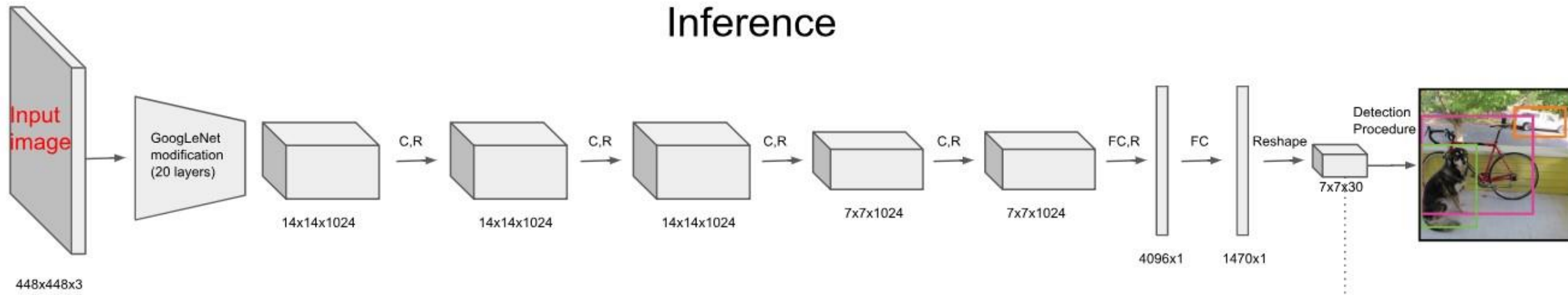
Inference



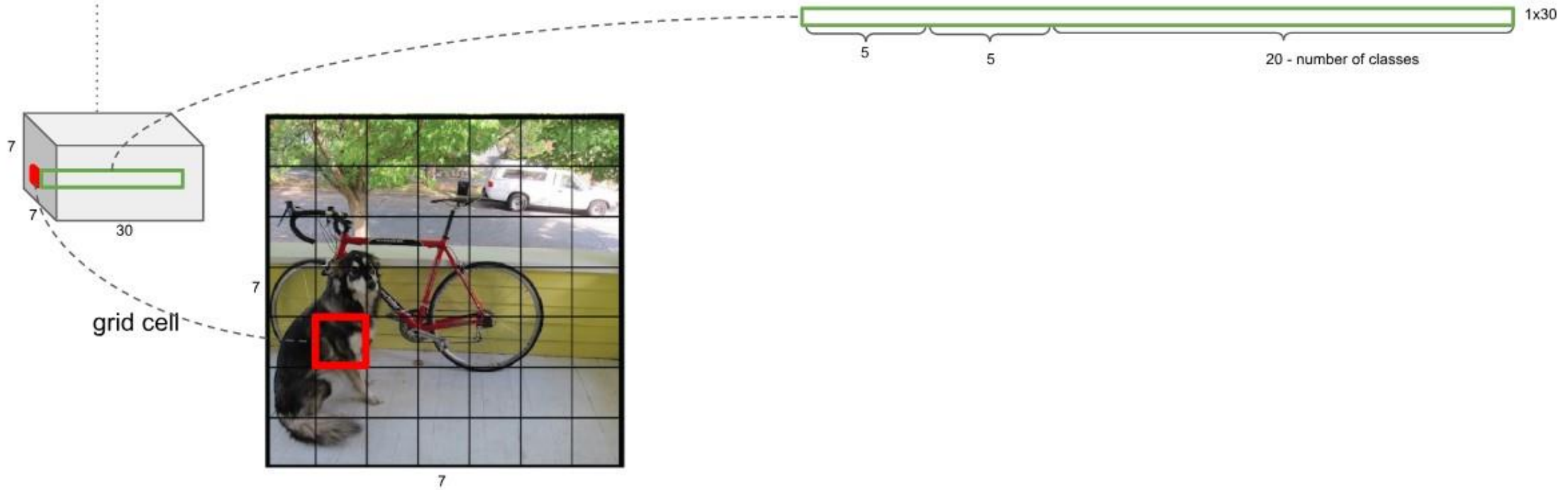
Tensor values interpretation



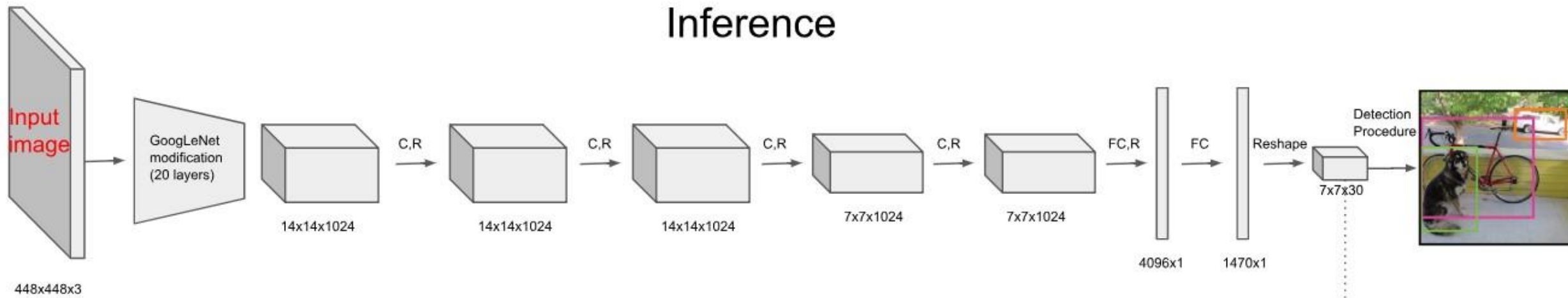
Inference



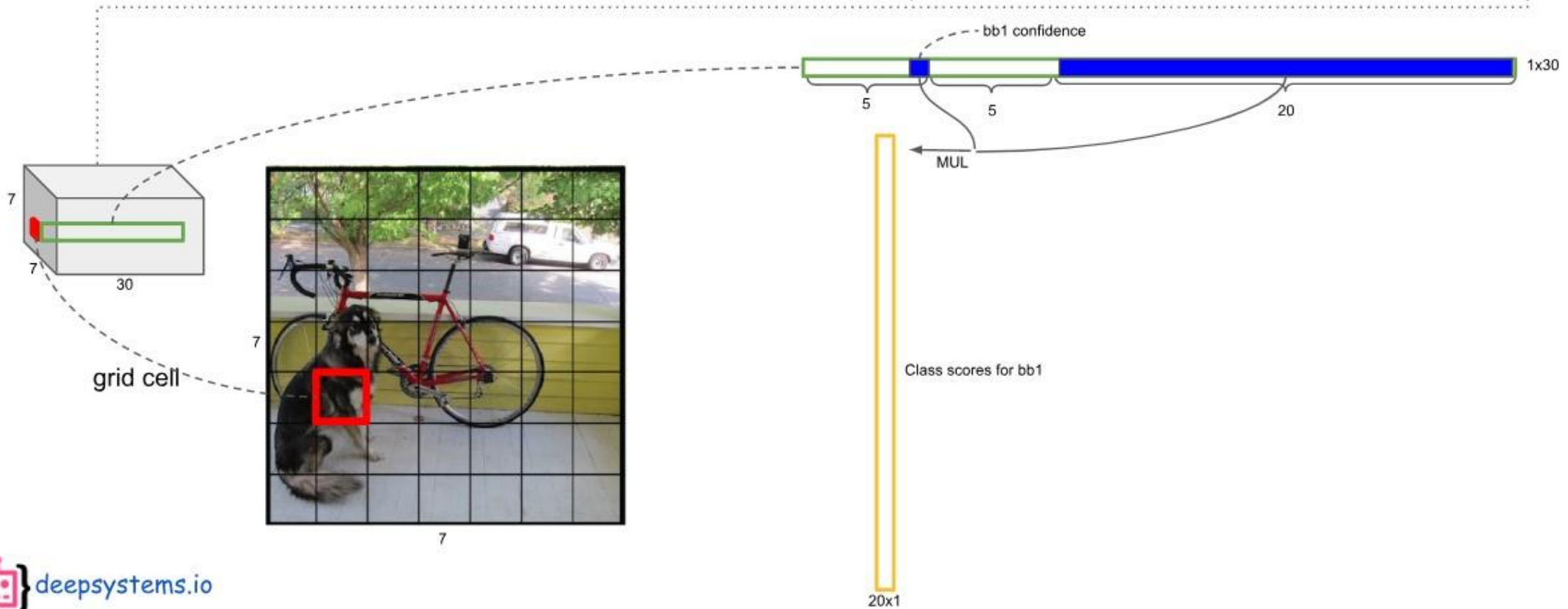
