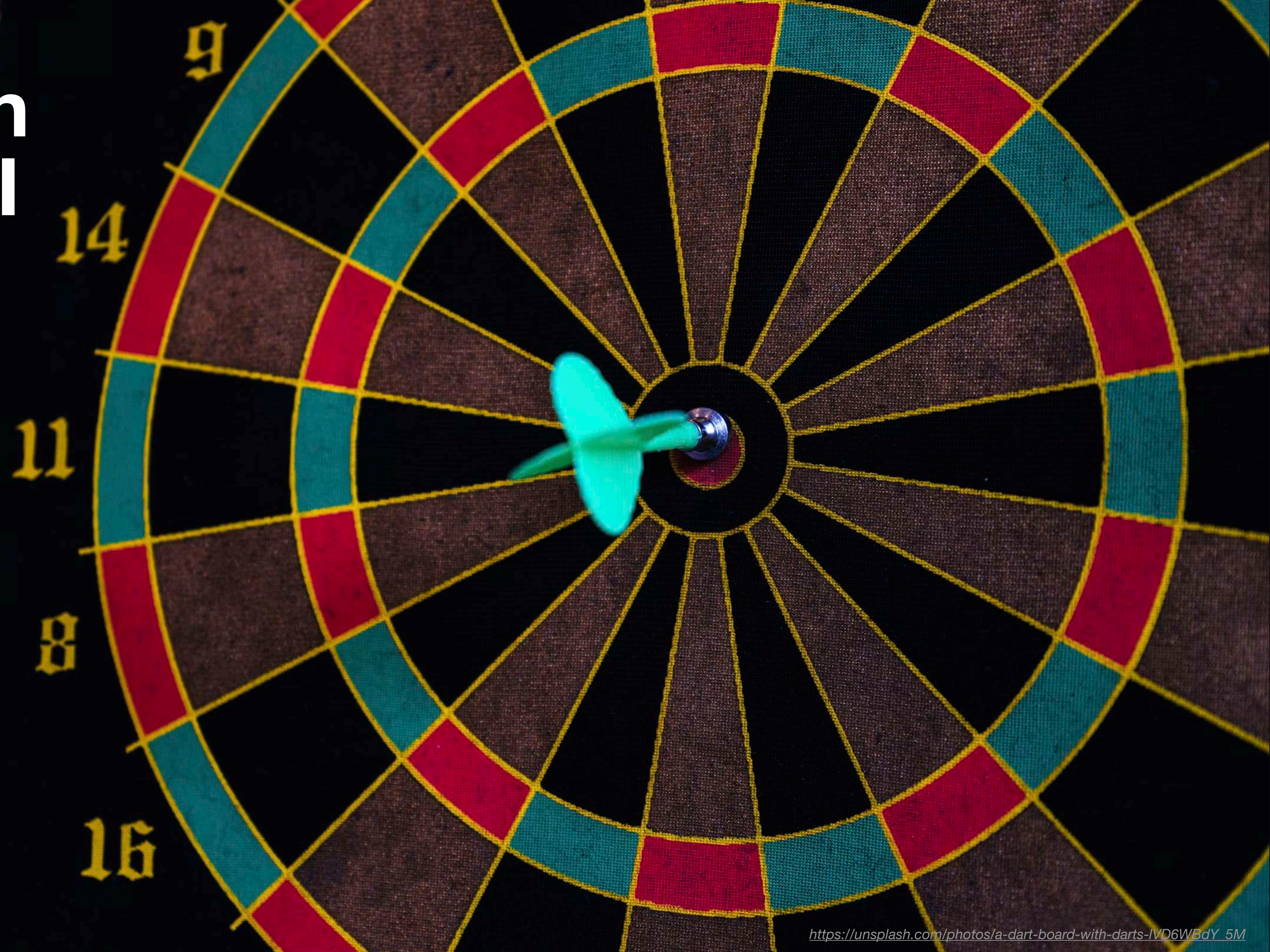
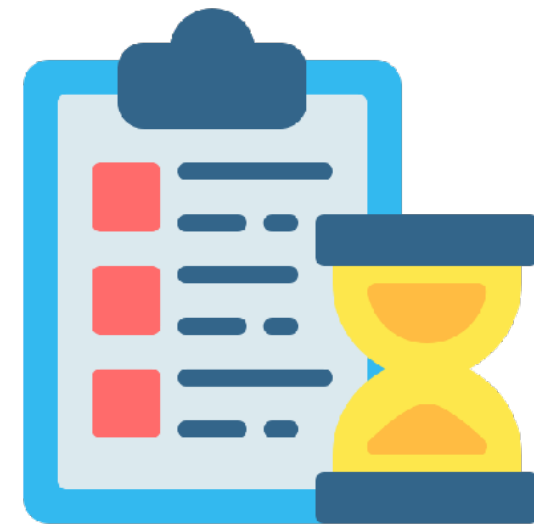


Detección Objetos II

Juan Carlos Arbeláez
jarbel16@eafit.edu.co



Contenido

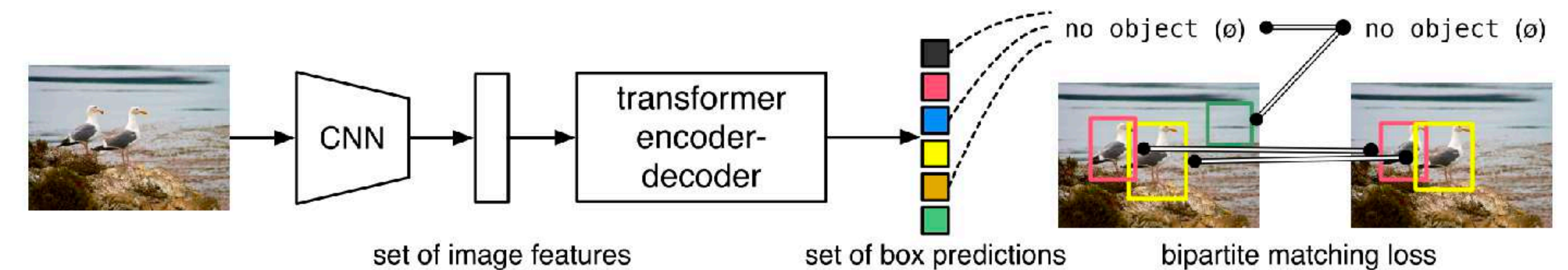


1. Introducción
2. Aplicaciones
3. Retos
4. Algoritmos
 1. R-CNN
 2. Faster R-CNN
 3. Single Shot Detector (SSD)
 4. You Only Look Once (Yolo)
 - 5. DETR**
5. Métricas de evaluación
6. Taller

- Presentada por Facebook en el 2020: combinación CNN + Transformer
- Elimina la necesidad de componentes “manuales” como: “anchor boxes” y “non-maximum suppression”
- End-to-end: realiza las predicciones directas sin necesidad de funciones adicionales

DETR

Detection Transformer

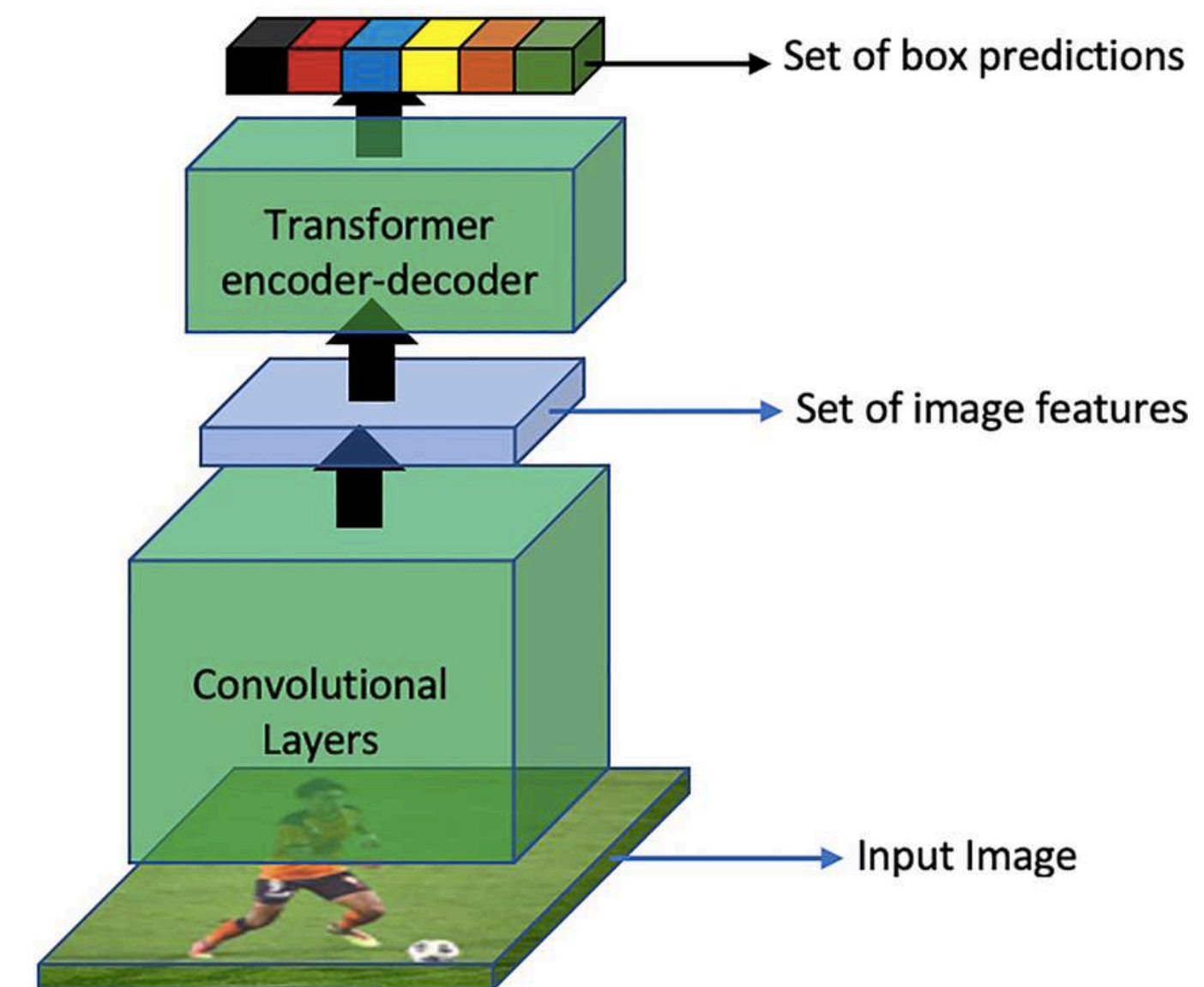


Predicciones

DETR

Dos elementos son necesarios para la predicción directa:

1. **Función de costo** que obligue la coincidencia única entre las predicciones y las etiquetas
2. **Arquitectura** que predice (en una sola pasada) un conjunto de objetos



Función de costo

DETR

DETR predice un set de tamaño fijo de N predicciones (clase y coordenadas) pero N es mucho más grande que los objetos de la imagen

$$\hat{y} = \{\hat{y}_i\}_{i=1}^N \text{ con } \hat{y}_i = (\hat{c}_i, \hat{b}_i)$$

\hat{y} usa un padding de no-objeto (\emptyset)

Se usa “Hungarian algorithm” para encontrar la asignación predicción-etiqueta con menor costo:

$$\hat{\sigma} = \arg \min_{\sigma \in \mathcal{S}_N} \sum_{i=1}^N L_{\text{match}}(y_i, \hat{y}_{\sigma(i)})$$

