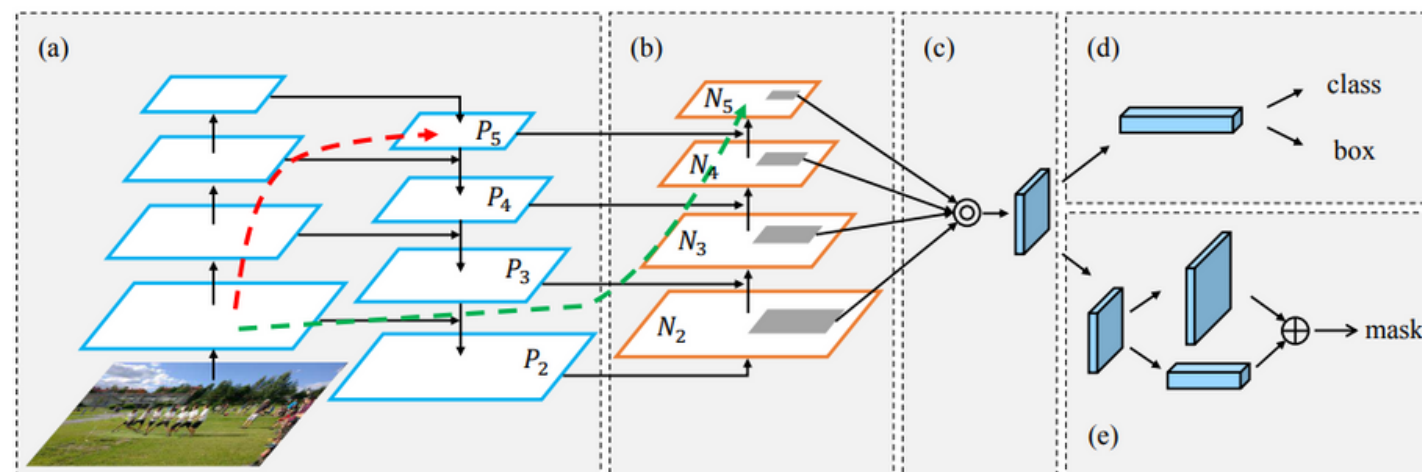


# Detección de Objetos

## [YOLO, FPN, Focal Loss, PAN]



**CLASE 7**

José M. Saavedra R.  
Profesor Asistente

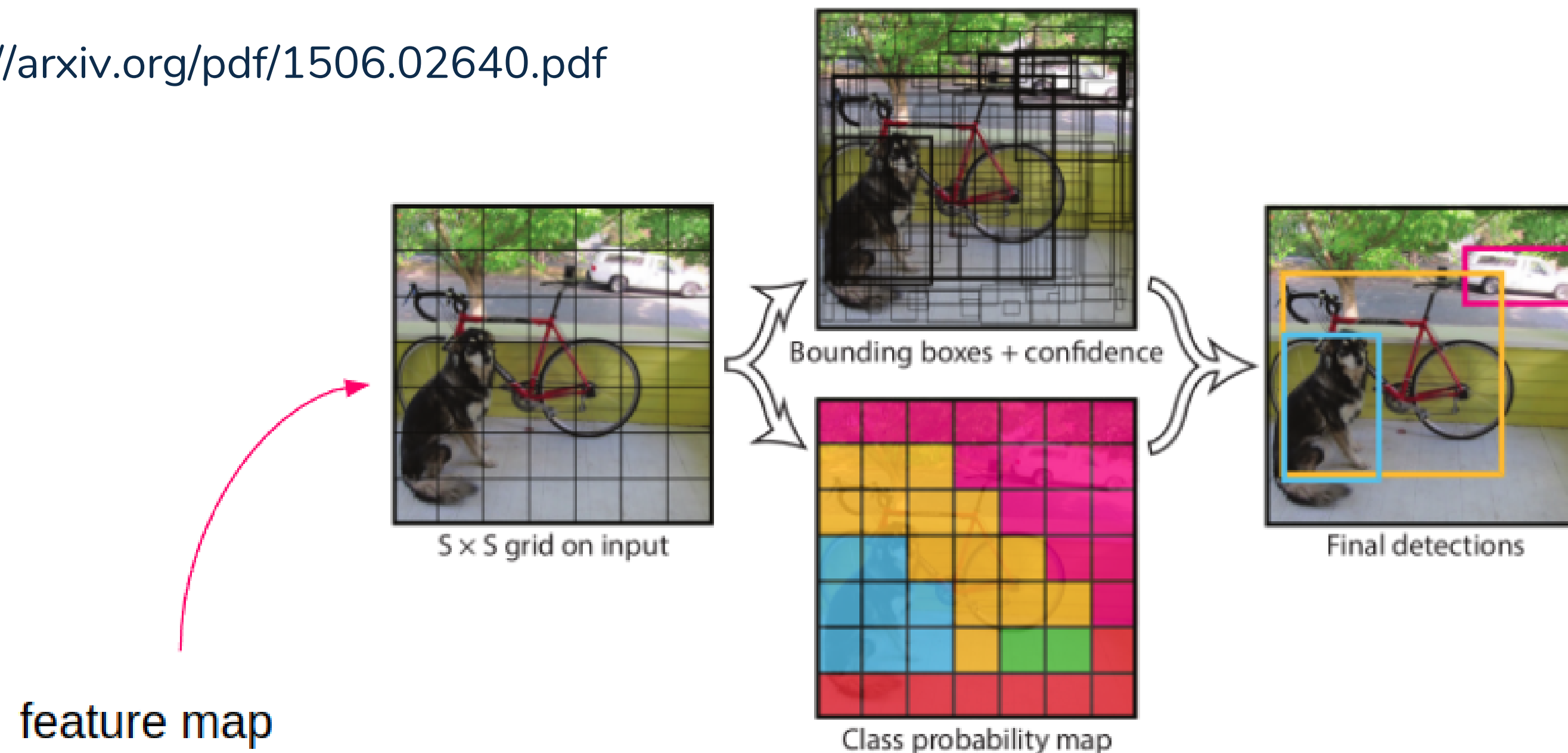
[jmsaavedrar@miuandes.cl](mailto:jmsaavedrar@miuandes.cl)

Ed. Ingeniería - Oficina 315

# YOLO [You Only Look Once]

Unlike Faster R-CNN, YOLO is a one-stage model performing region proposal and **classification in a single step**.

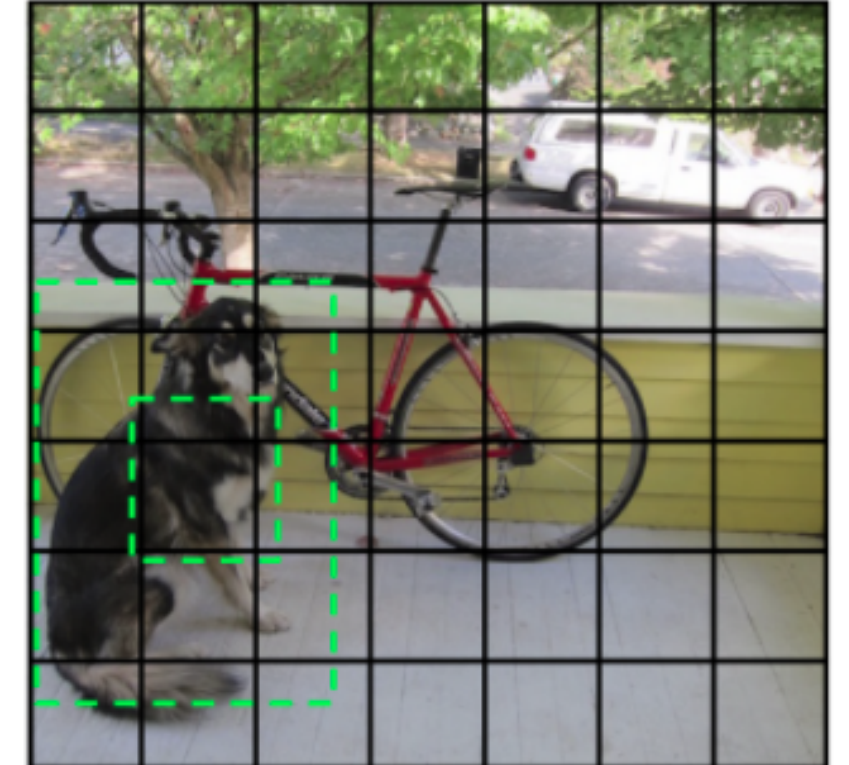
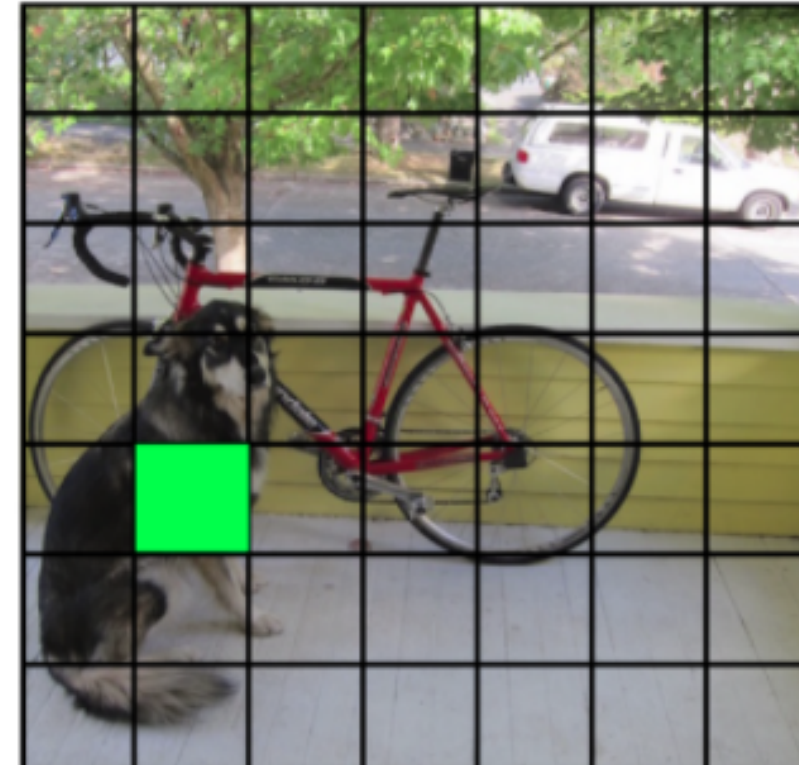
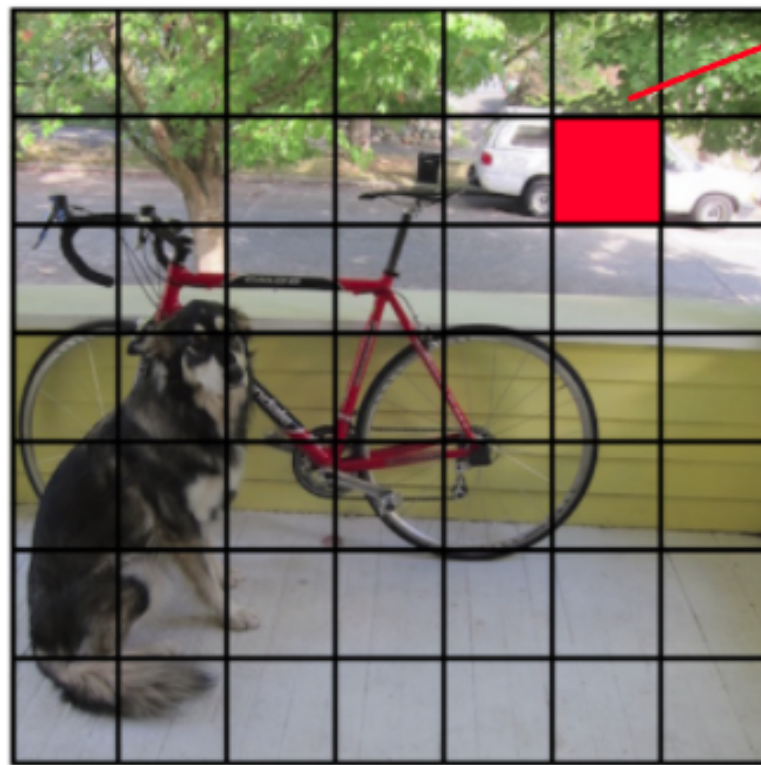
<https://arxiv.org/pdf/1506.02640.pdf>