Tensor values interpretation



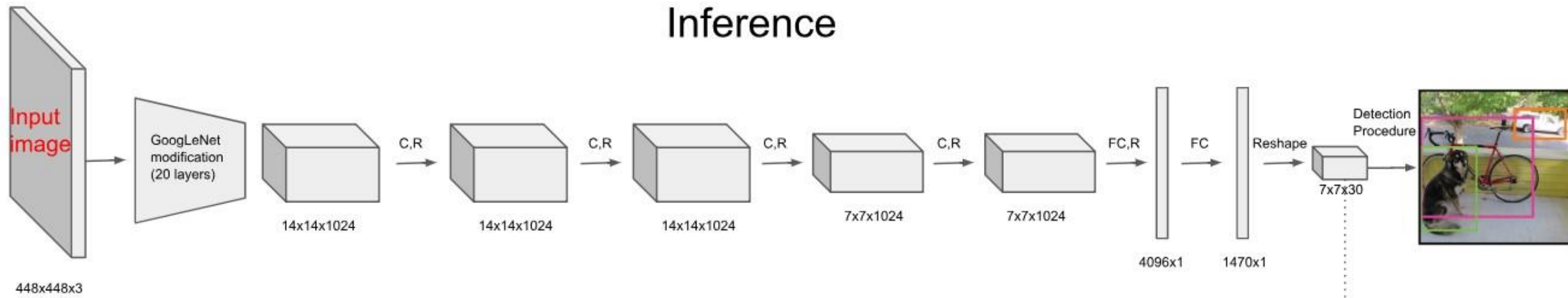
Inference



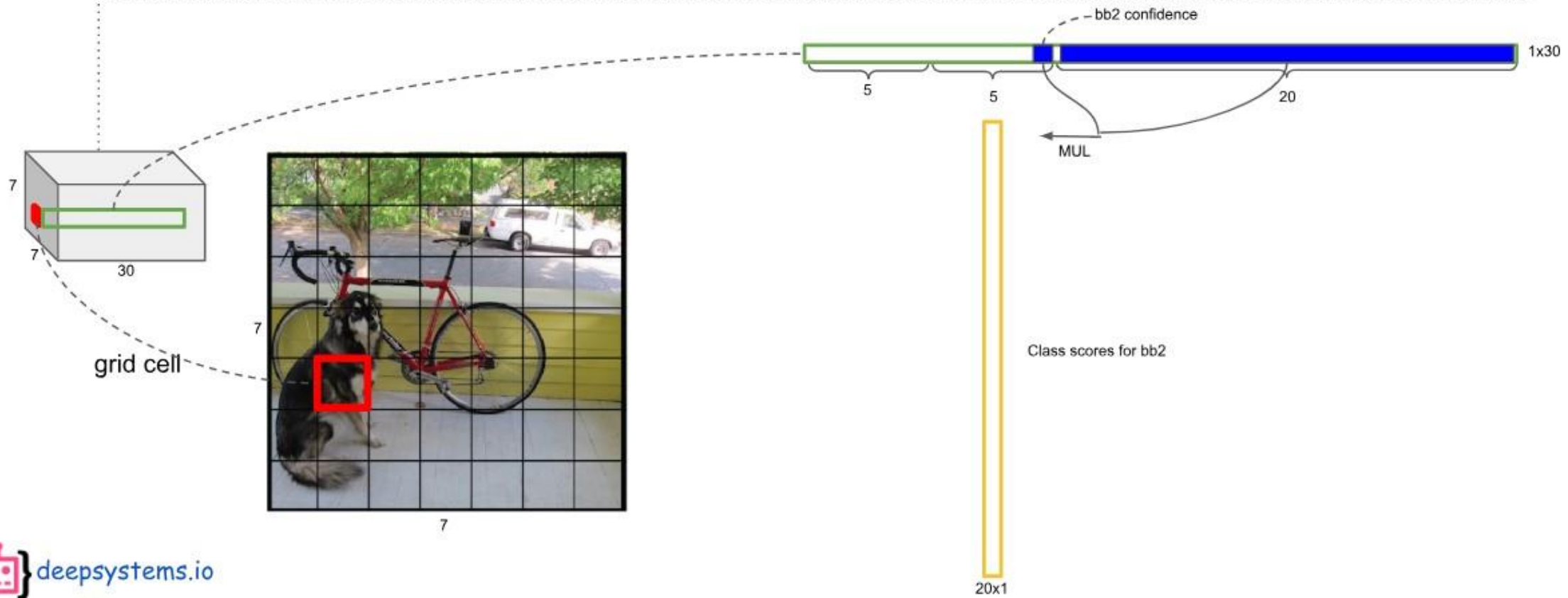
Tensor values interpretation



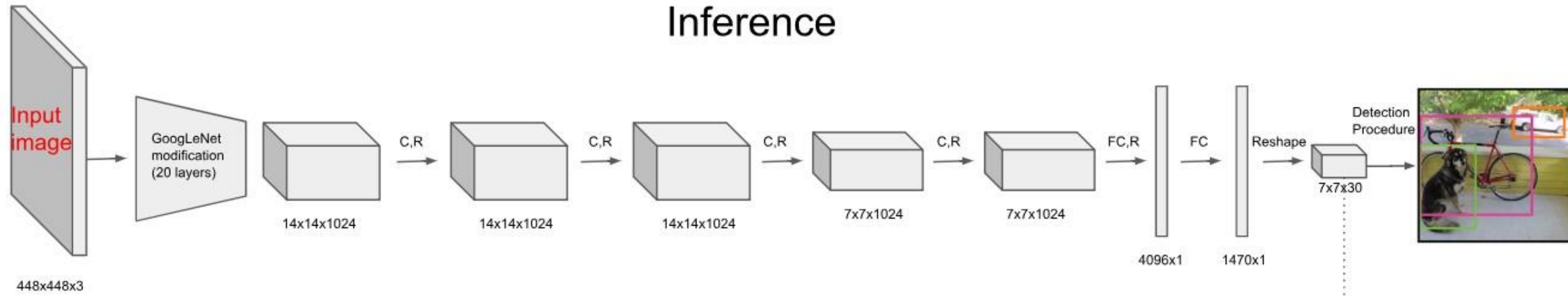
Inference