L_{match} penaliza la **clasificación** + **localización** donde el modelo predice positivos para objetos ($\mathbf{1}_{(c_i \neq \emptyset)}$)

$$L_{\text{match}}(y_i, \hat{y}_{\sigma(i)}) = -\mathbf{1}_{(c_i \neq \emptyset)} \hat{p}_{\sigma(i)}(c_i) + \mathbf{1}_{(c_i \neq \emptyset)} L_{\text{box}}(b_i, \hat{b}_{\sigma(i)})$$

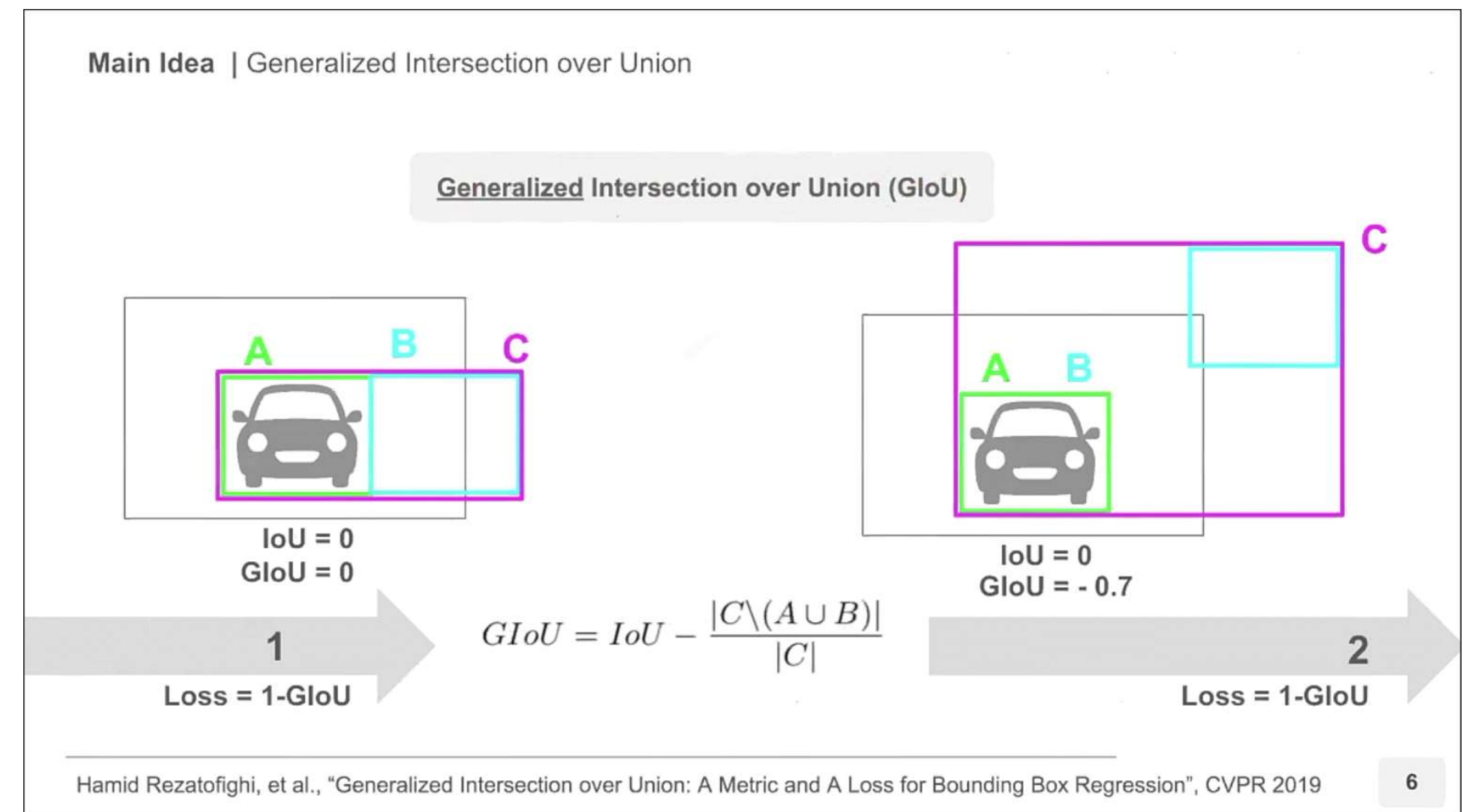
El costo de localización L_{box} es una combinación lineal de L_1 y IoU
Generalizado

$$L_{\text{box}}(b_i, \hat{b}_{\sigma(i)}) = \lambda_{\text{iou}} L_{\text{iou}}(b_i, \hat{b}_{\sigma(i)}) + \lambda_{L1} \|b_i - \hat{b}_{\sigma(i)}\|_1$$

GIOU

Generalized Intersection Over Union

- IOU tiene un problema principal: para todo $|A \cap B| = 0 \rightarrow IoU(A, B) = 0$ independiente de los lejos que se encuentren
- Usa el Convex-Hull C es la geometría más pequeña que envuelve A y B tener en cuenta la distancia cuando no se interceptan



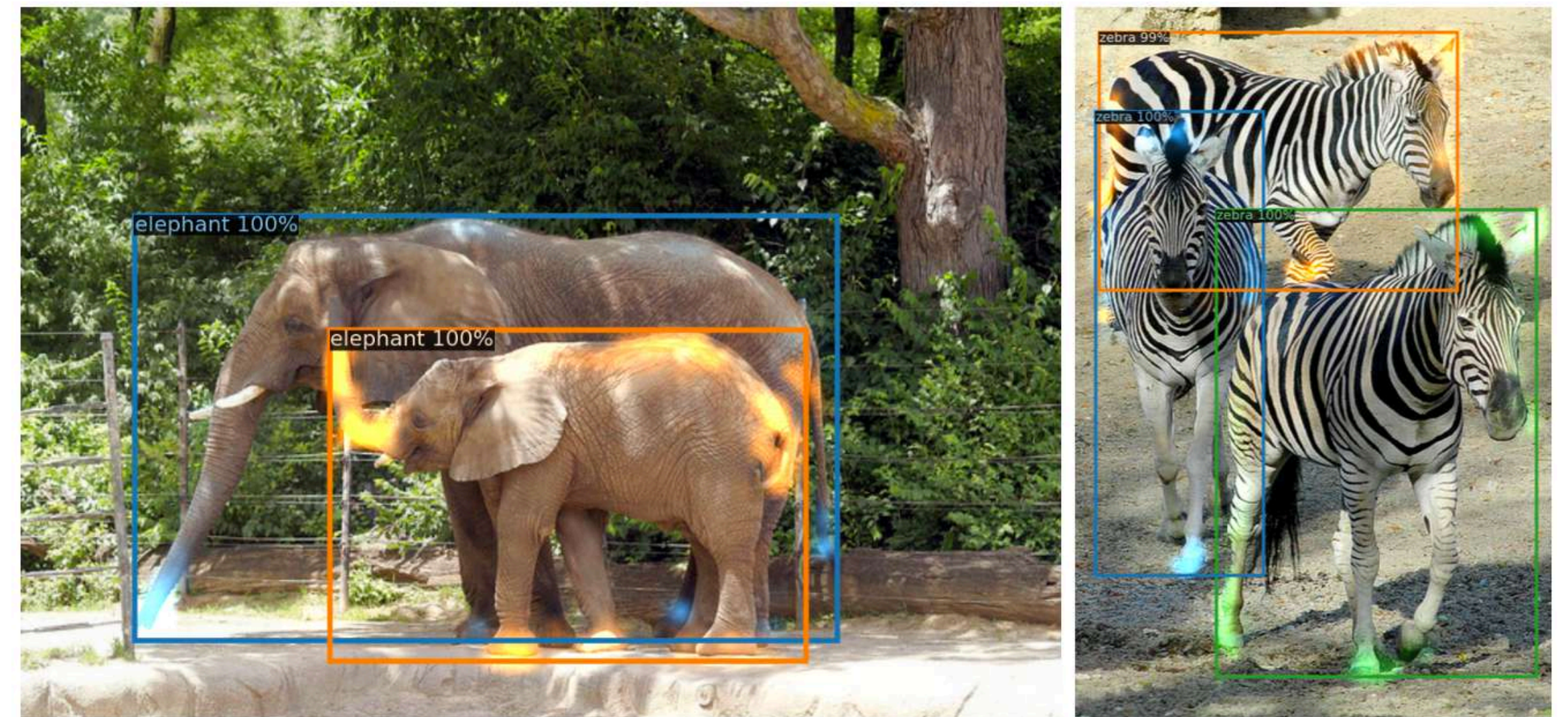
Hungarian loss

DETR

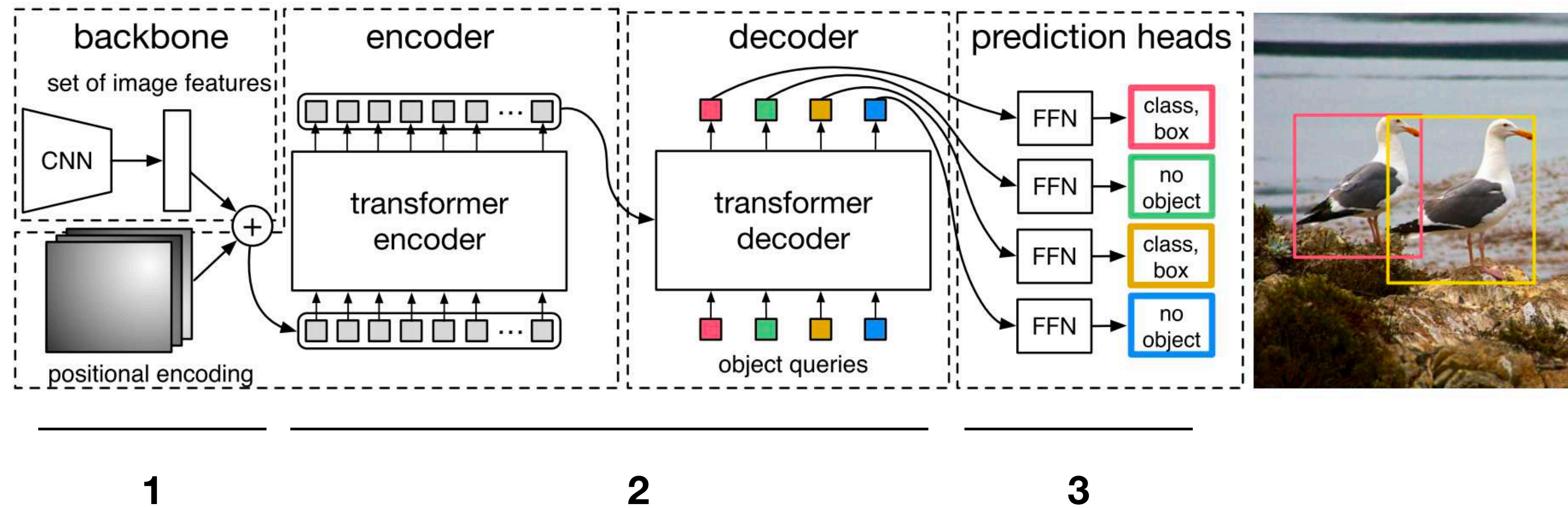
Función de costo a ser minimizada:

$$L_{hungarian}(y, \hat{y}) = \sum_{i=1}^N -\log \hat{p}_{\hat{\sigma}(i)}(c_i) + \mathbf{1}_{(c_i \neq \emptyset)} L_{box}(b_i, \hat{b}_{\hat{\sigma}(i)})$$

Clasificación + Detección de la mejor combinación encontrada con el algoritmo húngaro



