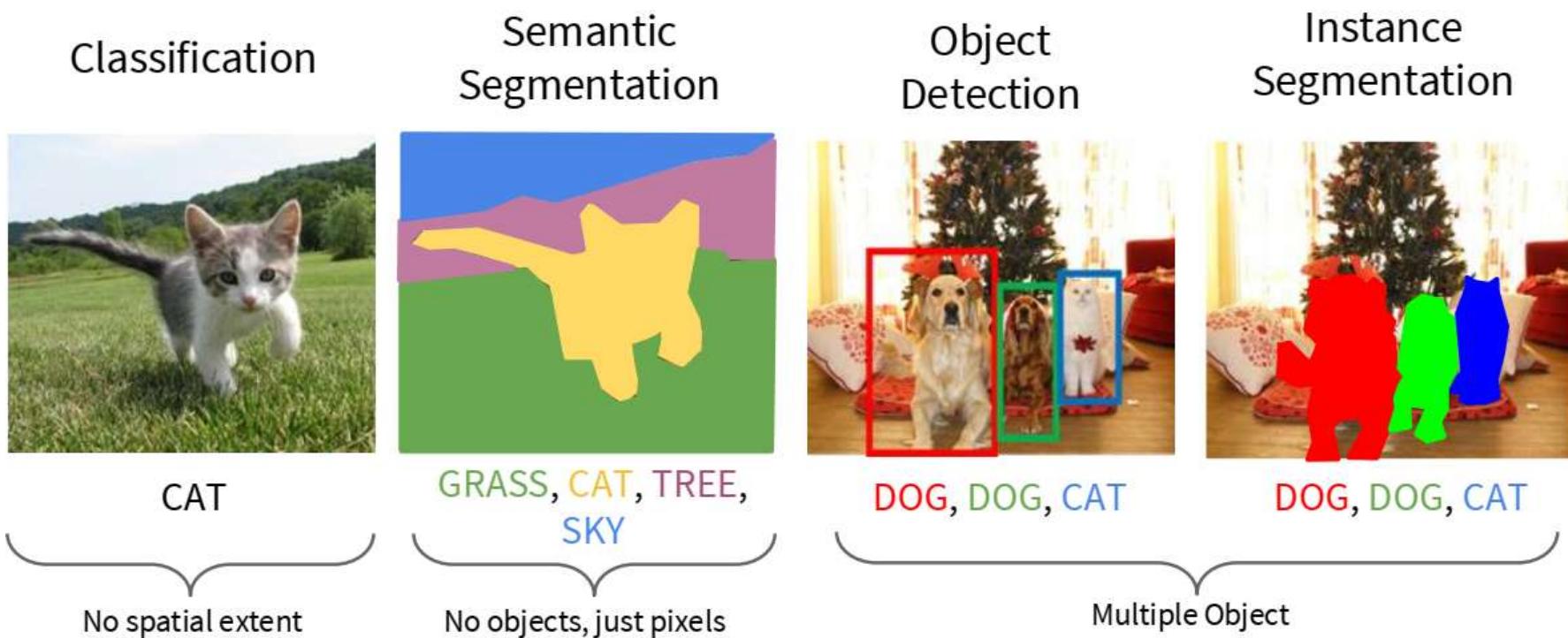


Visión Computacional

Ivan Sipiran

Tareas de Computer Vision



Clasificación

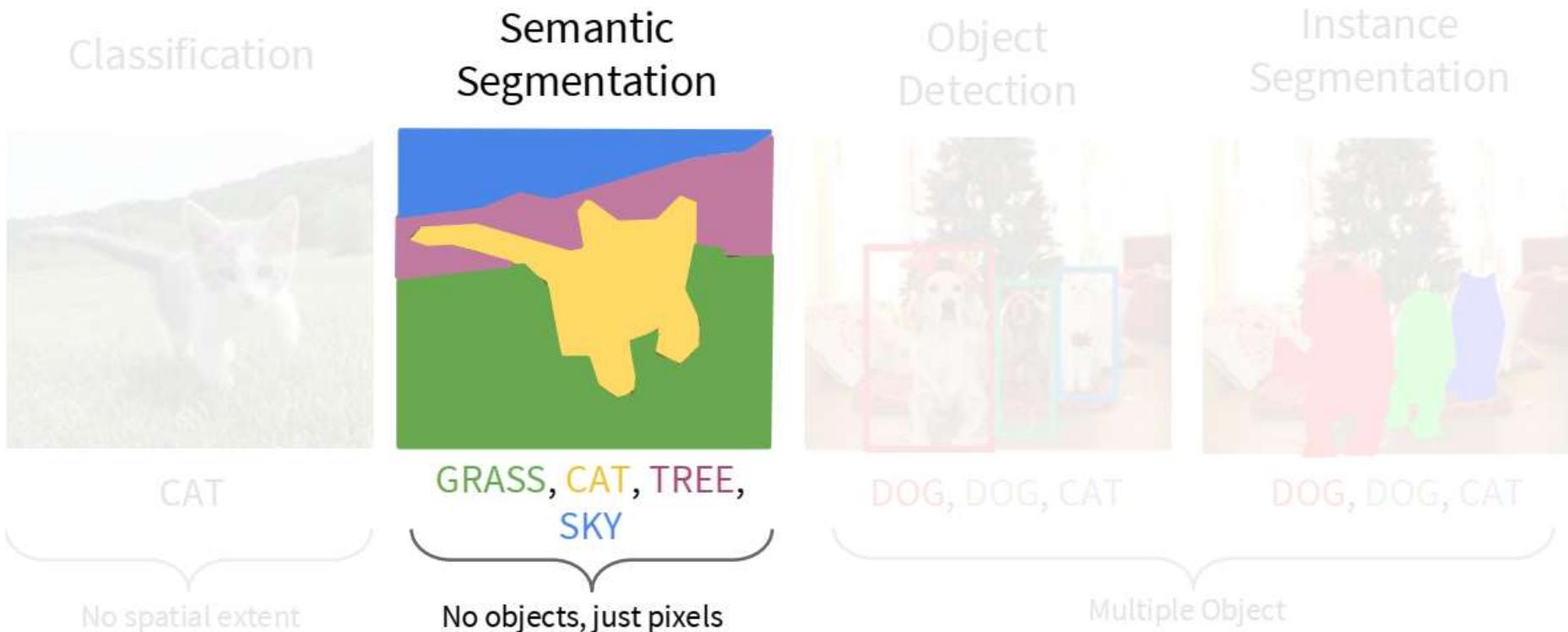
Tarea fundamental en CV



cat