# YOLO [You Only Look Once]

## Predicted Boxes

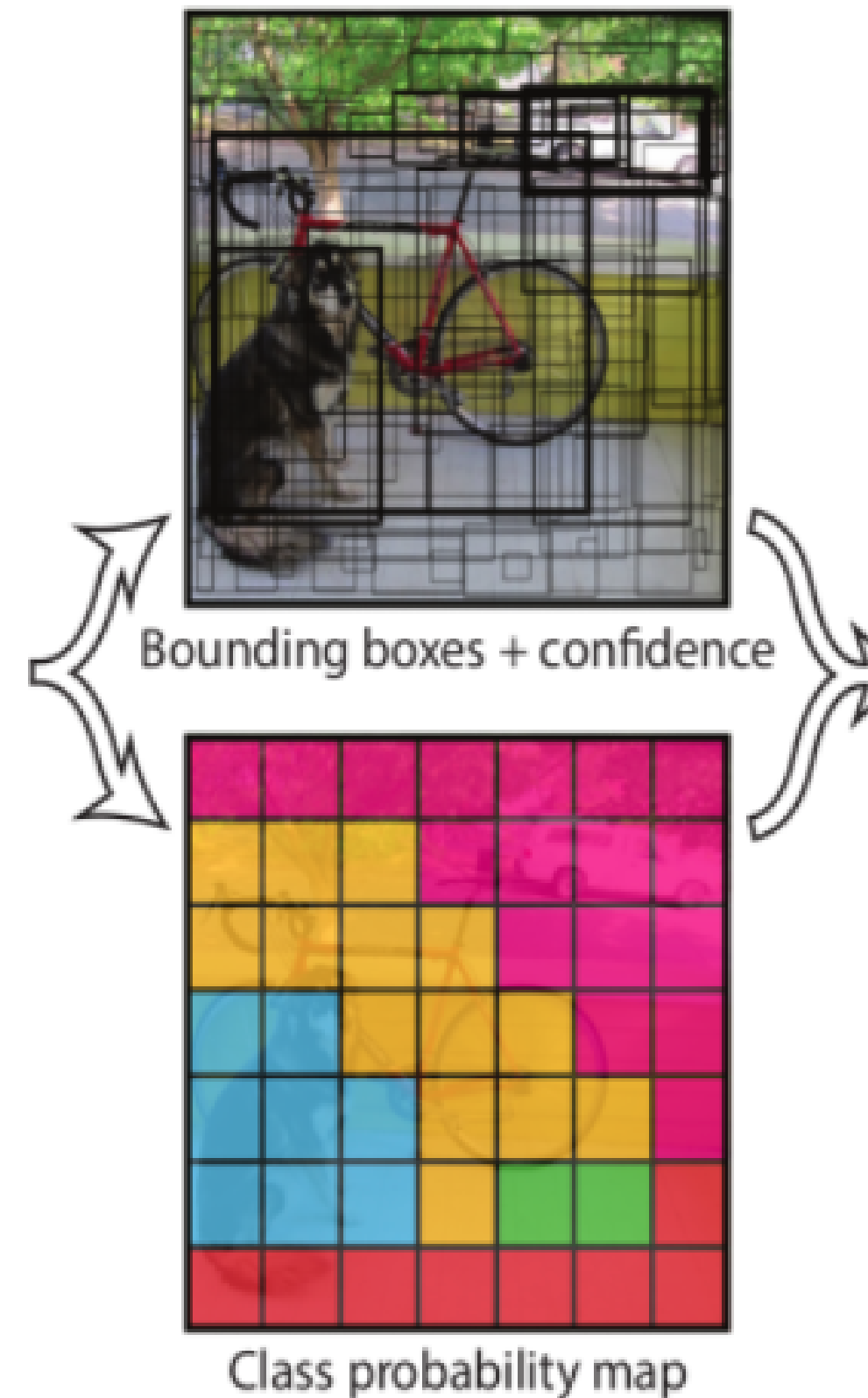
Each grid cell predicts  $B$  bounding boxes and confidence scores for those boxes under one-stage fashion.



# YOLO [You Only Look Once]

## Predicted Classes

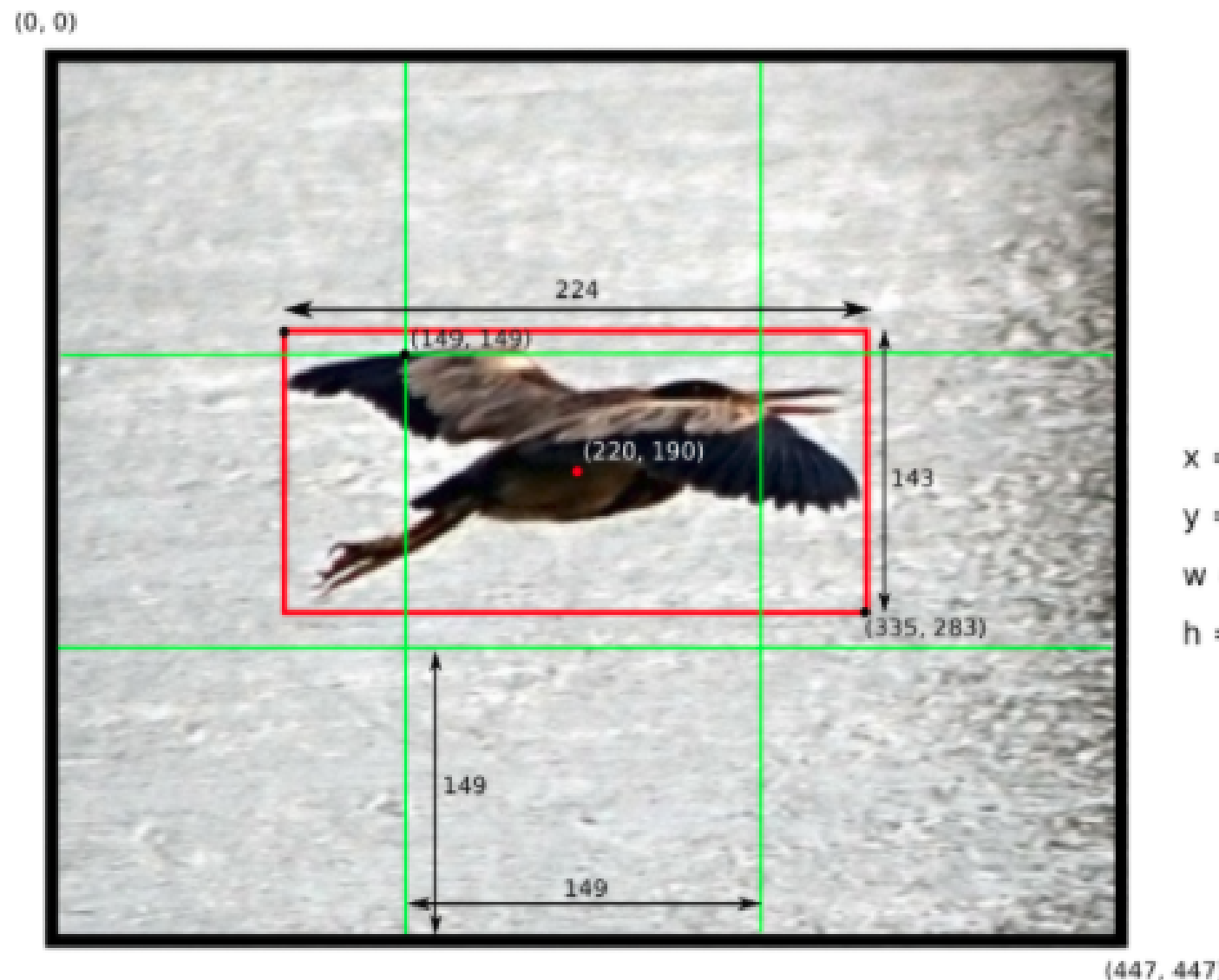
Each grid cell also predicts  $C$  conditional class probabilities,  $\Pr(\text{Class} | \text{Object})$ . These probabilities are conditioned on the grid cell containing an object





# YOLO [You Only Look Once]

- Coordinates of a box are relative to the whole image, varying in [0,1].
- If the center of an object falls into a grid cell, that grid cell is responsible for detecting that object.