Arquitectura DETR

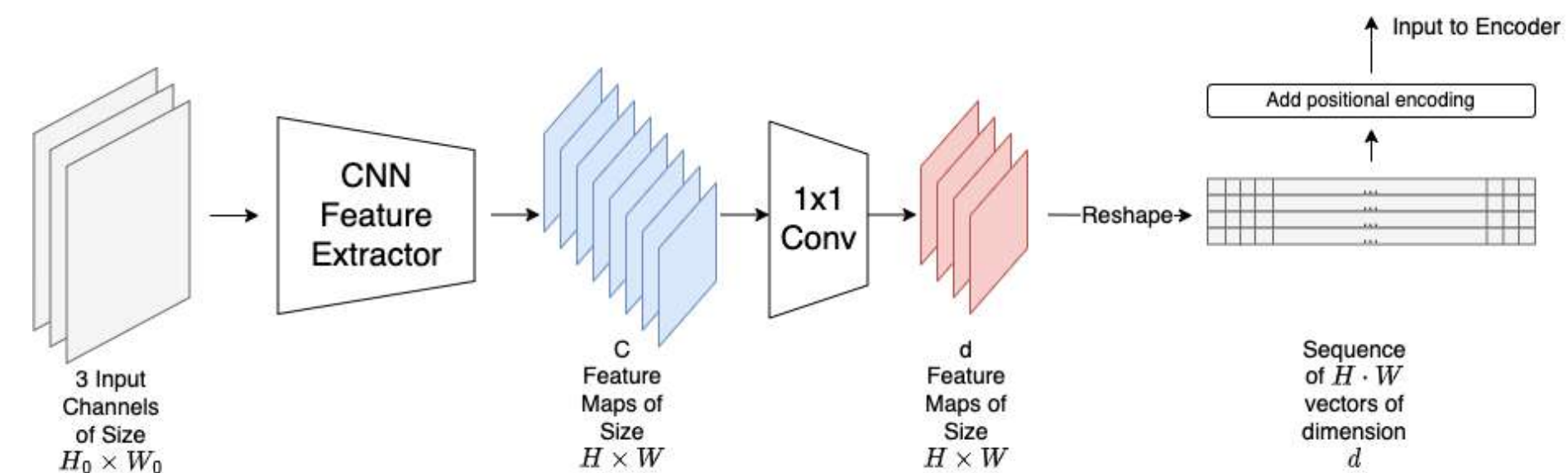


- Extractor de características CNN convencional (ResNet50)
- Mapa de características es transformado para pasar de un formato de secuencias

$$f \in \mathbb{R}^{C \times H \times W} \rightarrow z \in \mathbb{R}^{D \times HW}$$

1. Backbone

Arquitectura



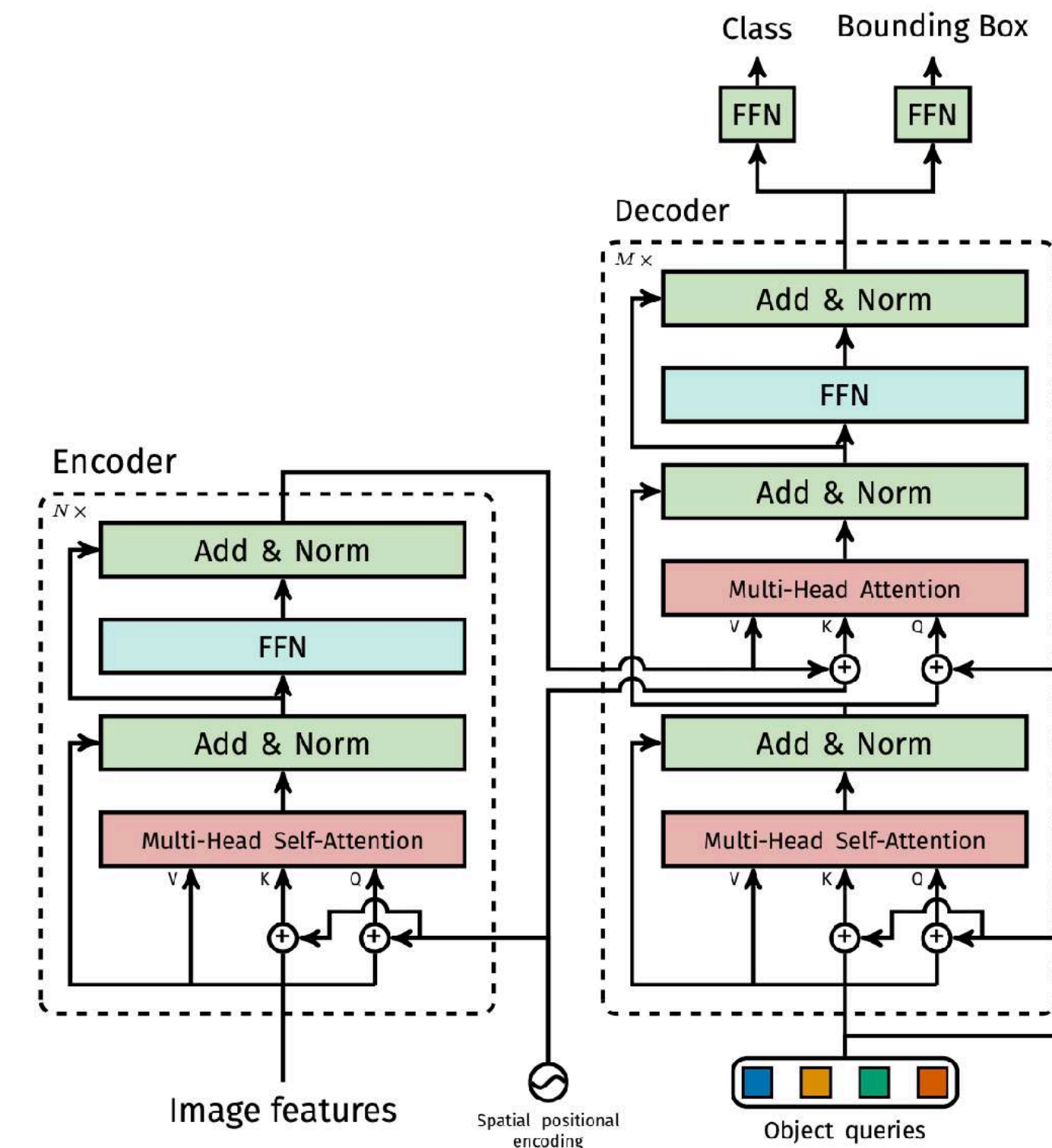
2. Transformer

Arquitectura

Transformer DETR consiste de multiples bloques de encoder y decoder (estándar con mínimas variaciones)

Encoder: “positional encodings” se agregan a K y Q . Ya que V representa características de la imagen

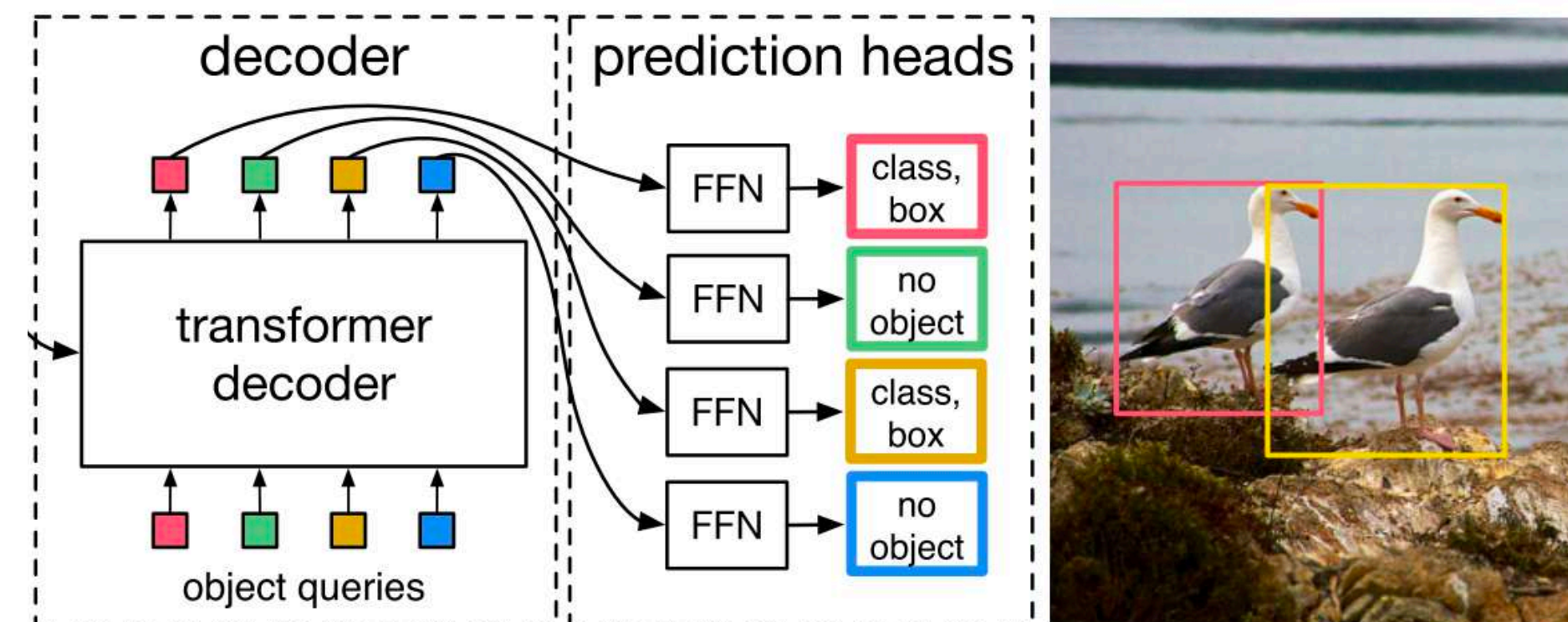
Decoder: Recibe N “object queries” (“learnt positional encodings”) que son transformados en un “output embedding” usado por la cabeza de predicción



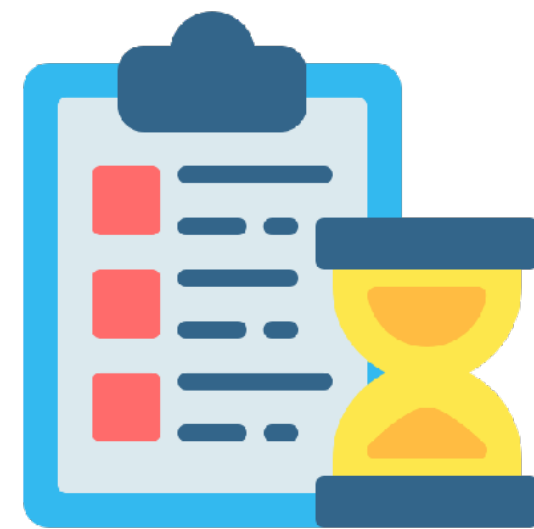
FC de 3 capas con activación ReLU que predice:

1. Las clases (softmax) más una clase especial de “no objeto” (\emptyset)
2. El centro normalizado (sigmoid) de las coordenadas,
3. El ancho y largo de los rectángulos

