

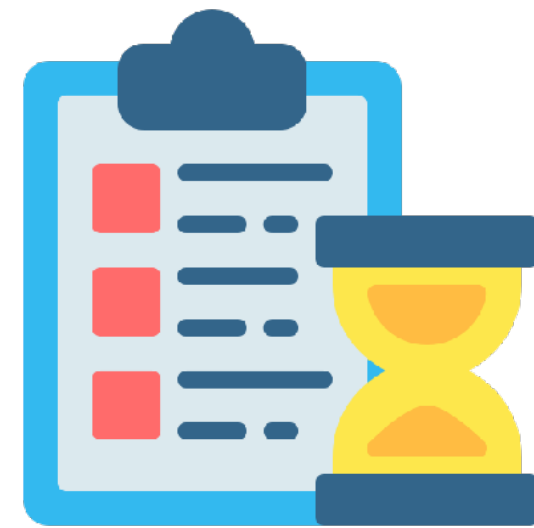
# Detección de objetos I

## Visión por computador

Juan Carlos Arbeláez  
jarbel16@eafit.edu.co



# Contenido

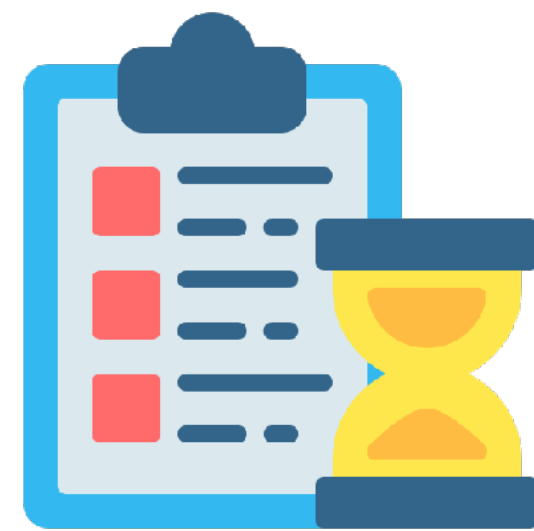


1. Introducción
2. Algoritmos
3. Métricas de evaluación
4. Aplicaciones
5. Retos
6. Taller

**Parte 1**

**Parte 2**

# Contenido

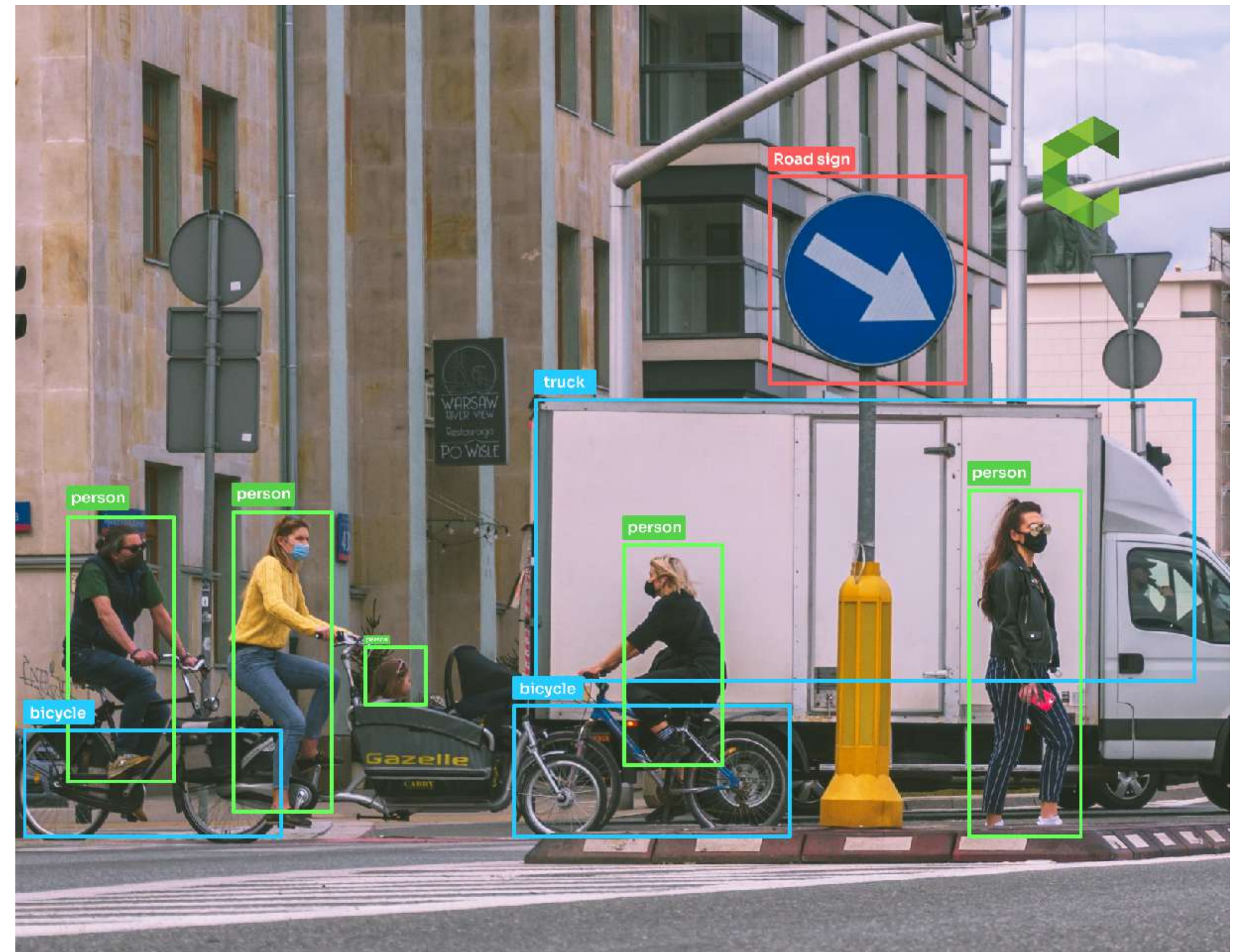


- 1. Introducción**
2. Algoritmos
3. Métricas de evaluación
4. Aplicaciones
5. Retos
6. Taller



# ¿Qué es Detección de Objetos?

- La detección de objetos es una tarea de visión artificial
- Se utiliza para detectar instancias de objetos visuales de determinadas clases
- El objetivo de la detección de objetos es desarrollar modelos que respondan a la pregunta: **¿Qué objetos y dónde?**





# ¿Qué es Detección de Objetos?

**La clasificación de imágenes** predecir la clase de un objeto en una imagen

**La localización de objetos** determinar el área (normalmente rectangular) que ocupa un objeto(s) en la imagen

**La detección de objetos** combina estas dos tareas y localiza y clasifica uno o más objetos en una imagen

*\*Dos términos que se usan para el reconocimiento de objetos son: “Object detection” y “Object recognition”*

**Classification**



CAT

**Classification  
+ Localization**

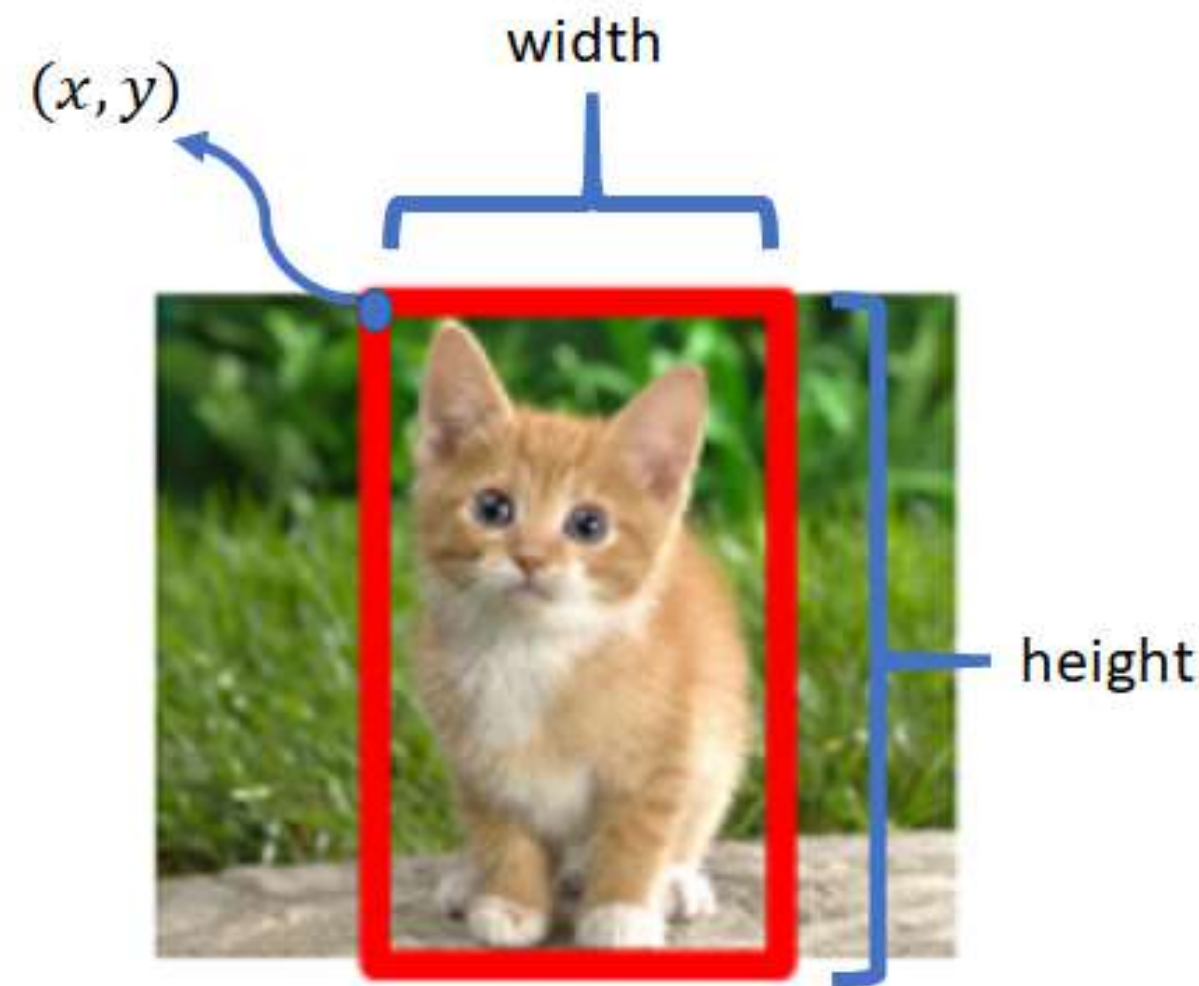


CAT

# ¿Qué es Detección de Objetos?



¿Cuántas variables por objeto debe predecir un modelo de detección?



Localization (4)

+



Clases (N)



# ¿Cómo son entrenados los modelos?

- El proceso de entrenamiento de un modelo para la *detección de objetos* es similar al *clasificación*
- Un conjunto de datos de detección de objetos vinculan una imagen con una lista de objetos que contiene y su ubicación
- El modelo acepta una imagen como entrada y devuelve una lista de predicciones por cada objeto con: ubicación (Coordenadas) y la clase

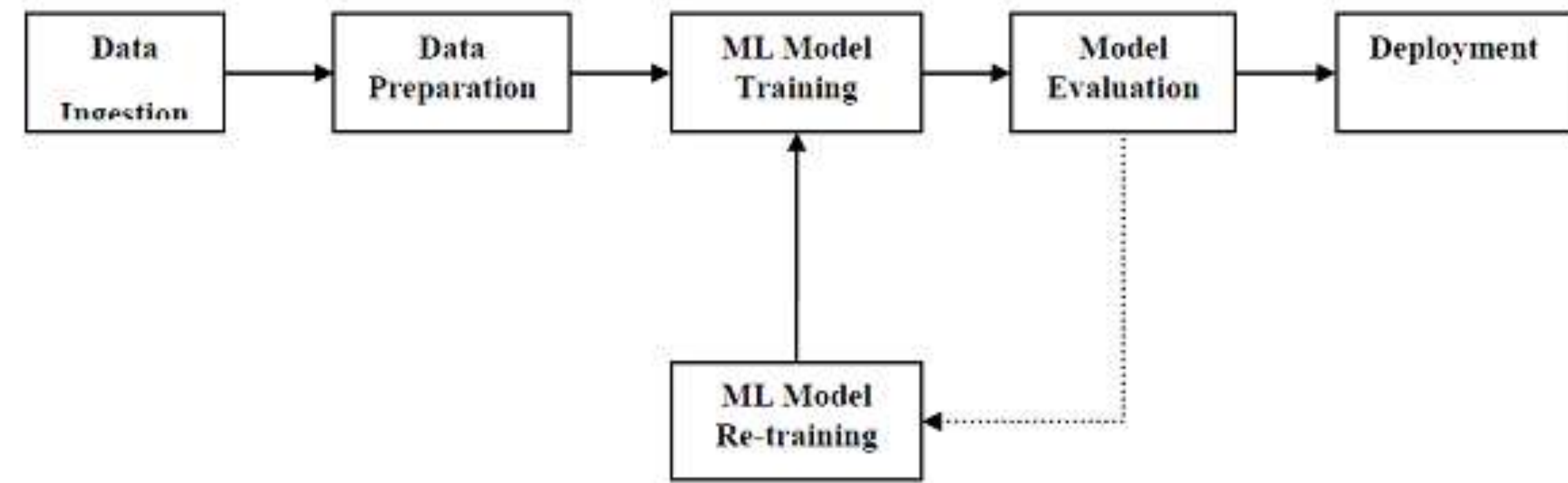
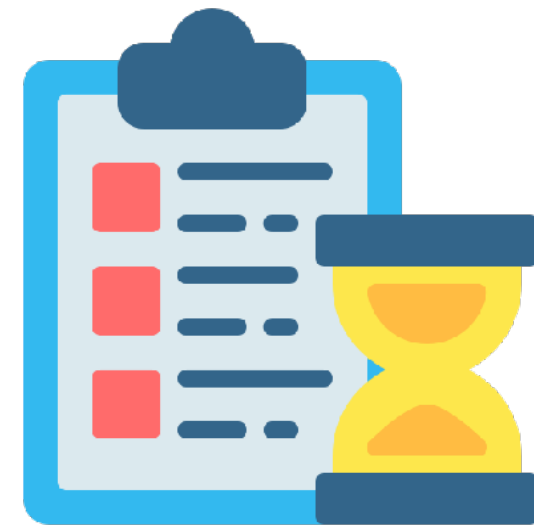


Diagrama de flujo de entrenamiento de modelos de ML

# Contenido



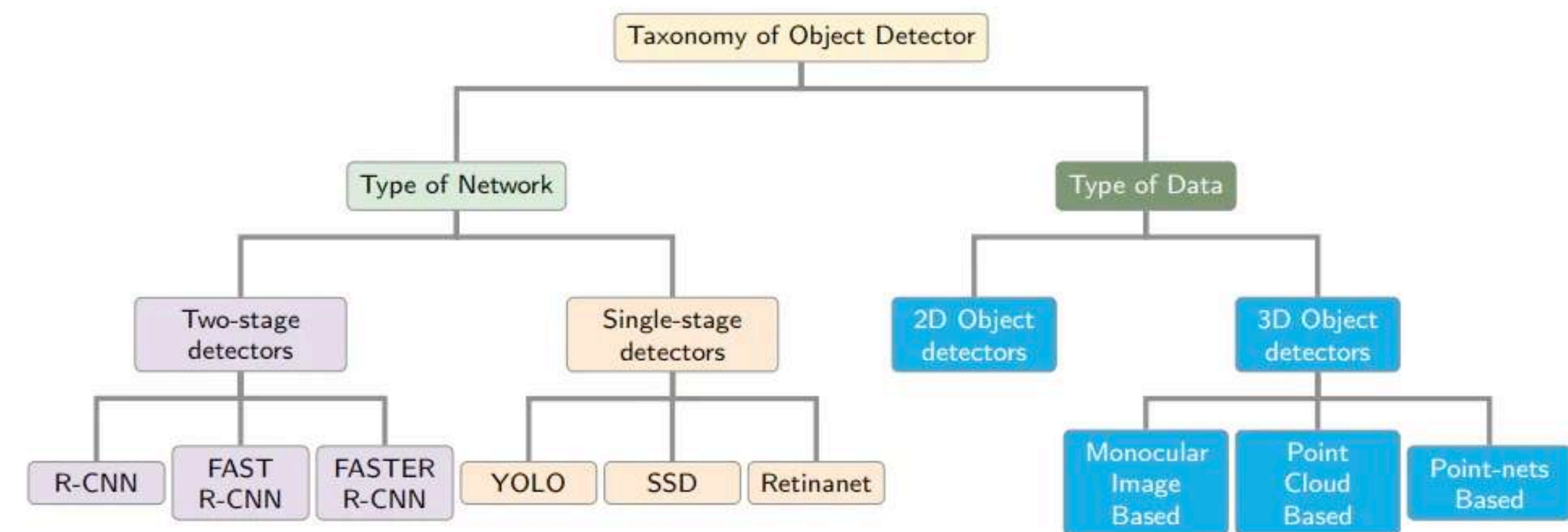
1. Introducción
- 2. Algoritmos**
3. Métricas de evaluación
4. Aplicaciones
5. Retos
6. Taller



# Algoritmos

Hay muchos algoritmos para la detección de objetos, cada uno tiene sus pros y contras. Principales que se utilizan en la industria:

1. Faster R-CNN
2. SSD
3. YOLO
4. DETR



# Sliding Window

- Solución convencional: ventana deslizante para buscar en cada posición de la imagen
- Diferentes objetos o incluso el mismo tipo pueden tener diferentes relaciones de aspecto y tamaños dependiendo del tamaño del objeto y la distancia desde la cámara.
- Extremadamente lento si utilizamos CNN para la clasificación de imágenes en cada ubicación

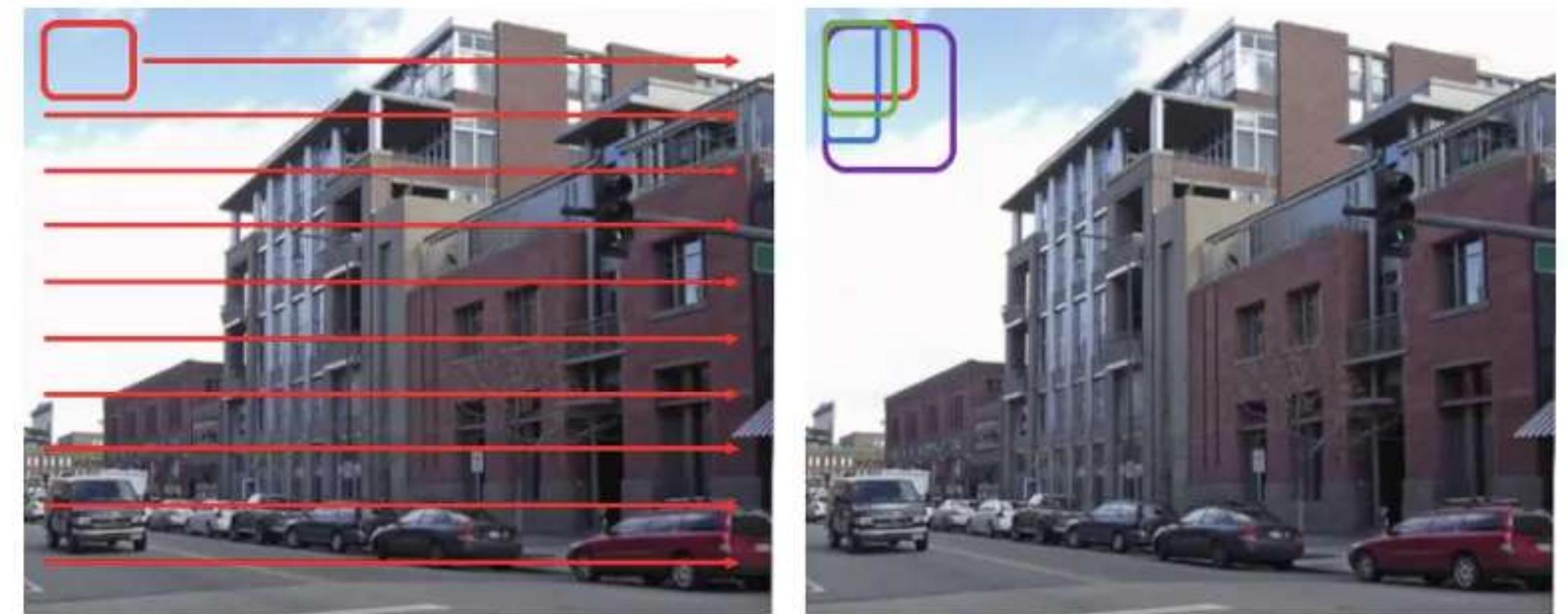


Illustration of Sliding Window (Left) with Different Aspect Ratios and Sizes (Right)