relative coordinates

$$x = (220 - 149) / 149 = 0.48$$

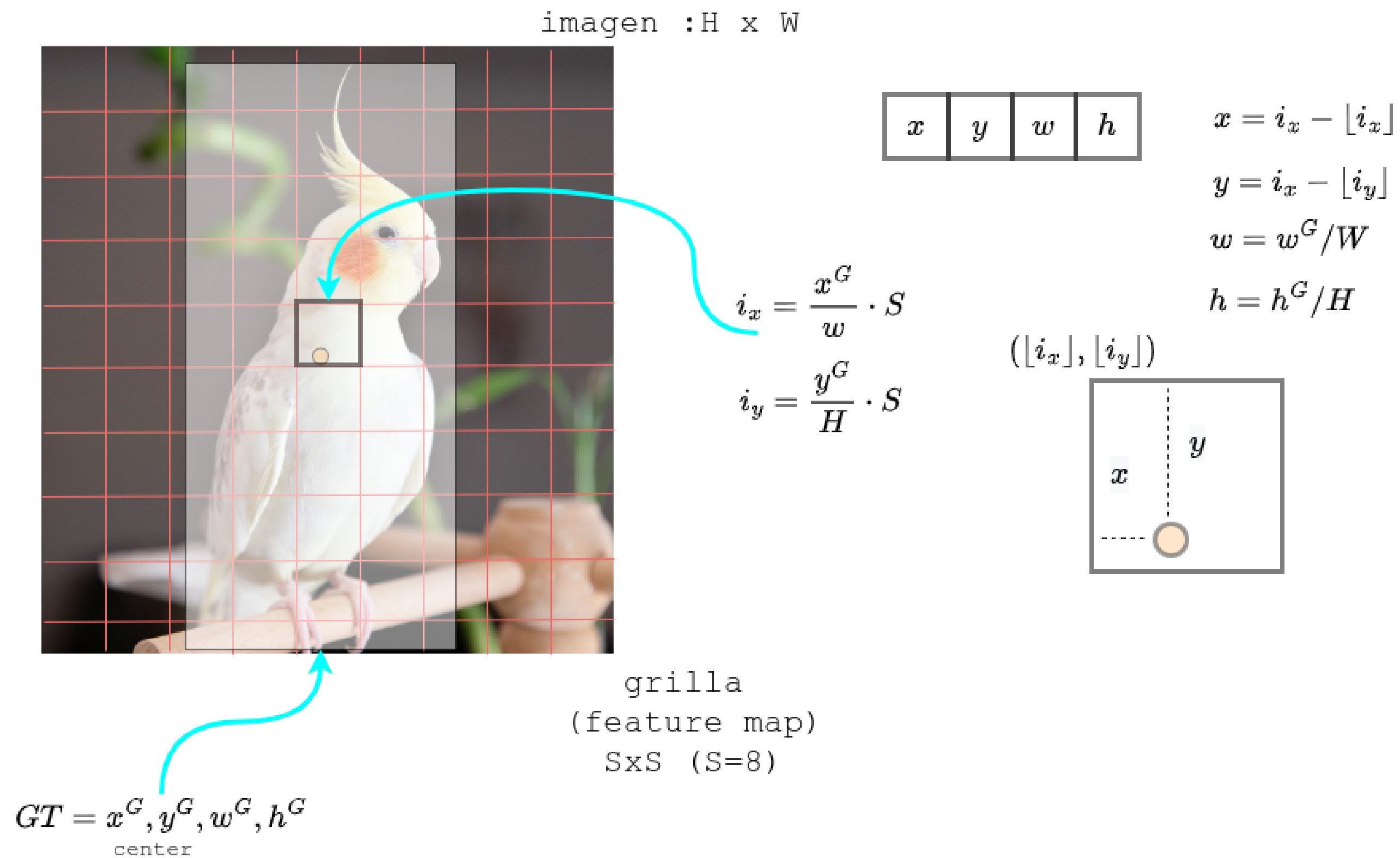
$$y = (190 - 149) / 149 = 0.28$$

$$w = 224 / 448 = 0.50$$

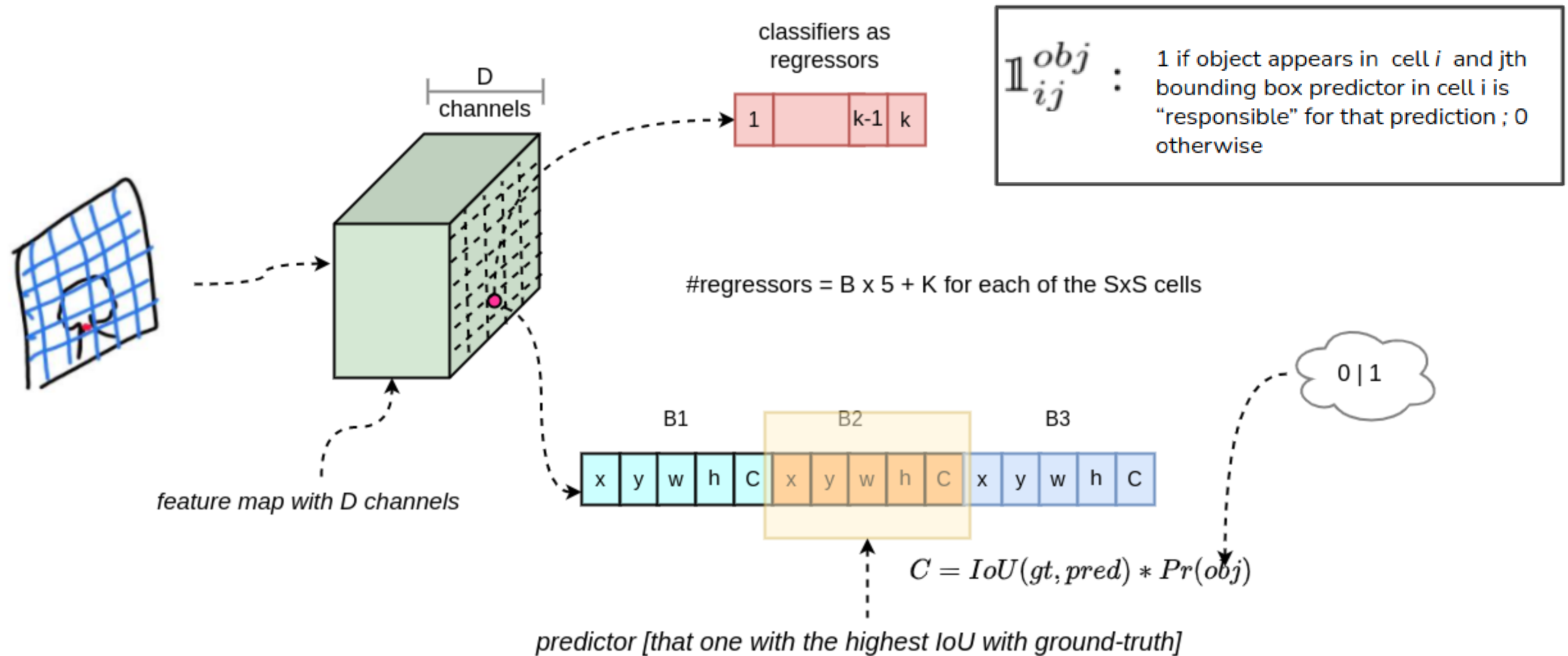
$$h = 143 / 448 = 0.32$$

Only one cell is responsible for an object.

# YOLO [You Only Look Once]



# YOLO [You Only Look Once]



YOLO predicts multiple bounding boxes per grid cell. At training time we only want one bounding box predictor to be responsible for each object. We assign one predictor to be "responsible" for predicting an object based on which prediction has the highest current IOU with the ground truth.

# YOLO [You Only Look Once]

## YOLO Loss

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

this term penalizes bad localization of center of cells

this term penalizes the bounding box with inaccurate height and width. The square root reflects that small deviations in large boxes matter less than in small boxes

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

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

this term tries to make the confidence score equal to the IOU between the object and the prediction when there is one object

this tries to make confidence score close to 0 when there is no object in the cell

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

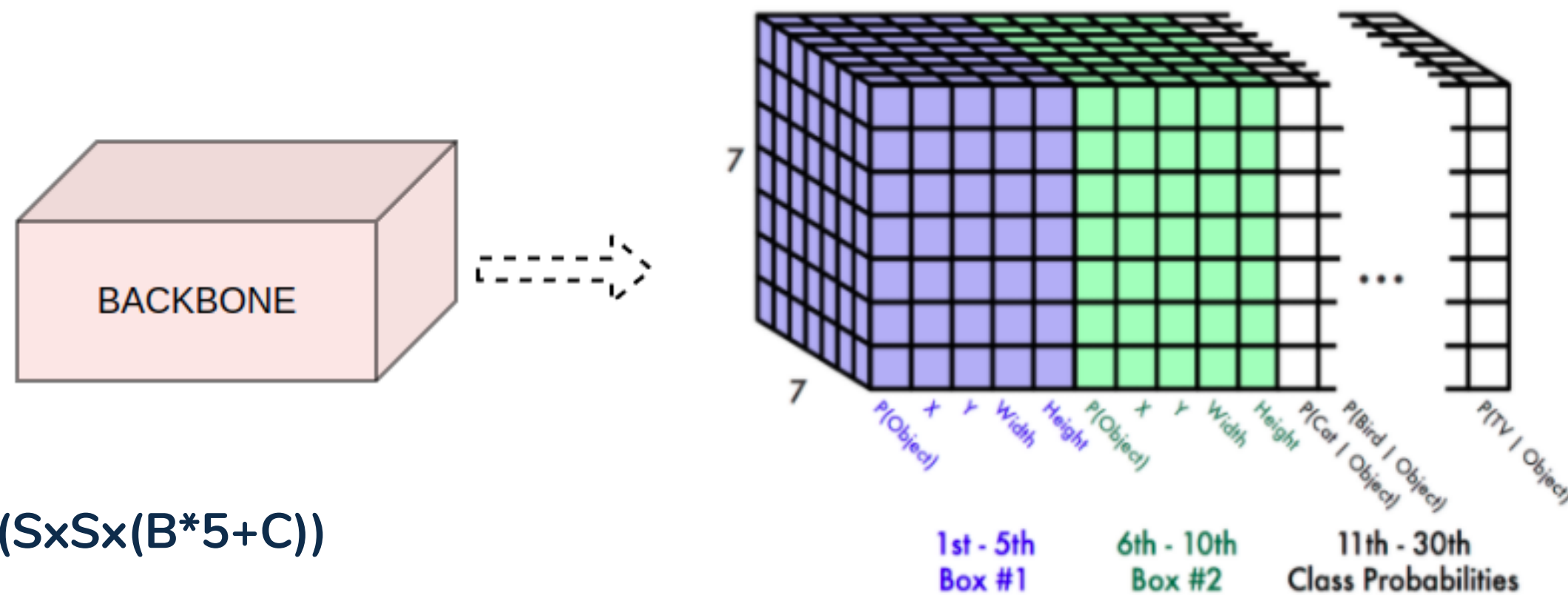
$$+ \sum_{i=0}^{S^2} \mathbb{1}_i^{\text{obj}} \sum_{c \in \text{classes}} (p_i(c) - \hat{p}_i(c))^2$$

This is a simple classification loss



# YOLO [You Only Look Once]

## YOLO Model



Head: regressor de  $(S \times S \times (B \times 5 + C))$

$S = 7$

$B = 2$

$C$  = número de clases, e.g. 20 en Pascal VOC

output  $\rightarrow 7 \times 7 \times 30$

# YOLO [You Only Look Once]

Real-Time Detectors	Train	mAP	FPS
100Hz DPM [31]	2007	16.0	100
30Hz DPM [31]	2007	26.1	30
Fast YOLO	2007+2012	52.7	<b>155</b>
YOLO	2007+2012	<b>63.4</b>	45
Less Than Real-Time			
Fastest DPM [38]	2007	30.4	15
R-CNN Minus R [20]	2007	53.5	6
Fast R-CNN [14]	2007+2012	70.0	0.5
Faster R-CNN VGG-16[28]	2007+2012	73.2	7
Faster R-CNN ZF [28]	2007+2012	62.1	18
YOLO VGG-16	2007+2012	66.4	21

**Table 1: Real-Time Systems on PASCAL VOC 2007.** Comparing the performance and speed of fast detectors. Fast YOLO is the fastest detector on record for PASCAL VOC detection and is still twice as accurate as any other real-time detector. YOLO is 10 mAP more accurate than the fast version while still well above real-time in speed.

# YOLO [You Only Look Once]

## YOLO 9000 (Yolo-v2)

An Improved version of YOLO (anchors are included)

<https://arxiv.org/pdf/1612.08242.pdf>



# YOLOv2 [You Only Look Once]

## Improvements

- Batch normalization
- A higher resolution classifier ( $224 \rightarrow 448$ ,  $S = 13$ )
- Anchors for bounding boxes
- Number of anchors is estimated by clustering on training data.



# YOLOv2 [You Only Look Once]

## Anchors

**Faster-RCNN:**  $t_x, t_y$  are not limited, and the center of a predicted box can fall anywhere in the image producing instability.

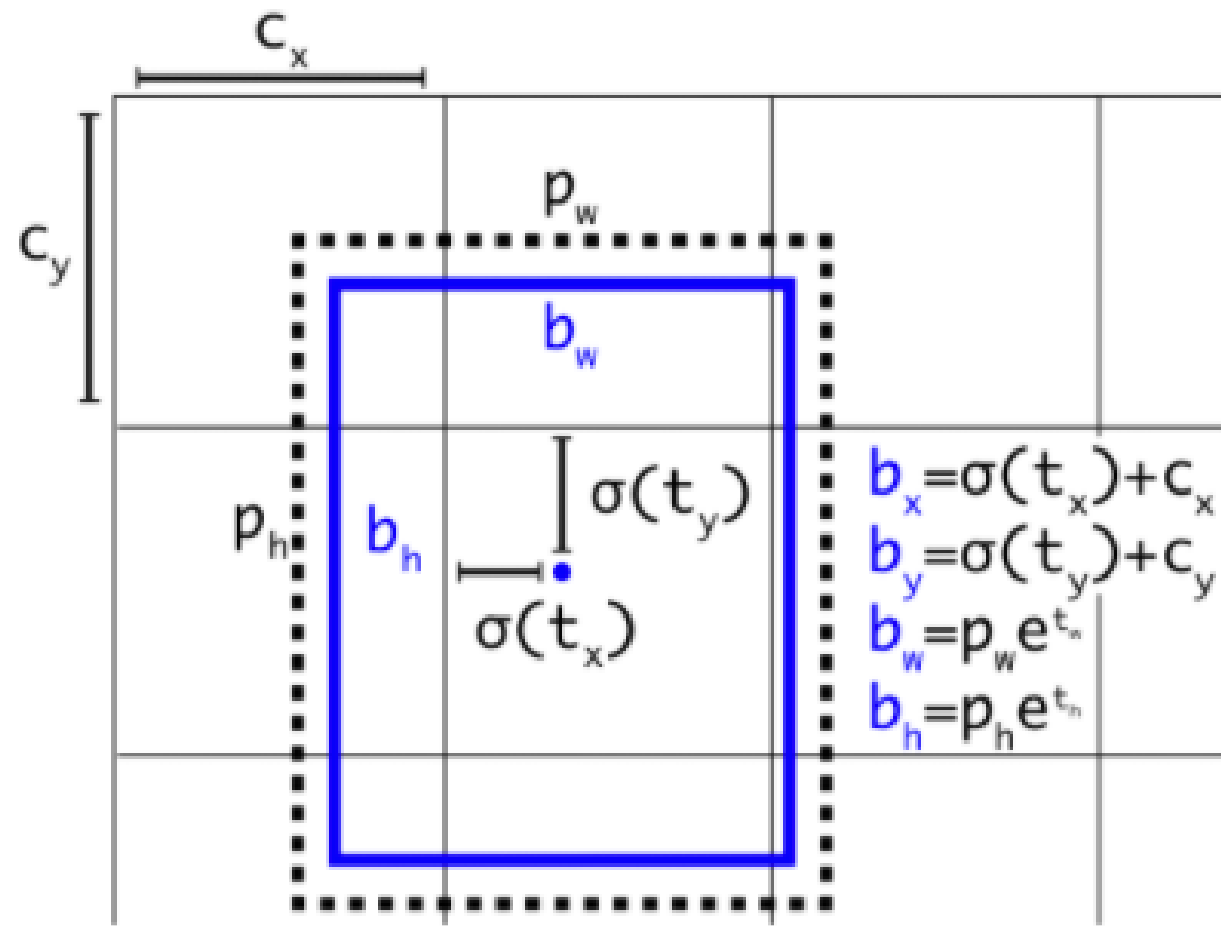
$$x = (t_x * w_a) + x_a$$

$$y = (t_y * h_a) + y_a$$

# YOLOv2 [You Only Look Once]

## Anchors

YOLOv2:  $x, y$  are constrained to fall within the corresponding cell



tx	ty	tw	th
----	----	----	----

sigma is the sigmoid function varying between 0 and 1

# YOLOv2 [You Only Look Once]

Detection Frameworks	Train	mAP	FPS
Fast R-CNN [5]	2007+2012	70.0	0.5
Faster R-CNN VGG-16[15]	2007+2012	73.2	7
Faster R-CNN ResNet[6]	2007+2012	76.4	5
YOLO [14]	2007+2012	63.4	45
SSD300 [11]	2007+2012	74.3	46
SSD500 [11]	2007+2012	76.8	19
YOLOv2 288 × 288	2007+2012	69.0	91
YOLOv2 352 × 352	2007+2012	73.7	81
YOLOv2 416 × 416	2007+2012	76.8	67
YOLOv2 480 × 480	2007+2012	77.8	59
YOLOv2 544 × 544	2007+2012	<b>78.6</b>	40

**Table 3: Detection frameworks on PASCAL VOC 2007.** YOLOv2 is faster and more accurate than prior detection methods. It can also run at different resolutions for an easy tradeoff between speed and accuracy. Each YOLOv2 entry is actually the same trained model with the same weights, just evaluated at a different size. All timing information is on a Geforce GTX Titan X (original, not Pascal model).