3. Cabeza de predicción Arquitectura



Contenido

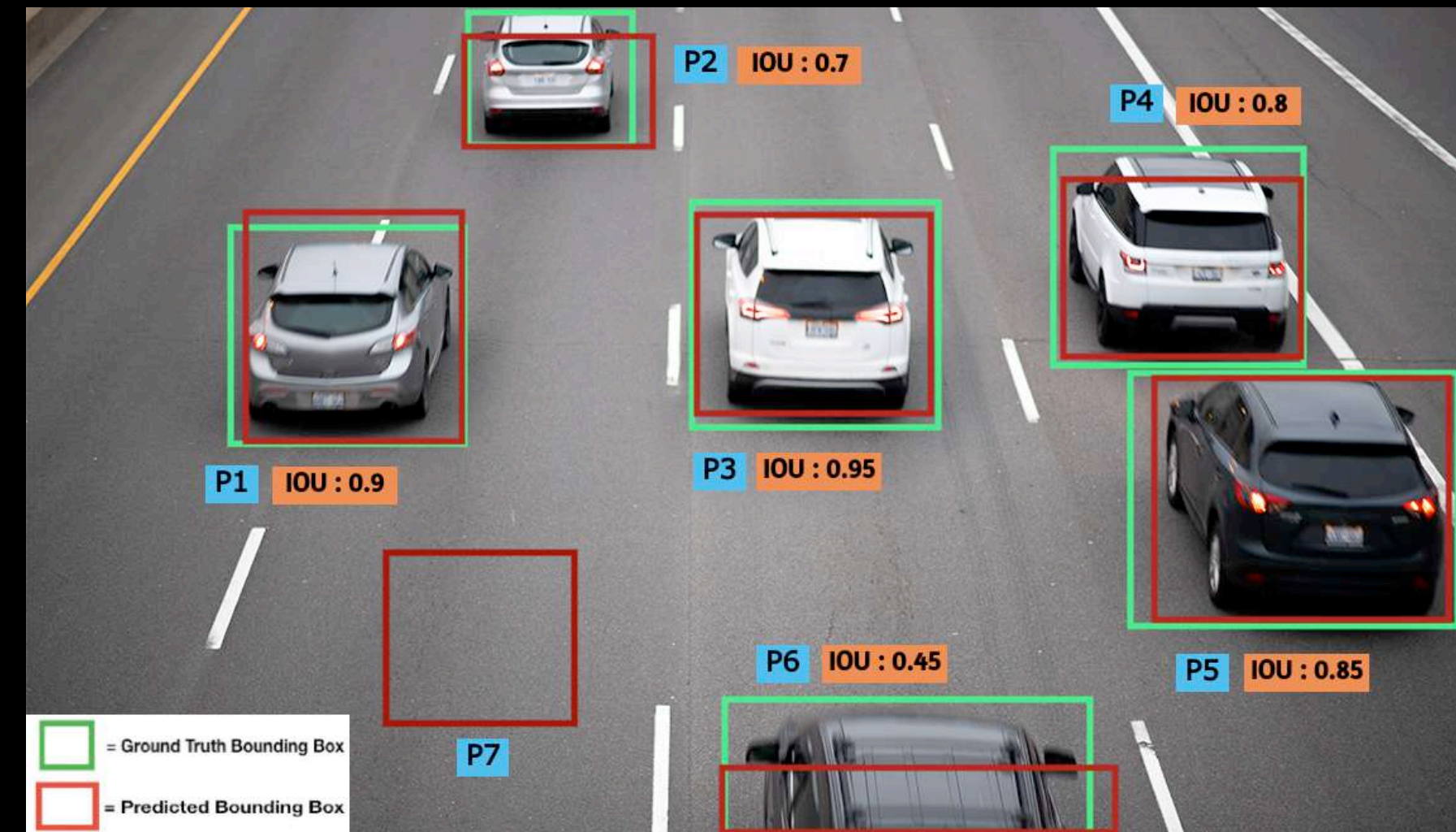


1. Introducción
2. Algoritmos
- 3. Métricas de evaluación**
4. Aplicaciones
5. Retos
6. Taller

mAP

Mean Average Precision

- El **mAP** es el promedio de **AP** de cada clase
- **AP** es el AUC de la curva precisión-sensibilidad



AP

Average Precision

AP No es el promedio de la precisión de las distintas clases

AP se calcula con la ayuda de otras métricas como IoU, matriz de confusión, precisión y la sensibilidad

$$\mathbf{mAP} = \frac{1}{n} \sum_{h=1}^n \mathbf{AP}_k$$

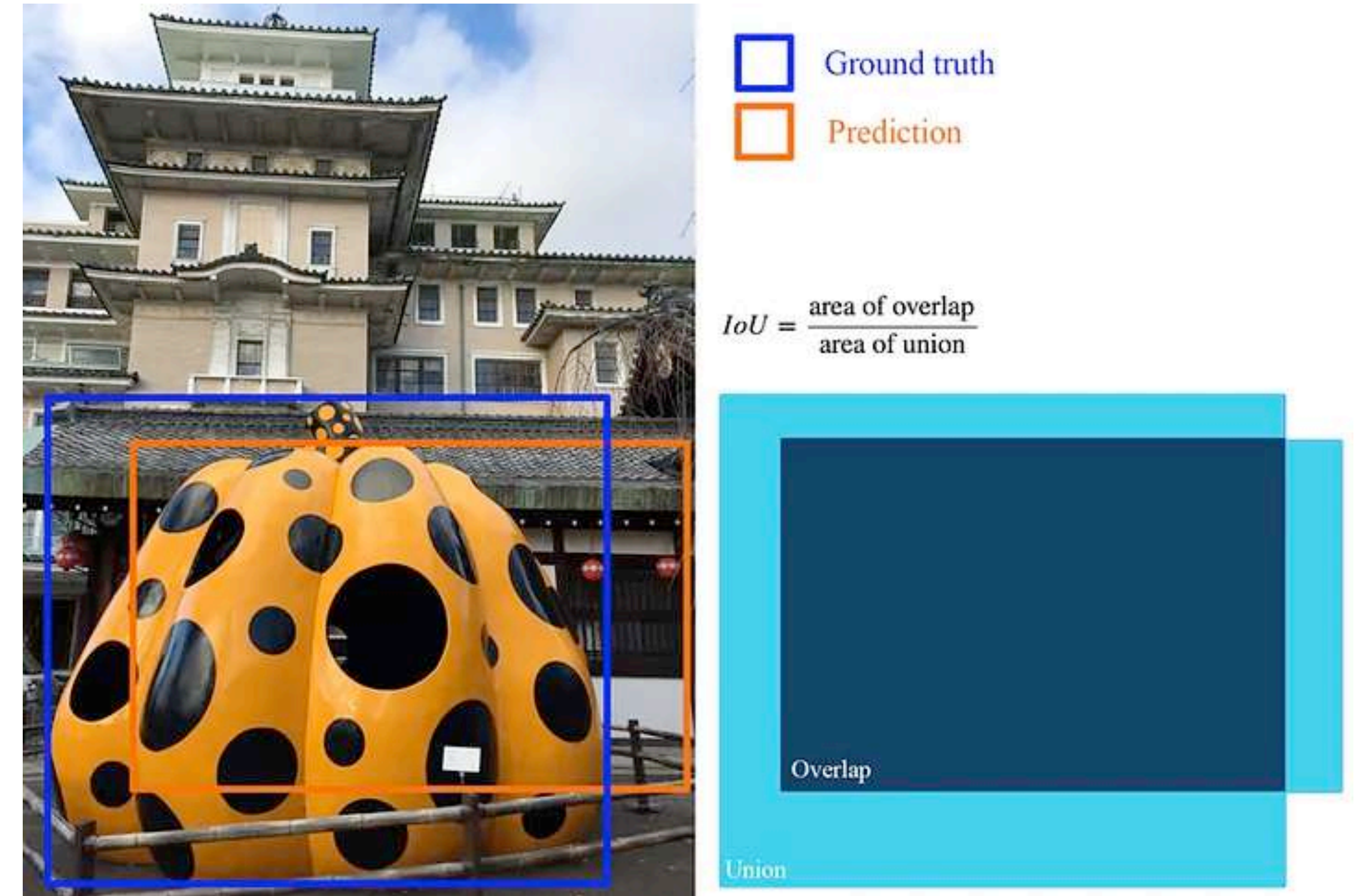
AP_k: Average Precision de la clase k

n: EL número de clases

IOU

Intersection Over Union

- Lo utilizamos para medir cuánto se solapa el rectángulo predicho con el real
- **Predefinimos un umbral de IoU** por ejemplo, 0,5 para clasificar si la predicción es un verdadero positivo o un falso positivo

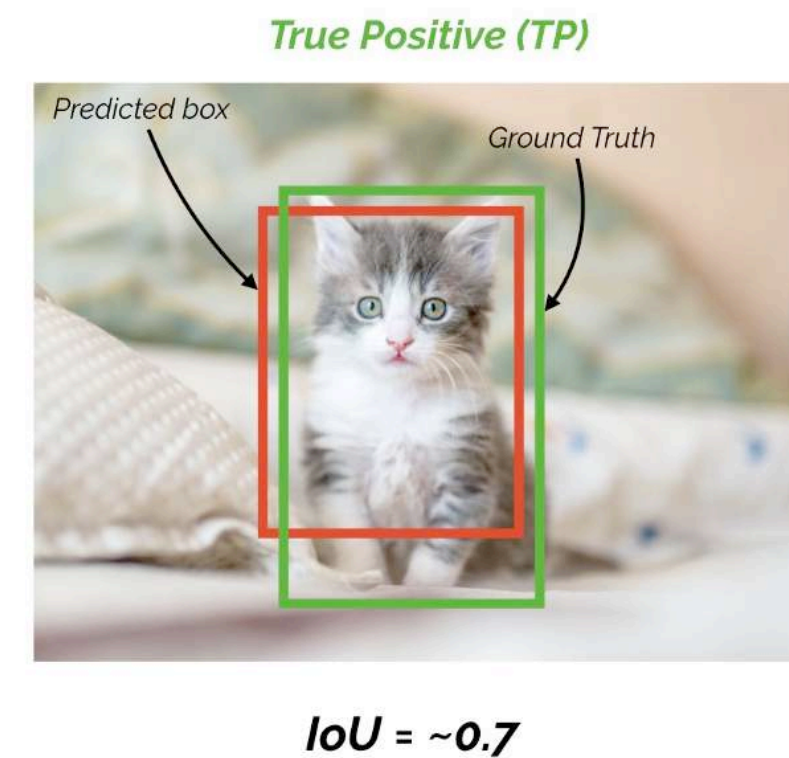
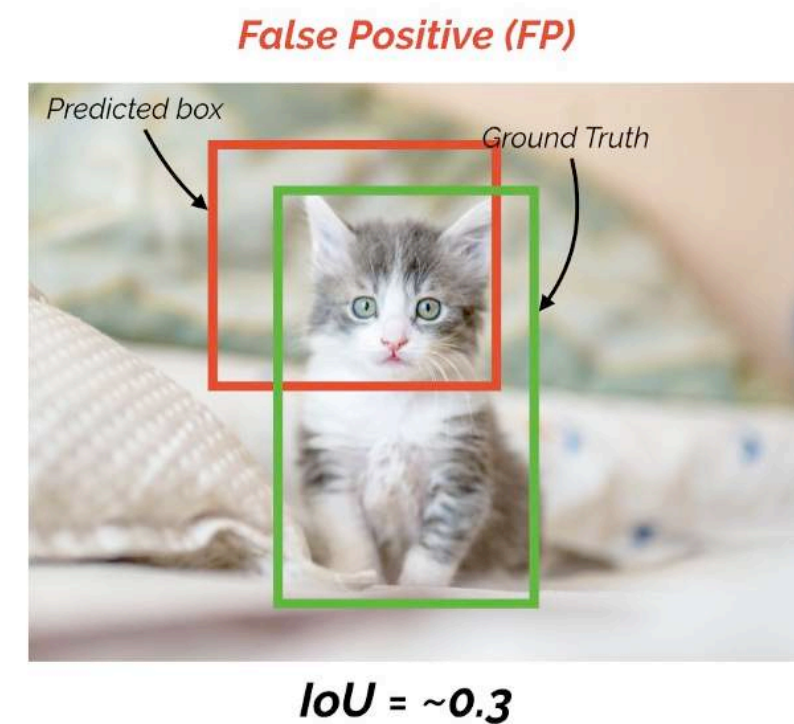