De entre un conjunto posible de etiquetas

Segmentación semántica

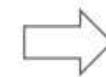


Segmentación semántica: Problema



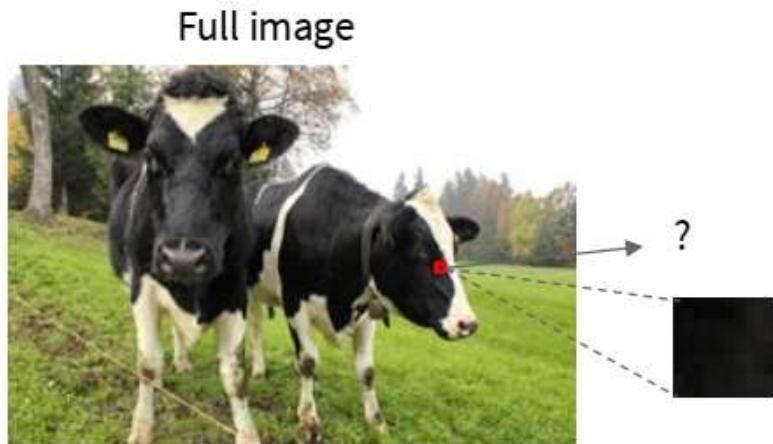
GRASS, CAT, TREE,
SKY, ...