# YOLOv2 [You Only Look Once]

		0.5:0.95	0.5	0.75	S	M	L	1	10	100	S	M	L
Fast R-CNN [5]	train	19.7	35.9	-	-	-	-	-	-	-	-	-	-
Fast R-CNN[1]	train	20.5	39.9	19.4	4.1	20.0	35.8	21.3	29.5	30.1	7.3	32.1	52.0
Faster R-CNN[15]	trainval	21.9	42.7	-	-	-	-	-	-	-	-	-	-
ION [1]	train	23.6	43.2	23.6	6.4	24.1	38.3	23.2	32.7	33.5	10.1	37.7	53.6
Faster R-CNN[10]	trainval	24.2	45.3	23.5	7.7	26.4	37.1	23.8	34.0	34.6	12.0	38.5	54.4
SSD300 [11]	trainval35k	23.2	41.2	23.4	5.3	23.2	39.6	22.5	33.2	35.3	9.6	37.6	56.5
SSD512 [11]	trainval35k	<b>26.8</b>	<b>46.5</b>	<b>27.8</b>	<b>9.0</b>	<b>28.9</b>	<b>41.9</b>	<b>24.8</b>	<b>37.5</b>	<b>39.8</b>	<b>14.0</b>	<b>43.5</b>	<b>59.0</b>
YOLOv2 [11]	trainval35k	21.6	44.0	19.2	5.0	22.4	35.5	20.7	31.6	33.3	9.8	36.5	54.4

Table 5: Results on COCO test-dev2015. Table adapted from [11]



# COCO

<https://cocodataset.org>



info@cocodataset.org

[Home](#) [People](#) [Dataset](#) [Tasks](#) [Evaluate](#)

## News

- We are pleased to announce the [LVIS 2021 Challenge and Workshop](#) to be held at ICCV.
- Please note that there will not be a COCO 2021 Challenge, instead, we encourage people to participate in the LVIS 2021 Challenge.
- We have partnered with the team behind the open-source tool [FiftyOne](#) to make it easier to download, visualize, and evaluate COCO
- [FiftyOne](#) is an open-source tool facilitating visualization and access to COCO data resources and serves as an evaluation tool for model analysis on COCO.

## What is COCO?



COCO is a large-scale object detection, segmentation, and captioning dataset. COCO has several features:

- ✓ Object segmentation
- ✓ Recognition in context
- ✓ Superpixel stuff segmentation
- ✓ 330K images (>200K labeled)
- ✓ 1.5 million object instances
- ✓ 80 object categories
- ✓ 91 stuff categories
- ✓ 5 captions per image
- ✓ 250,000 people with keypoints

## Collaborators

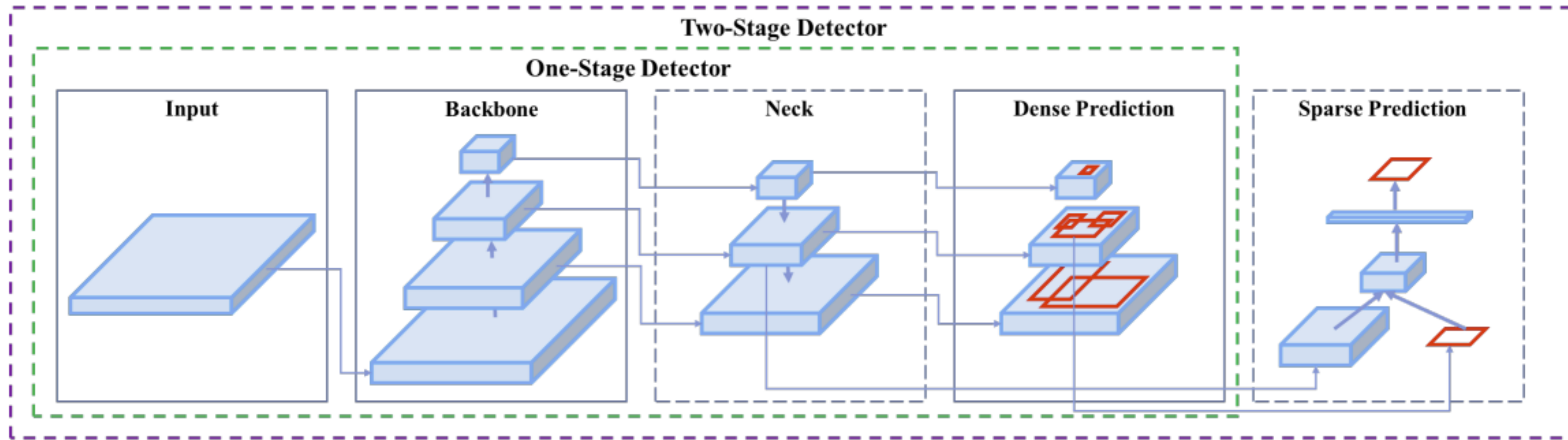
Tsung-Yi Lin Google Brain  
Genevieve Patterson MSR, Trash TV  
Matteo R. Ronchi Caltech  
Yin Cui Google  
Michael Maire TTI-Chicago  
Serge Belongie Cornell Tech  
Lubomir Bourdev WaveOne, Inc.  
Ross Girshick FAIR  
James Hays Georgia Tech  
Pietro Perona Caltech  
Deva Ramanan CMU  
Larry Zitnick FAIR  
Piotr Dollár FAIR

## Sponsors



## Dataset examples





**Input:** { Image, Patches, Image Pyramid, ... }

**Backbone:** { VGG16 [68], ResNet-50 [26], ResNeXt-101 [86], Darknet53 [63], ... }

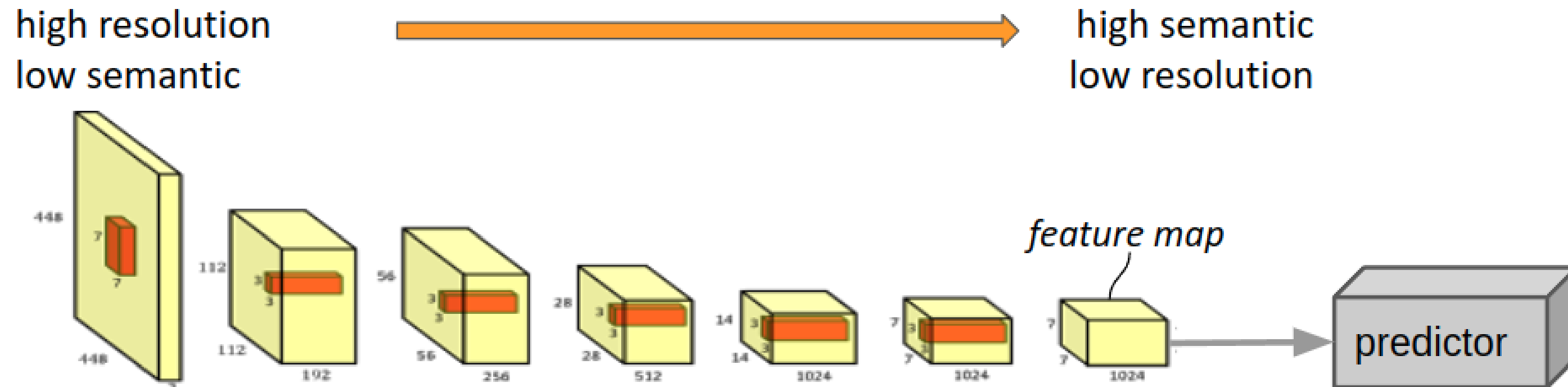
**Neck:** { FPN [44], PANet [49], Bi-FPN [77], ... }

**Head:**

**Dense Prediction:** { RPN [64], YOLO [61, 62, 63], SSD [50], RetinaNet [45], FCOS [78], ... }

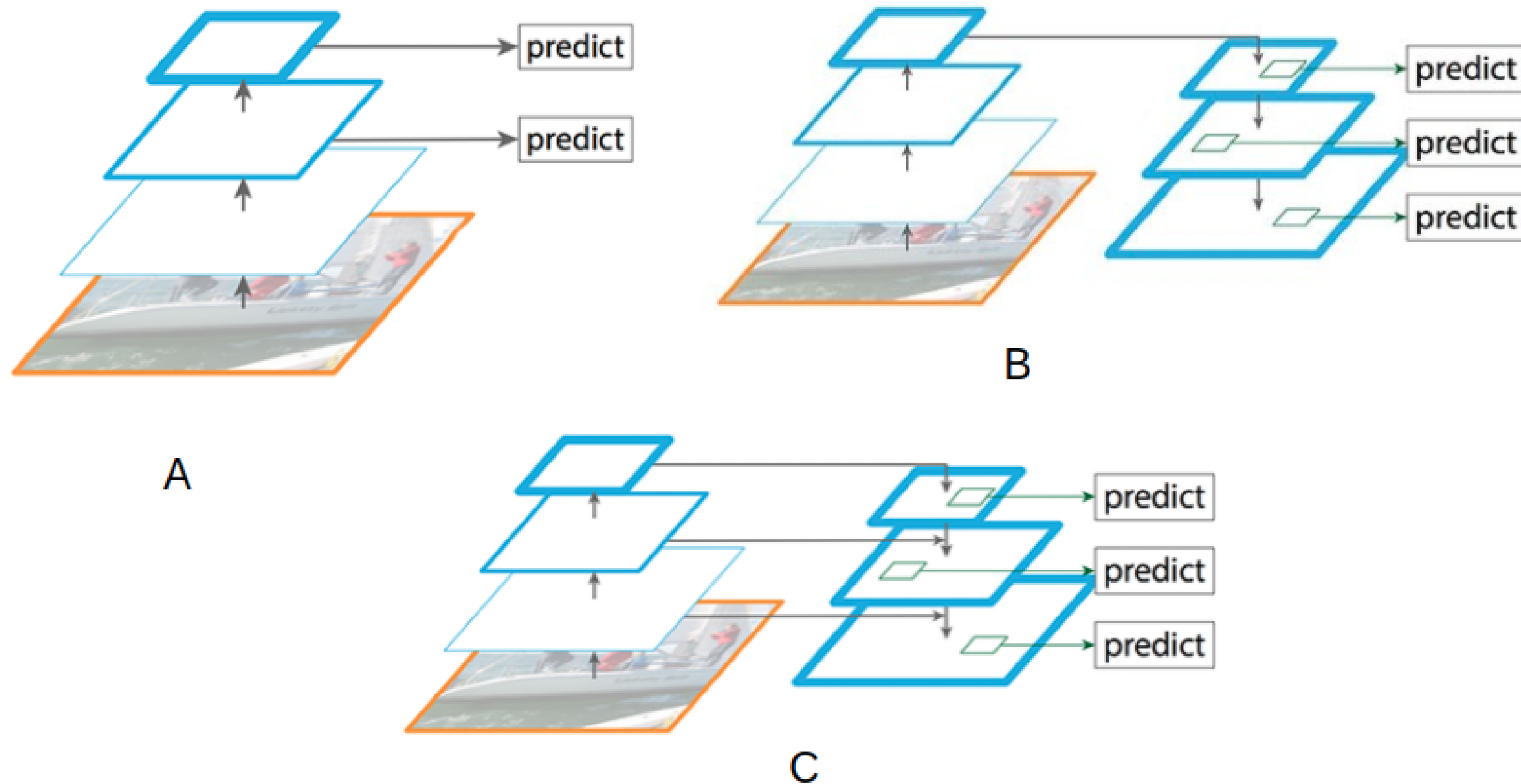
**Sparse Prediction:** { Faster R-CNN [64], R-FCN [9], ... }