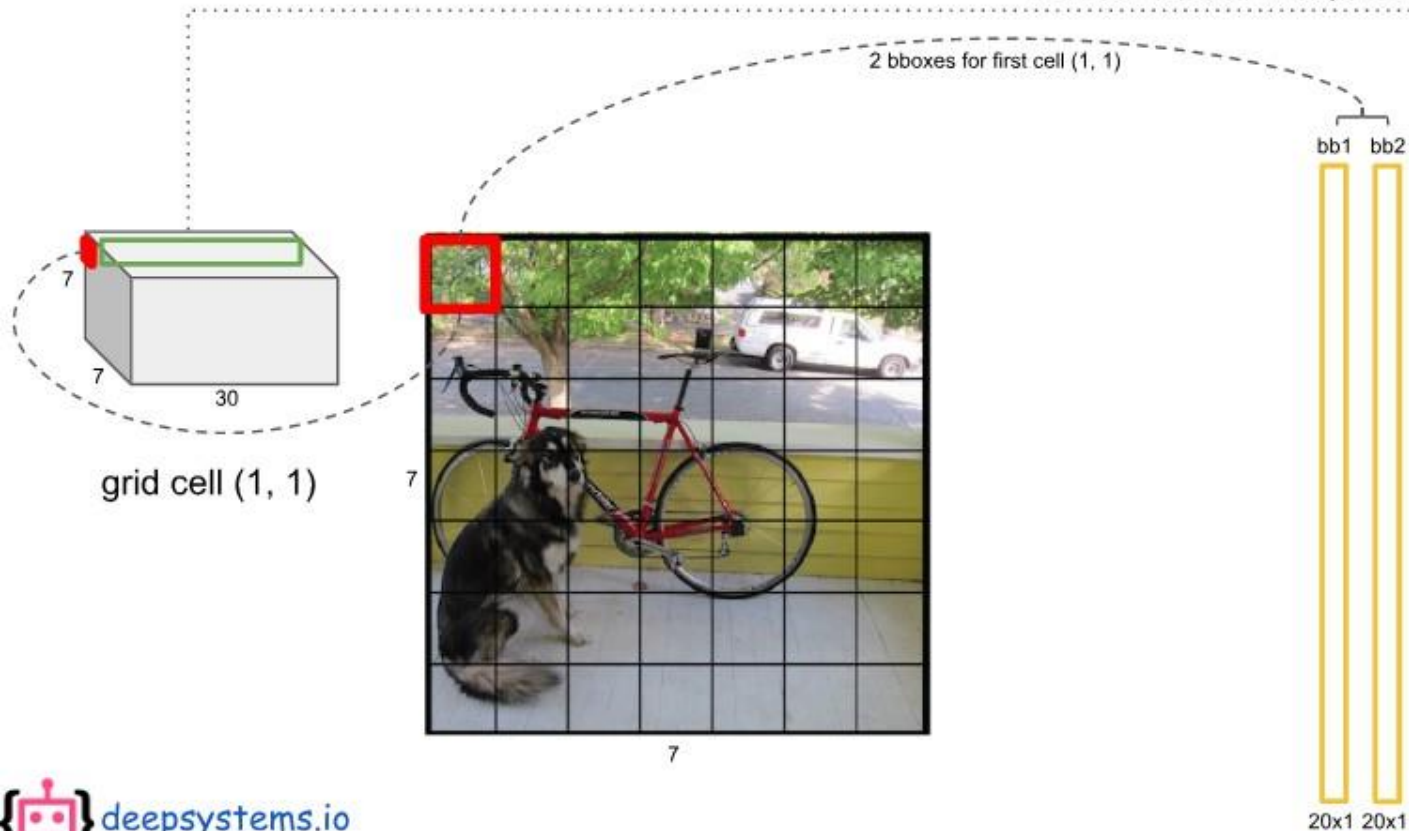
Tensor values interpretation



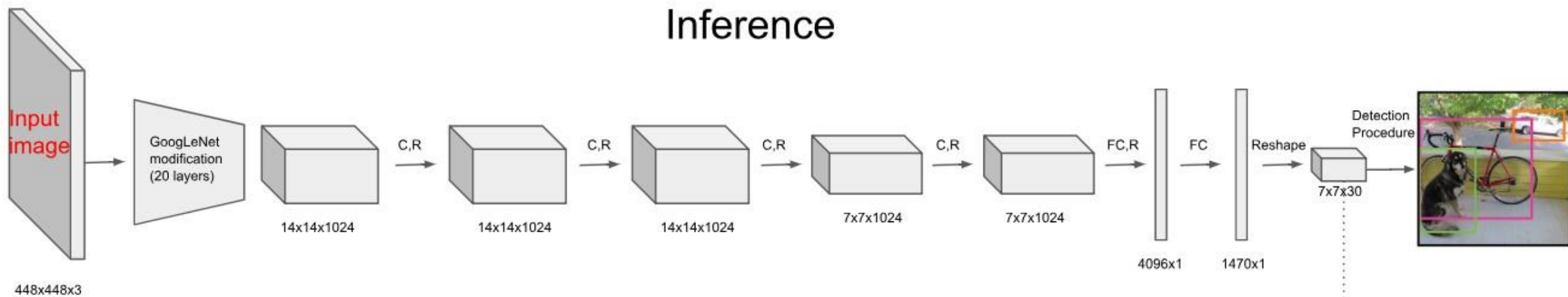
Inference



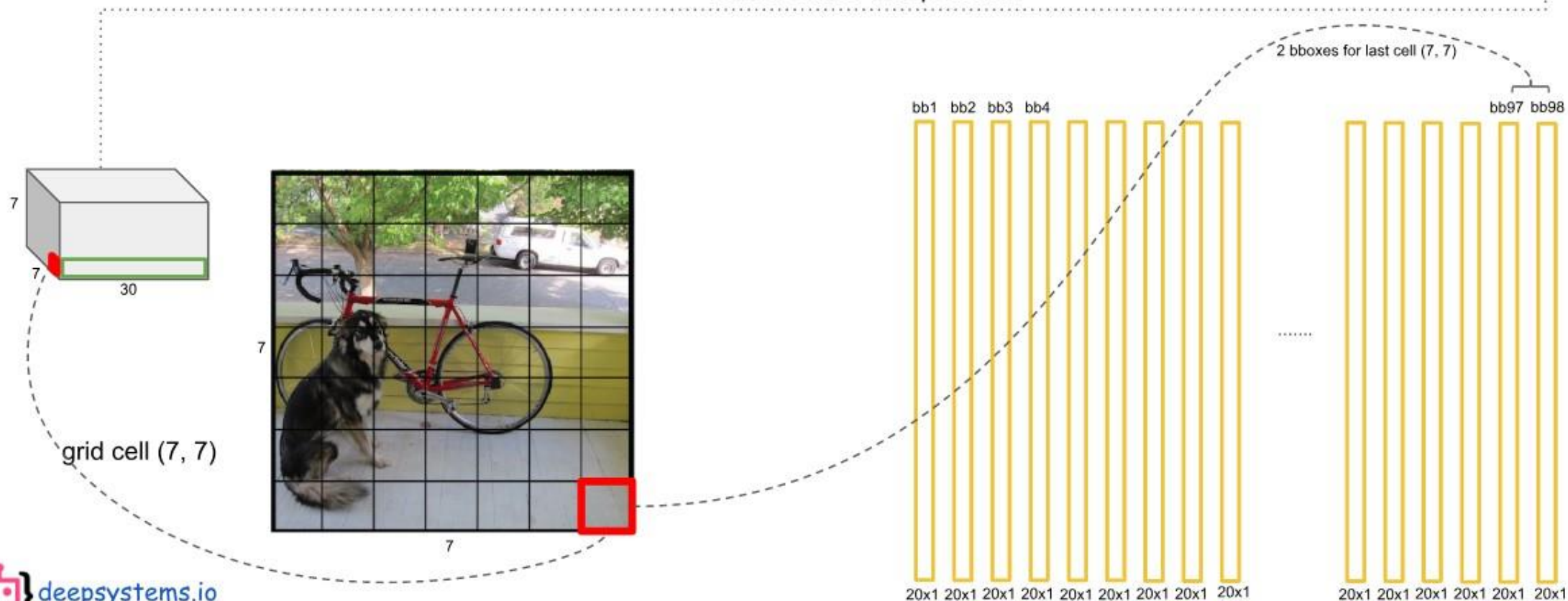
Tensor values interpretation