# Contenido

1. Introducción

2. Algoritmos

**1. R-CNN**

2. Faster R-CNN

3. Single Shot Detector (SSD)

4. You Only Look Once (Yolo)

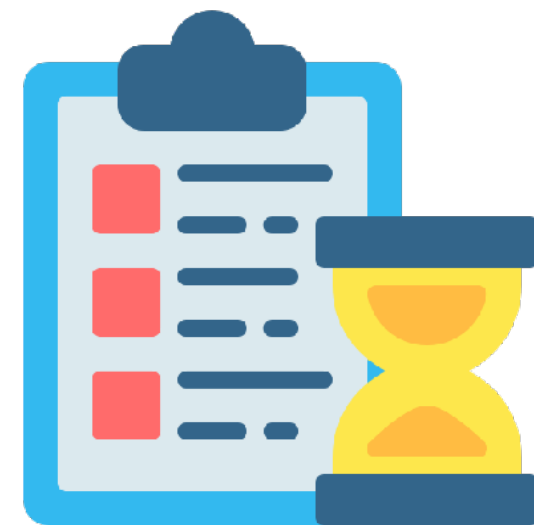
5. DETR

3. Métricas de evaluación

4. Aplicaciones

5. Retos

6. Taller





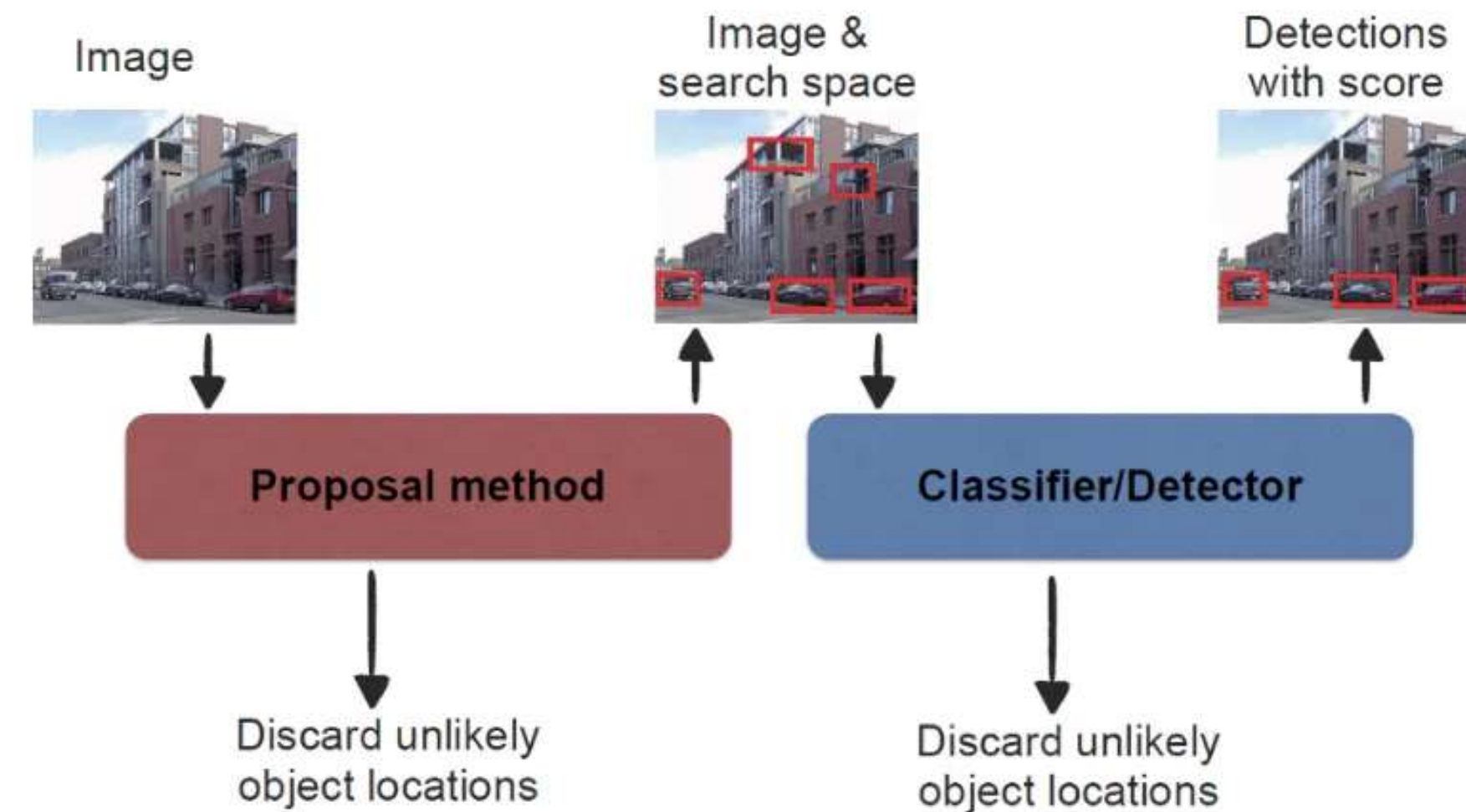
# R-CNN

## Region-based Convolutional Neural Network

Es un detector de dos etapas

1. Localizar regiones de interés
2. Clasificación de regiones

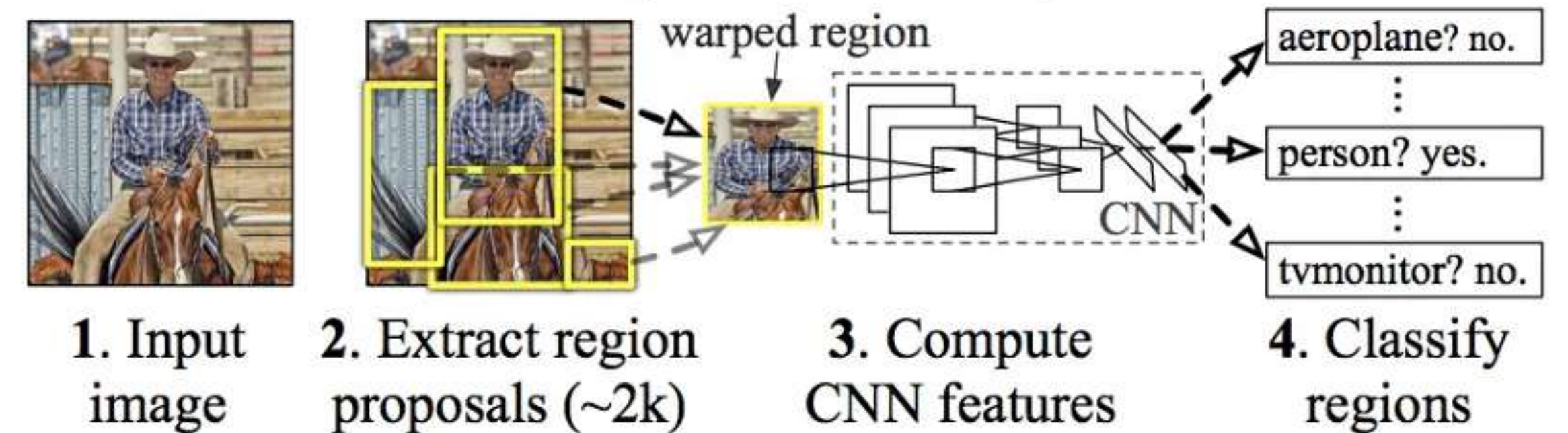
Dos implementaciones comúnmente usadas: Fast RCNN y Faster RCNN



R-CNN Flowchart

# R-CNN

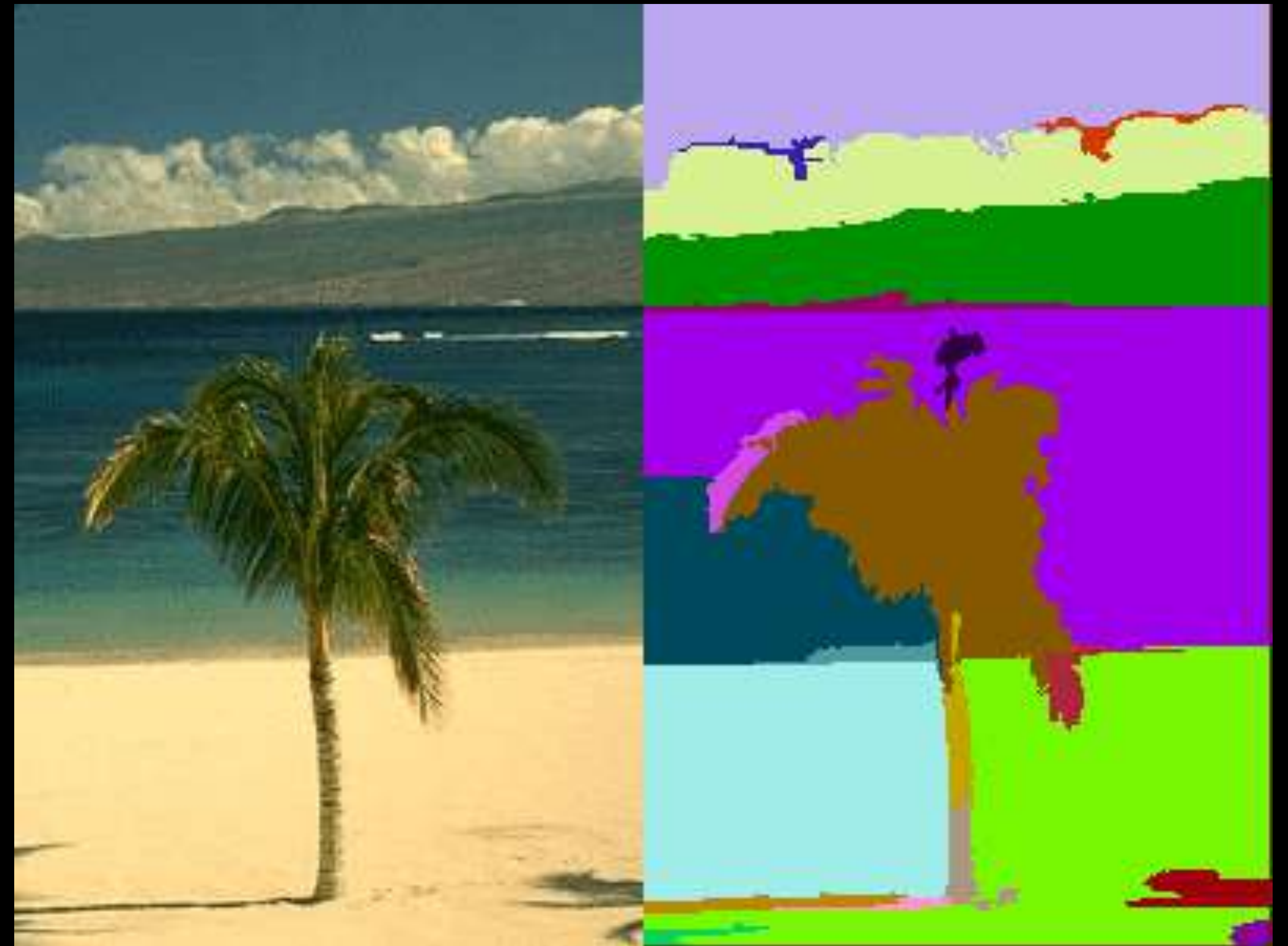
- A. Proponer regiones de interés:  
Extraemos las regiones con un algoritmo (búsqueda selectiva)
- B. Clasificación: Redimensionar todos los regiones y pasarlas por una CNN





# Selective Search

Es un algoritmo que se usa para la propuesta de regiones que agrupa regiones de forma jerárquica basada su similitud





# Selective Search

1. Sobre-segmentar la imagen usando el algoritmo *Felzenszwalb & Huttenlocher*
2. Añade los recuadros de las partes segmentadas a la lista de propuestas de regiones
3. Agrupa los segmentos adyacentes en función de la similitud

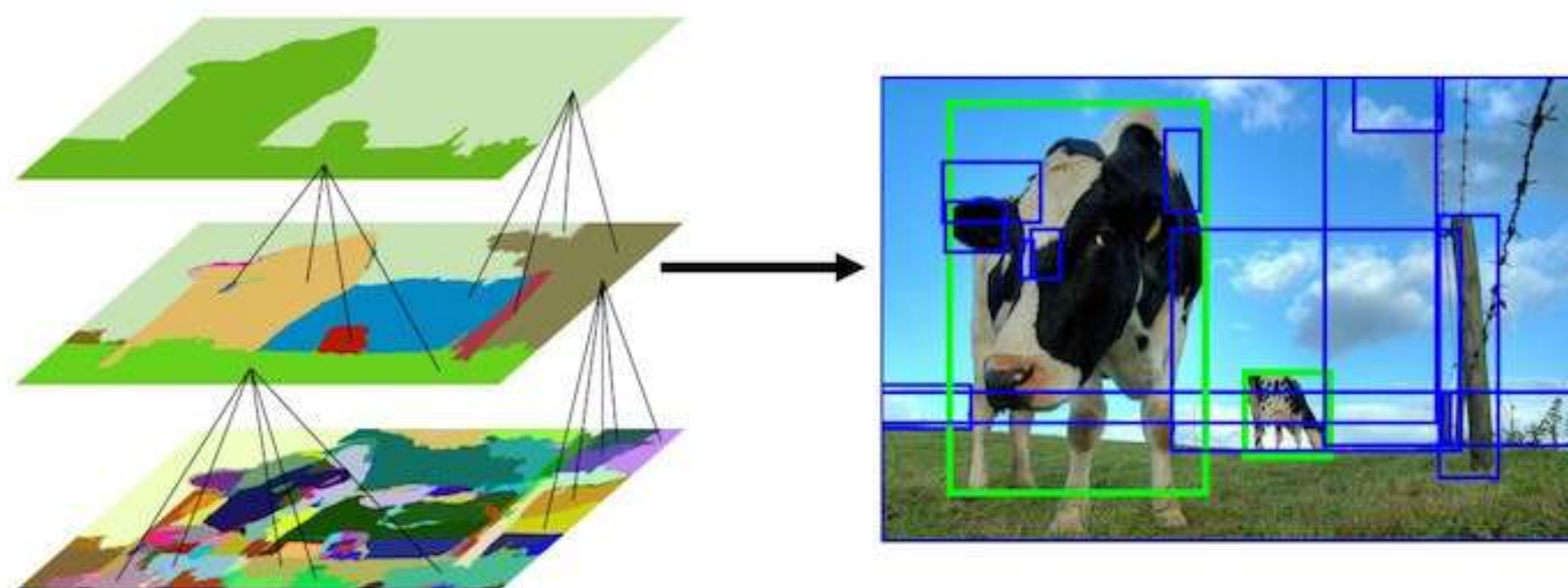


Imagen original



Over-segmented

# Selective Search



- En cada iteración se forman segmentos más grandes y se añaden a la lista de propuestas de regiones
- Se crean propuestas de regiones de segmentos más pequeños a segmentos más grandes

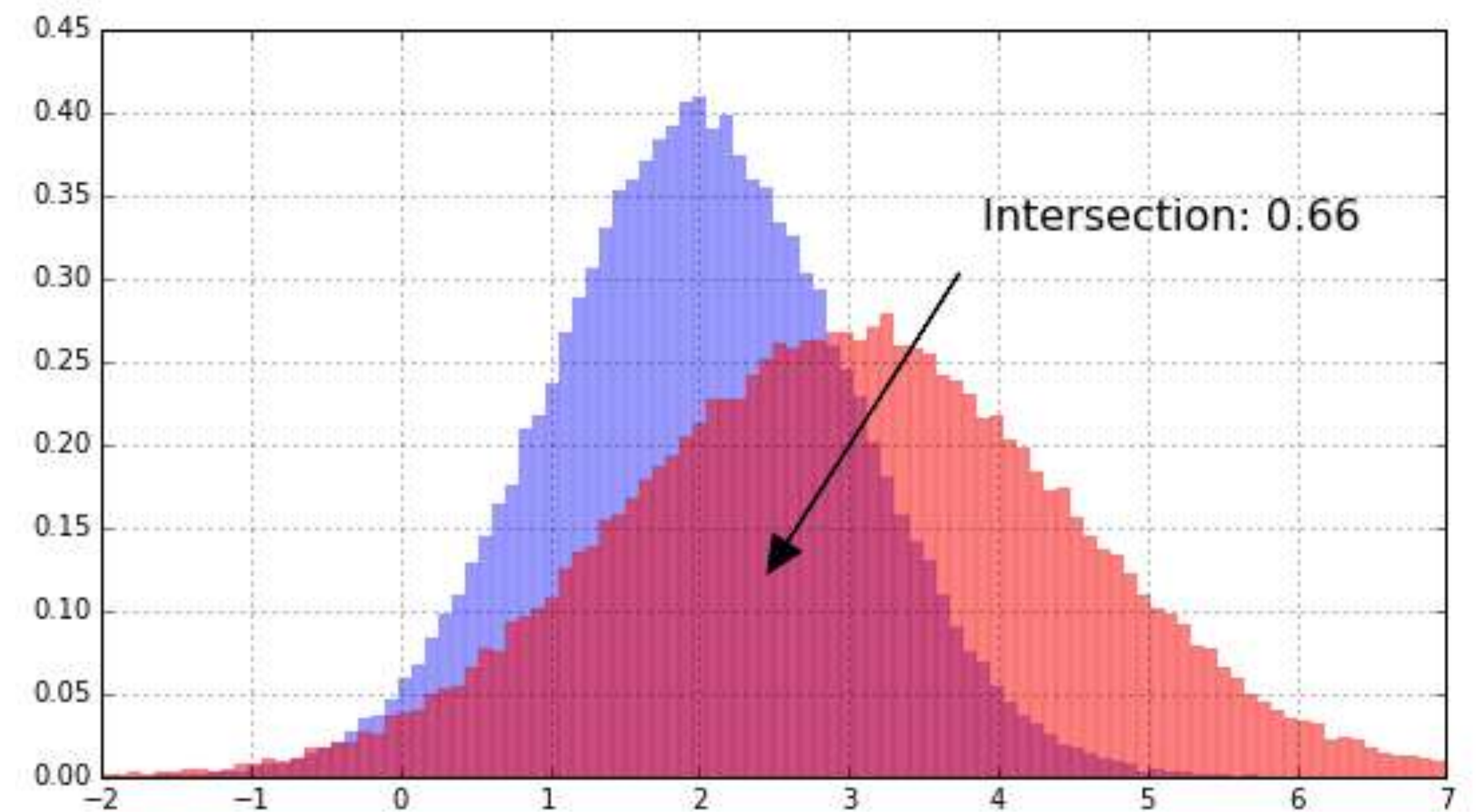


# Similitud entre recuadros

1. Intersección de histogramas de:
  - Color
  - Textura (derivados de color)
2. Tamaño: anima a las regiones más pequeñas a fusionarse antes
3. Compatibilidad de forma: lo bien que encajan dos regiones entre sí

**Similitud final:**

$$s(r_i, r_j) = \alpha_1 S_{color} + \alpha_2 S_{texture} + \alpha_3 S_{size} + \alpha_4 S_{shape}$$

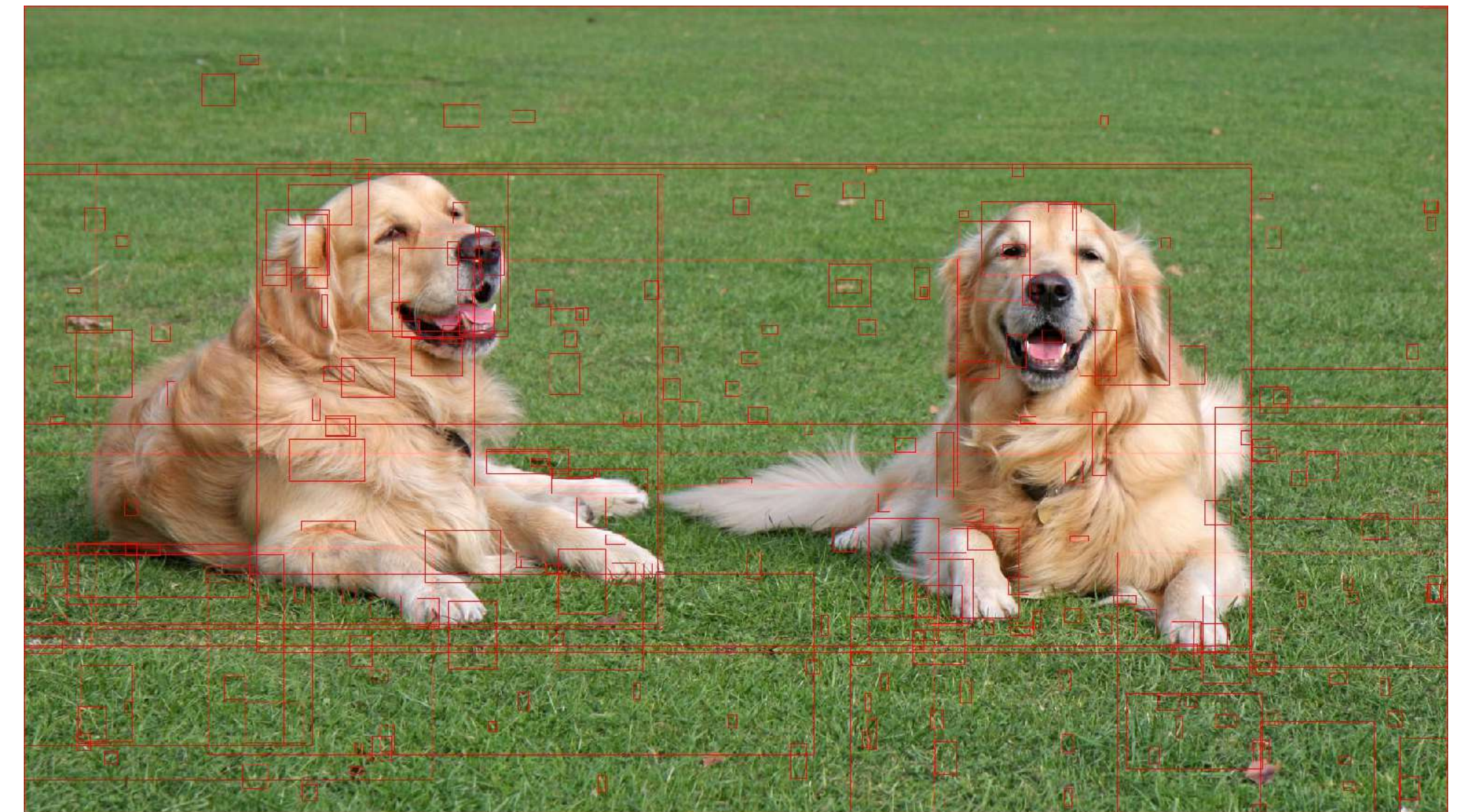


Intersección de histogramas



# Selective Search

- Es necesario etiquetar cada una de las regiones propuestas para entrenar y medir el modelo de clasificación
- Asignar a qué clase pertenece dependiendo de su intersección con la etiqueta manual (IOU)