# A convnet [single-scale]



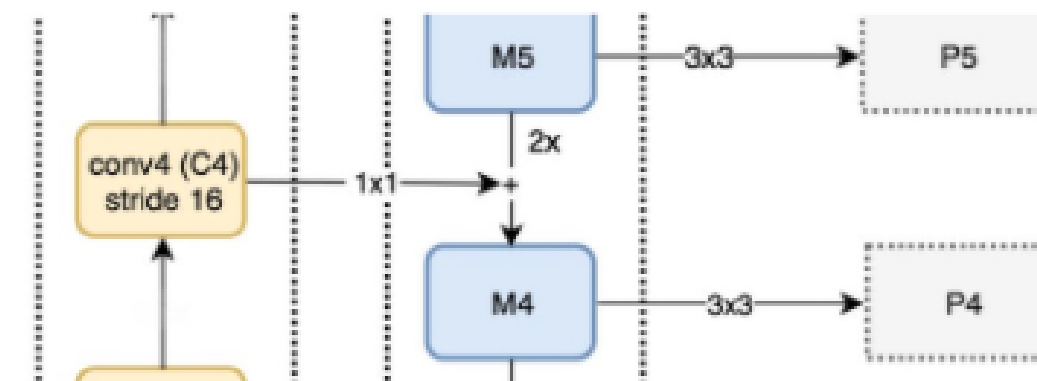
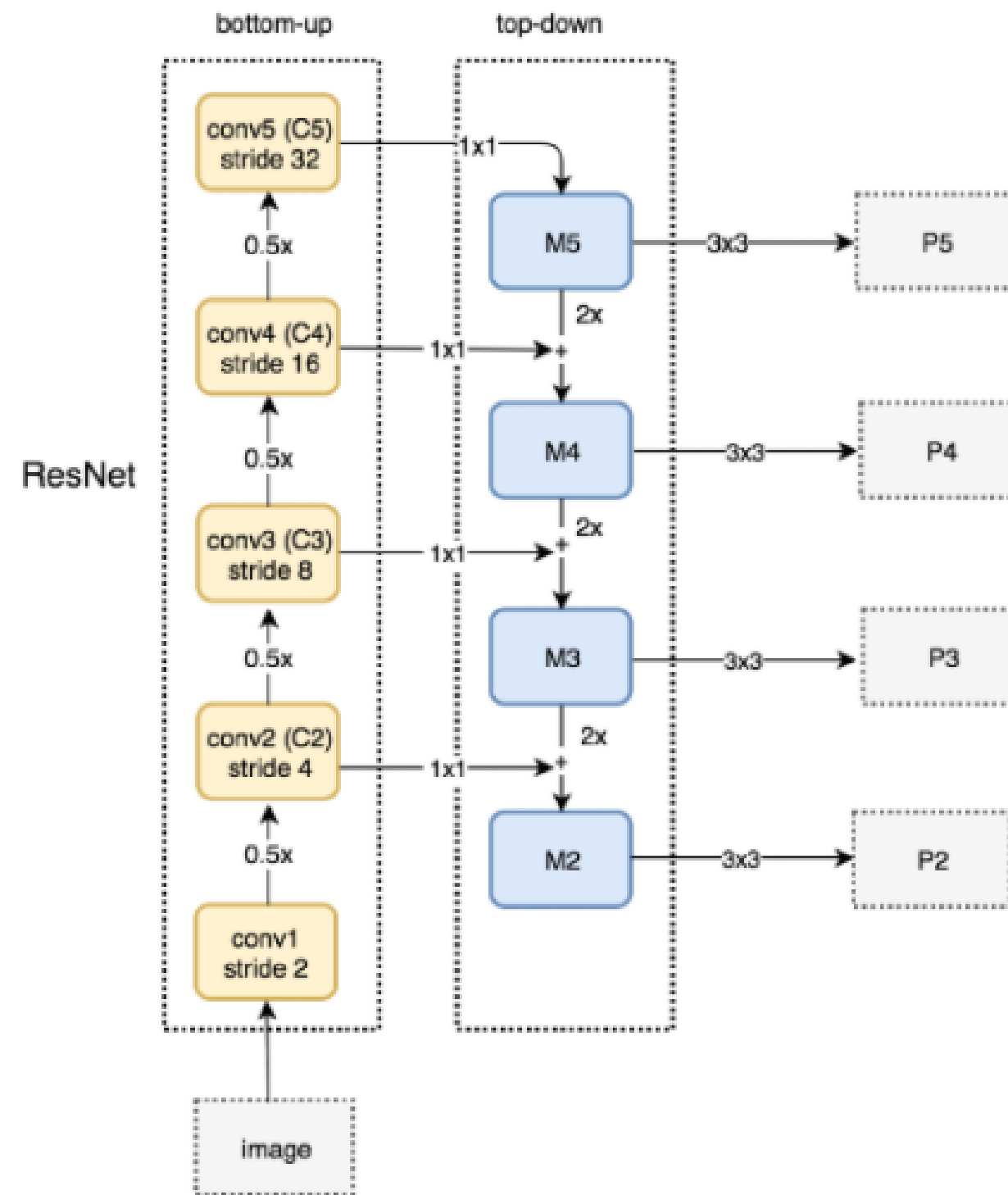
prediction head processes a single-scale feature map

# Feature Pyramid Network



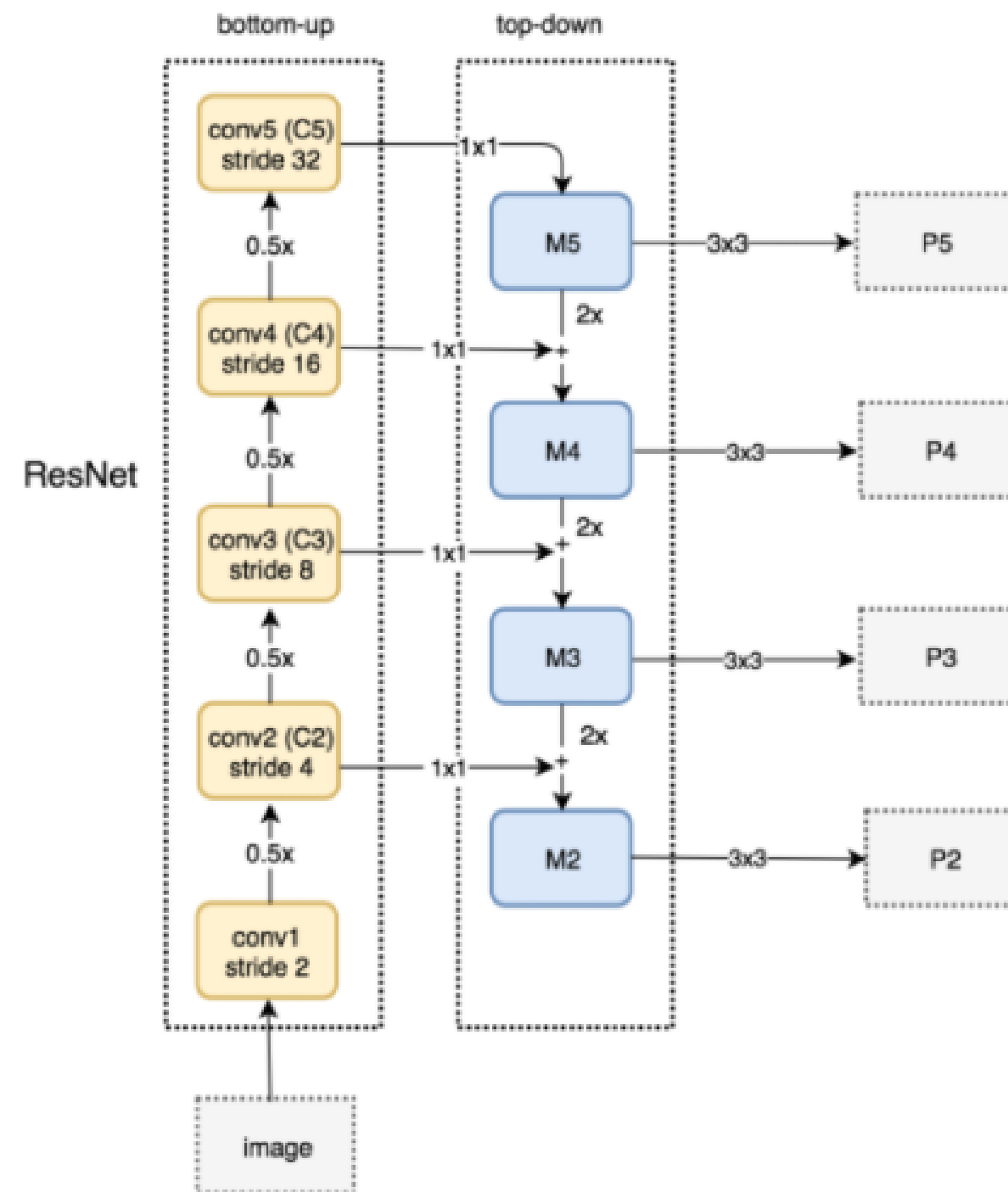


# Feature Pyramid Network

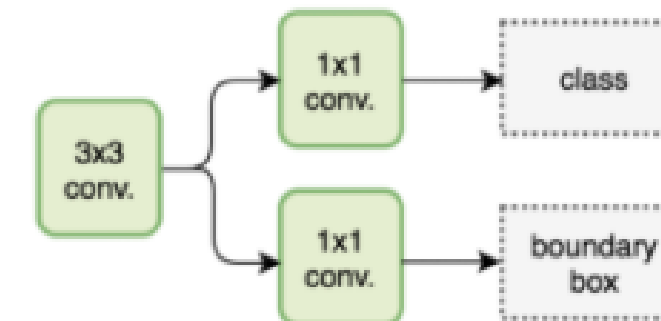


<https://arxiv.org/pdf/1612.03144.pdf>

# Feature Pyramid Network



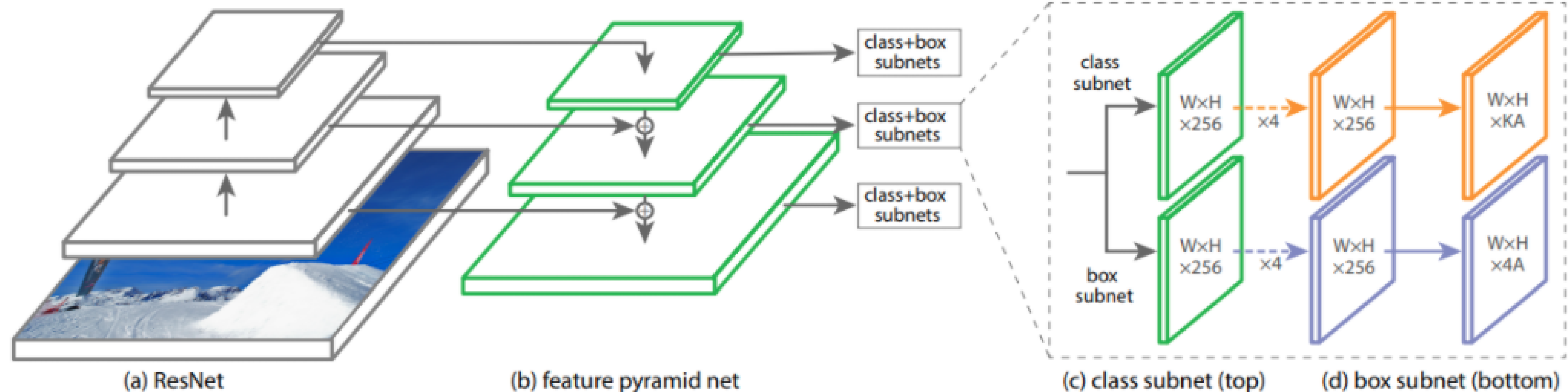
We note that the parameters of the heads are shared across all feature pyramid levels; we have also evaluated the alternative without sharing parameters and observed similar accuracy. The good performance of sharing parameters indicates that all levels of our pyramid share similar semantic levels. This advantage is analogous to that of using a featurized image pyramid, where a common head classifier can be applied to features computed at any image scale.



predictor head

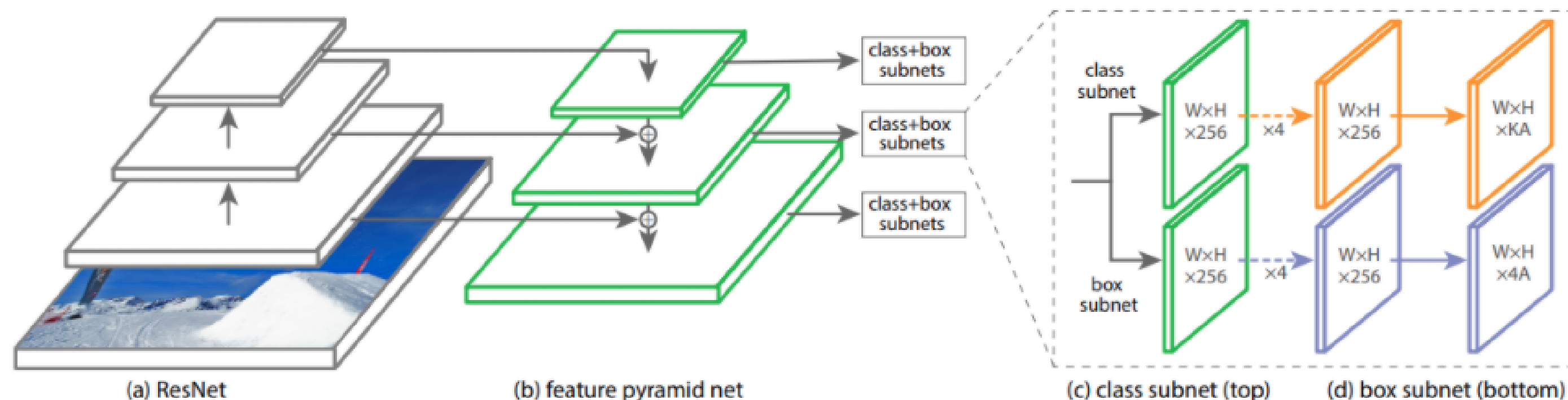
shared weights

# Feature Pyramid Network



Because all levels of the pyramid use shared classifiers/regressors as in a traditional featurized image pyramid, we fix the feature dimension (numbers of channels, denoted as  $d$ ) in all the feature maps. We set  $d = 256$  in this paper and thus all extra convolutional layers have 256-channel outputs. There are no non-linearities in these extra layers, which we have empirically found to have minor impacts.