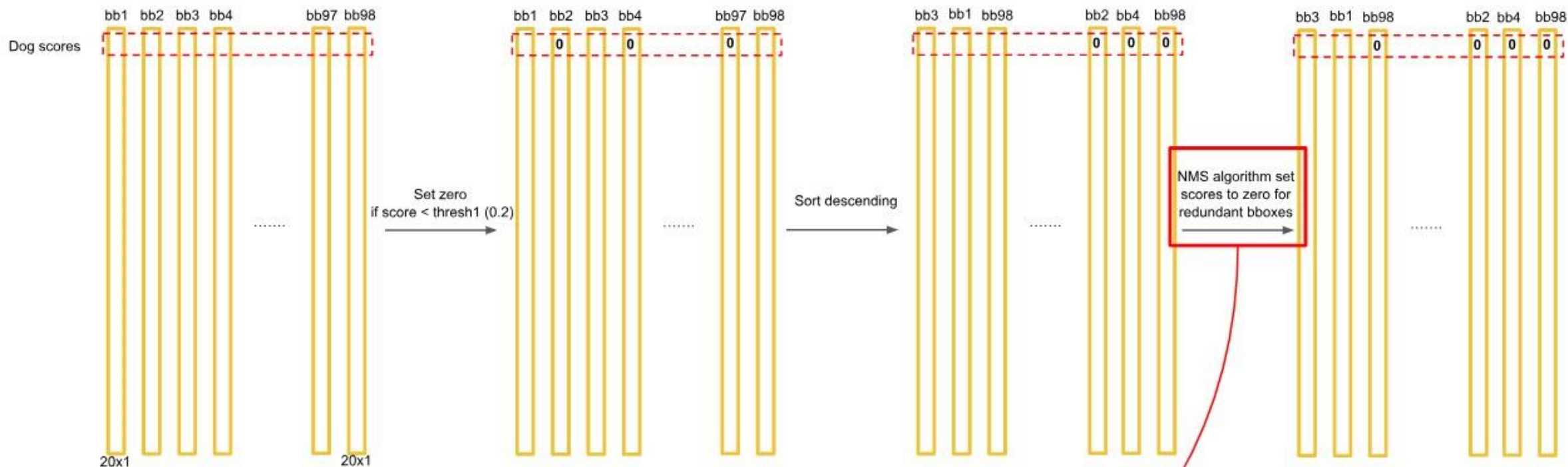


Inference



Tensor values interpretation



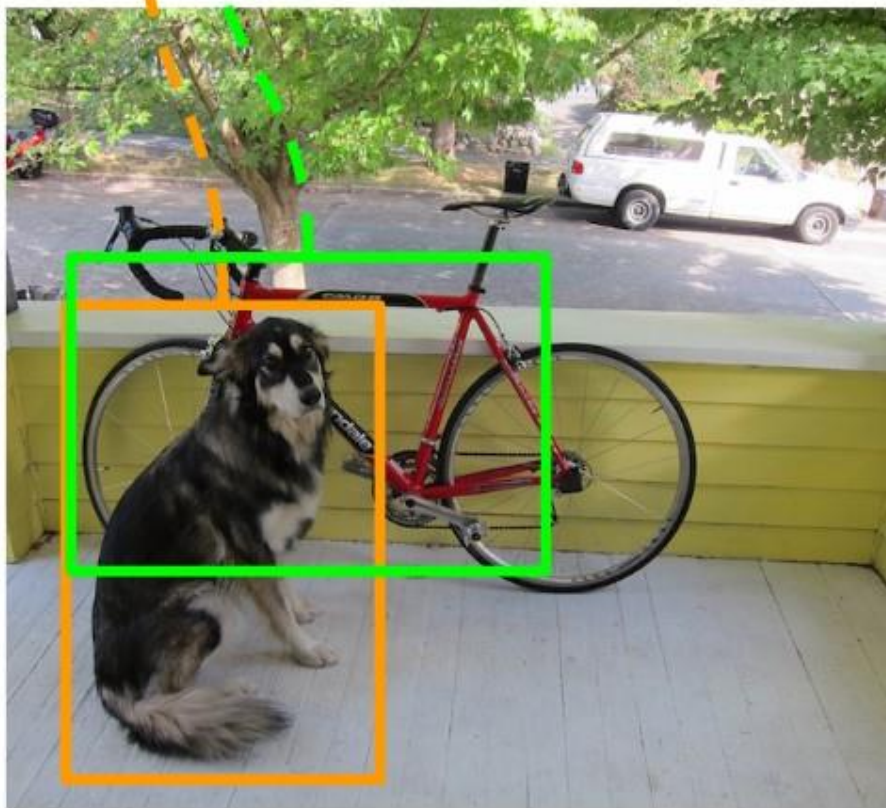


How it works

Non-Maximum Suppression: intuition

class (dog) scores for each bbox