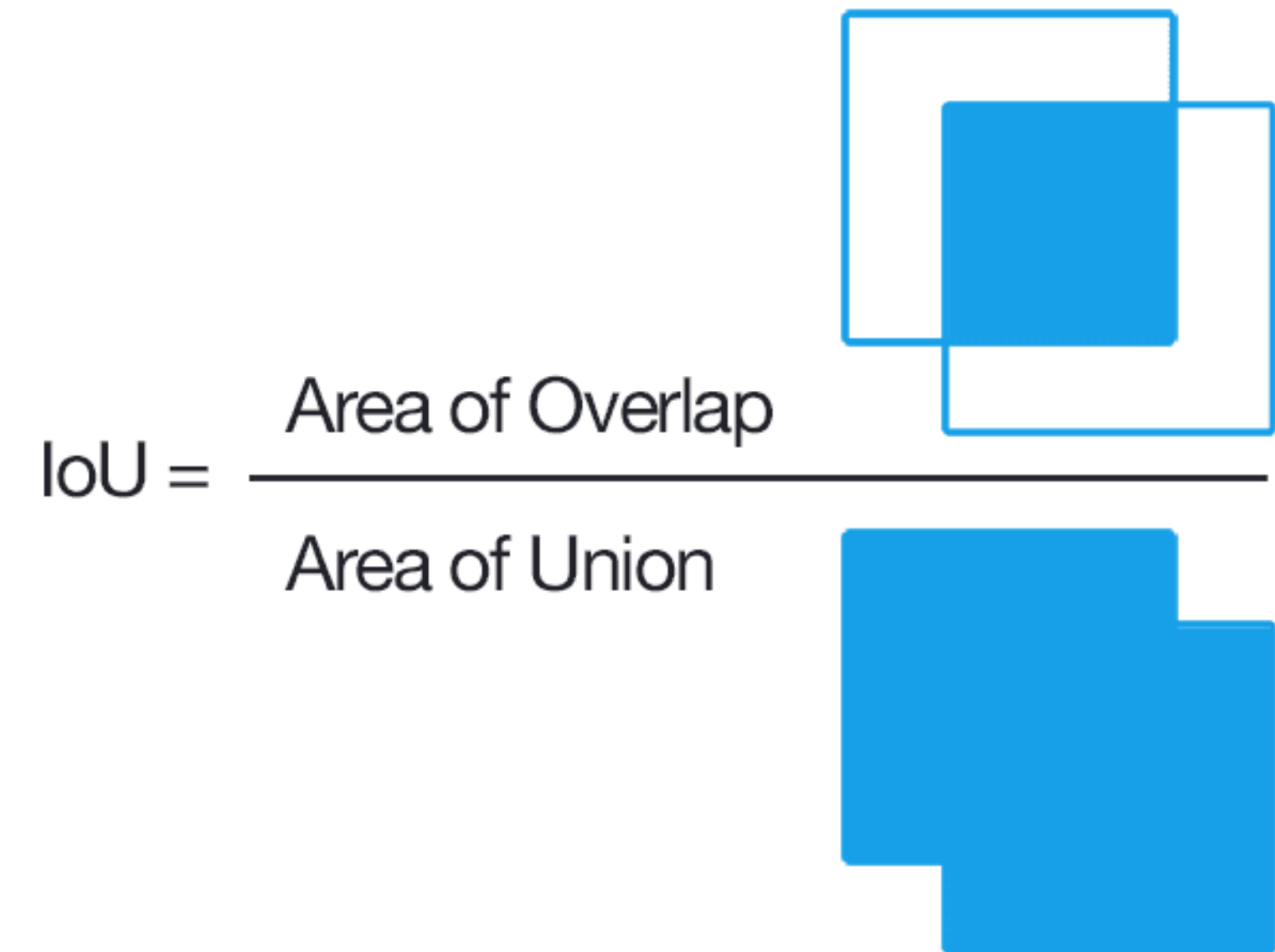
Dogs: top 250 region proposals

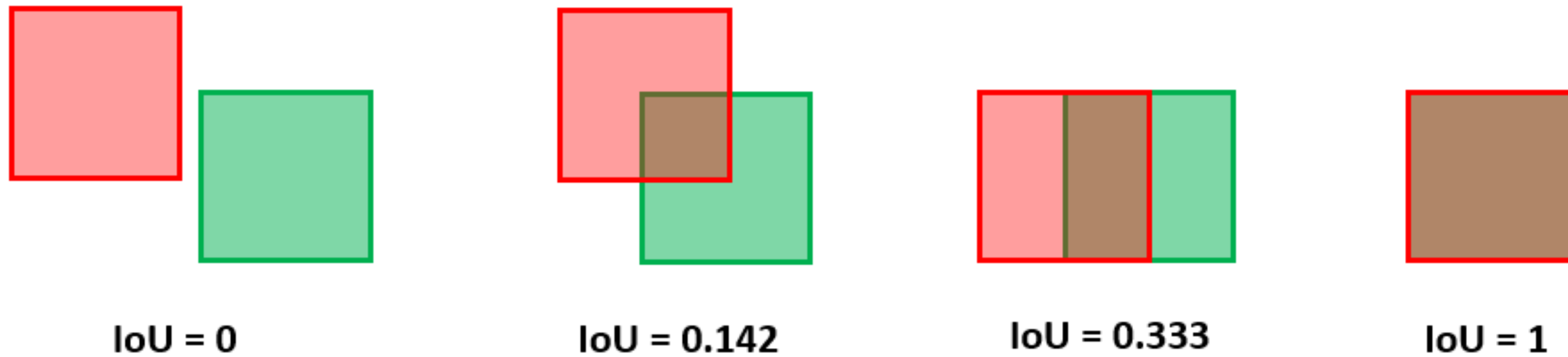


# IOU

## Intersection over Union

- Describe el grado de superposición de rectángulos
- Cuanto mayor sea la región de superposición: mayor será el IOU
- Se utiliza principalmente en aplicaciones relacionadas con la detección de objetos



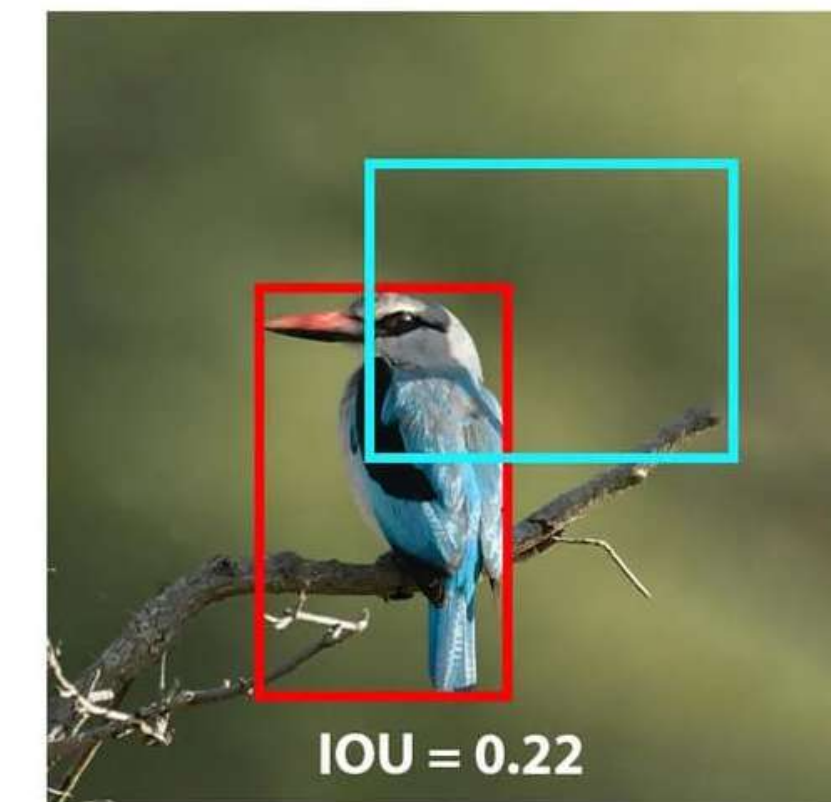
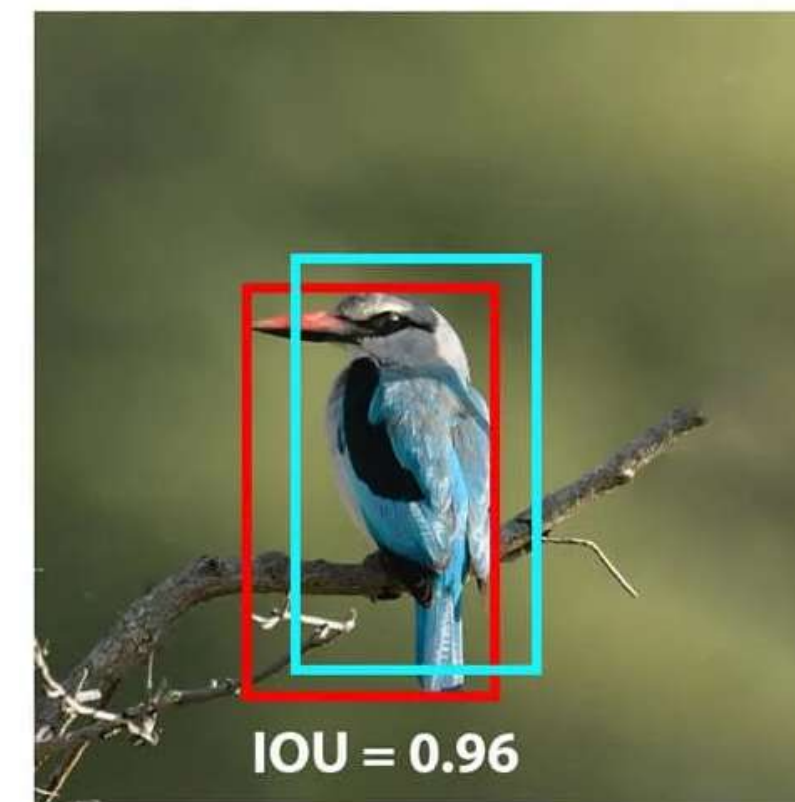


# IOU (Intersection over Union)



# Etiquetado de regiones

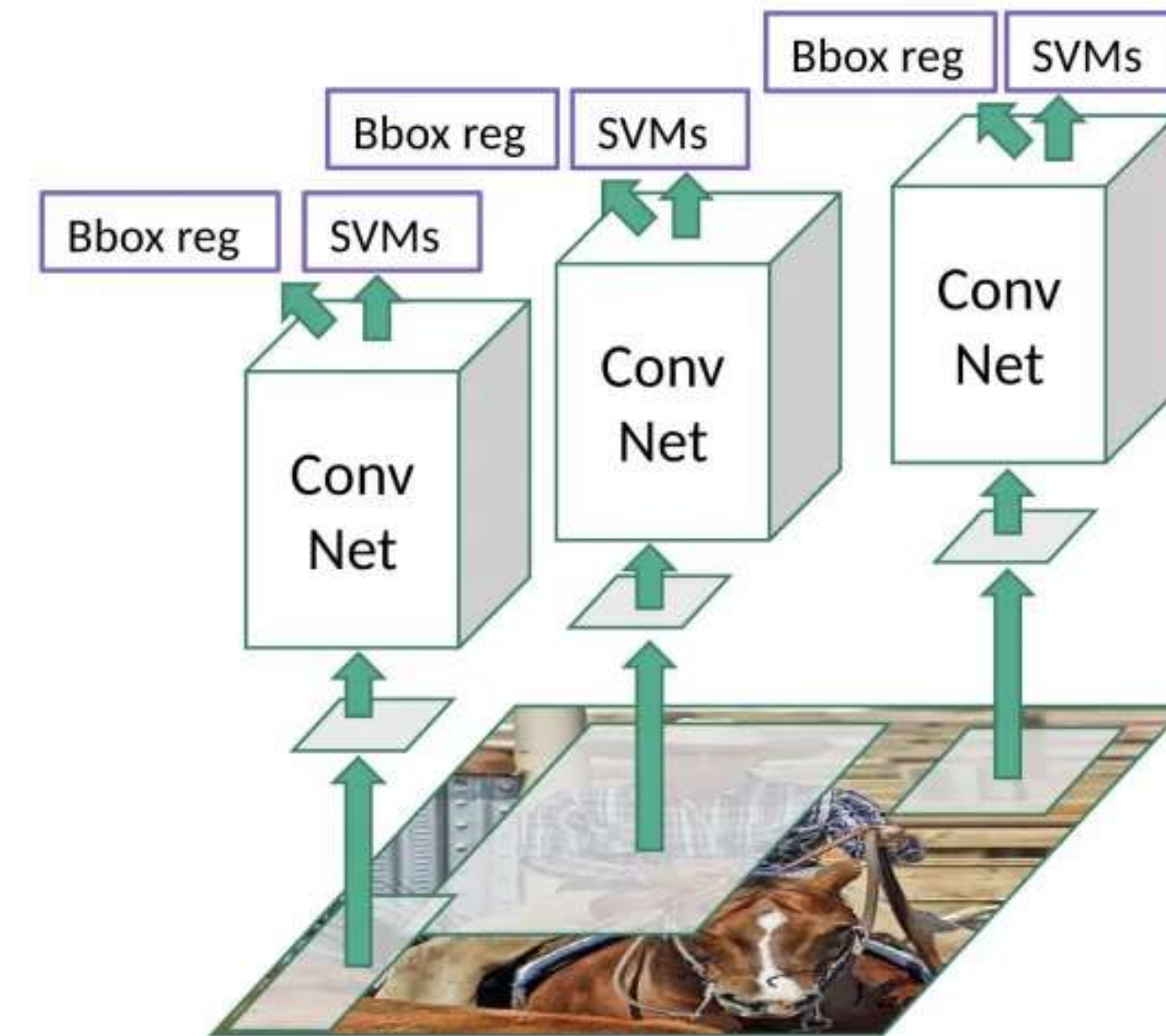
- Es necesario definir el valor del threshold sobre el cual considerar un recuadro como perteneciente a una clase
- En el artículo original usan “grid search” para definir este valor (*IOU* inferior a 0.3 se etiquetan como fondo)



# Arquitectura

## Clasificación de regiones

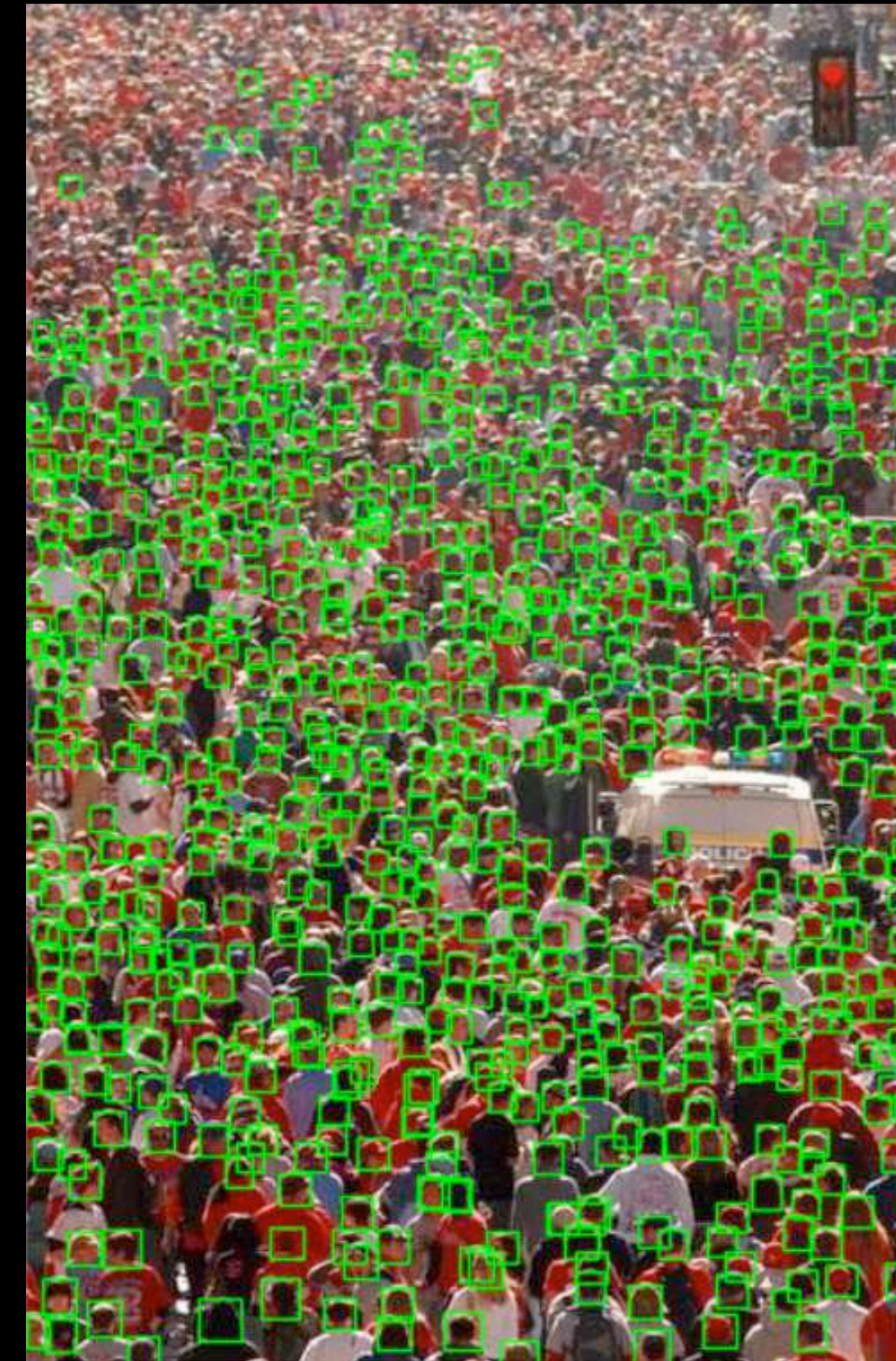
- Las regiones candidatas se deforman en un cuadrado y se introducen en una CNN que produce un vector de características como salida
- Las características extraídas se introducen en una *SVM* para clasificar la presencia del objeto
- También se predicen 4 valores de desplazamiento de las coordenadas del rectángulo para aumentar la precisión





# Problemas con R-CNN

- Alto tiempo de entrenamiento: clasificar cada región (~2000)
- No puede aplicarse en tiempo real ( ~47 segundos por cada imagen)
- No es posible mejorar el algoritmo de proposición de regiones, desacoplado del clasificador y sin ML