# Feature Pyramid Network

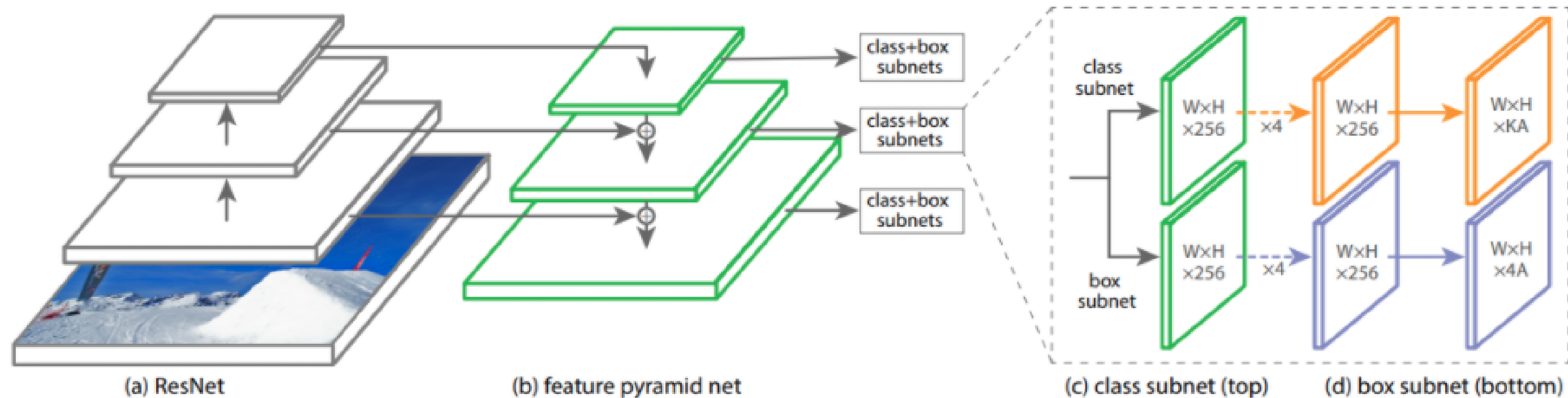


For region proposal (RPN)

anchors on a specific level. Instead, we assign anchors of a single scale to each level. Formally, we define the anchors to have areas of  $\{32^2, 64^2, 128^2, 256^2, 512^2\}$  pixels on  $\{P_2, P_3, P_4, P_5, P_6\}$  respectively.<sup>1</sup> As in [29] we also use anchors of multiple aspect ratios  $\{1:2, 1:1, 2:1\}$  at each level. So in total there are 15 anchors over the pyramid.



# Feature Pyramid Network

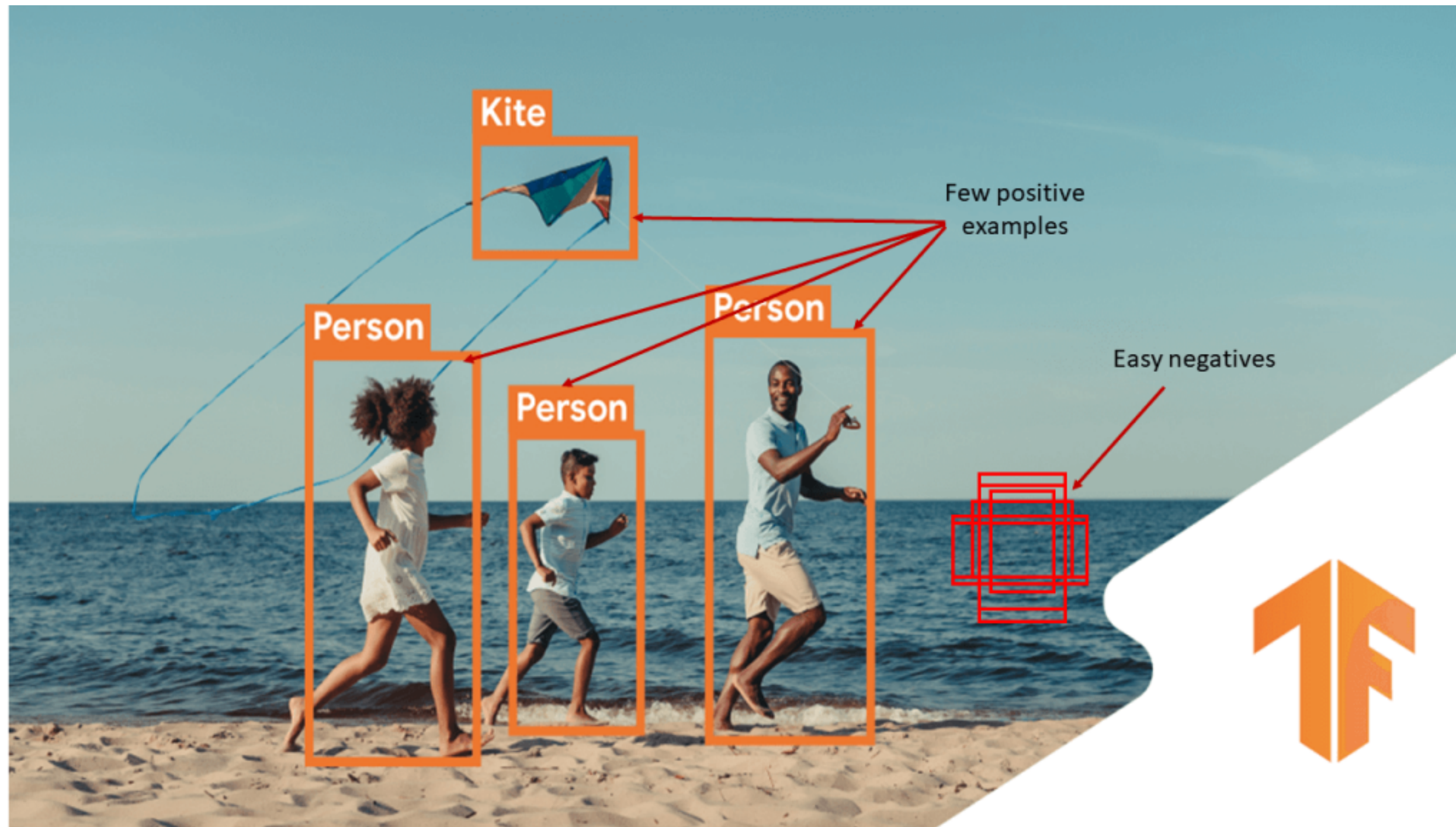


For classifier (Fast R-CNN)

We view our feature pyramid as if it were produced from an image pyramid. Thus we can adapt the assignment strategy of region-based detectors [15, 11] in the case when they are run on image pyramids. Formally, we assign an RoI of width  $w$  and height  $h$  (on the input image to the network) to the level  $P_k$  of our feature pyramid by:

$$k = \lfloor k_0 + \log_2(\sqrt{wh}/224) \rfloor. \quad (1)$$

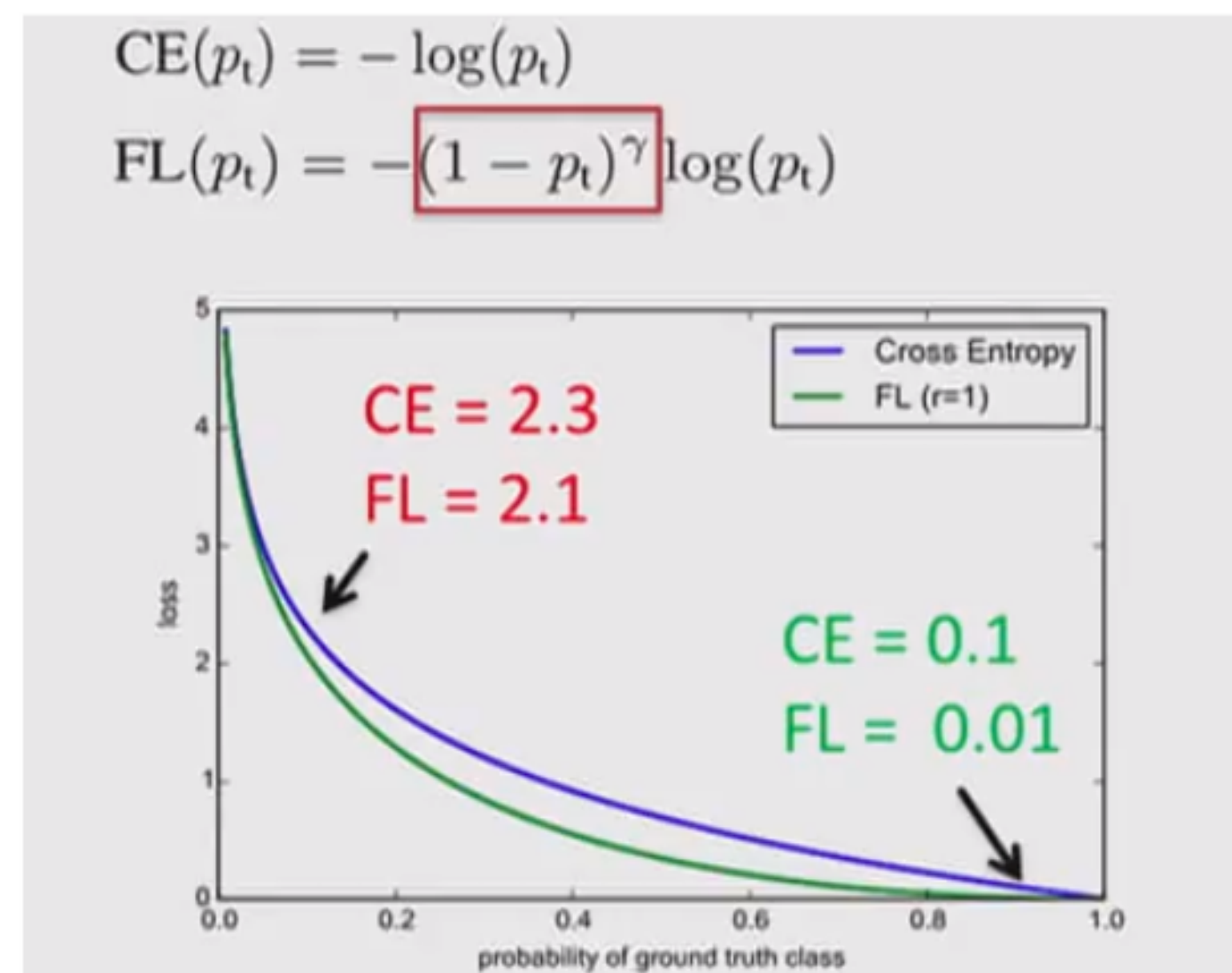
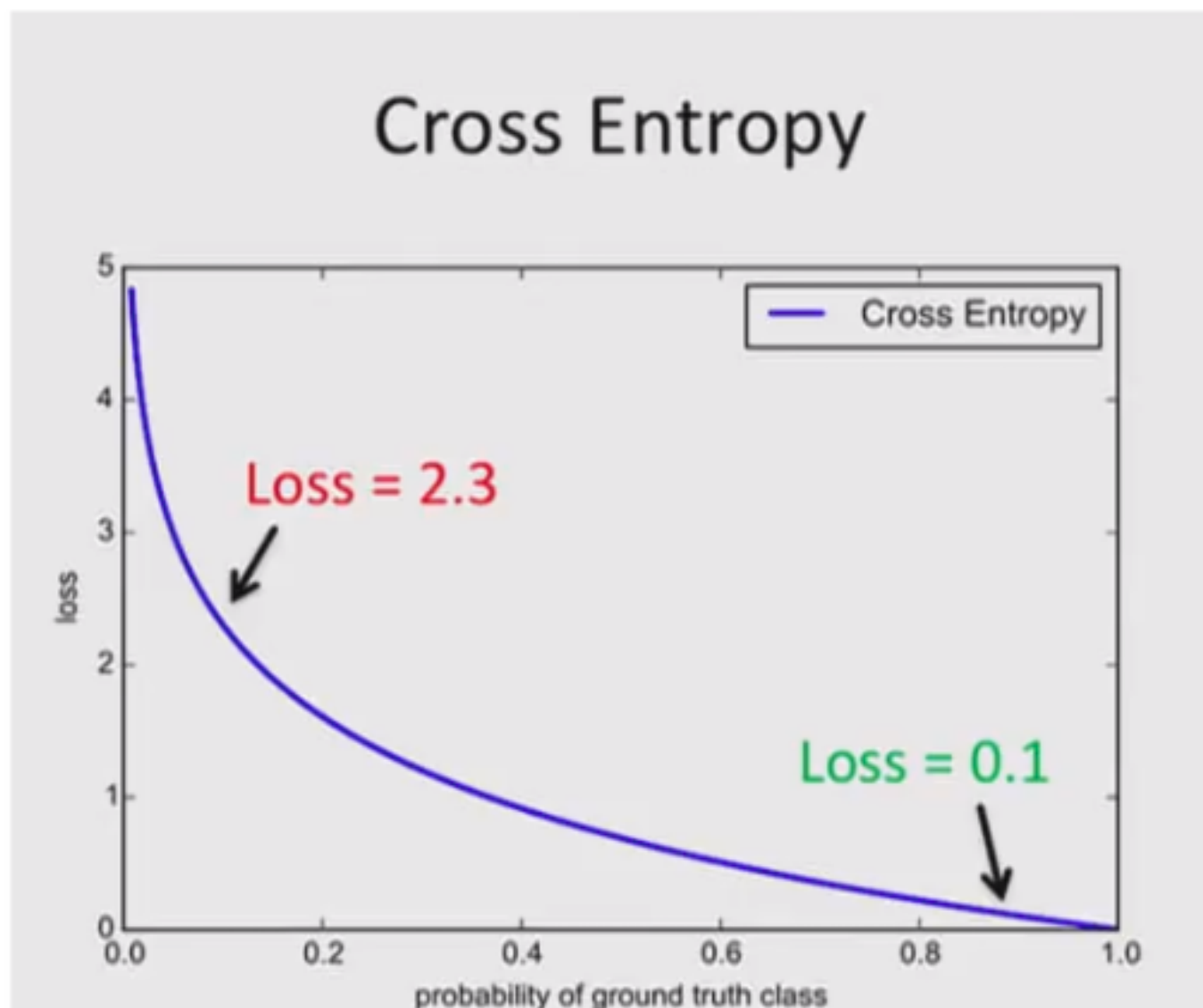
# Focal Loss