class: dog	bb47	bb20	bb15	bb7									bb1	bb4	bb8	bb98	1x98
	0.5	0.3	0.2	0.1									0	0	0	0	

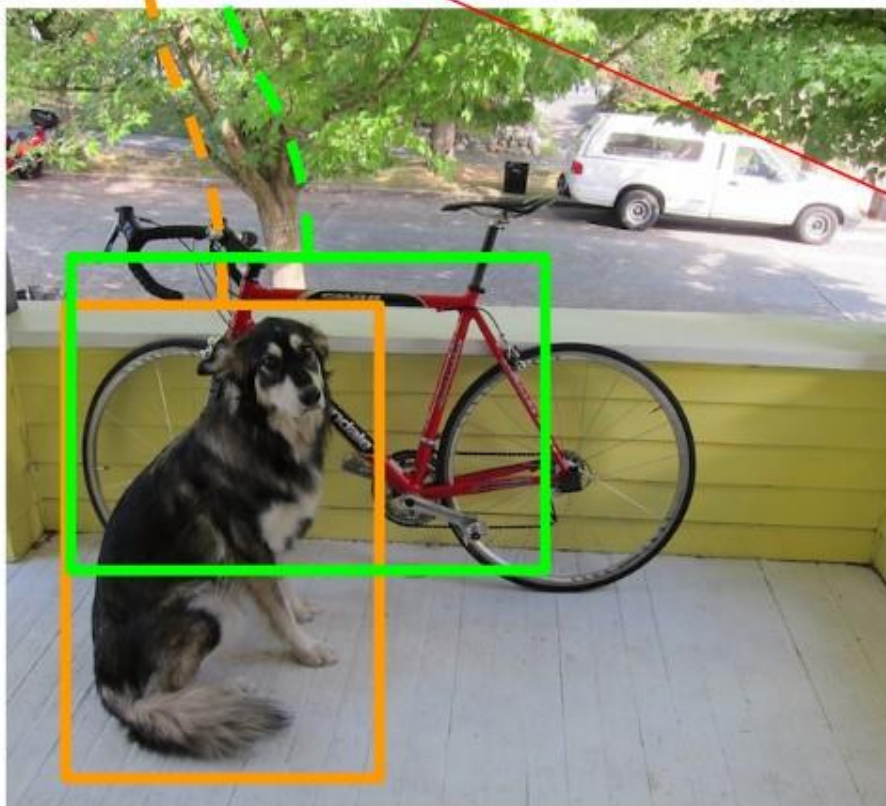


If $\text{IoU}(\text{bbox_max}, \text{bbox_cur}) > 0.5$ then set 0 score to bbox_cur .