# Contenido

1. Introducción

2. Algoritmos

1. R-CNN

**2. Faster R-CNN**

3. Single Shot Detector (SSD)

4. You Only Look Once (Yolo)

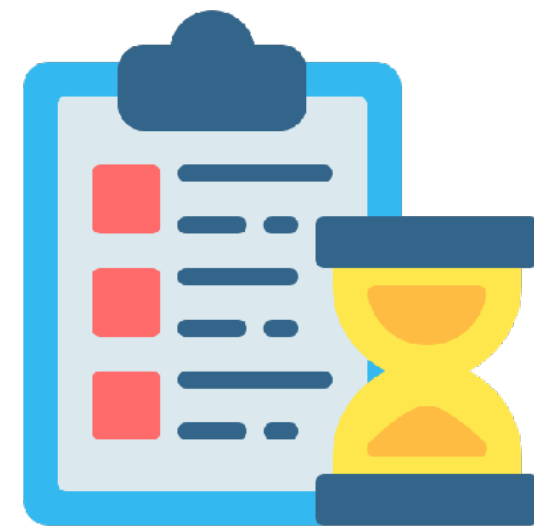
5. DETR

3. Métricas de evaluación

4. Aplicaciones

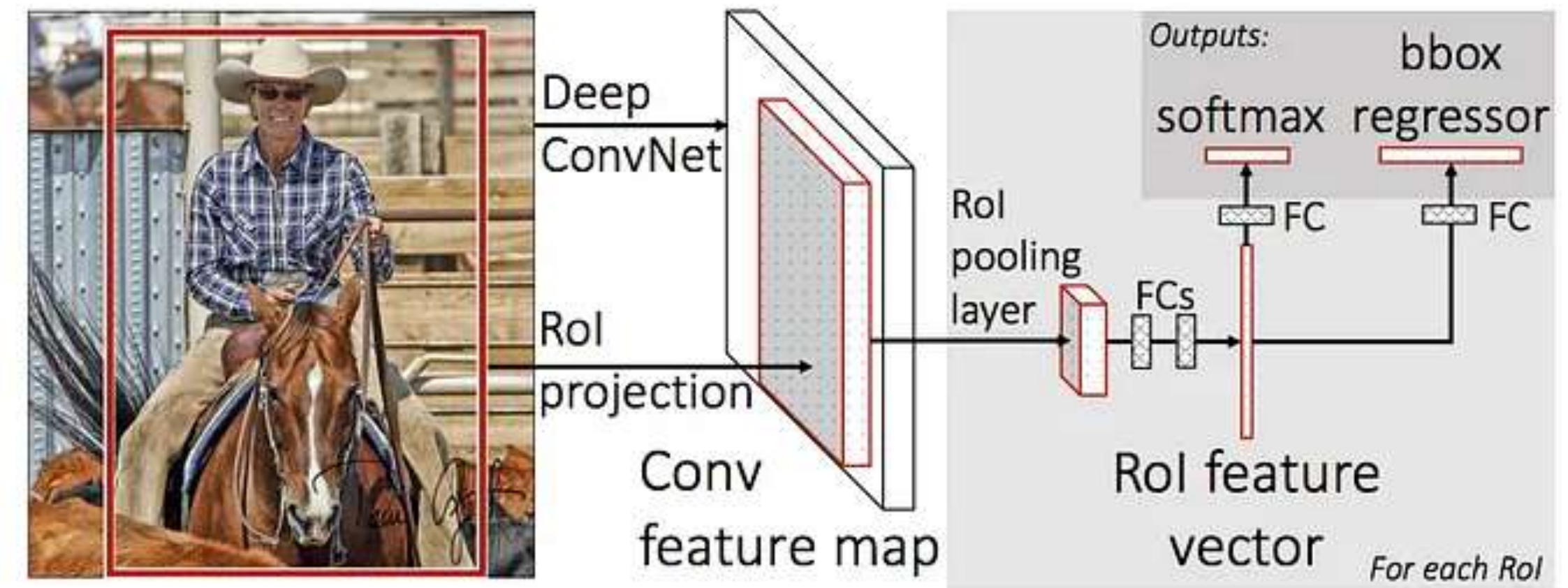
5. Retos

6. Taller



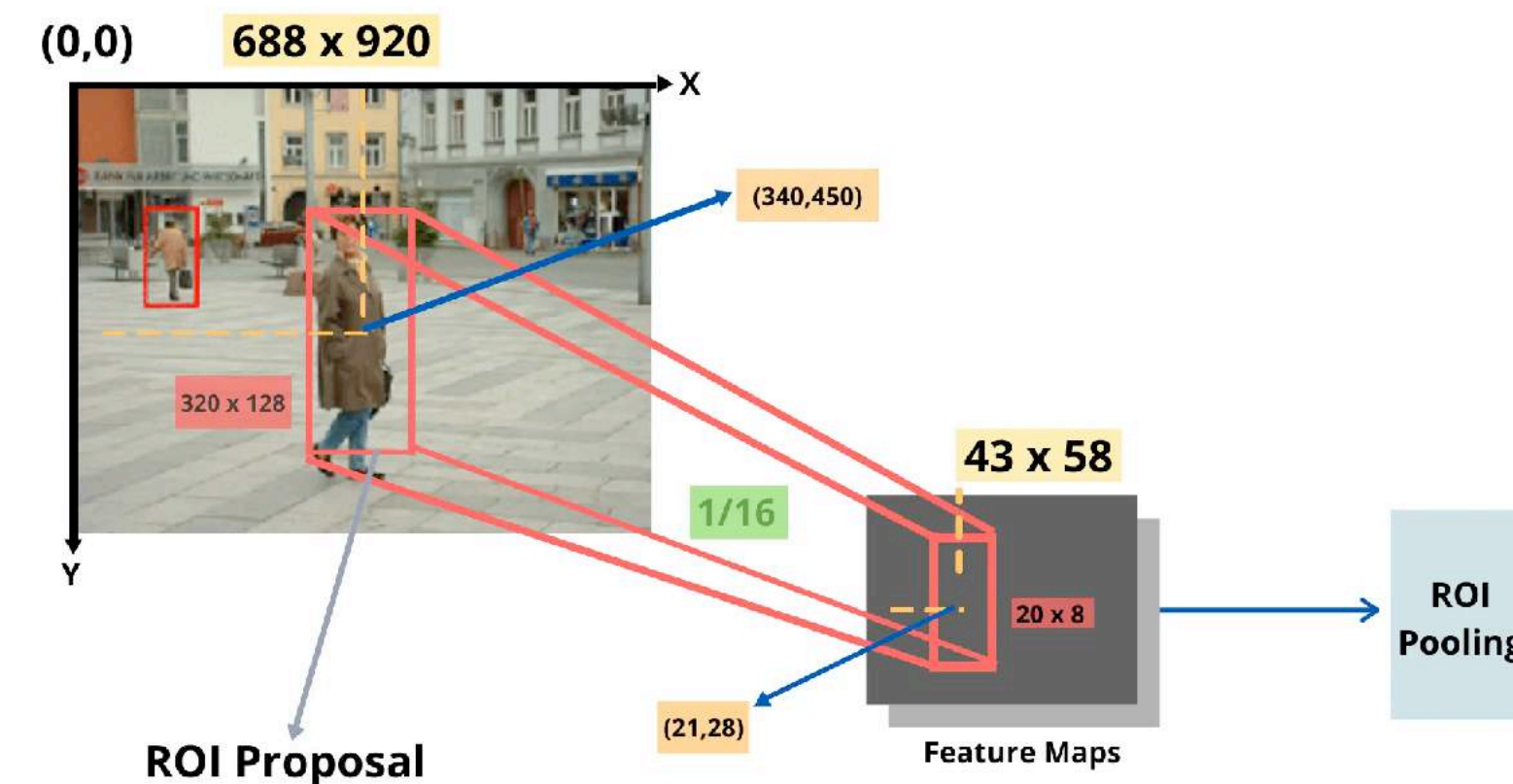
# Fast-RCNN

- El mismo autor resolvió algunos de los problemas para construir un algoritmo de detección de objetos más rápido y lo llamó *Fast R-CNN*
- Similar a R-CNN pero se usa un mapa de características de toda la imagen generado por una CNN



# Proyección de la región de interés

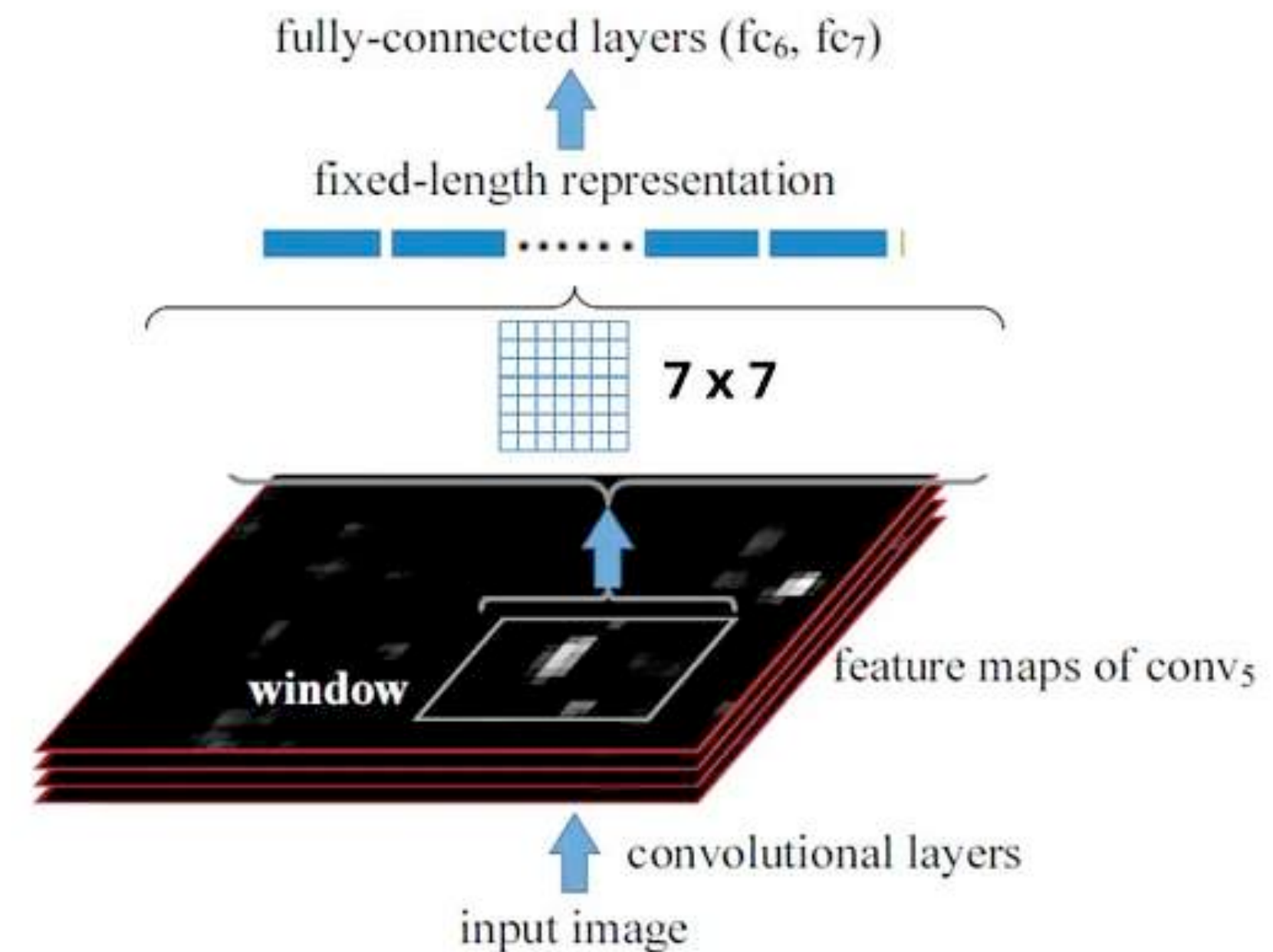
- Proyectar las regiones propuestas sobre el mapa de características
- Una imagen de  $688 \times 920$  se alimenta a una CNN con “subsampling” de  $1/16$ , el mapa de características resultante tiene un tamaño de  $43 \times 58$
- Coordenadas son escaladas usando la transformación geométrica de escalado





# Prediction Head

- Las capas FC de la cabeza de predicción esperan vectores de tamaño fijo pero las proyecciones ROI son de tamaños variables
- Para resolver este problema, los autores proponen “ROI Pooling”
- ROI Pooling: transforma en dimensiones fijas de las proyección de ROI del mapa de características



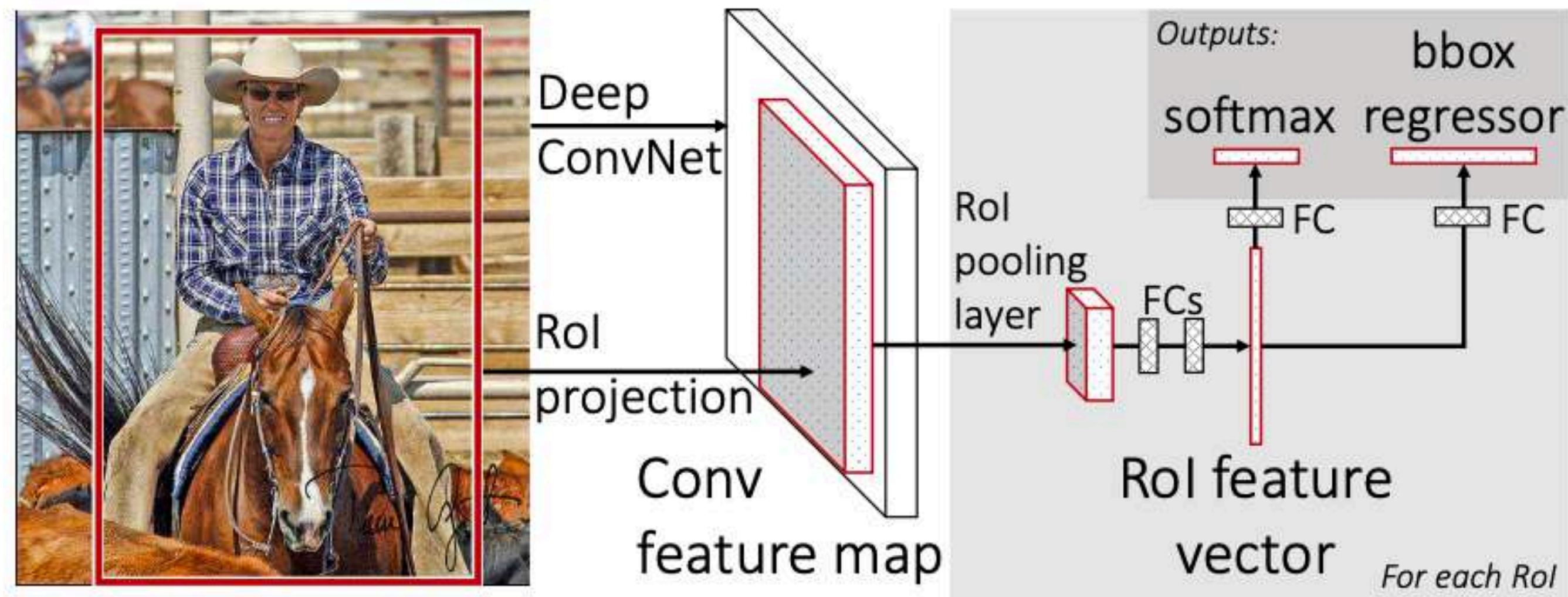
# ROI Pooling

## Region of Interest Pooling

- Produce mapas de características de tamaño fijo a partir de entradas no uniformes
- El mapa de características es dividido en dimensiones fijas y se aplica *max-pooling*

8 x 8

0.88	0.44	0.16	0.14	0.37	0.77	0.96	0.27
0.19	0.45	0.16	0.57	0.63	0.29	0.71	0.70
0.66	0.36	0.64	0.82	0.54	0.73	0.59	0.25
0.85	0.24	0.84	0.76	0.29	0.75	0.62	0.24
0.32	0.74	0.39	0.31	0.34	0.03	0.33	0.48
0.20	0.69	0.13	0.16	0.73	0.65	0.96	0.32
0.19	0.14	0.86	0.09	0.88	0.07	0.01	0.48
0.83	0.24	0.04	0.97	0.34	0.35	0.50	0.91



# Arquitectura Fast RCNN



# Contenido

1. Introducción

2. Algoritmos

1. R-CNN

2. Faster R-CNN

**3. Single Shot Detector (SSD)**

4. You Only Look Once (Yolo)

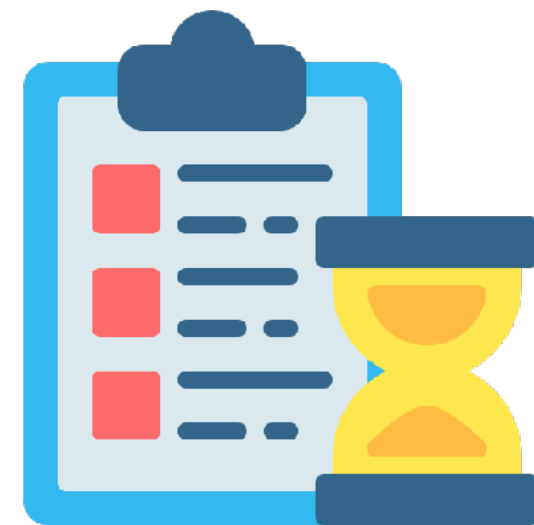
5. DETR

3. Métricas de evaluación

4. Aplicaciones

5. Retos

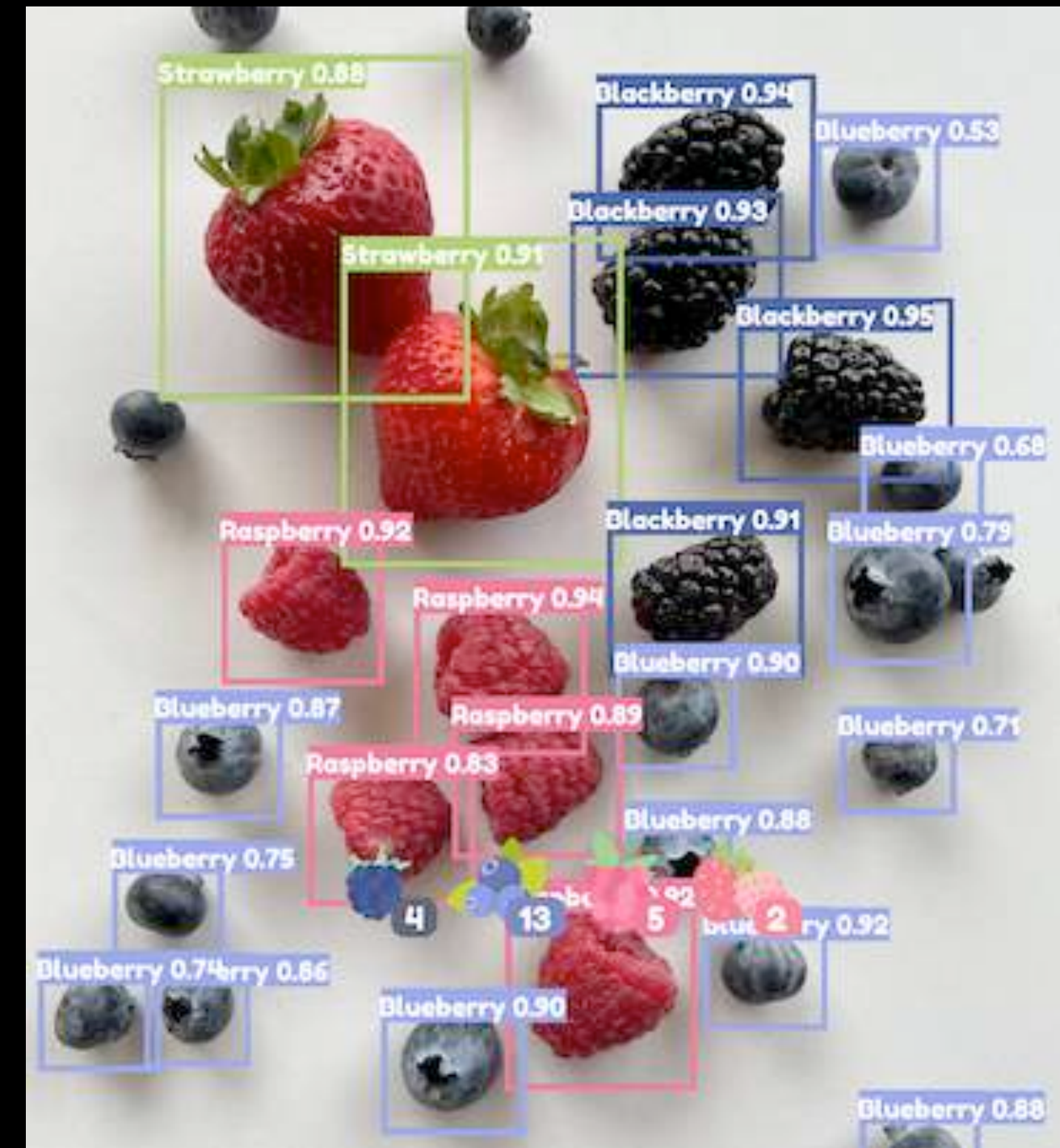
6. Taller



# SSD

## Single Shot Detector

- SSD está diseñado para la detección de objetos en tiempo real
- Es un detector de una sola etapa. (Fast R-CNN utiliza un “*region proposal*” y un clasificador CNN)
- SSD acelera el proceso eliminando el “*region proposal*”

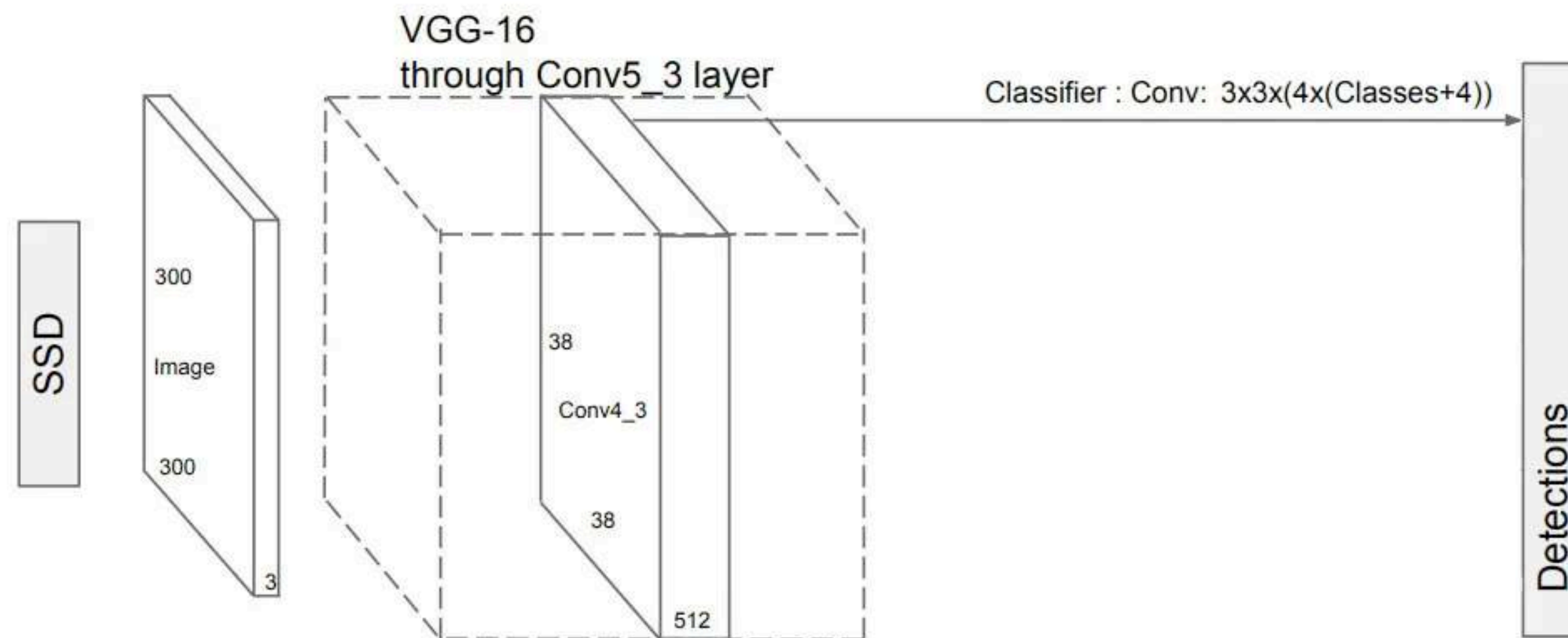


# SSD

## Single Shot Detector

Componentes del modelo SSD:

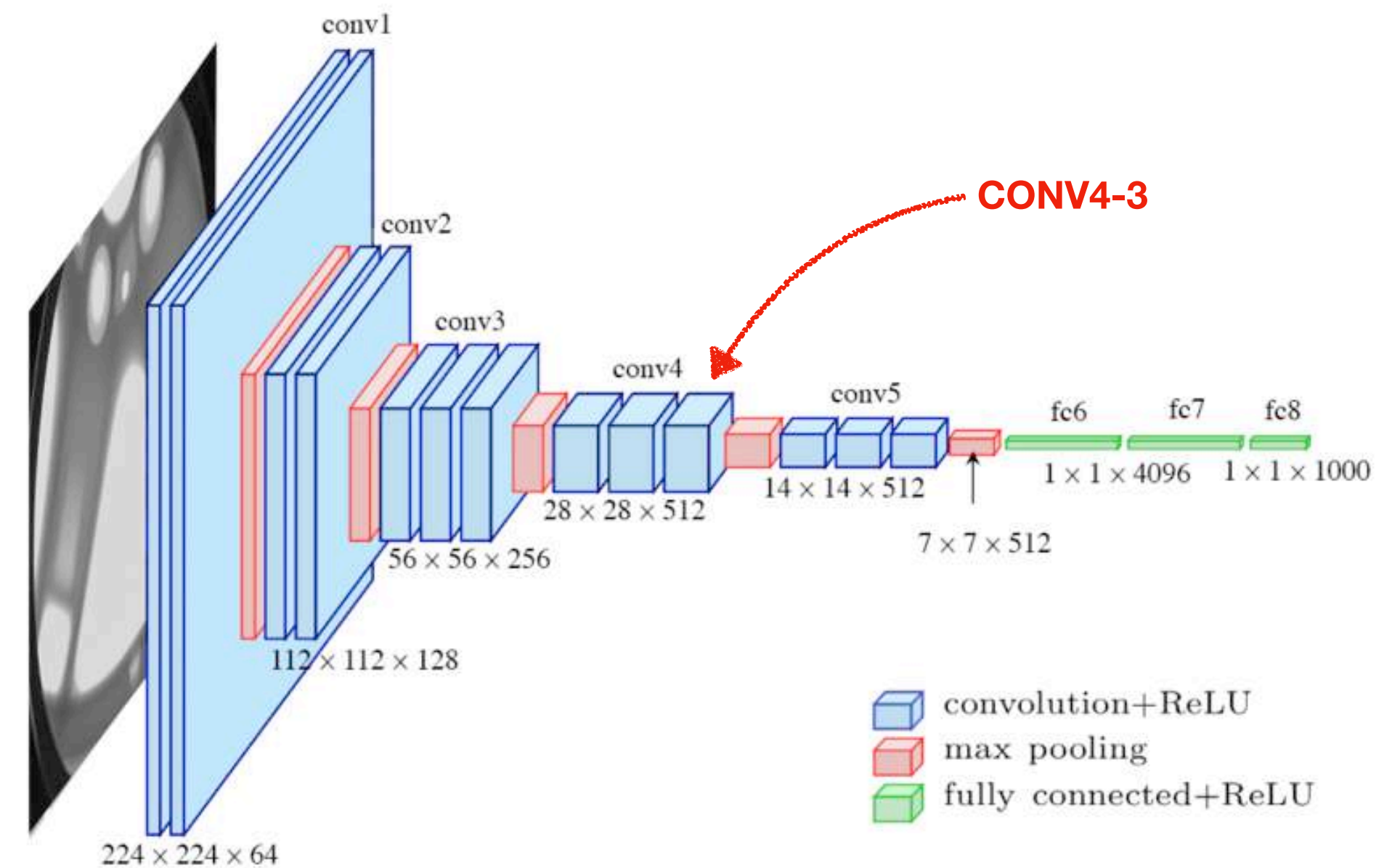
1. Extractor de características
2. CNN de detección