AP

Average Precision

@IOU = 0.5

Falso positivo: El modelo predijo que existía un rectángulo en una posición determinada pero se equivocó



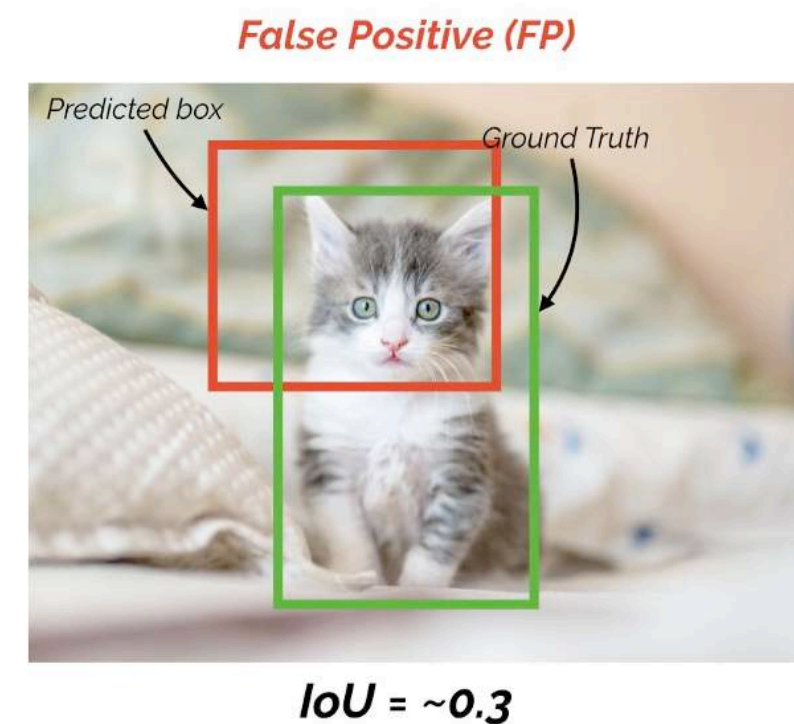
Verdadero positivo: El modelo predijo que existía un rectángulo en una posición y acertó

AP

Average Precision

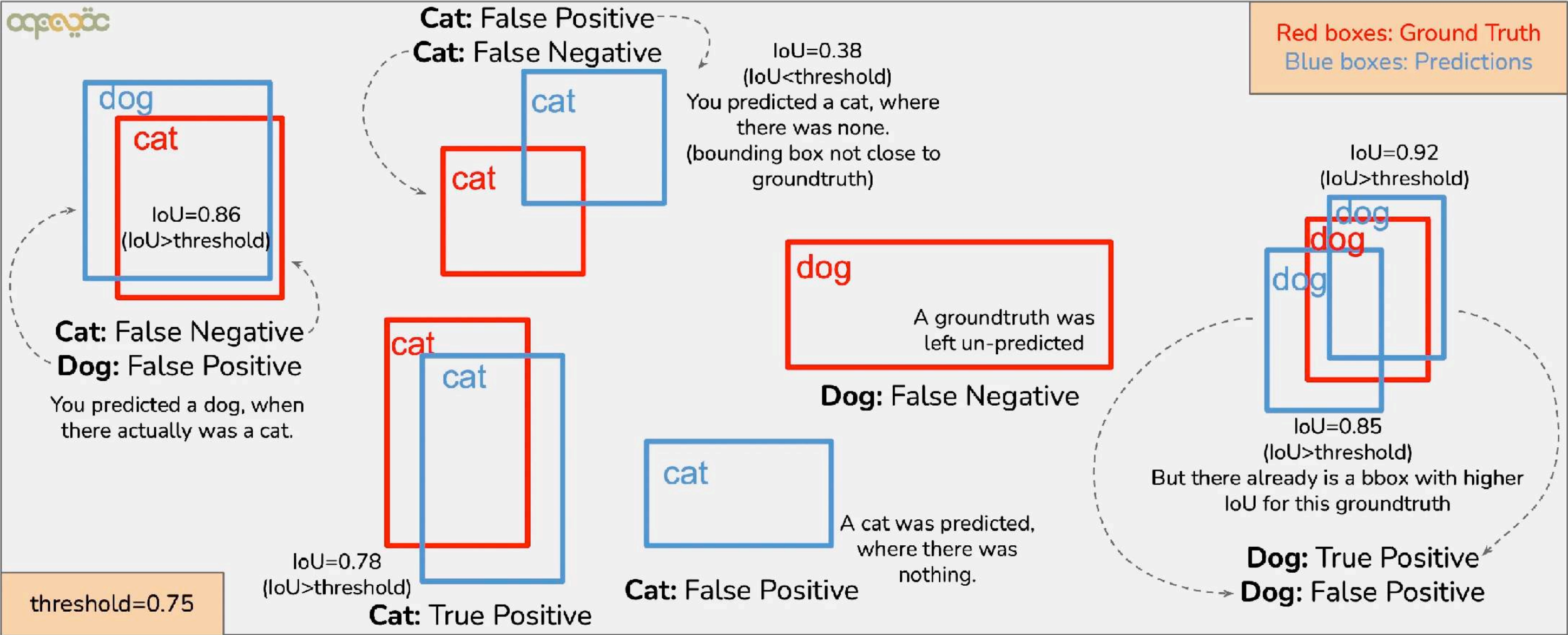
@IOU = 0.5

Falso negativo: El modelo no predijo un rectángulo en una posición determinada y se equivocó



Verdadero negativo: Corresponde al fondo, la zona sin rectángulos, y **no se utiliza para calcular las métricas finales**

Object Detection and Localization - IoU, True Positive, False Positive, False Negative



$$IoU = \frac{\text{Area of Intersection}}{\text{Area of Union}}$$

@_aqeelanwar
aqeelanwarmalik

Threshold	Class	# GroundTruth	# predictions	TP	FP	FN	Precision	Recall
0.75	Cat	3	3	1	2	2	1/3	1/3
	Dog	2	3	1	2	1	1/3	1/2
0.35	Cat	3	3	2	1	1	2/3	2/3
	Dog	2	3	1	2	1	1/3	1/2

Sensibilidad y precisión

Sensibilidad: ¿Qué proporción de positivos reales se identificó correctamente?

$$\frac{TP}{TP + FN}$$

Precisión: ¿Qué proporción de predicciones positivas fue realmente correcta?

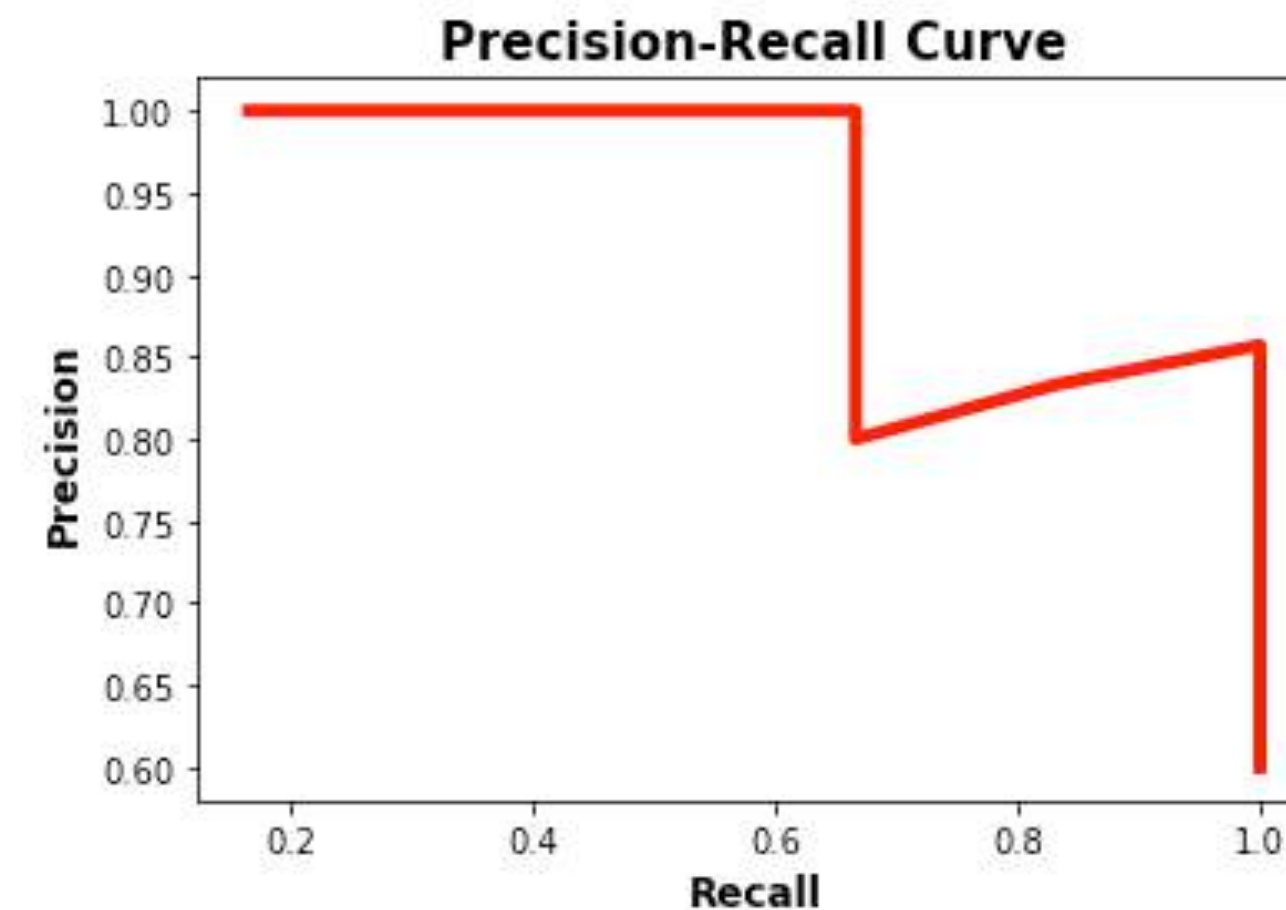
$$\frac{TP}{TP + FP}$$

		True Class	
		Positive	Negative
Predicted Class	Positive	TP	FP
	Negative	FN	TN

Matriz de confusión para la clasificación binaria

Precision-Recall Curve

- Para evaluar un modelo, hay que examinar tanto la precisión como la sensibilidad
- La precisión y la recuperación suelen estar en “tensión”
- Una curva de precisión-sensibilidad representa el valor de la precisión frente a la sensibilidad para distintos umbrales de confianza.