Non-Maximum Suppression: intuition

class (dog) scores for each bbox

class: dog	bb47	bb20	bb15	bb7											bb1	bb4	bb8	bb98	1x98
	0.5	0	0.2	0.1											0	0	0	0	



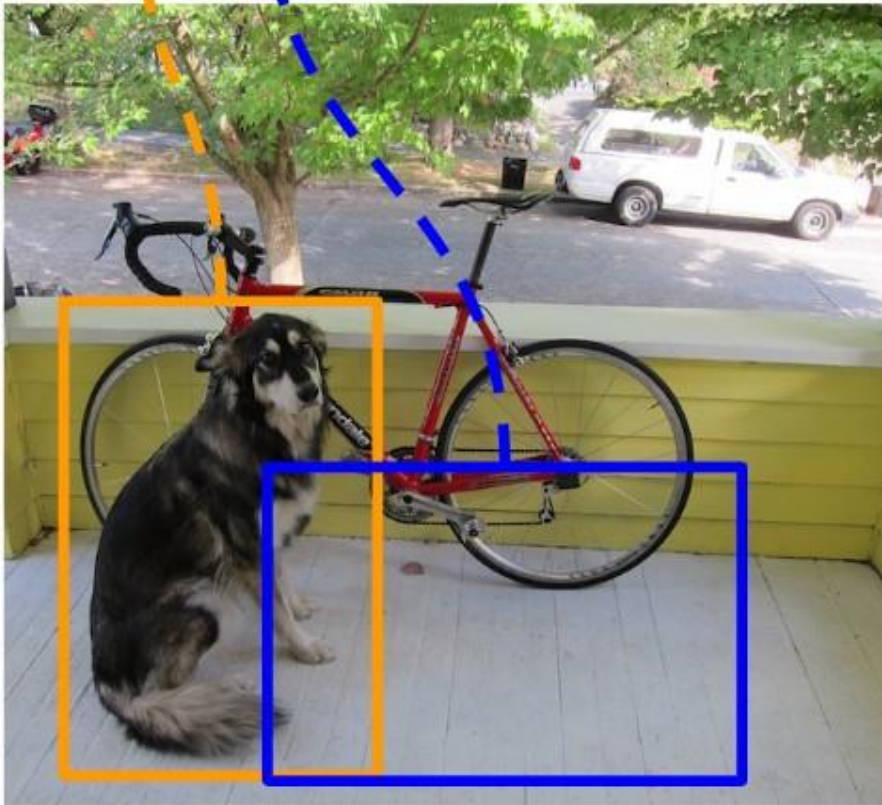
If $\text{IoU}(\text{bbox_max}, \text{bbox_cur}) > 0.5$ then set 0 score to bbox_cur .

In this case: set to 0.

Non-Maximum Suppression: intuition

class (dog) scores for each bbox

class: dog	bb47	bb20	bb15	bb7									bb1	bb4	bb8	bb98	1x98
	0.5	0	0.2	0.1									0	0	0	0	



Go to next `bbox_cur`.

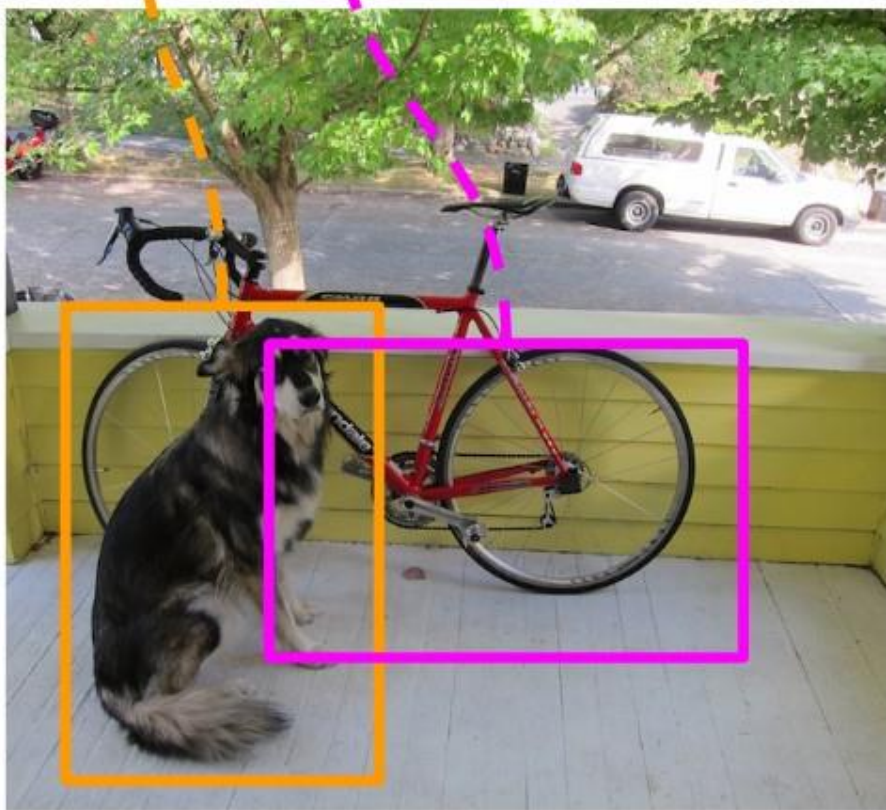
If $\text{IoU}(\text{bbox_max}, \text{bbox_cur}) > 0.5$ then set 0 score to `bbox_cur`.

In this case: continue.

Non-Maximum Suppression: intuition

class (dog) scores for each bbox

class: dog	bb47	bb20	bb15	bb7									bb1	bb4	bb8	bb98	1x98
	0.5	0	0.2	0.1									0	0	0	0	



Go to next **bbox_cur**.

If $\text{IoU}(\text{bbox_max}, \text{bbox_cur}) > 0.5$ then set 0 score to **bbox_cur**.

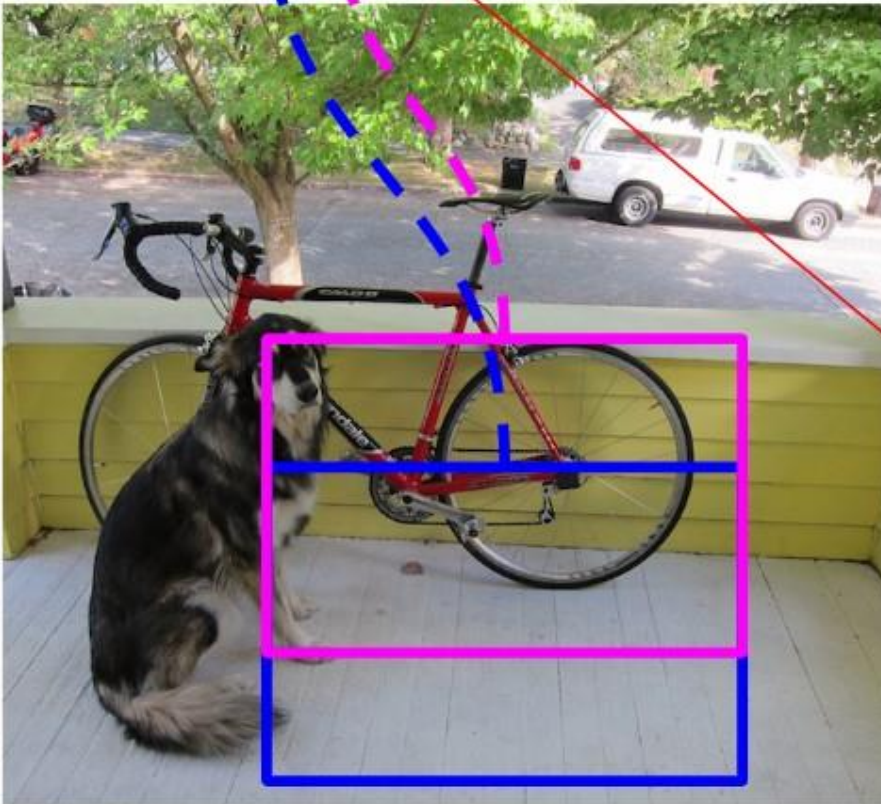
In this case: continue.