# Extractor de Características

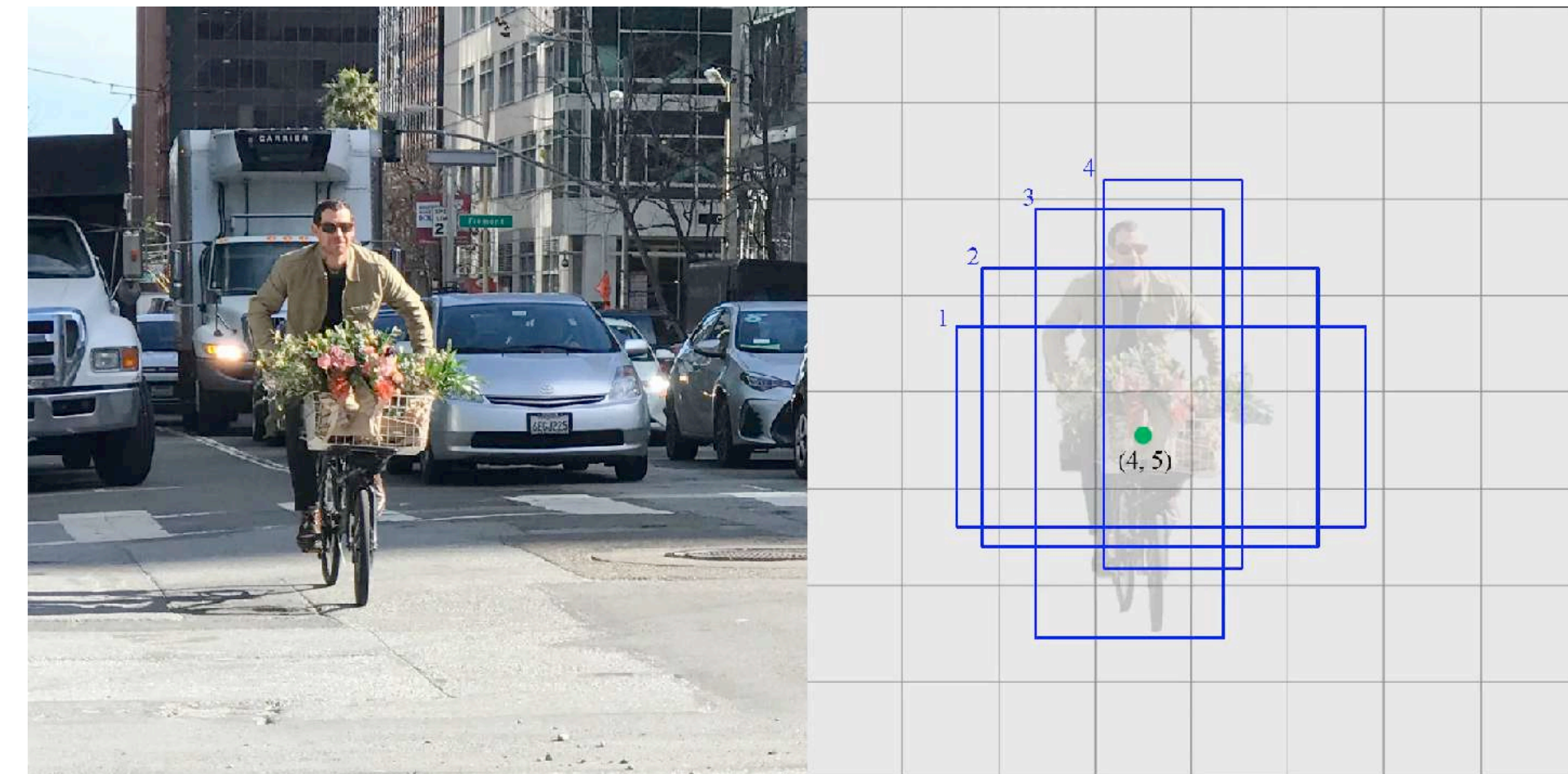
SSD Usa VGG16 para generar el mapa de características. Luego identifica los objetos usando los mapas desde CONV4-3



# Predicción

## Single Shot Detector

- Por cada celda del mapa de características **predice  $k$  objetos**
- Cada predicción se compone de:
  - A. Rectángulo (Bounding Box)
  - B. El vector correspondiente a la clase del BB + una clase de no objeto

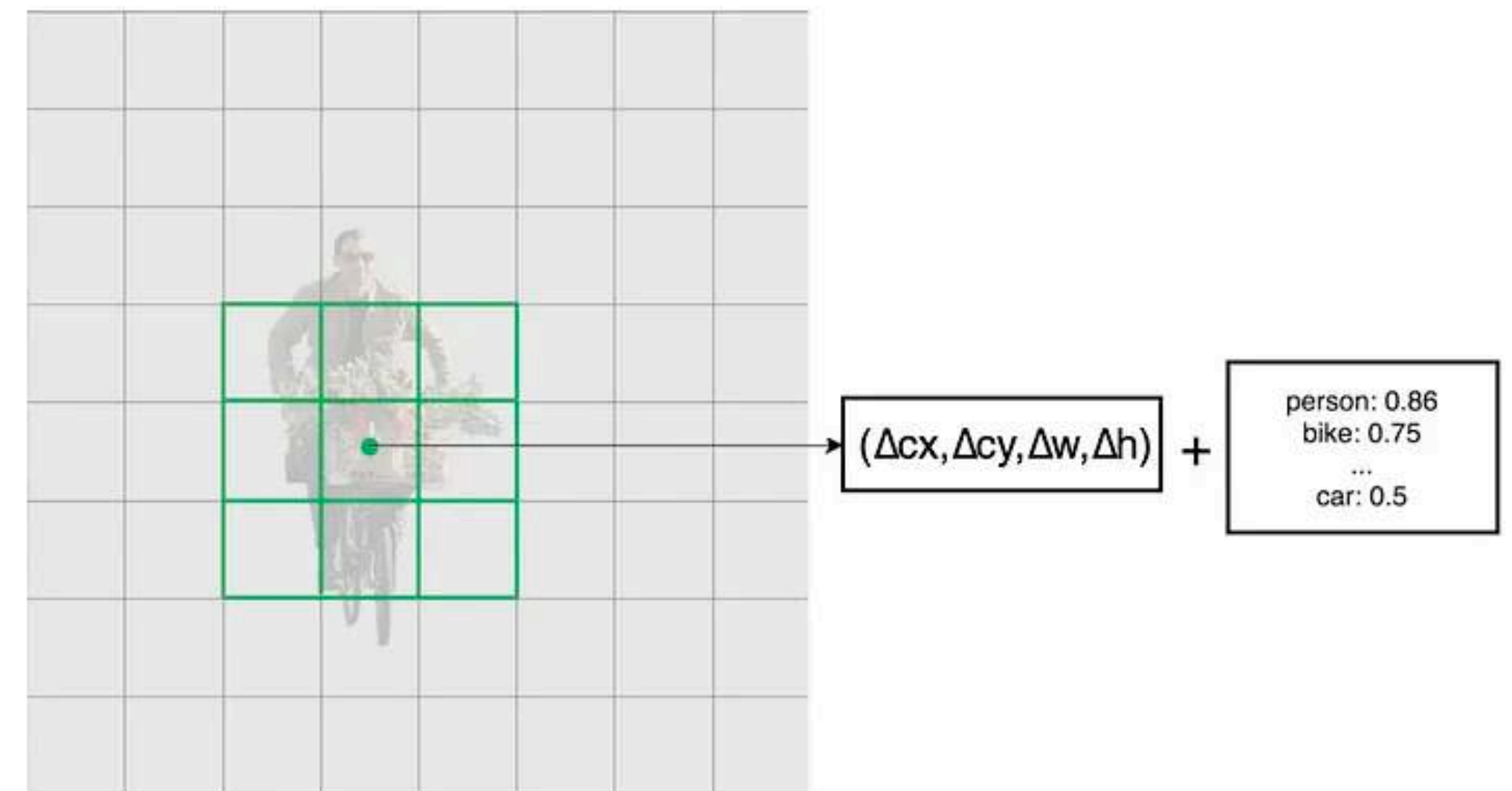


Predicción de detección por cada celda del mapa de características

# Predicción

## Single Shot Detector

- Por cada “celda” del mapa predice  $k$  objetos (rectángulos más la clase)
- Cada  $k$  corresponden a rectángulos de distintos ratios y tamaños
- Un rectángulo vertical es más adecuado para una persona y uno horizontal para un vehículo
- Aplica una capa de filtros convolucionales de 3x3 con stride de 1 sobre el mapa de características



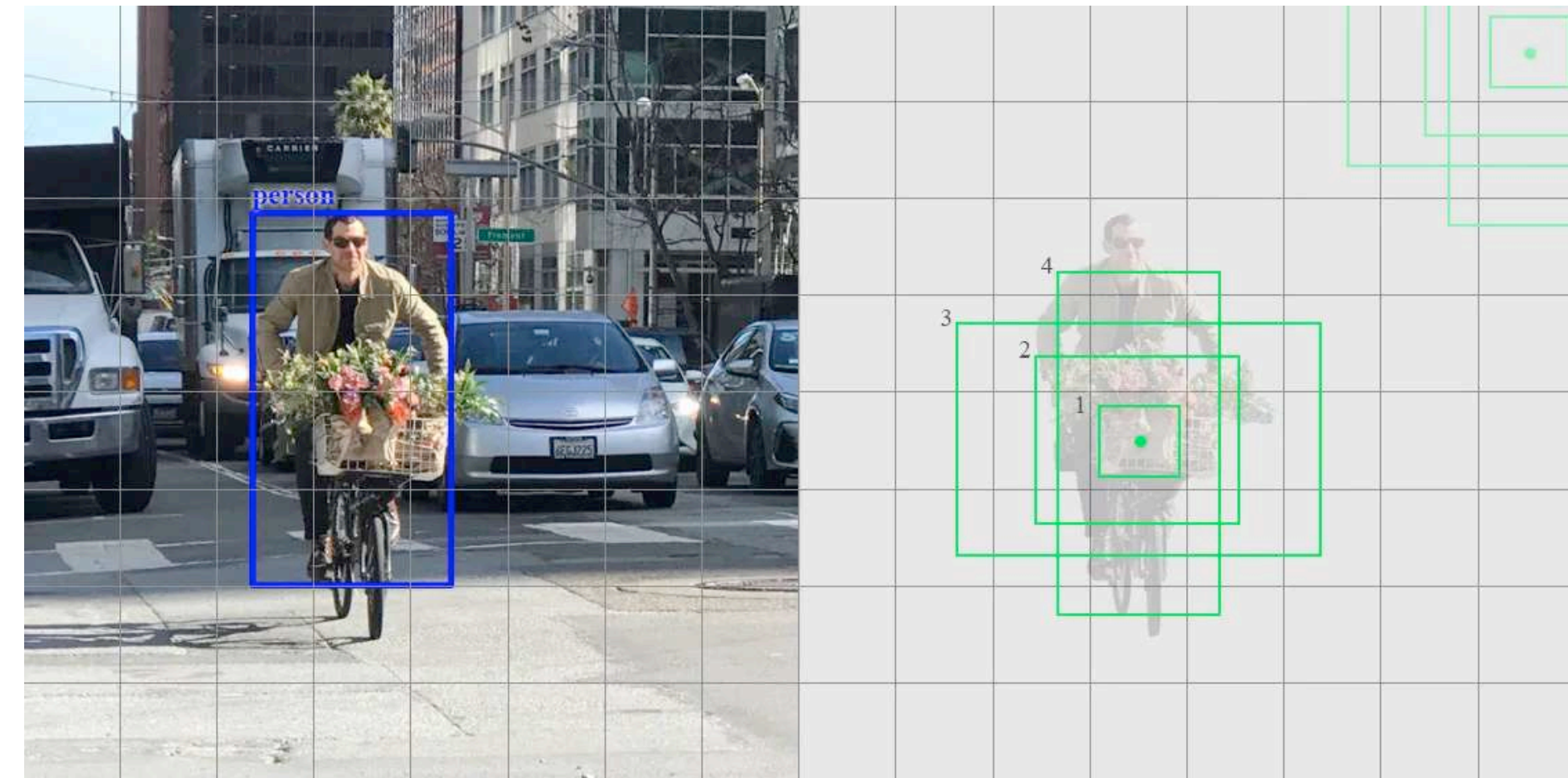
Predicción de detección por cada celda del mapa de características



# Predicción

## Single Shot Detector

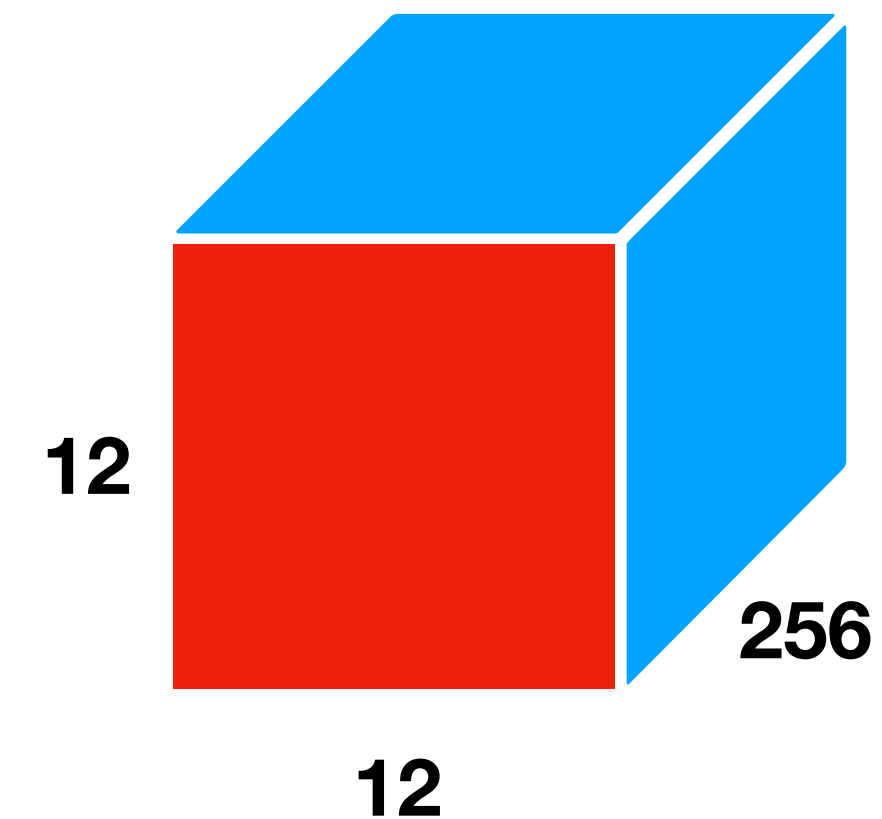
- El mapa de características es de  $m \times n$
- Por cada celda predice  $k$  número de objetos (aspect ratios)
- Cada  $k$  está compuesto de:
  - $c$  puntajes de la clase
  - 4 offsets relativos al rectángulo base



# Predicción

## Single Shot Detector

- Si tenemos un mapa de características de 12x12
- Vamos a predecir 10 clases
- Y tenemos 2 rectángulos base (2 aspect ratio)
- ¿Cuántos filtros debe tener la capa convolucional de detección?



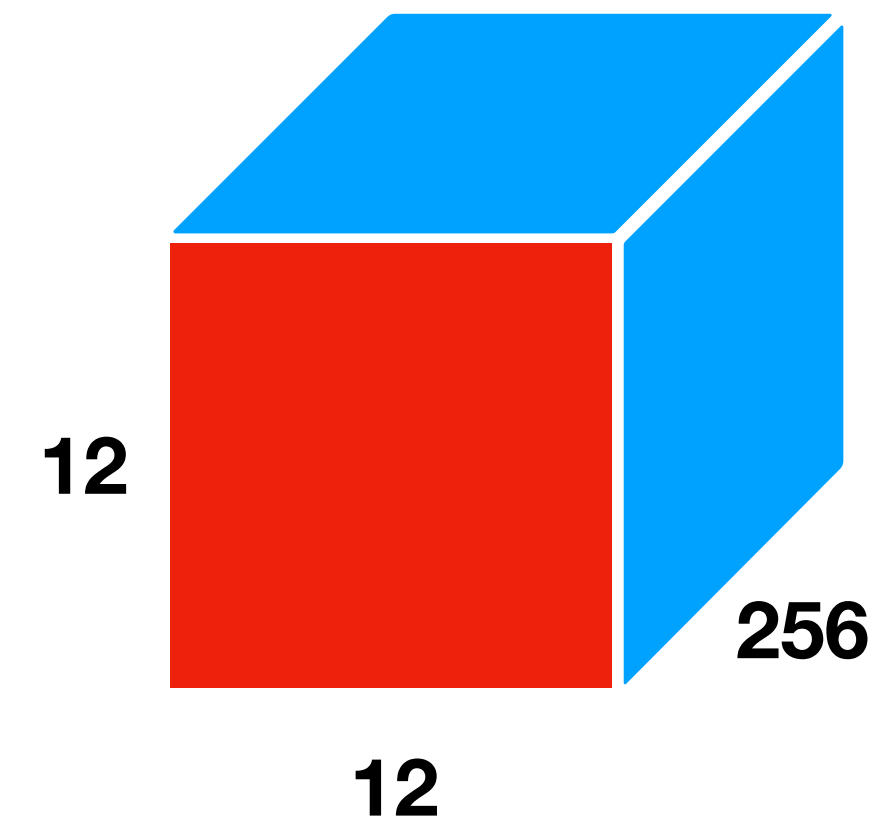
# Predicción

## Single Shot Detector

Número de filtros de la capa de predicción

$$(c + 5) \times k$$

$$(10 + 5) \times 2 = 30$$

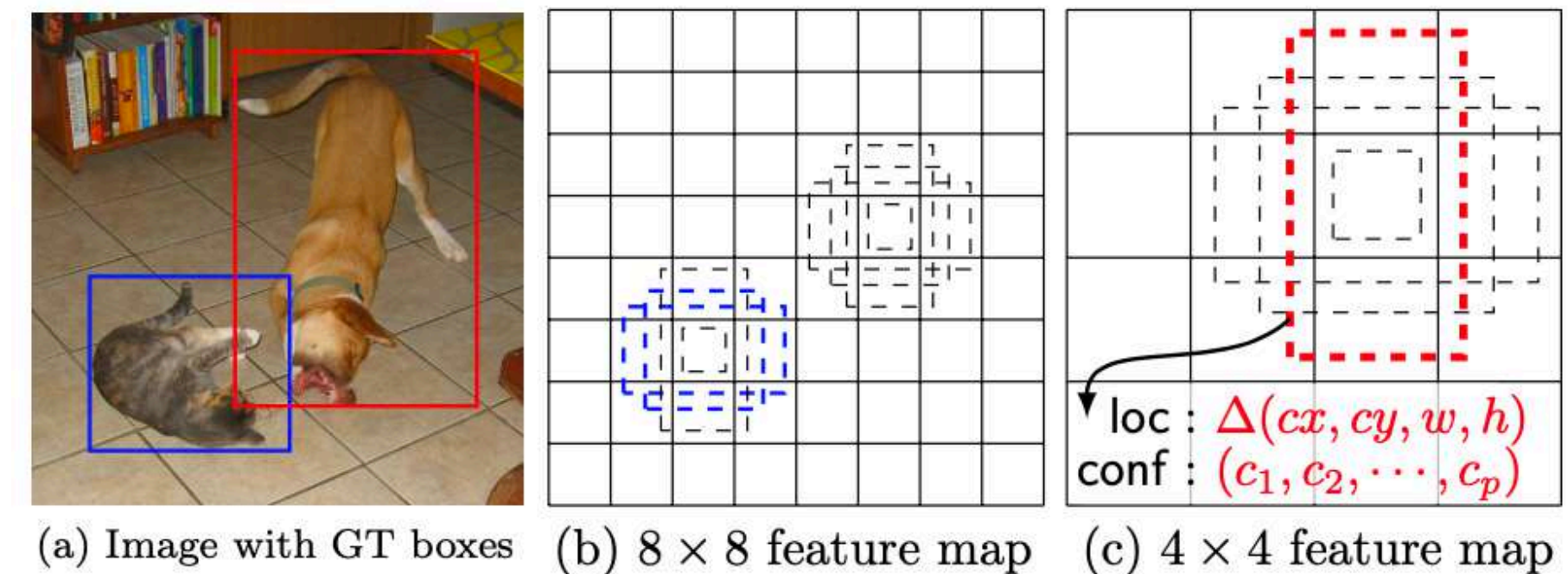




# Múltiples escalas

## Single Shot Detector

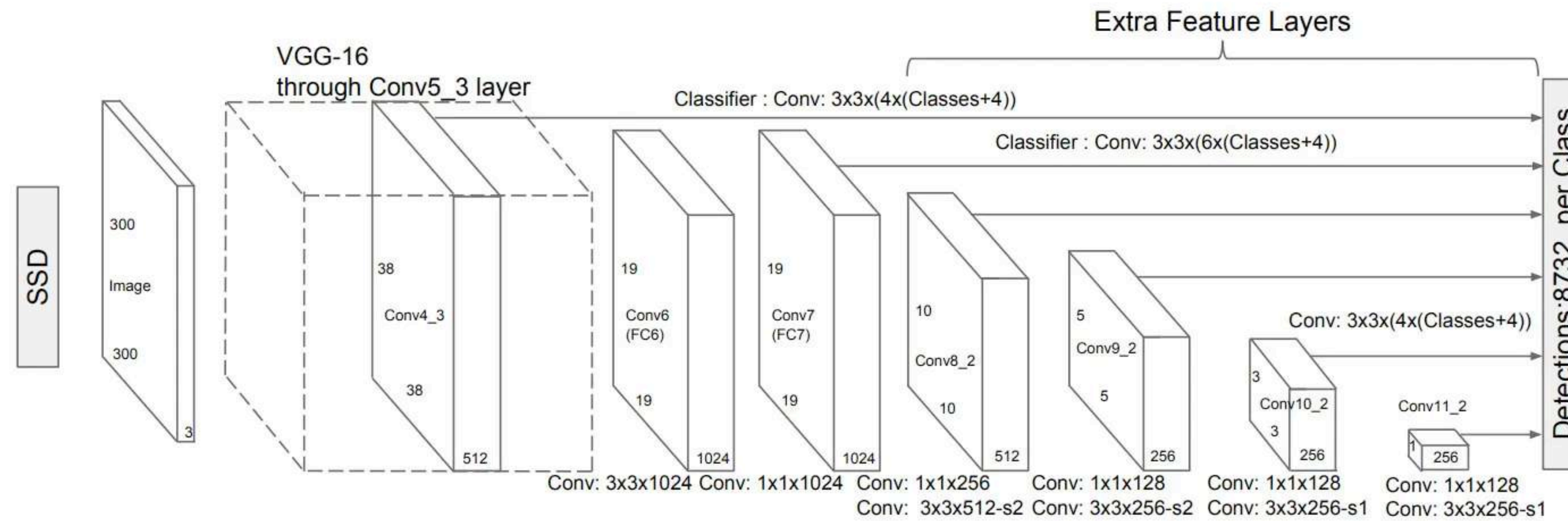
- VGG16 reduce gradualmente la dimensión espacial, la resolución de los mapas de características también disminuye
- SSD utiliza capas de menor resolución para detectar objetos a mayor escala



Los mapas de características de menor resolución  
(derecha) detectan objetos de mayor escala.