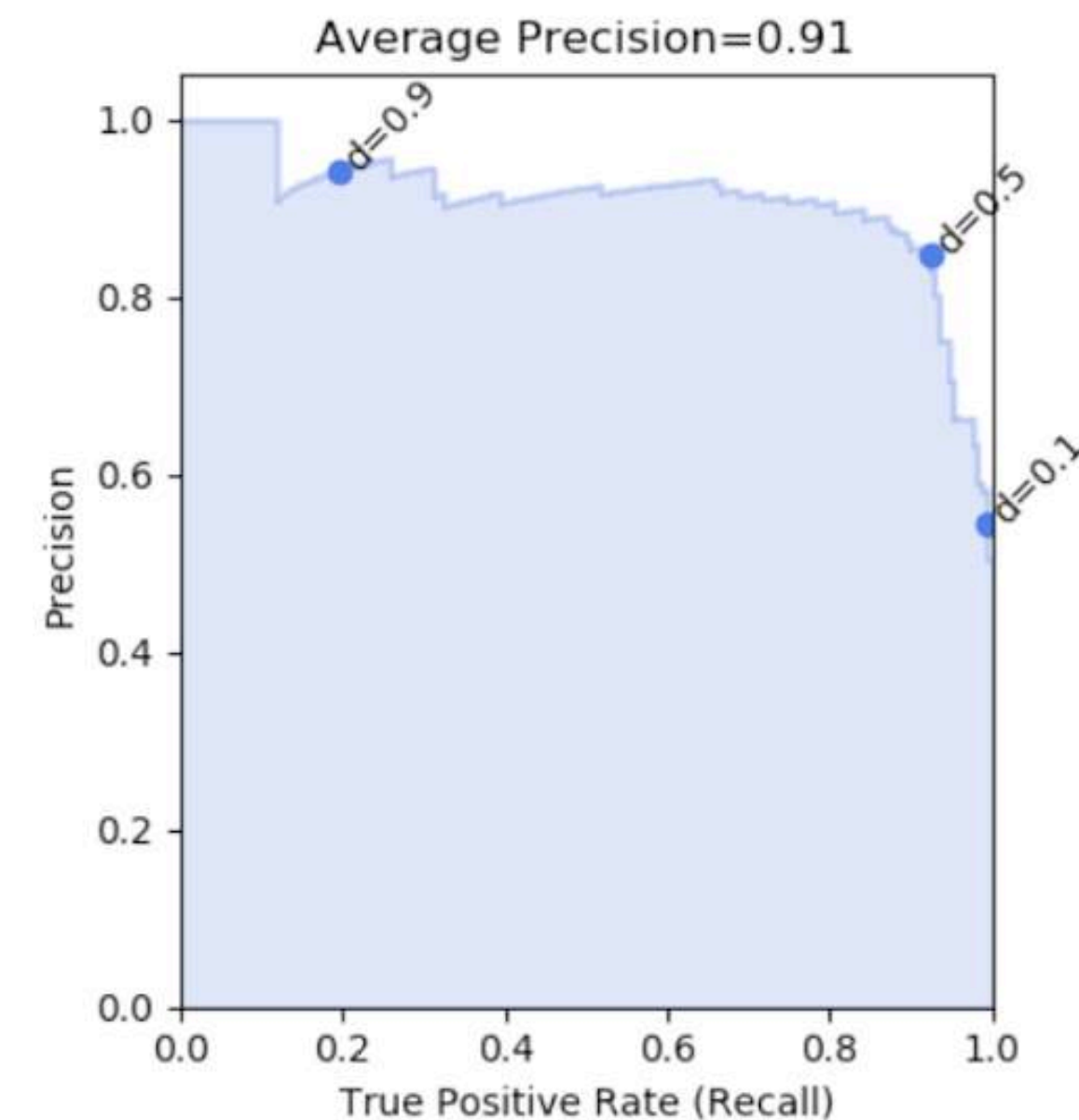
Average Precision

AP: Área bajo la curva PR de una clase

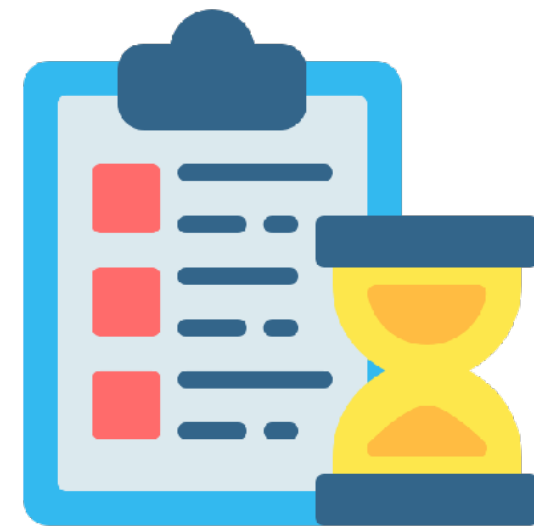
$$AP = \int_{r=0}^1 p(r)$$

mAP: Promedio de AP de las clases

$$mAP = \frac{1}{n} \sum_{h=1}^n AP_k$$



Contenido



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

Retail

- Conteo de personas
- Información sobre cómo pasan el tiempo los clientes y su afluencia
- Comprender la interacción y la experiencia del cliente
- Optimizar la distribución de la tienda y hacer más eficientes las operaciones
- Reducir el tiempo de espera en las tiendas minoristas



Robótica

- Los robots autónomos dependen de la detección de objetos para reconocer peatones, señales de tráfico, otros vehículos, etc.
- Por ejemplo, el *autopilot* de *Tesla* utiliza en gran medida la detección de objetos para percibir el entorno, como vehículos que circulan en sentido contrario u obstáculos.



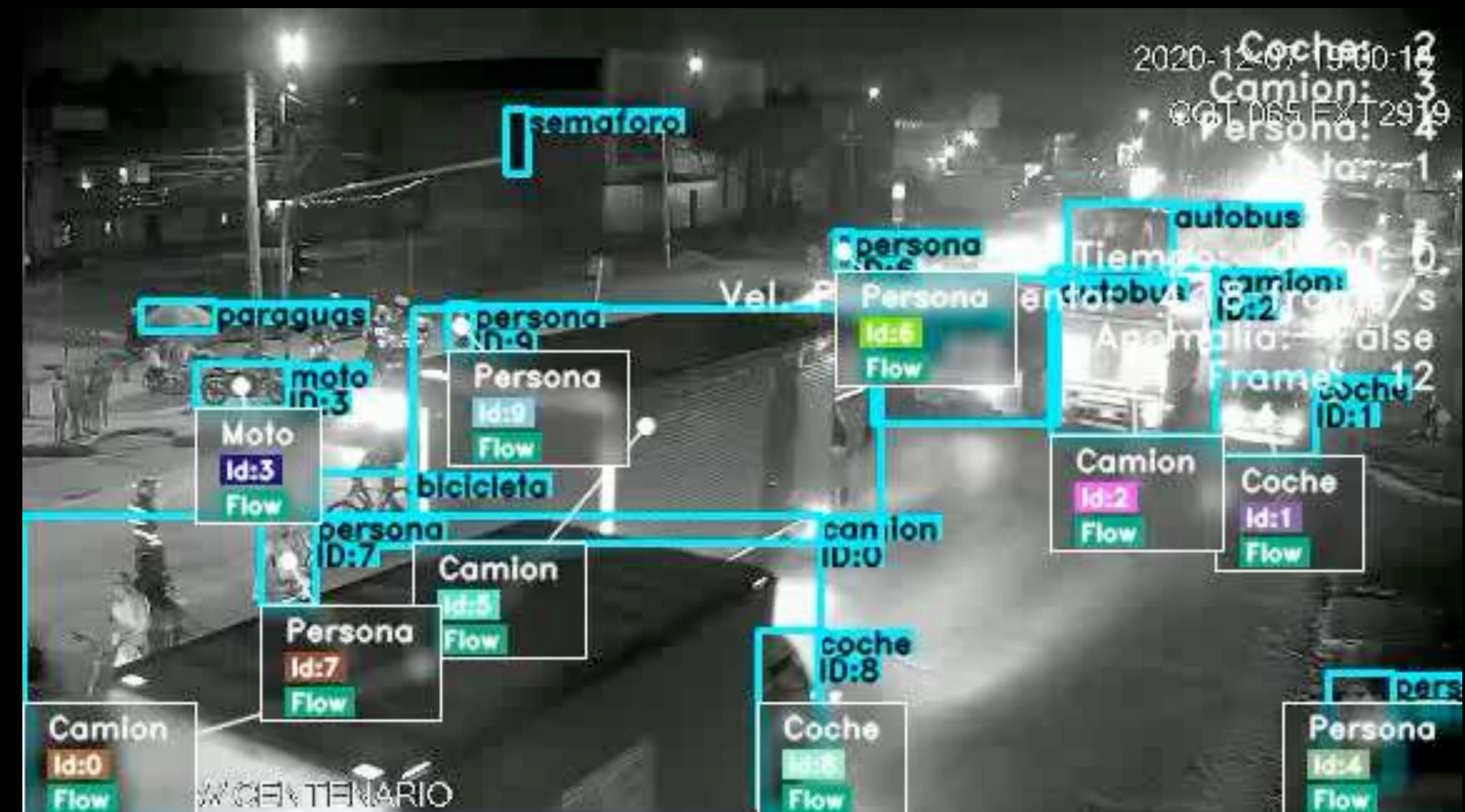


Video vigilancia

- Seguimiento de personas a la vez en tiempo real (punto de partida para tracking)
- Rendimiento y seguridad de los trabajadores, tráfico peatonal
- Detección en cámaras de CCTV de personas en zonas peligrosas, objetos en bordes, personas en zonas no autorizadas, ...

Monitoreo de tráfico

- Análisis del tráfico (punto de partida para tracking)
- Tipos de vehículos (motos, peatones, taxis, ...), afluencia y comportamiento en las vías (giros, paros, velocidades)
- Accidentes, detención en zonas peligrosas



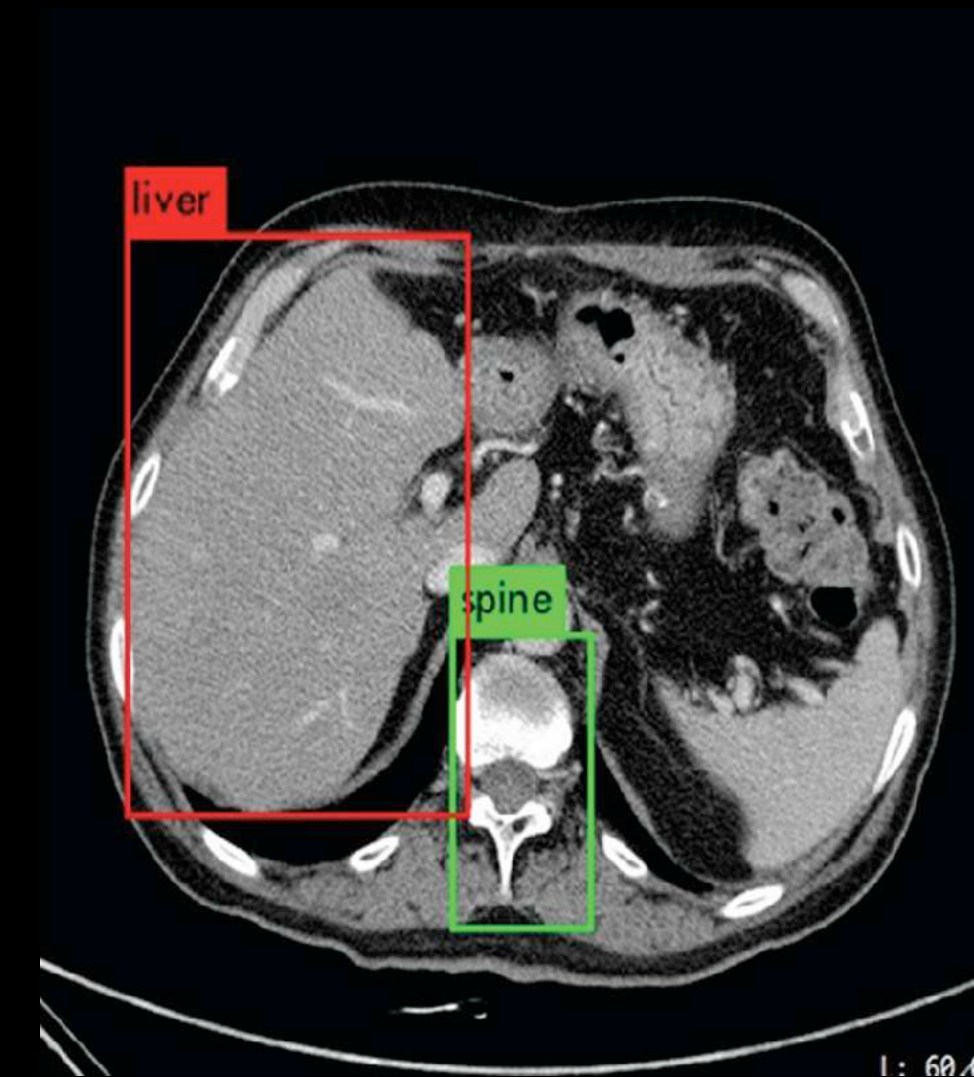
Agricultura

- Conteo de células, frutas, animales, flores, hojas
- Monitoreo de animales, comportamiento basado en movimientos
- Evaluación de la calidad de los productos agrícolas
- Agricultura de precisión
- Detección de enfermedades en plantas y animales