Do this procedure for other "bbox_cur". After that ...

Non-Maximum Suppression: intuition

class (dog) scores for each bbox

class: dog	bb47	bb20	bb15	bb7									bb1	bb4	bb8	bb98	1x98
	0.5	0	0.2	0									0	0	0	0	

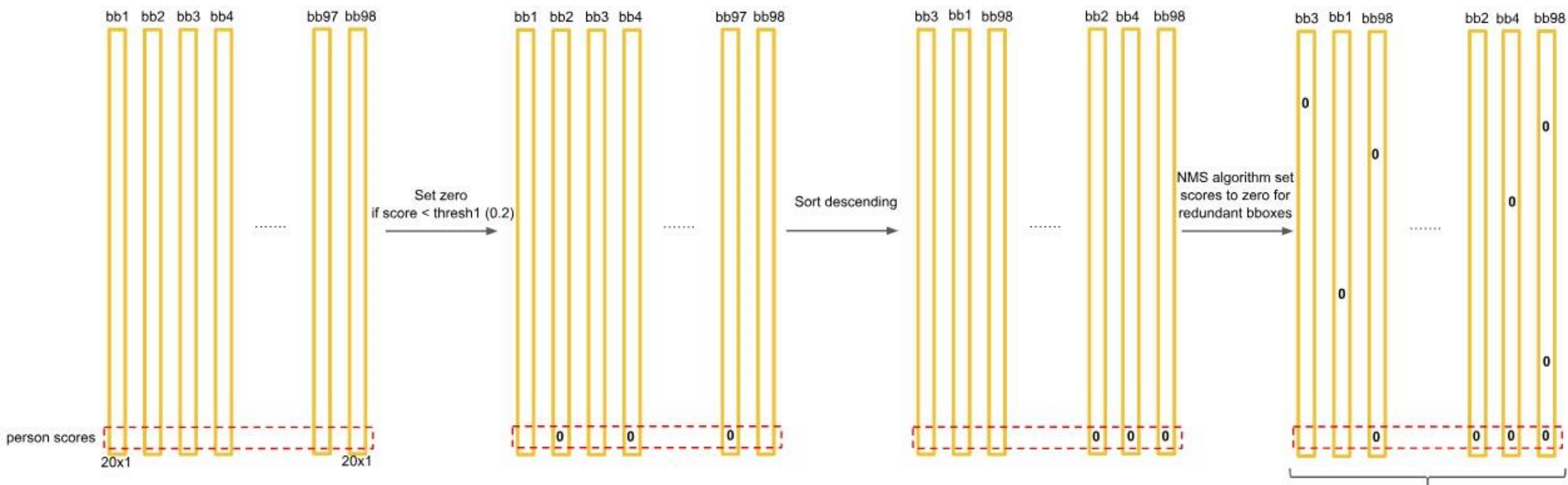


Go to next **bbox_cur**.

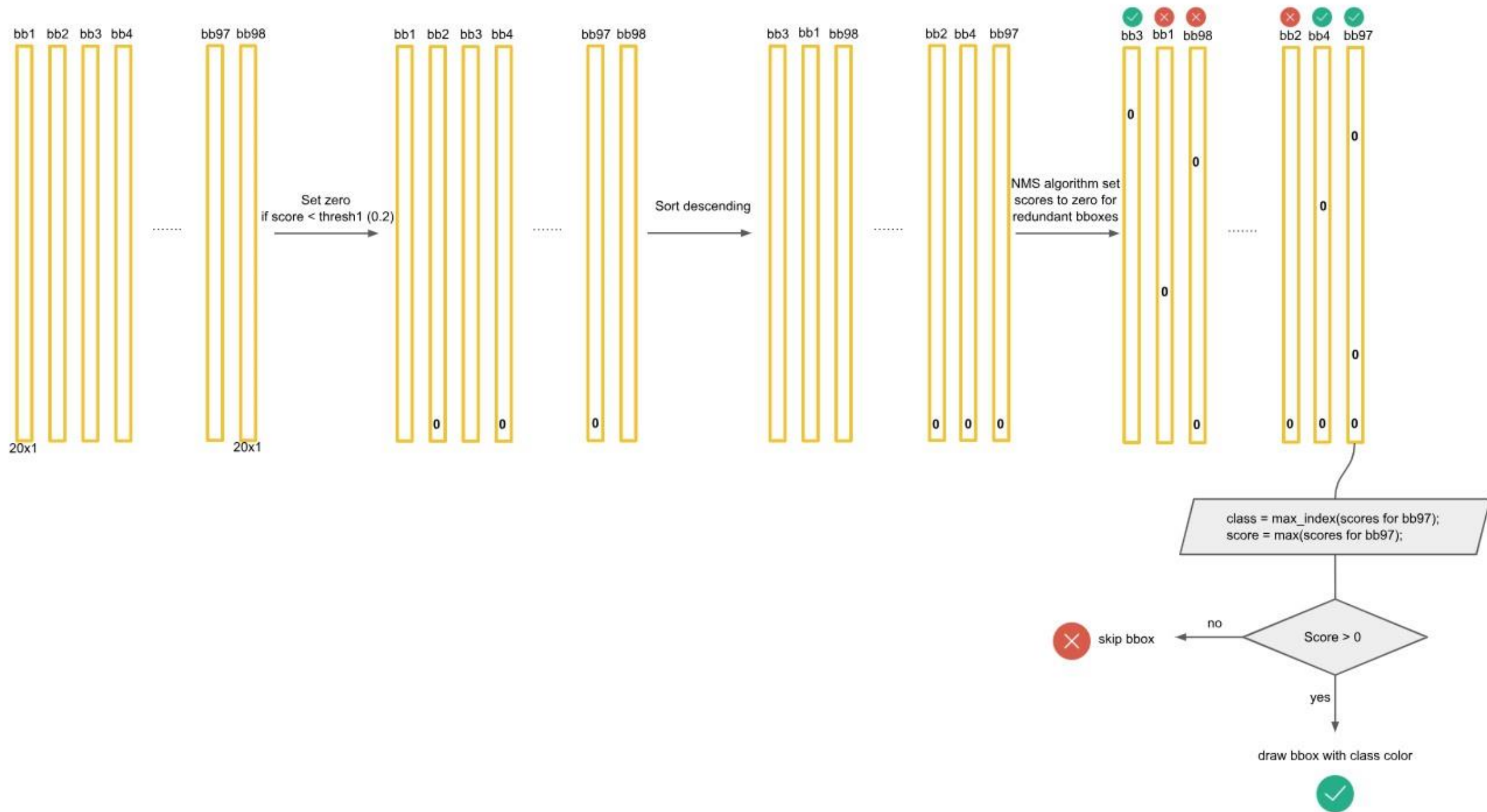
If $\text{IoU}(\text{bbox_max}, \text{bbox_cur}) > 0.5$ then set 0 score to **bbox_cur**.