Para cada imagen del dataset, cada píxel
tiene una etiqueta con una categoría semántica



Durante inferencia, clasificar cada píxel
De la imagen

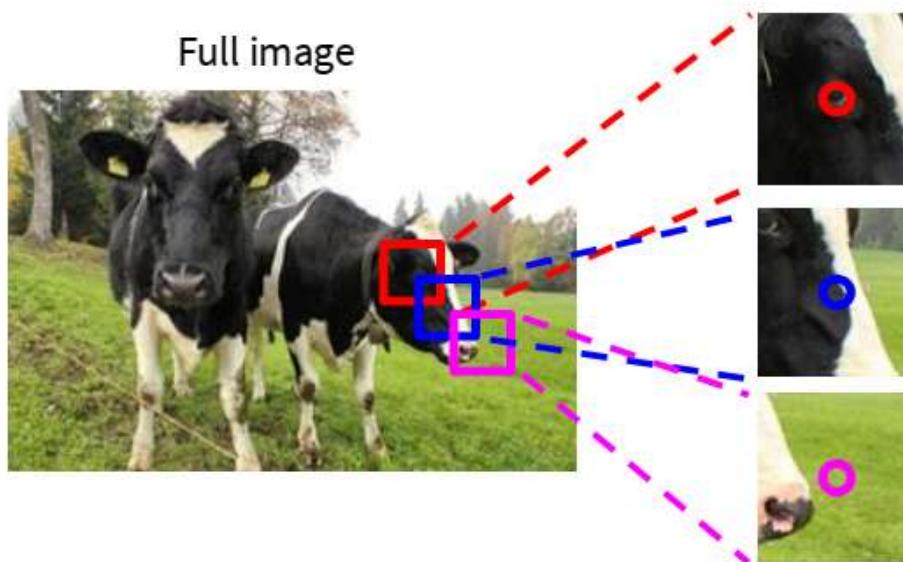
Segmentación semántica: Idea básica



Impossible clasificar sin un contexto

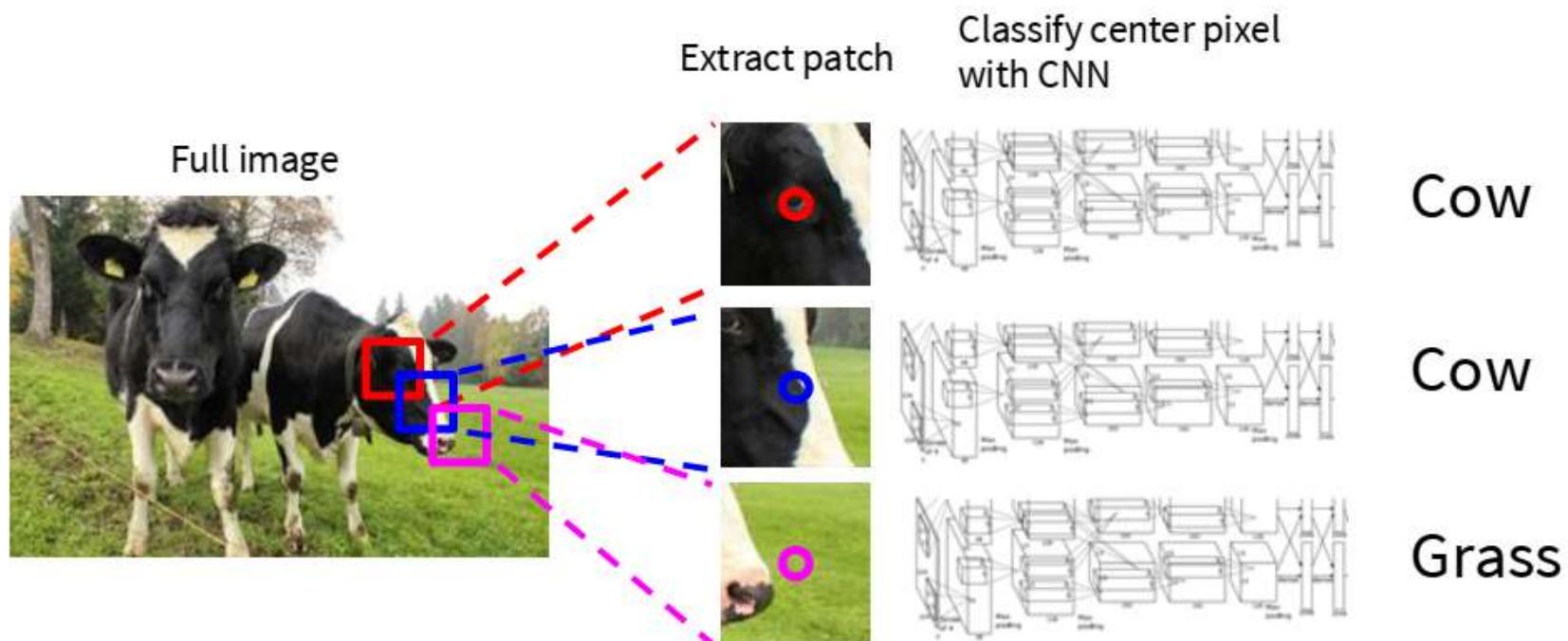
Cómo podemos incorporar contexto?