# SSD

## Arquitectura

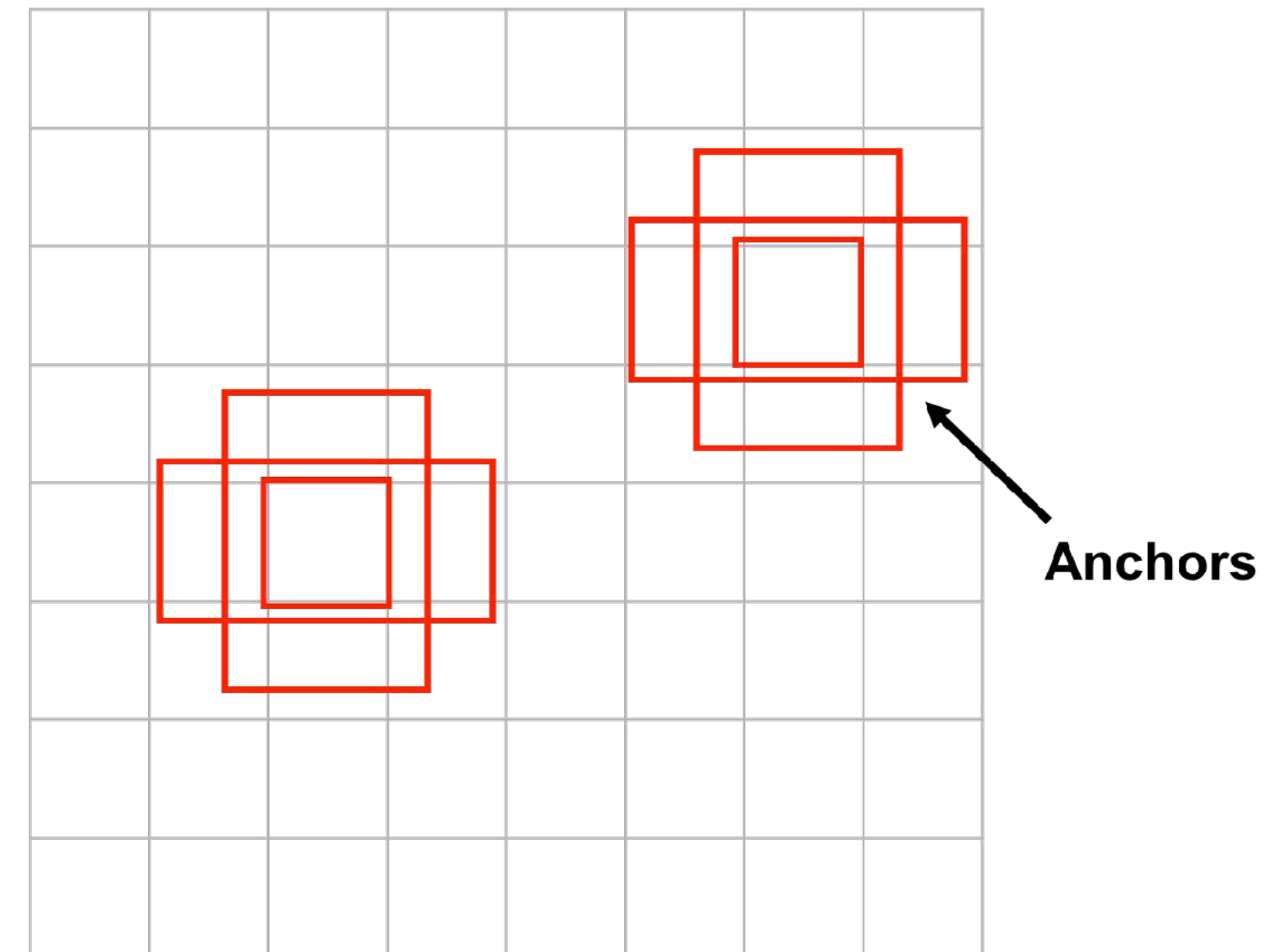


SSD añade 6 capas más a la VGG16 que produce  $k$  predicciones de objetos por capa

# Default boxes & aspect ratios

## Single Shot Detector

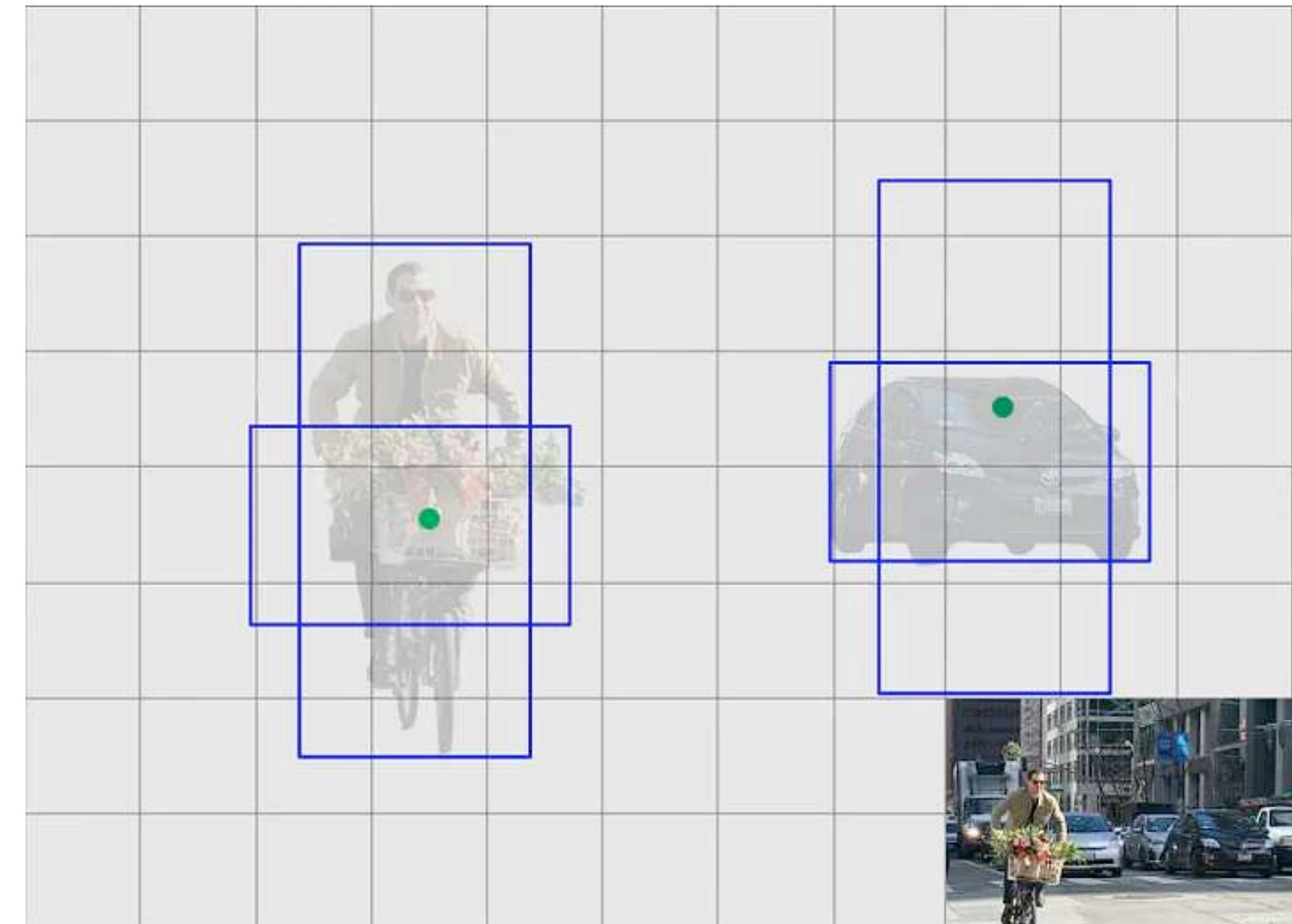
- Durante el entrenamiento el modelo puede luchar entre sí para determinar qué formas (personas o vehículos) deben optimizarse
- El entrenamiento inicial puede ser muy inestable. Las predicciones de los rectángulos que funcionan bien para una categoría para otras no
- Se quiere que las predicciones iniciales sean diversas y no se parezcan





# Default boxes and aspect ratios

- Se usan recuadros predeterminados para cubrir el espectro de los objetos a predecir y reducir la complejidad computacional
- SSD sugiere 4 o 6 rectángulos por defecto
- Las diferentes capas utilizan diferentes rectángulos por defecto para detecciones de diferentes tamaños



# Default boxes and aspect ratios

- SSD define un valor de escala para cada capa del mapa de características
- Se calcula el tamaño del rectángulo usando la escala y los aspect ratios pre definidos
- Las escalas van desde 0.2 a 0.9 (última capa)
- SSD propone 5 aspect ratios: 1, 2, 3, 1/2, y 1/3

**Cálculo del ancho**

$$w = s\sqrt{a}$$

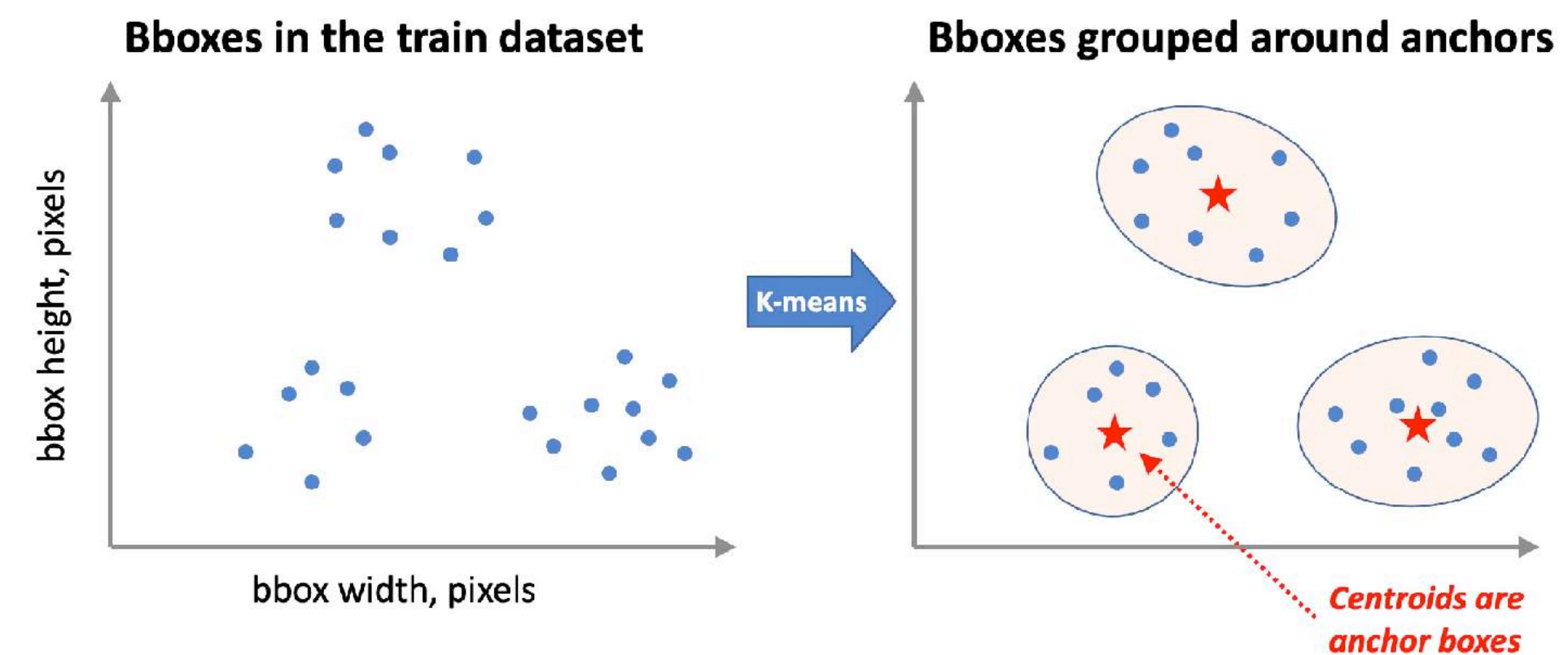
**Calculo del alto**

$$h = \frac{s}{\sqrt{a}}$$

En donde  $s$  es la escala y  $a$  el “*aspect ratio*”

# Default boxes and aspect ratios

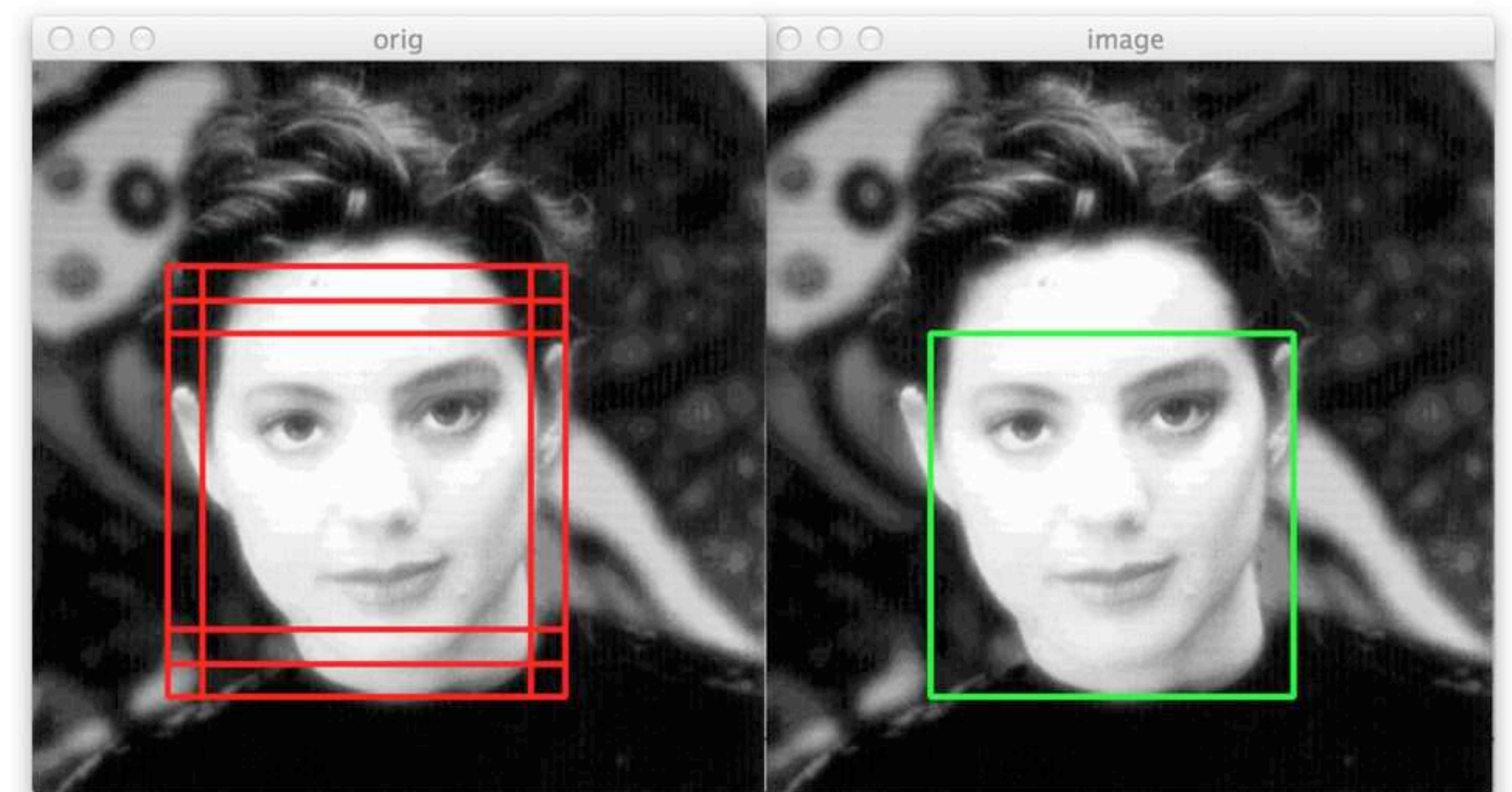
- Los *Bounding Box* también pueden ser definidos por agrupamiento (k-means)
- Se usan las etiquetas de entrenamiento para el calculo de los BB
- [Como calcular los BB con K-means](#)





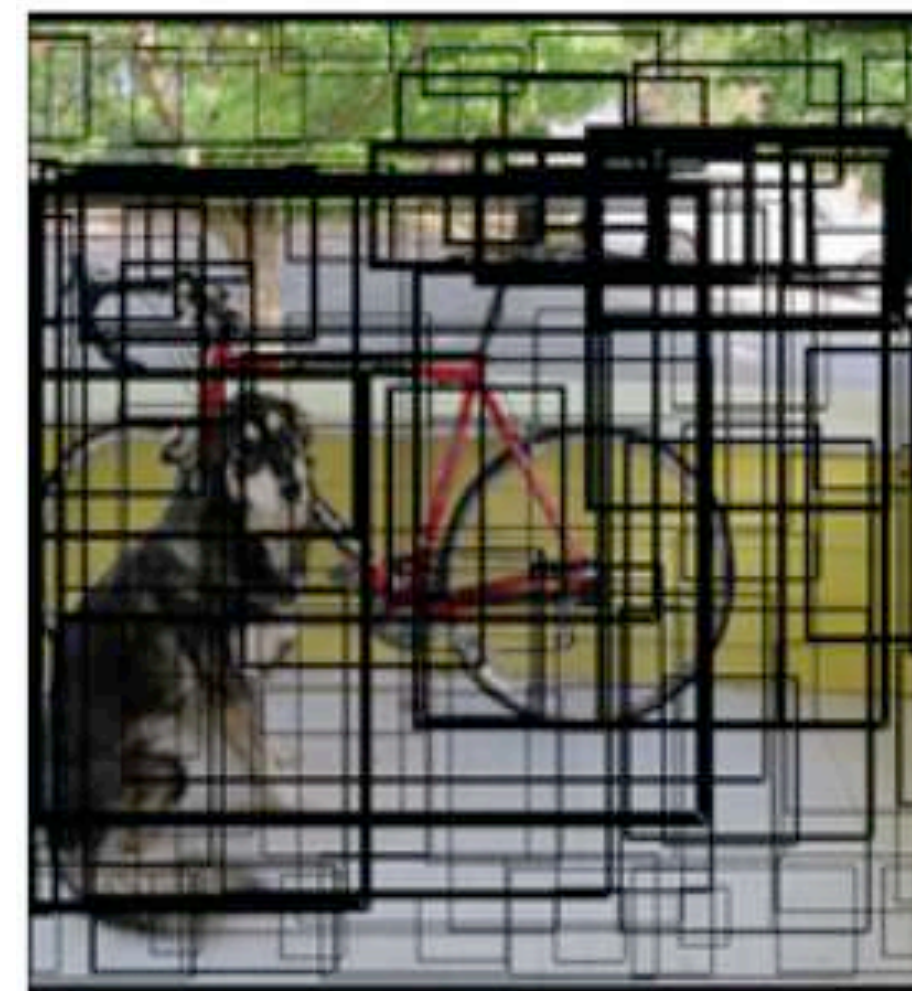
# Predicción final

- Luego de la predicción se tienen múltiples rectángulos que se solapan
- Se realiza un paso de *non-maximum suppression* para producir el resultado final