<https://arxiv.org/pdf/1708.02002.pdf>

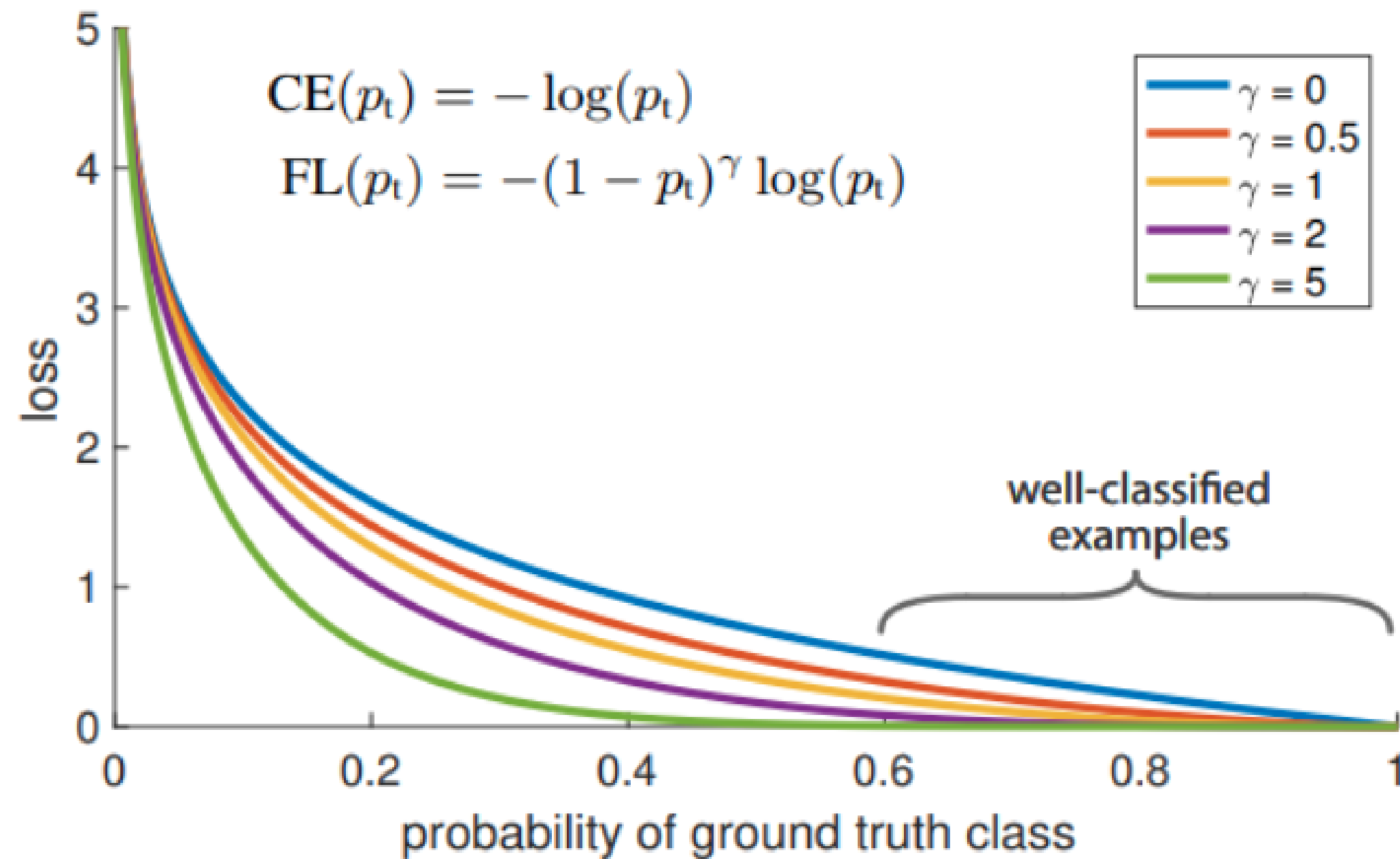
# Focal Loss

100000 easy : 100 hard examples



Loss: centrado en ejemplos fáciles

# Focal Loss



# Focal Loss

- $\alpha$ -balanced Cross entropy

$$\text{CE}(p_t) = -\alpha_t \log(p_t)$$

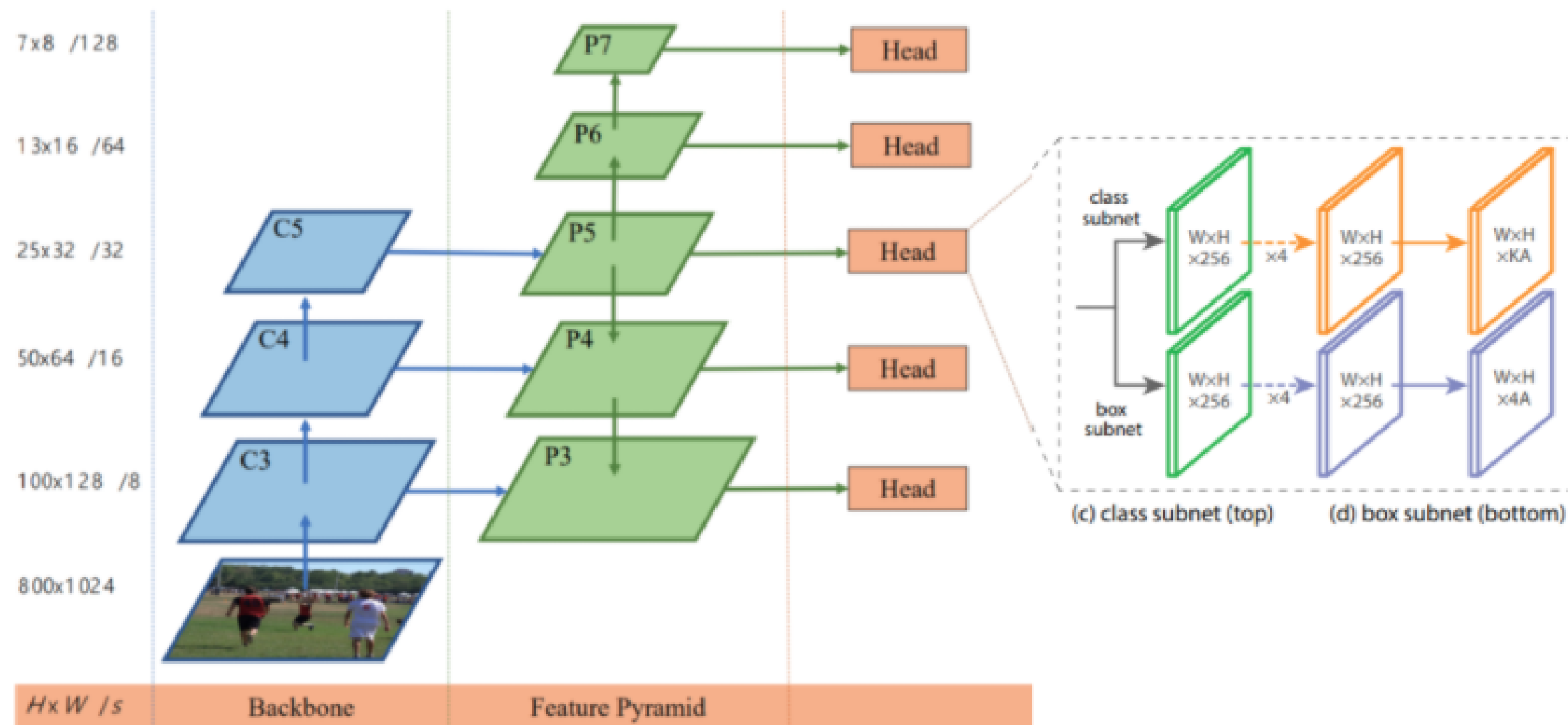
- $\alpha$ -balanced Focal Loss

$$\text{FL}(p_t) = -\alpha_t (1 - p_t)^\gamma \log(p_t)$$

- $\gamma$ : focus more on hard examples
- $\alpha$ : offset class imbalance of number of examples

# RetinaNet

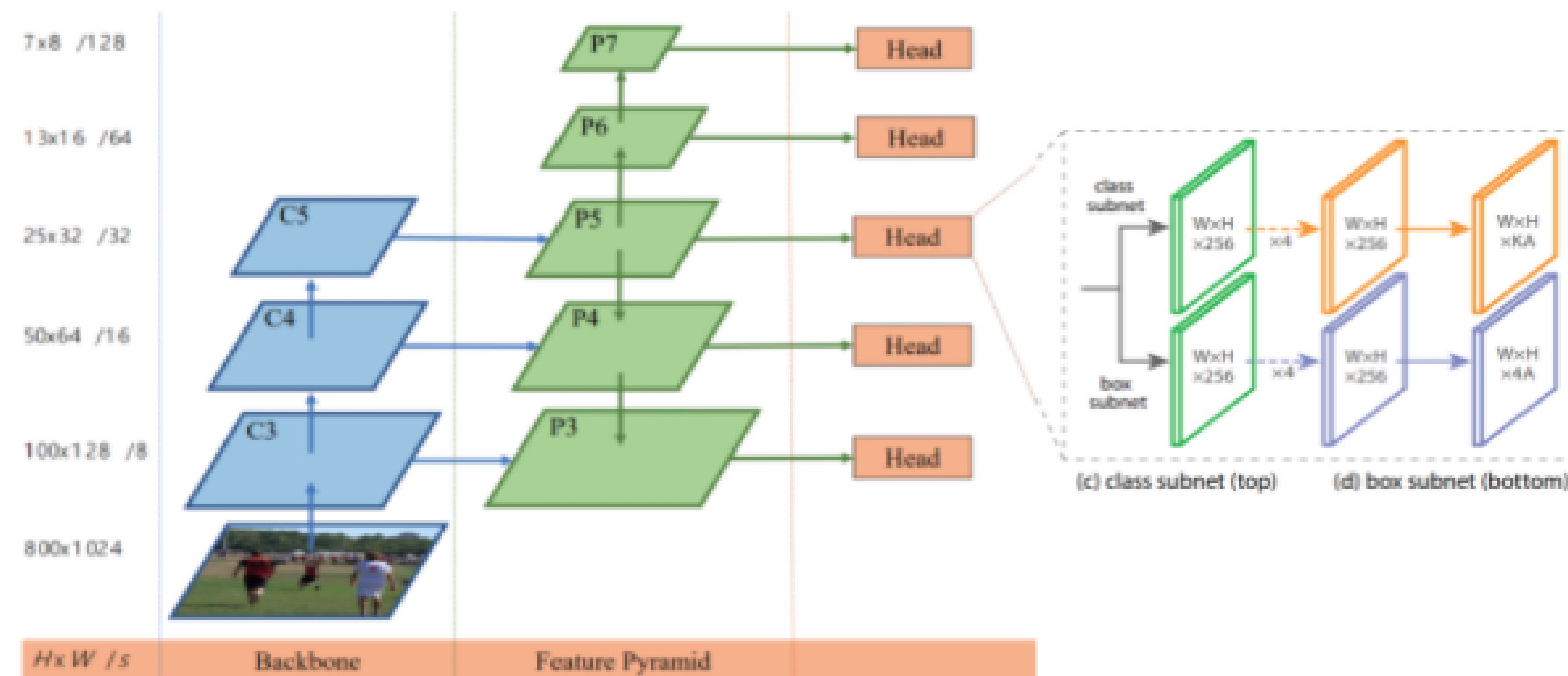
One-stage object detector combining FPN + Focal Loss