Segmentación semántica: Idea básica

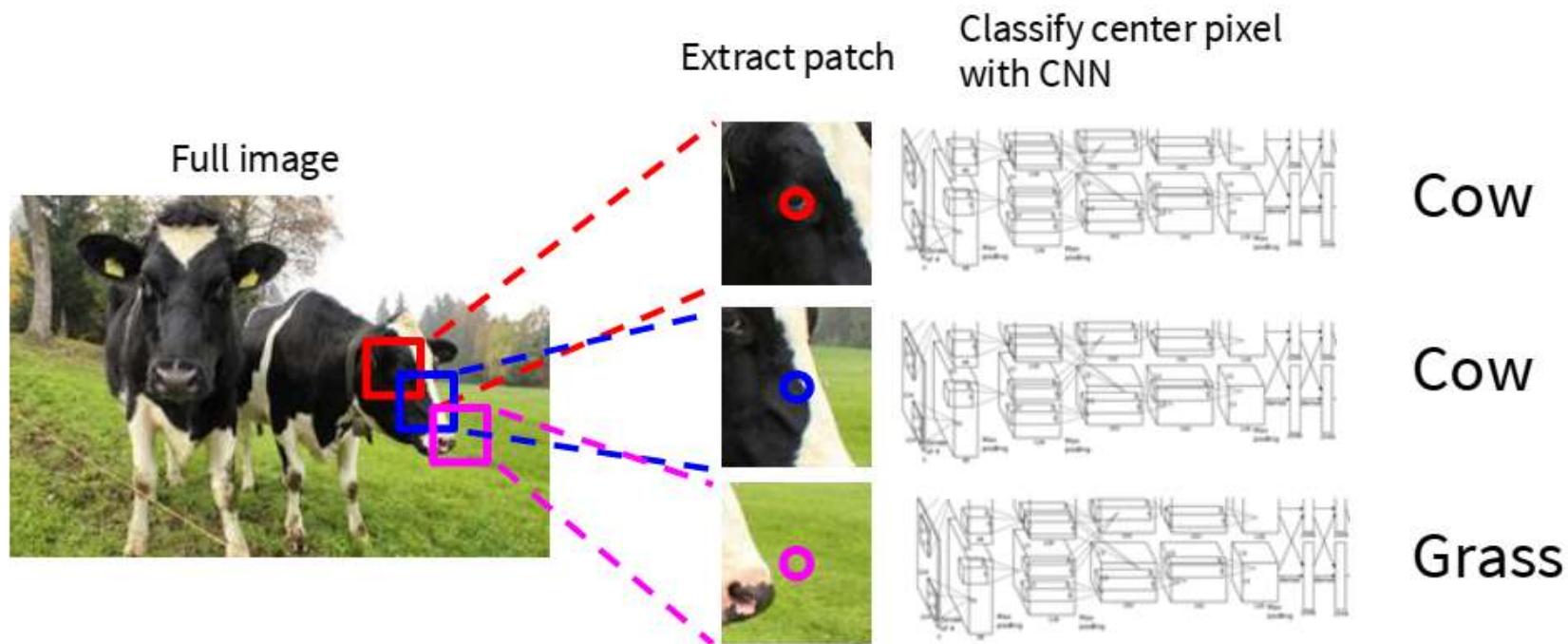


Sliding window

Segmentación semántica: Idea básica



Segmentación semántica: Idea básica

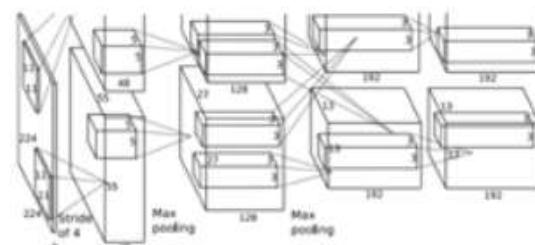


Problema: muy ineficiente!

No se reusan features compartidas entre parches

Segmentación semántica: Convolución

Full image