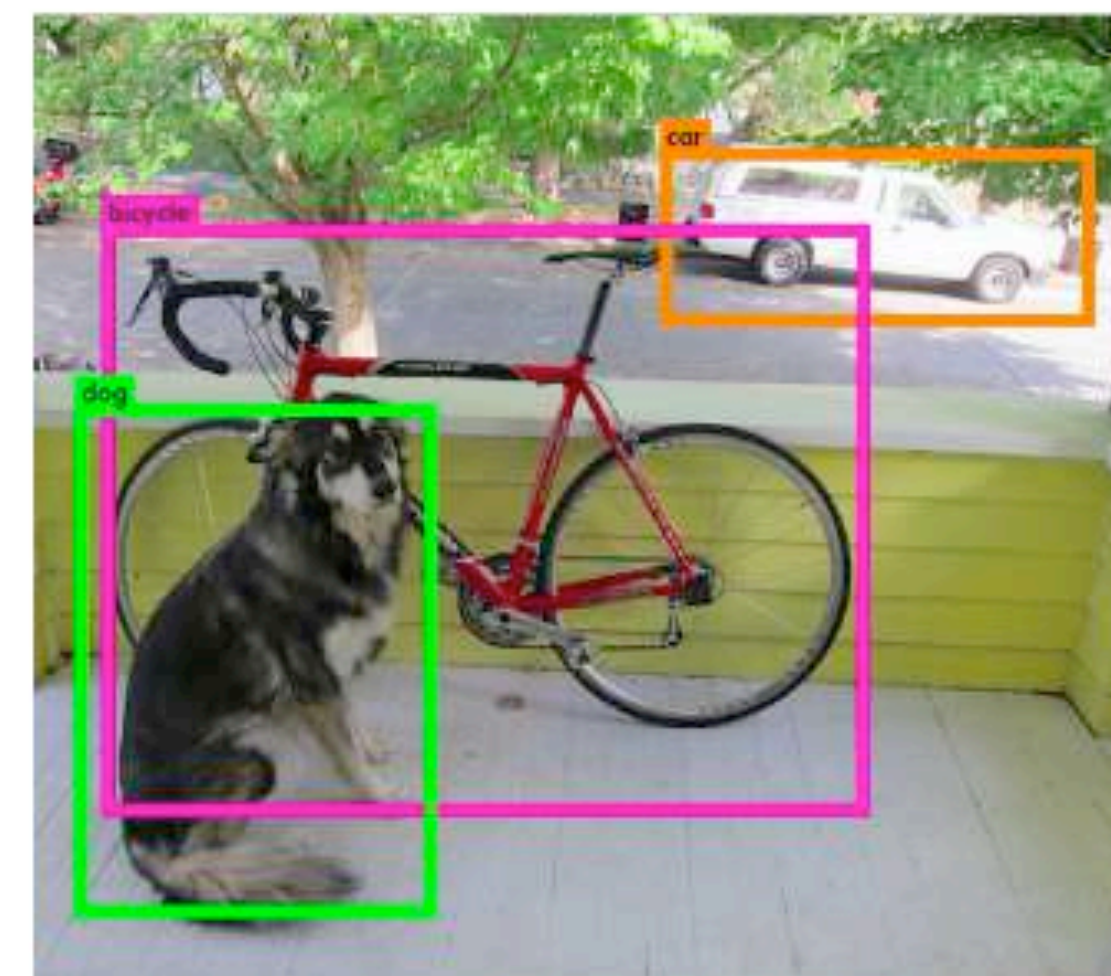


# Non-maximum Suppression

- La primera etapa tiene una alta sensibilidad para predecir todas las posibles regiones, esto da paso a tener cientos de regiones similares
- Non-maximum Suppression es una técnica para “fusionar” rectángulos del mismo objeto



Multiple Bounding Boxes



Final Bounding Boxes

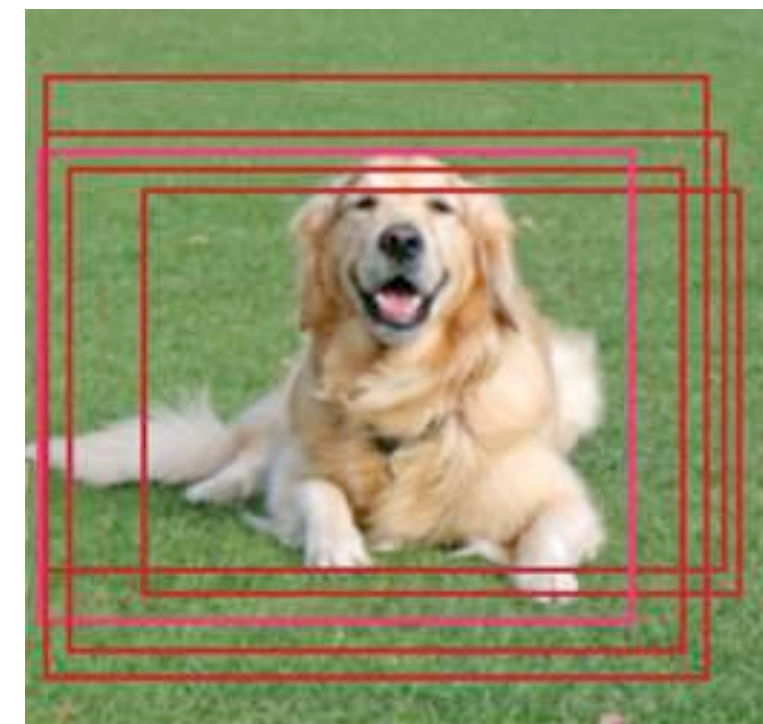


# Non-maximum Suppression

1. Seleccione la propuesta con mayor score, elimínela de  $B$  y añádala a la lista final de propuestas  $D$  (Inicialmente  $D$  está vacía)
2. Calcula el IOU de la propuesta con los otros rectángulos en  $B$ . Si el IOU es mayor que el umbral  $N$  (\*Y es la misma clase) se elimina
3. Este proceso se repite hasta que no queden más propuestas en  $B$



After applying non-maximum suppression

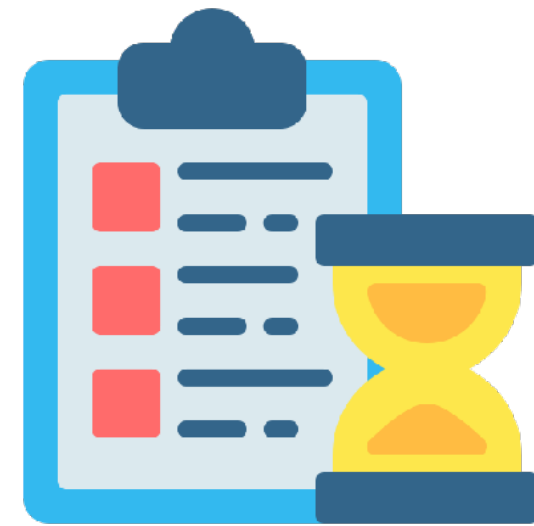


Predictions before NMS



# Contenido

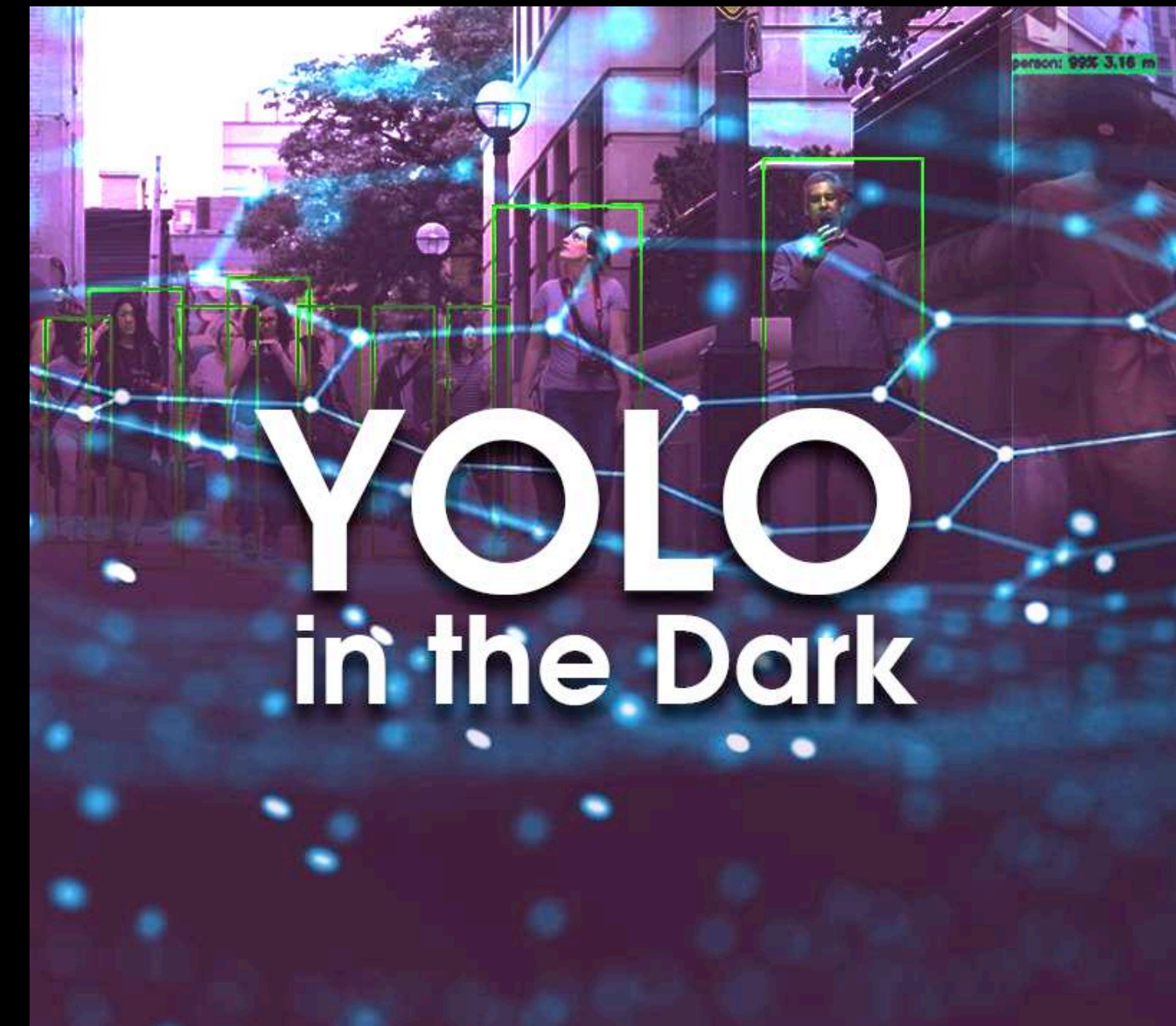
1. Introducción
2. Algoritmos
  1. R-CNN
  2. Faster R-CNN
  3. Single Shot Detector (SSD)
  - 4. You Only Look Once (Yolo)**
  5. DETR
3. Métricas de evaluación
4. Aplicaciones
5. Retos
6. Taller



# Yolo

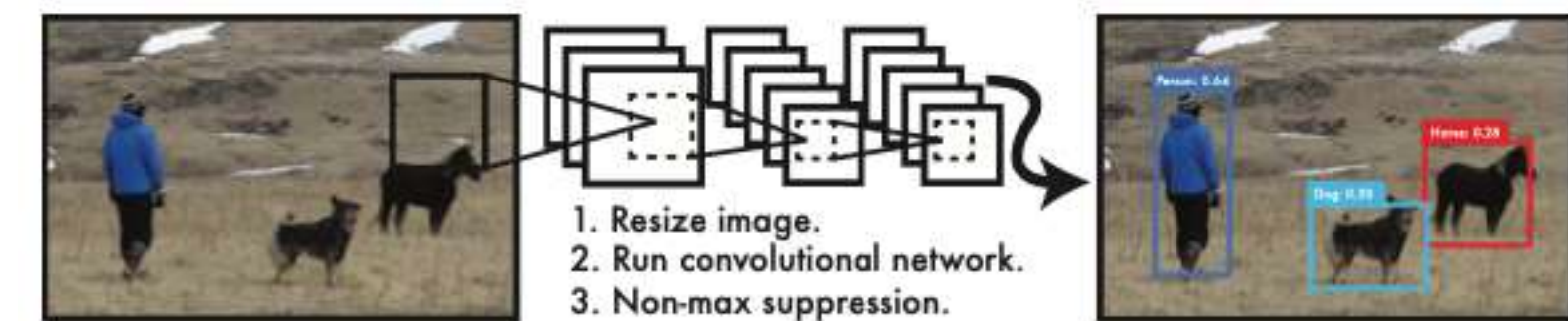
You Only Look Once

- Es uno de los algoritmos de detección más conocidos y usados
- Es un detector de una sola etapa
- Arquitectura basada en los clasificadores CNN



# Introducción

- YOLOv1 (2016): la primera de las muchas iteraciones por las que ha pasado esta arquitectura
- La idea básica de la arquitectura sigue siendo la misma. YOLOv1 se denomina simplemente YOLO (Actual v12)
- La clave es velocidad, predice a 45 frames por segundo: es una buena opción para aplicaciones en tiempo real



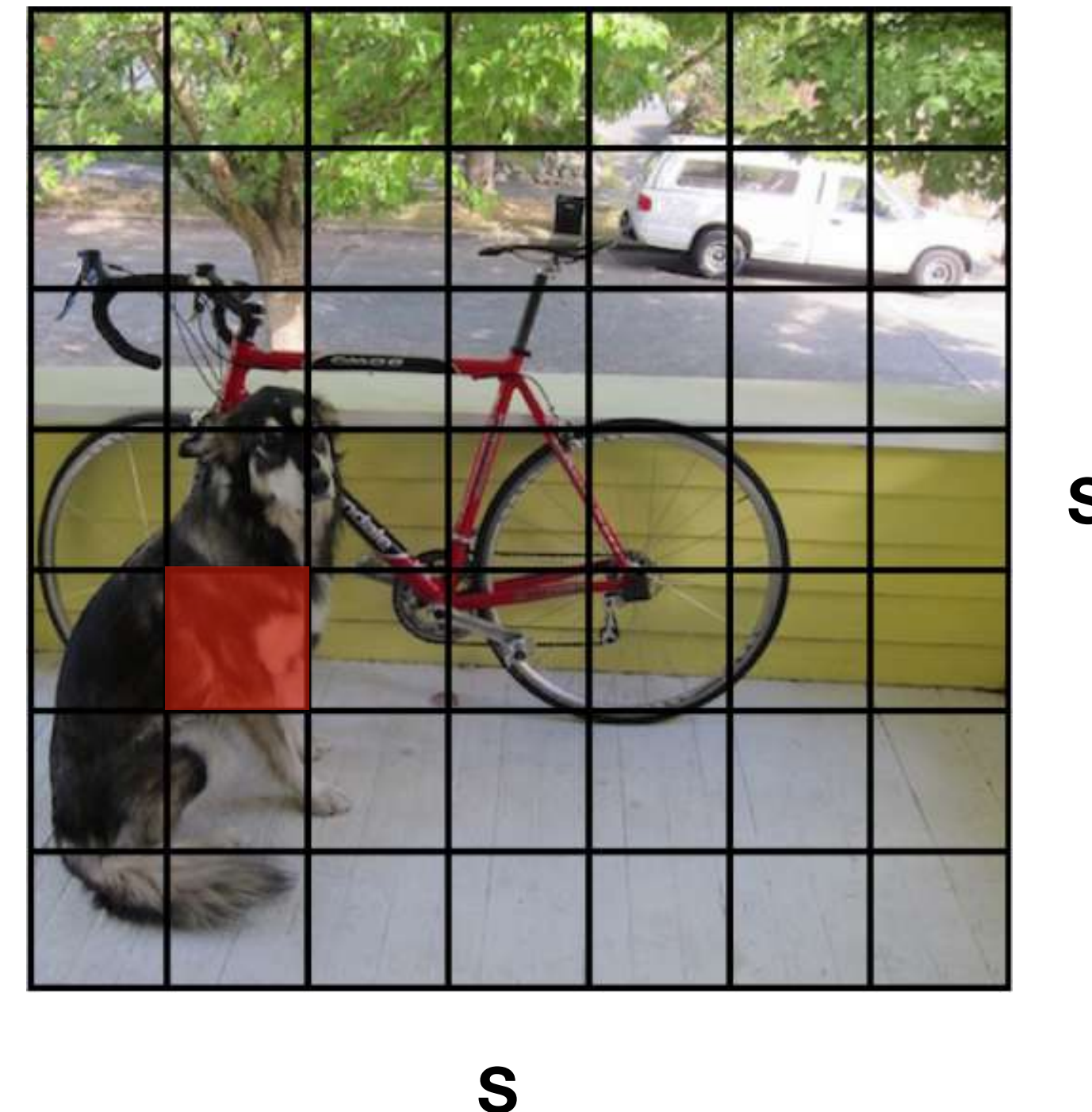
**Figure 1: The YOLO Detection System.** Processing images with YOLO is simple and straightforward. Our system (1) resizes the input image to  $448 \times 448$ , (2) runs a single convolutional network on the image, and (3) thresholds the resulting detections by the model's confidence.



# Introducción

## Yolo

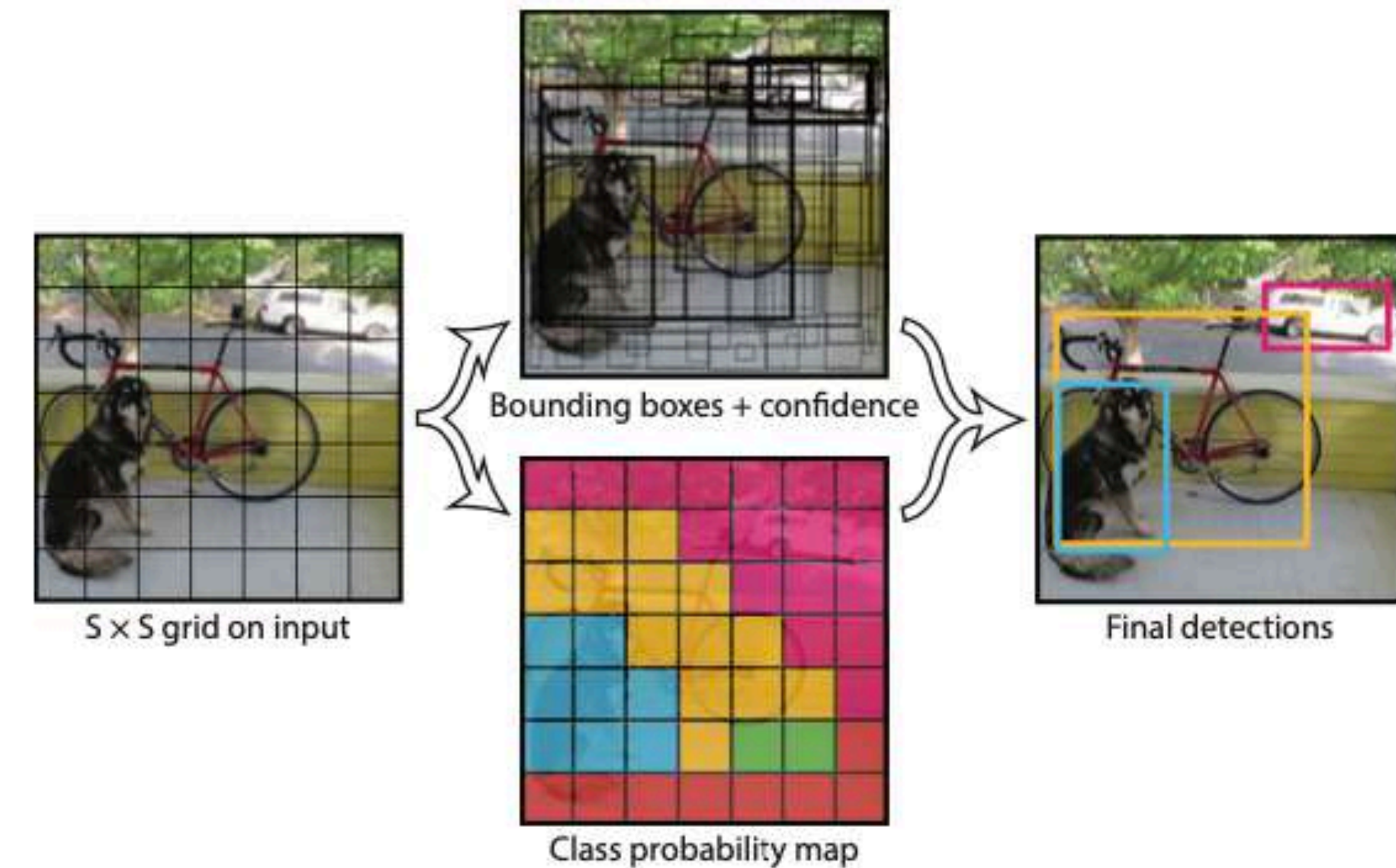
- YOLO se basa en la idea de segmentar una imagen en partes más pequeñas
- La imagen se divide en una cuadrícula de dimensiones **S×S**
- La celda en la que se encuentra el centro de un objeto es la responsable de detectar ese objeto (centro del perro)



# Predicciones

## Yolo

- Cada celda predecirá **K** rectángulos y una puntuación de confianza para cada uno
- Estos niveles de confianza reflejan la certeza del modelo de que existe un objeto en esa celda
- La confianza representa la diferencia entre el IOU de la predicción y el rectángulo verdadero