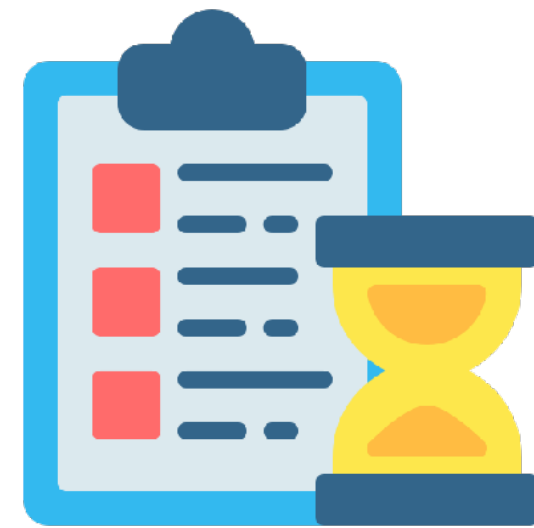


Cuidado de la salud

- “Medical object detection” es la tarea de identificar objetos médicos en una imagen
- Los diagnósticos médicos se basan en gran medida en el estudio de imágenes, escáneres y fotografías
- La detección de objetos en tomografías y resonancias magnéticas es una herramienta para el diagnóstico de enfermedades



Contenido



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

Retos

- Doble prioridad, clasificación y localización de objetos: La detección de objetos tiene dos objetivos, la clasificación y la localización
- Se utiliza una función de pérdida multitarea que penaliza tanto los errores de clasificación como los de localización

Función de costo (weighted sum λ)

Clasificación impone una pérdida logarítmica

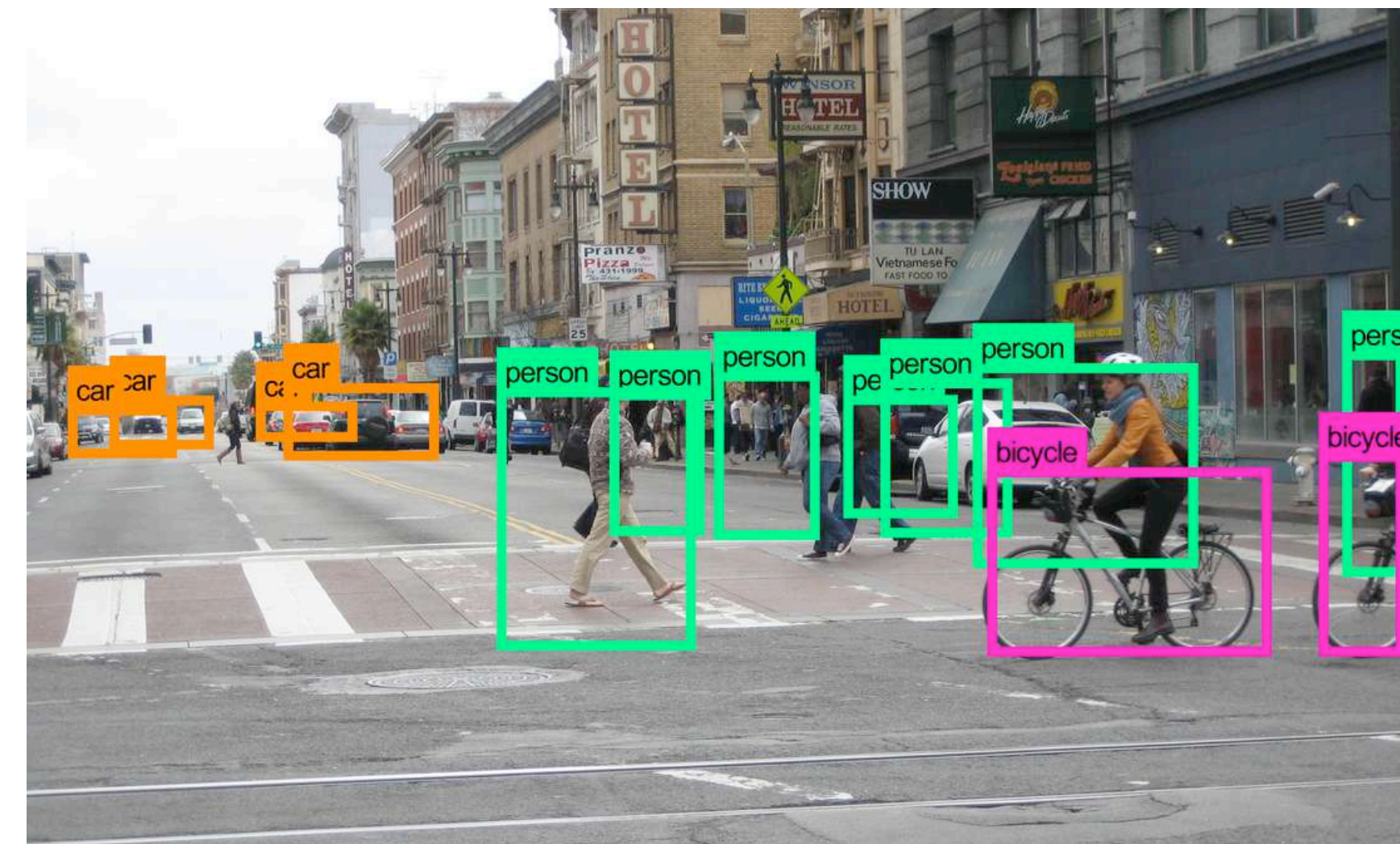
Localización es una pérdida *Smooth L₁*-loss para los cuatro componentes que definen el rectángulo

$$\ell(p, u, t^u, v) = \ell_c(pu) + \lambda_{[u \geq 1]} \ell_l(t^u, v)$$

Retos

Detección en tiempo real

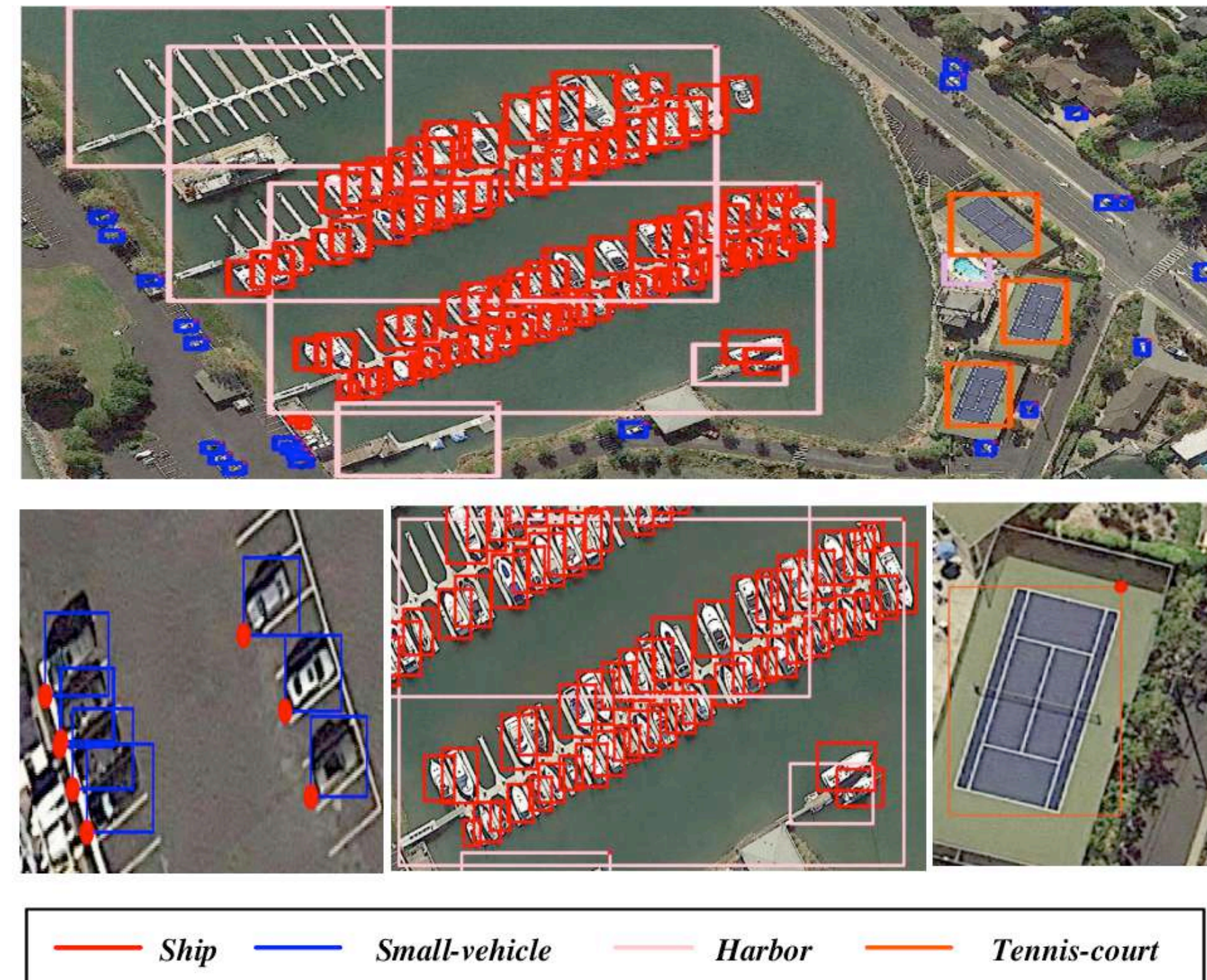
- Muchas de las aplicaciones vistas requieren análisis en tiempo real
- Los vídeos suelen grabarse a un mínimo de 24 fps
- los algoritmos actuales de detección de objetos intentan encontrar un equilibrio entre velocidad y precisión



Retos

Múltiples escalas y relaciones de aspecto

- Los elementos de interés pueden aparecer en una amplia gama tamaños y relación entre ancho y alto (*Aspect ratio*)
- Las imágenes con gran diferencia de tamaños son un reto hoy en día