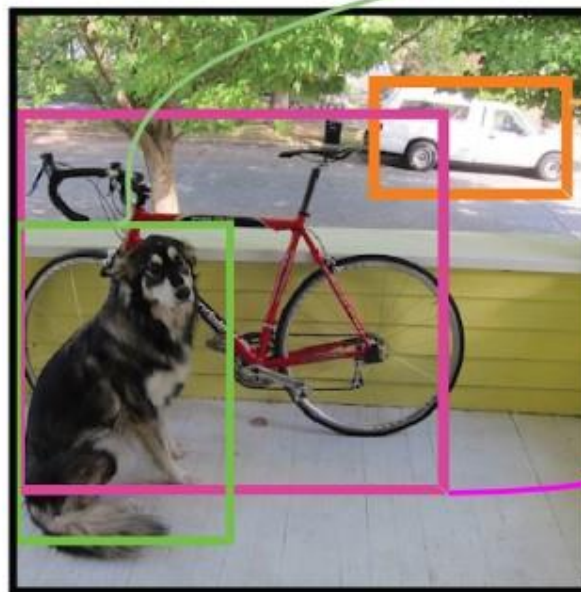
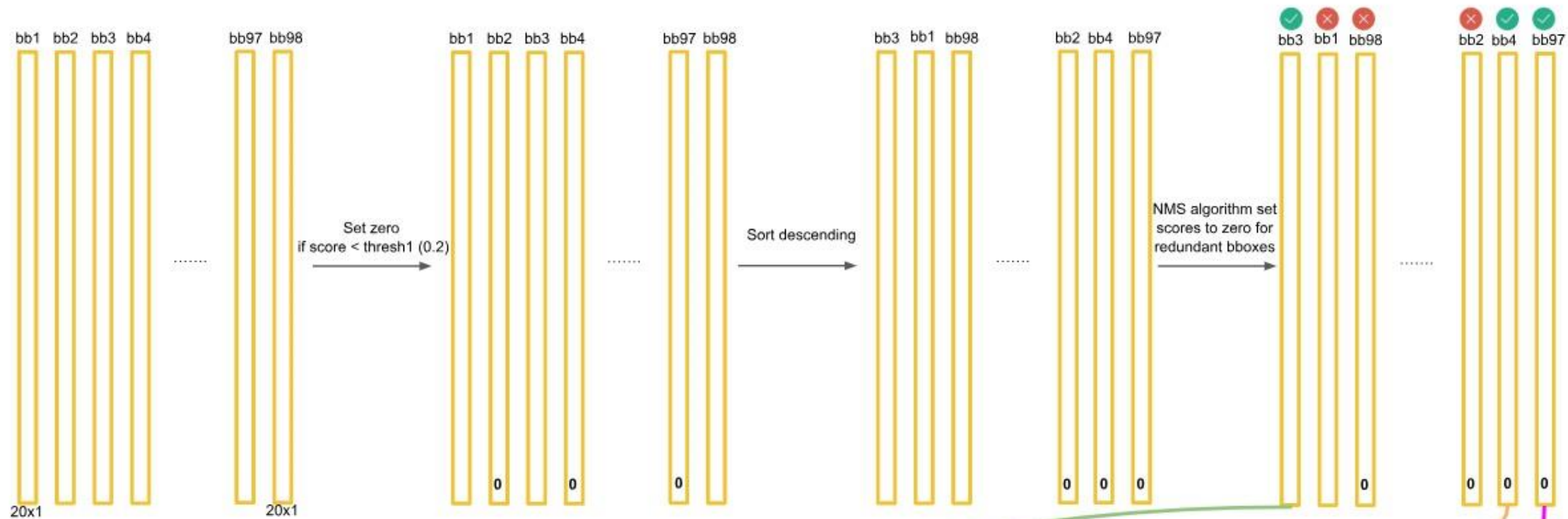
In this case: set to 0.

Do this procedure for other “bbox_max” and for other corresponding “bbox_cur”.



After this procedure -
a lot of zeros





$$\begin{aligned}
& \lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} \left[(x_i - \hat{x}_i)^2 + (y_i - \hat{y}_i)^2 \right] \\
& + \lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} \left[\left(\sqrt{w_i} - \sqrt{\hat{w}_i} \right)^2 + \left(\sqrt{h_i} - \sqrt{\hat{h}_i} \right)^2 \right] \\
& + \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} (C_i - \hat{C}_i)^2 \\
& + \lambda_{\text{noobj}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{noobj}} (C_i - \hat{C}_i)^2 \\
& + \sum_{i=0}^{S^2} \mathbb{1}_i^{\text{obj}} \sum_{c \in \text{classes}} (p_i(c) - \hat{p}_i(c))^2
\end{aligned}$$

Regression problem : SSE

$\lambda_{\text{coord}} = 5$ and $\lambda_{\text{noobj}} = .5$.

Classification based cell