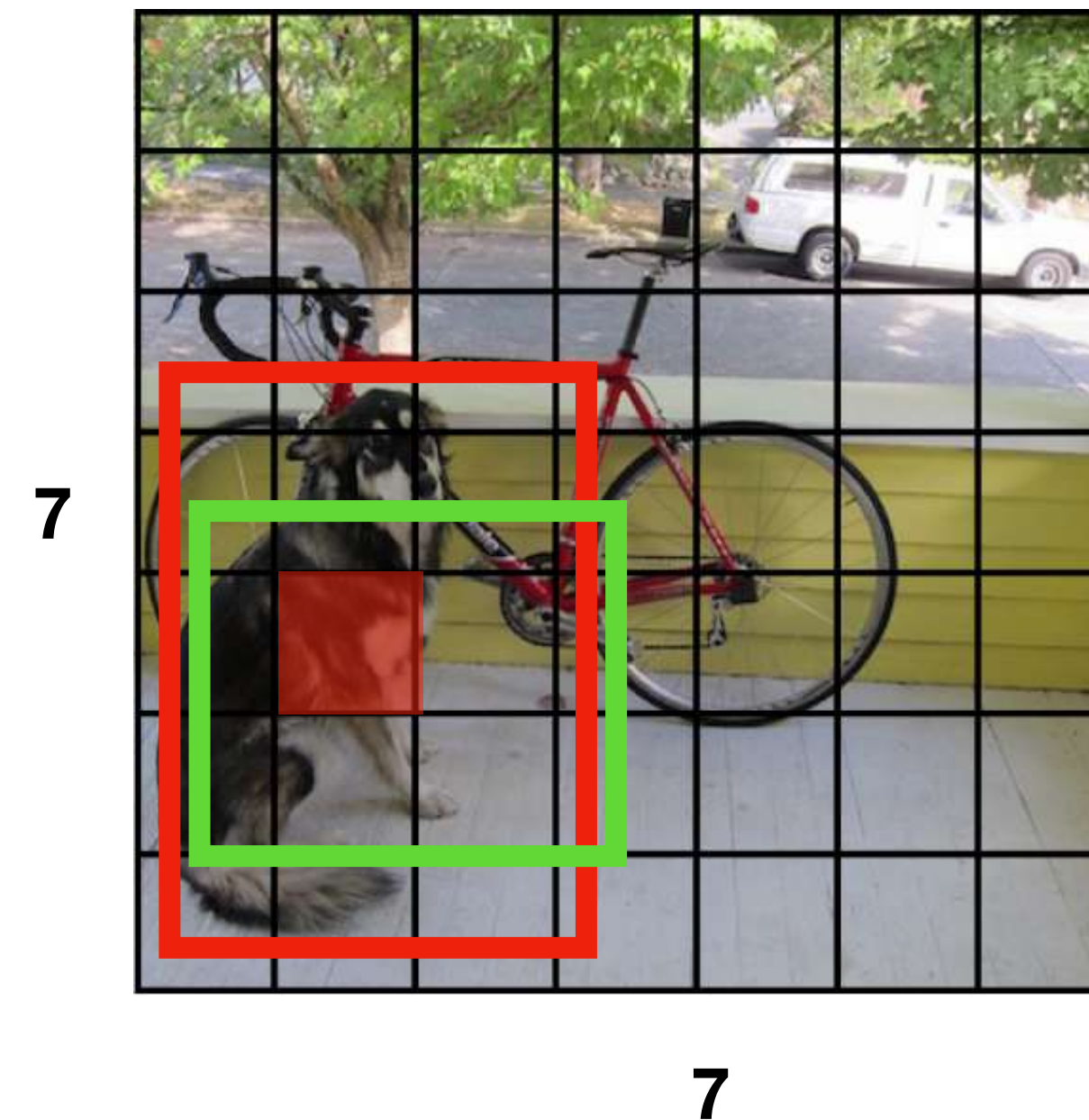


# Predicciones

## Yolo

Además de rectángulos y la puntuación de confianza, cada celda predice la clase del objeto (on-hot)

**¿Cuántas predicciones hace un modelo si tiene 5 clases y una cuadrícula de 7x7 y predice 2 rectángulos por cuadro?**



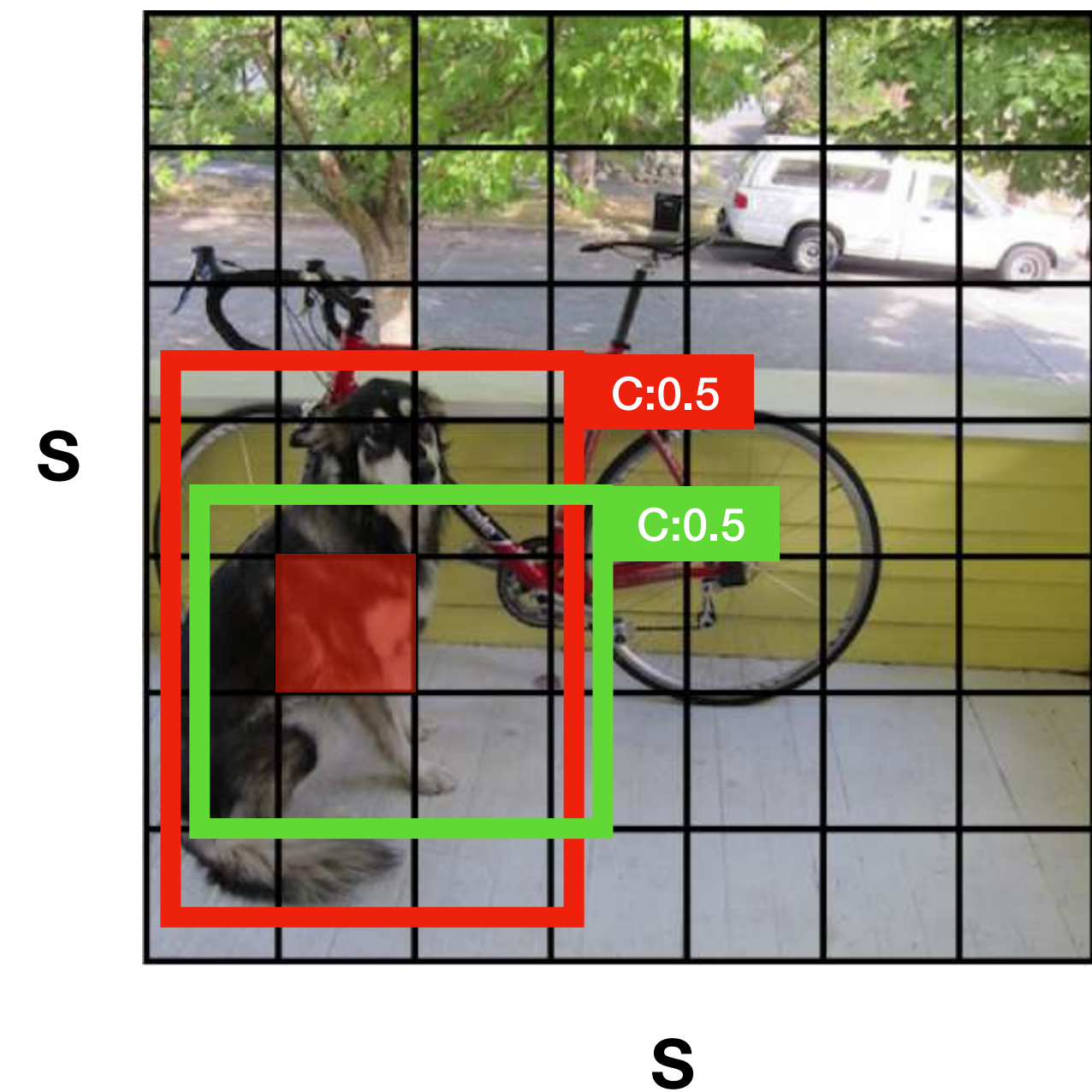


# Predicciones

## Yolo

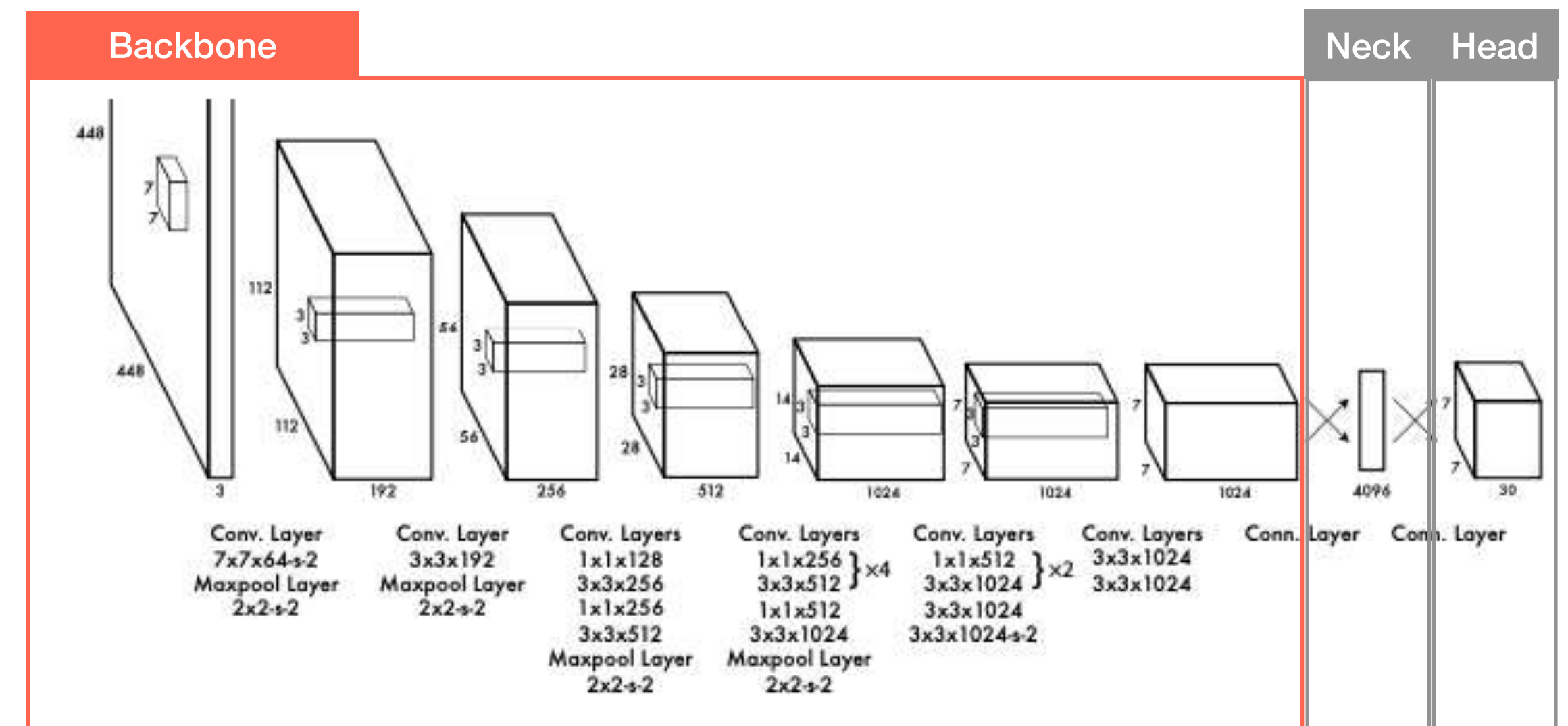
$$\begin{array}{ccccc} 7 \times 7 & \times & (1+4+5) & \times & 2 \\ \text{CUADRÍCULA} & & \text{BOUNDING BOX + CONFIABILIDAD + CLASES} & & \text{PREDICCIONES} \\ & & = 980 & & \end{array}$$

\*\* El modelo predice el centro del recuadro con el ancho y alto en lugar de las posiciones de las esquinas superior izquierda e inferior derecha.



# Arquitectura Backbone

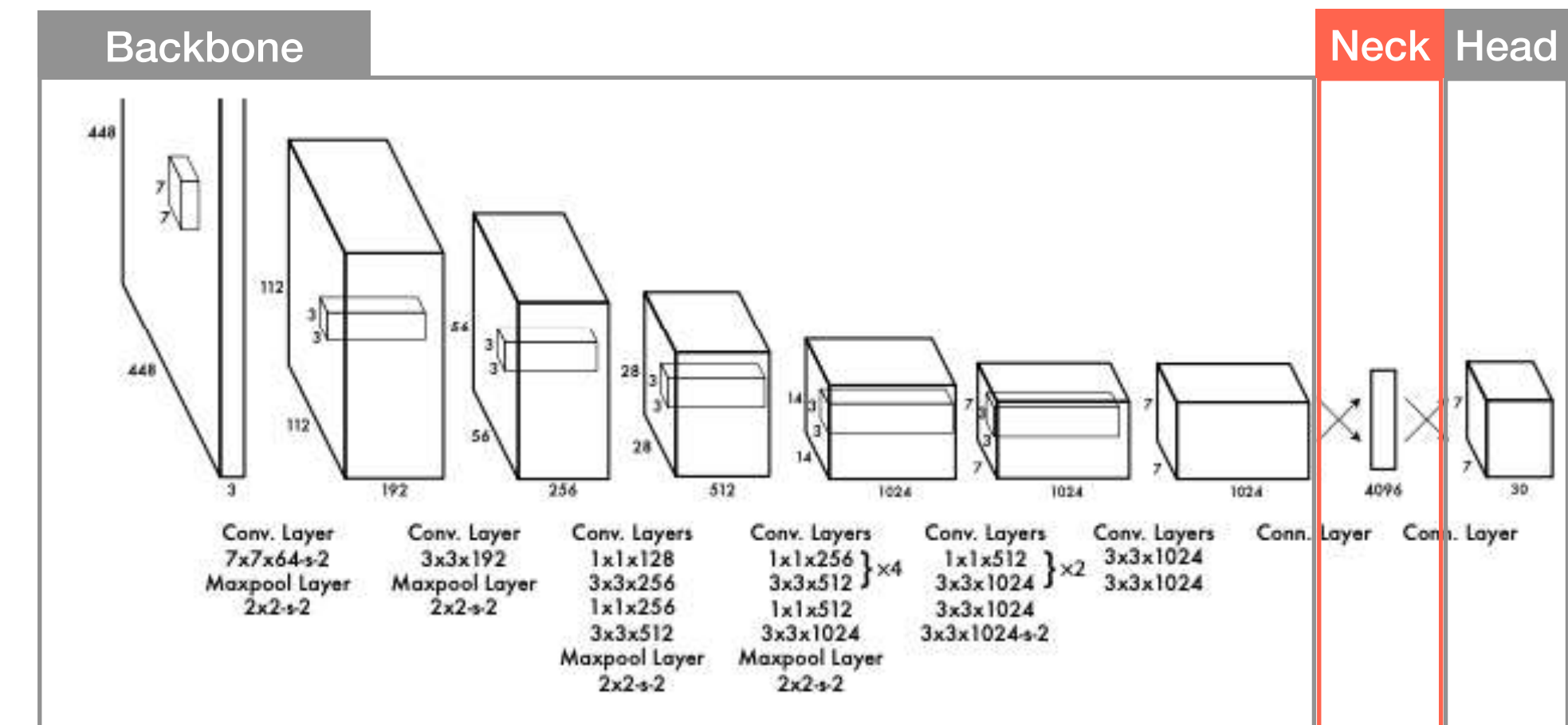
- Arquitectura de YOLO tiene tres componentes (“head”, “neck”, y “backbone”) básicamente un modelo de clasificación
- El “**backbone**” encarga de la extracción de características y se entrena primero en un conjunto de datos de clasificación (ImageNet)
- Los “parches” corresponden a la reducción espacial del extractor de características



# Arquitectura

## Neck

- El mapa de características se aplanan y comprime en una representación latente mediante una capa totalmente conectada (FC)
- Produce un vector que contiene la información necesaria para las predicciones que posteriormente es utilizado por la cabeza

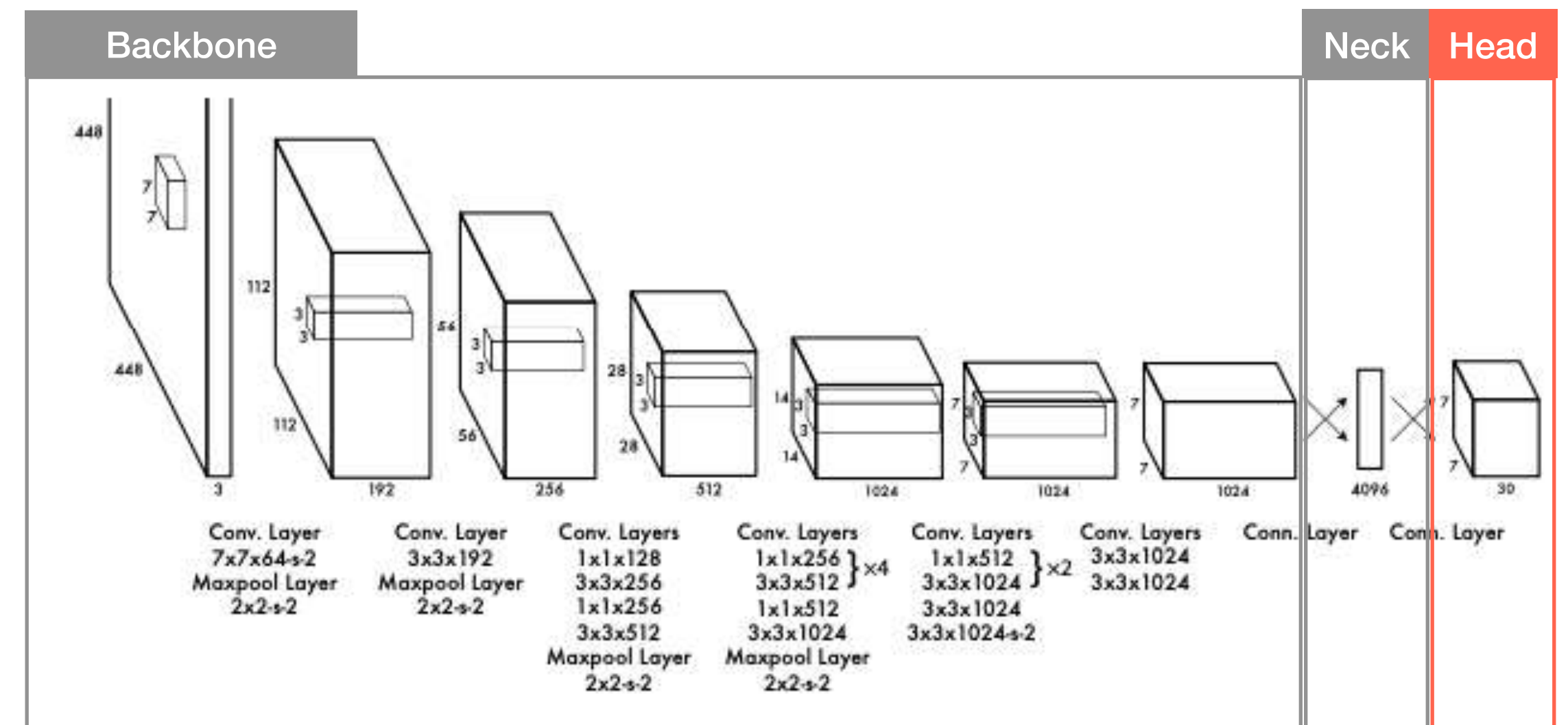




# Arquitectura

## Head

- La capa de salida final de la red que realiza las predicciones
- Una segunda capa (FC) toma este vector latente y lo expande en una salida estructurada correspondiente a la malla  $S \times S$



# Función de costo

## Yolo

- Es la suma del error cuadrático ponderado de todas sus predicciones
- $\lambda_{coord}$  y  $\lambda_{noobj}$  balancean el costo dandole mas fuerza a la localización y disminuyendo el impacto de confidence de no objeto
- $\mathbb{1}_{ij}^{obj}$  Sólo penaliza el error si un objeto está presente en esa celda “responsable” de su predicción

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

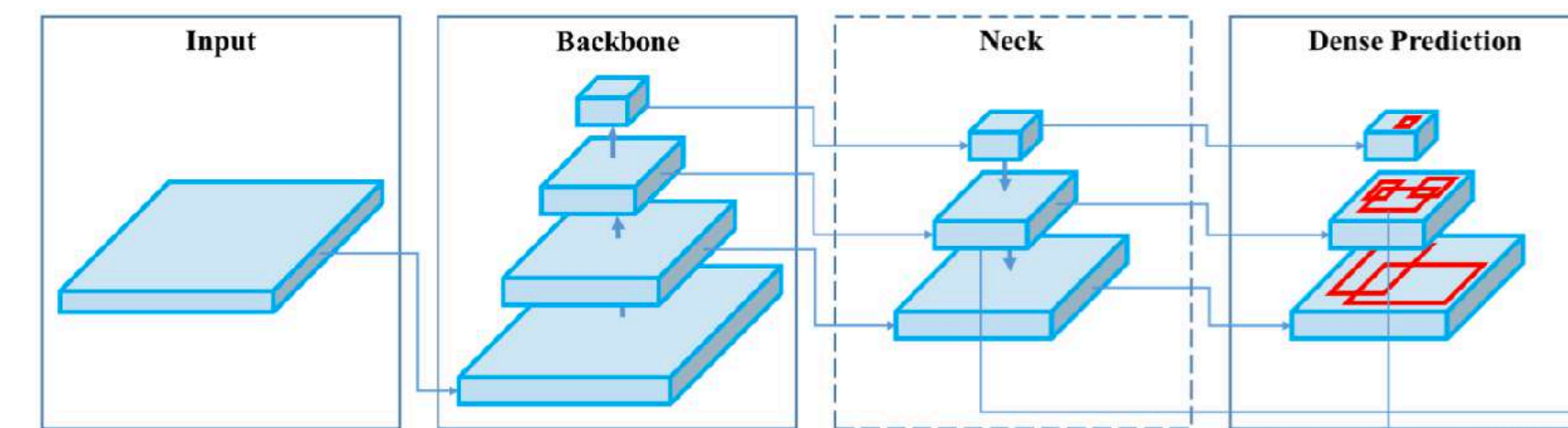
# Actualizaciones

## Yolo

**Anchor Boxes (YOLOv2):** Adoptado de Faster R-CNN para mejorar la localización

**Multi-Scale Detection (YOLOv3):**  
inspirado en Feature Pyramid Networks,  
mejora la detección de objetos a diferentes  
escalas

**Data Augmentation (YOLOv4):** Uso de  
“Mosaic augmentation” mejorar la  
generalización



Feature Pyramid



Mosaic augmentation



# Referencias

- 2014 - Rich feature hierarchies for accurate object detection and semantic segmentation
- 2004 - Efficient Graph-Based Image Segmentation
- 2013 - Selective Search for Object Recognition
- 2015 - Fast R-CNN
- 2016 - SSD: Single Shot MultiBox Detector
- 2016 - You Only Look Once Unified Real-Time Object Detection
- 2017 - YOLO9000 Better, Faster, Stronger
- 2018 - YOLOv3 An Incremental Improvement
- 2020 - YOLOv4 Optimal Speed and Accuracy of Object Detection
- 2020 - End-to-End Object Detection with Transformers
- 2019 - Generalized Intersection over Union: A Metric and A Loss for Bounding Box Regression