Retos

Datos limitados

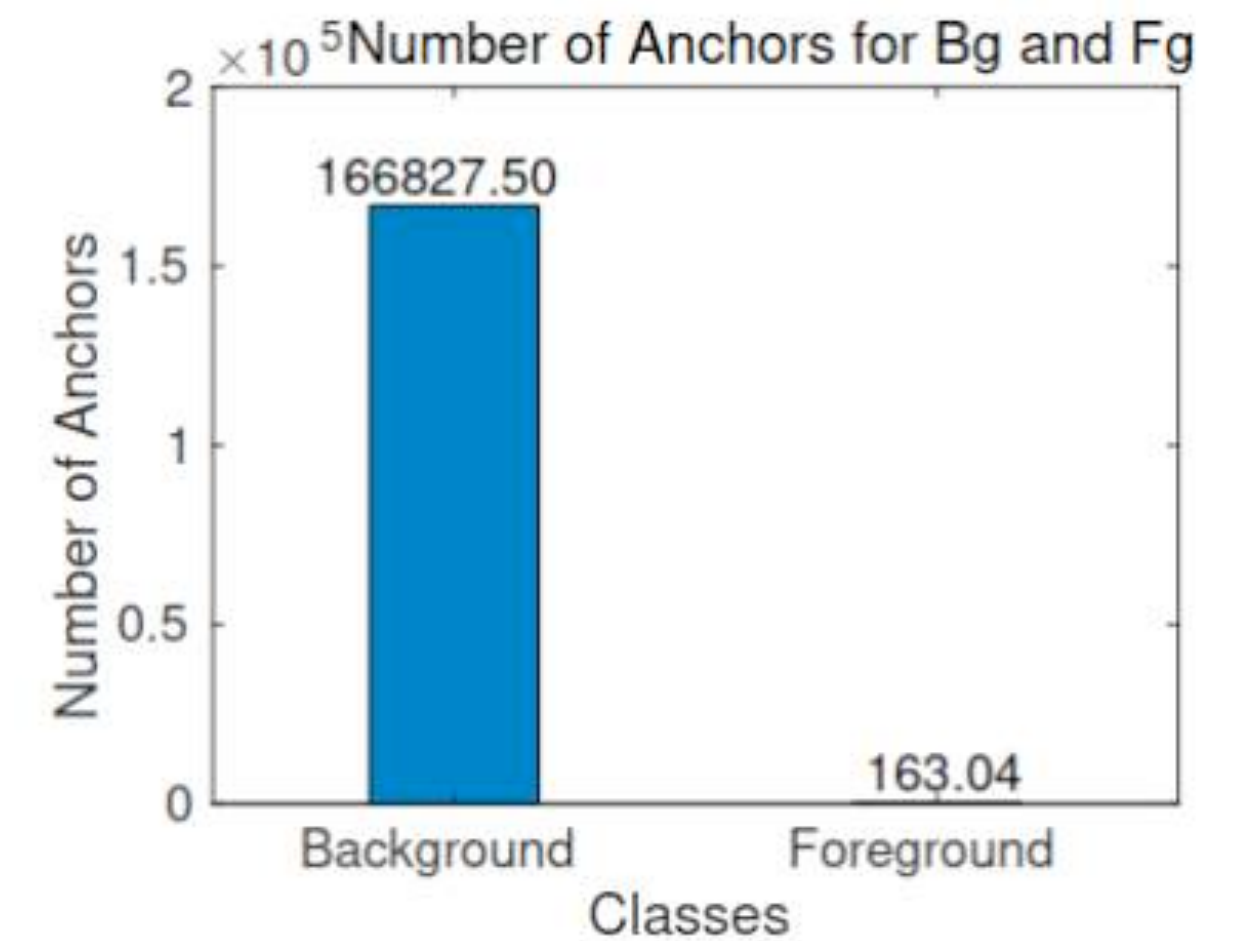
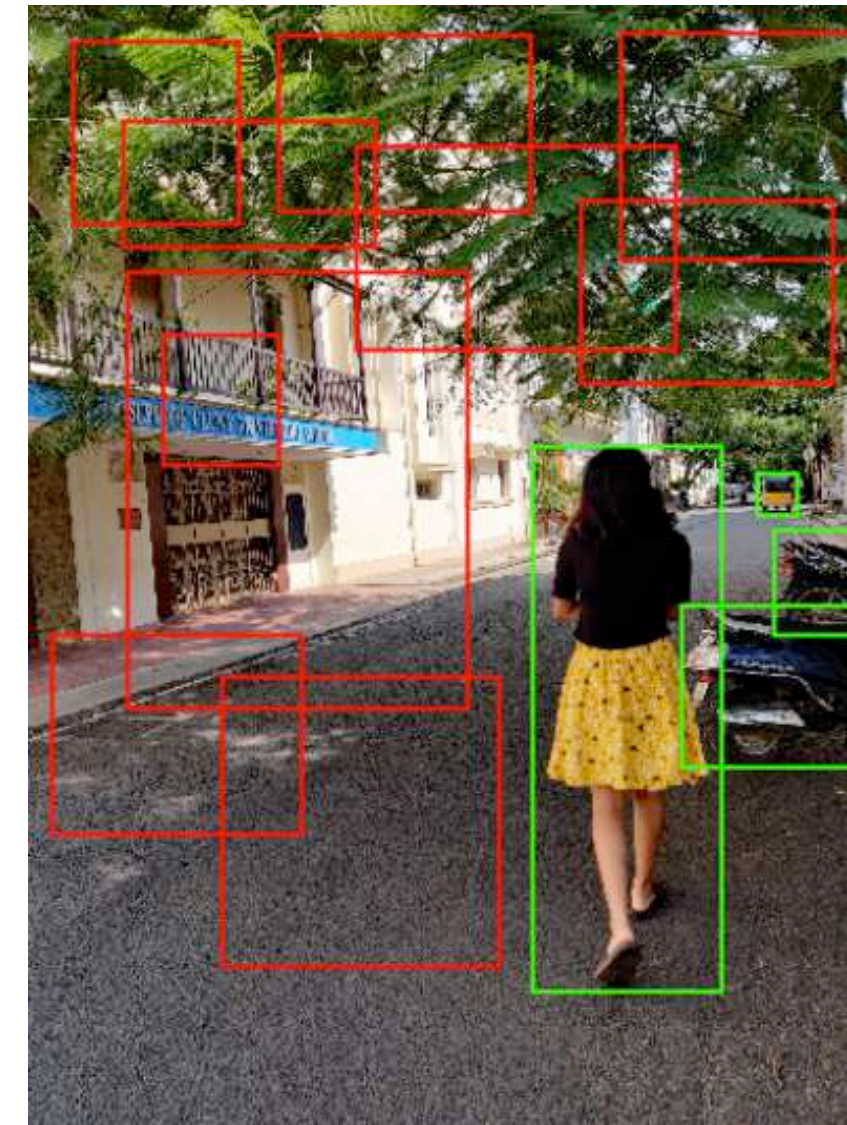
- La limitada cantidad de datos anotados que se dispone para la detección de objetos es un obstáculo importante
- Se requiere una alta cantidad de datos para la correcta localización de los objetos



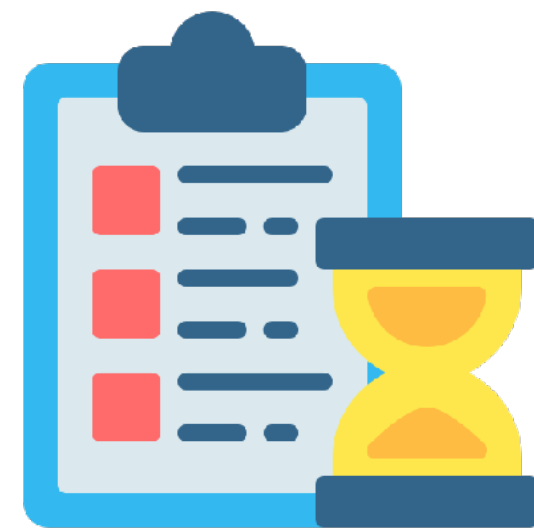
Retos

Class imbalance

En una fotografía típica lo mas probable es que contenga unos pocos objetos principales y que el resto de la imagen esté llena de fondo



Contenido



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

Ejercicio

Realizar un modelo para la localización de objetos
(Regresión de las coordenadas de los Bounding Box)

Pasos:

1. Descargar dataset y etiquetas en imágenes
2. Crear matriz de etiquetas (y)
3. Modelo: Extractor de características
(Convoluciones) + Cabeza (regresión de 4 variables)
4. Entrenar con función de pérdida MSE (Mean Squared Error)

(x_{top} , y_{top})



(x_{bottom} , y_{bottom})

Ejercicio

1. Descargar dataset de Teams
2. Graficar etiquetas en imágenes: Las etiquetas están dadas por dos puntos

Inicial (x_{top} , y_{top})

Final (x_{bottom} , y_{bottom})

(x_{top} , y_{top})



(x_{bottom} , y_{bottom})

Ejercicio

Generar Etiquetas:

Normalizar por el tamaño de la imagen, de esta forma el modelo predice valores entre 0-1 por cada coordenada

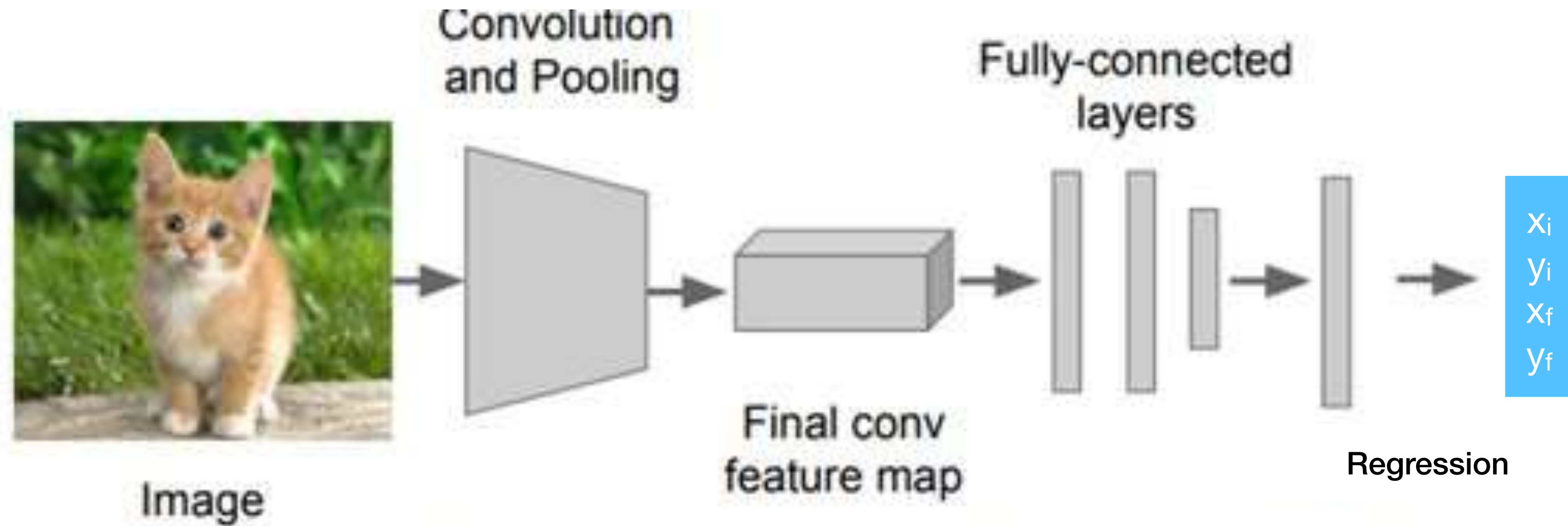
$$\left[\frac{x_{...}}{w}, \frac{y_{...}}{h} \right]$$

(x_top, y_top)



(x_bottom, y_bottom)

Ejercicio



Lecturas

2022 - CenterNet++ for Object Detection

2022 - BEIT: Bert pre-training of image transformers

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