# RetinaNet

One-stage object detector combining FPN + Focal Loss



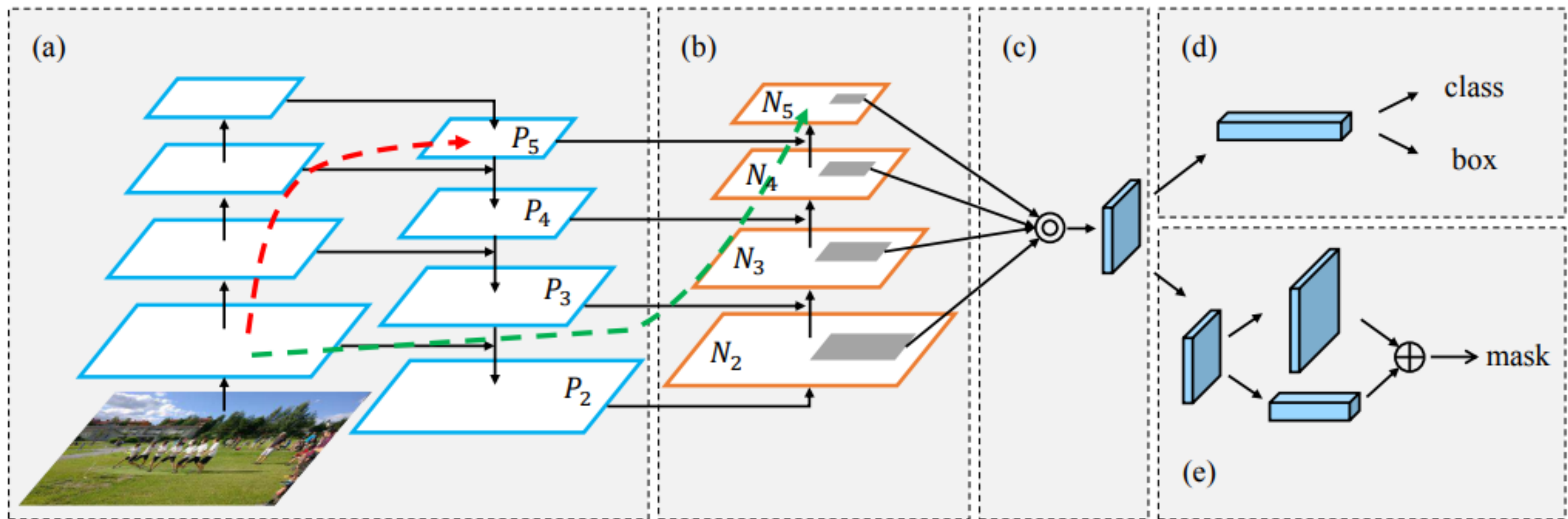
**Anchors:** We use translation-invariant anchor boxes similar to those in the RPN variant in [20]. The anchors have areas of  $32^2$  to  $512^2$  on pyramid levels  $P_3$  to  $P_7$ , respectively. As in [20], at each pyramid level we use anchors at three aspect ratios  $\{1:2, 1:1, 2:1\}$ . For denser scale coverage than in [20], at each level we add anchors of sizes  $\{2^0, 2^{1/3}, 2^{2/3}\}$  of the original set of 3 aspect ratio anchors. This improves AP in our setting. In total there are  $A = 9$  anchors per level and across levels they cover the scale range 32 - 813 pixels with respect to the network's input image.

# RetinaNet

	backbone	AP	AP <sub>50</sub>	AP <sub>75</sub>	AP <sub>S</sub>	AP <sub>M</sub>	AP <sub>L</sub>
<i>Two-stage methods</i>							
Faster R-CNN+++ [16]	ResNet-101-C4	34.9	55.7	37.4	15.6	38.7	50.9
Faster R-CNN w FPN [20]	ResNet-101-FPN	36.2	59.1	39.0	18.2	39.0	48.2
Faster R-CNN by G-RMI [17]	Inception-ResNet-v2 [34]	34.7	55.5	36.7	13.5	38.1	52.0
Faster R-CNN w TDM [32]	Inception-ResNet-v2-TDM	36.8	57.7	39.2	16.2	39.8	<b>52.1</b>
<i>One-stage methods</i>							
YOLOv2 [27]	DarkNet-19 [27]	21.6	44.0	19.2	5.0	22.4	35.5
SSD513 [22, 9]	ResNet-101-SSD	31.2	50.4	33.3	10.2	34.5	49.8
DSSD513 [9]	ResNet-101-DSSD	33.2	53.3	35.2	13.0	35.4	51.1
<b>RetinaNet</b> (ours)	ResNet-101-FPN	39.1	59.1	42.3	21.8	42.7	50.2
<b>RetinaNet</b> (ours)	ResNeXt-101-FPN	<b>40.8</b>	<b>61.1</b>	<b>44.1</b>	<b>24.1</b>	<b>44.2</b>	51.2

Table 2. **Object detection** *single-model* results (bounding box AP), vs. state-of-the-art on COCO test-dev. We show results for our RetinaNet-101-800 model, trained with scale jitter and for  $1.5\times$  longer than the same model from Table 1e. Our model achieves top results, outperforming both one-stage and two-stage models. For a detailed breakdown of speed versus accuracy see Table 1e and Figure 2.

# Path Aggregation Network (PAN)



Path Aggregation Network for Instance Segmentation

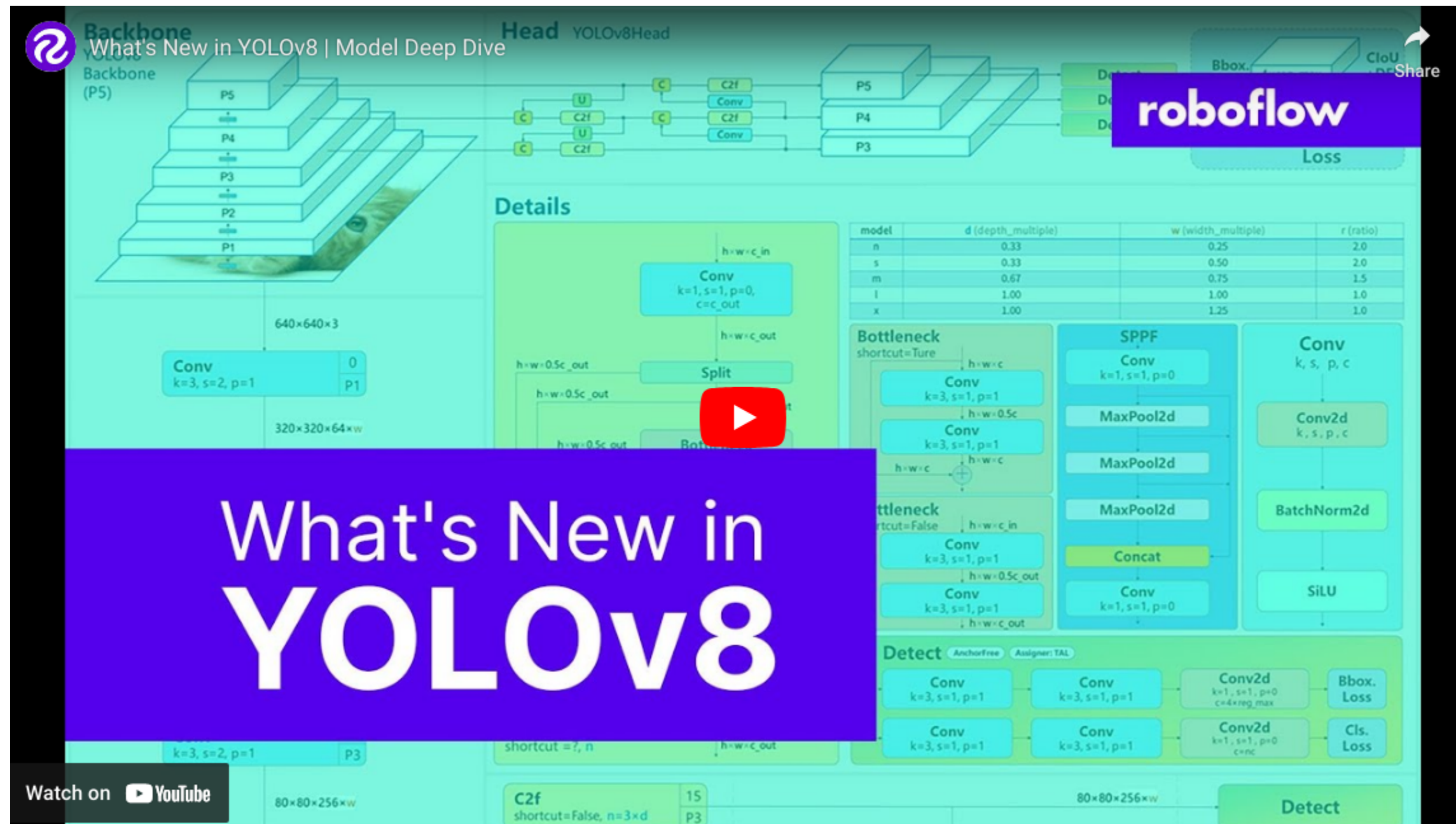
# Path Aggregation Network

	$AP^{bb}$	$AP_{50}^{bb}$	$AP_{75}^{bb}$	$AP_S^{bb}$	$AP_M^{bb}$	$AP_L^{bb}$
Champion 2015 [23]	37.4	59.0	40.2	18.3	41.7	52.9
Champion 2016 [27]	41.6	62.3	45.6	24.0	43.9	55.2
Our Team 2017	<b>51.0</b>	<b>70.5</b>	<b>55.8</b>	<b>32.6</b>	<b>53.9</b>	<b>64.8</b>

Table 7. Box AP of COCO Object Detection Challenge in different years on *test-dev*.

Path Aggregation Network for Instance Segmentation

# YOLO8



<https://blog.roboflow.com/whats-new-in-yolov8/>



# YOLO: A Brief History

**YOLO** (You Only Look Once), a popular object detection and image segmentation model, was developed by Joseph Redmon and Ali Farhadi at the University of Washington. Launched in 2015, YOLO quickly gained popularity for its high speed and accuracy.

- **YOLOv2**, released in 2016, improved the original model by incorporating batch normalization, anchor boxes, and dimension clusters.
- **YOLOv3**, launched in 2018, further enhanced the model's performance using a more efficient backbone network, multiple anchors and spatial pyramid pooling.
- **YOLOv4** was released in 2020, introducing innovations like Mosaic data augmentation, a new anchor-free detection head, and a new loss function.
- **YOLOv5** further improved the model's performance and added new features such as hyperparameter optimization, integrated experiment tracking and automatic export to popular export formats.
- **YOLOv6** was open-sourced by **Meituan** in 2022 and is in use in many of the company's autonomous delivery robots.
- **YOLOv7** added additional tasks such as pose estimation on the COCO keypoints dataset.
- **YOLOv8** is the latest version of YOLO by Ultralytics. As a cutting-edge, state-of-the-art (SOTA) model, YOLOv8 builds on the success of previous versions, introducing new features and improvements for enhanced performance, flexibility, and efficiency. YOLOv8 supports a full range of vision AI tasks, including **detection**, **segmentation**, **pose estimation**, **tracking**, and **classification**. This versatility allows users to leverage YOLOv8's capabilities across diverse applications and domains.



# Muchas Gracias

José Manuel Saavedra Rondo  
[jmsaavedrar@miuandes.cl](mailto:jmsaavedrar@miuandes.cl)