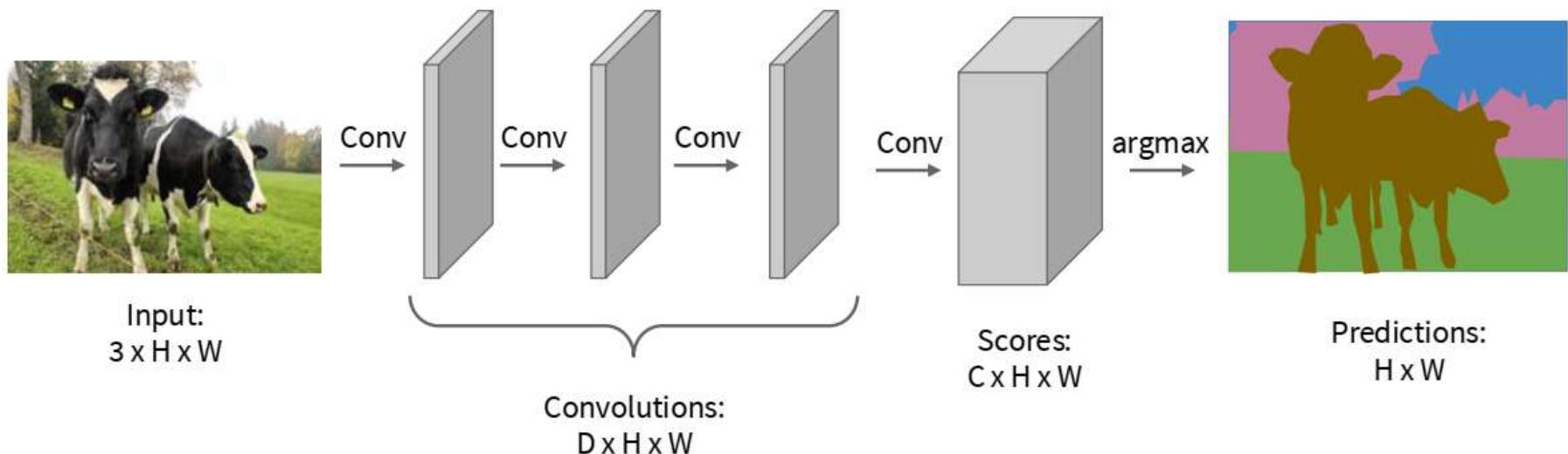
Idea intuitiva: codificar la imagen entera con ConvNet y hacer segmentación semántica directamente

Problema: arquitecturas de clasificación reducen las características espaciales de acuerdo a la profundidad de la red.

Segmentación semántica requiere que el tamaño de la salida sea igual que la entrada.

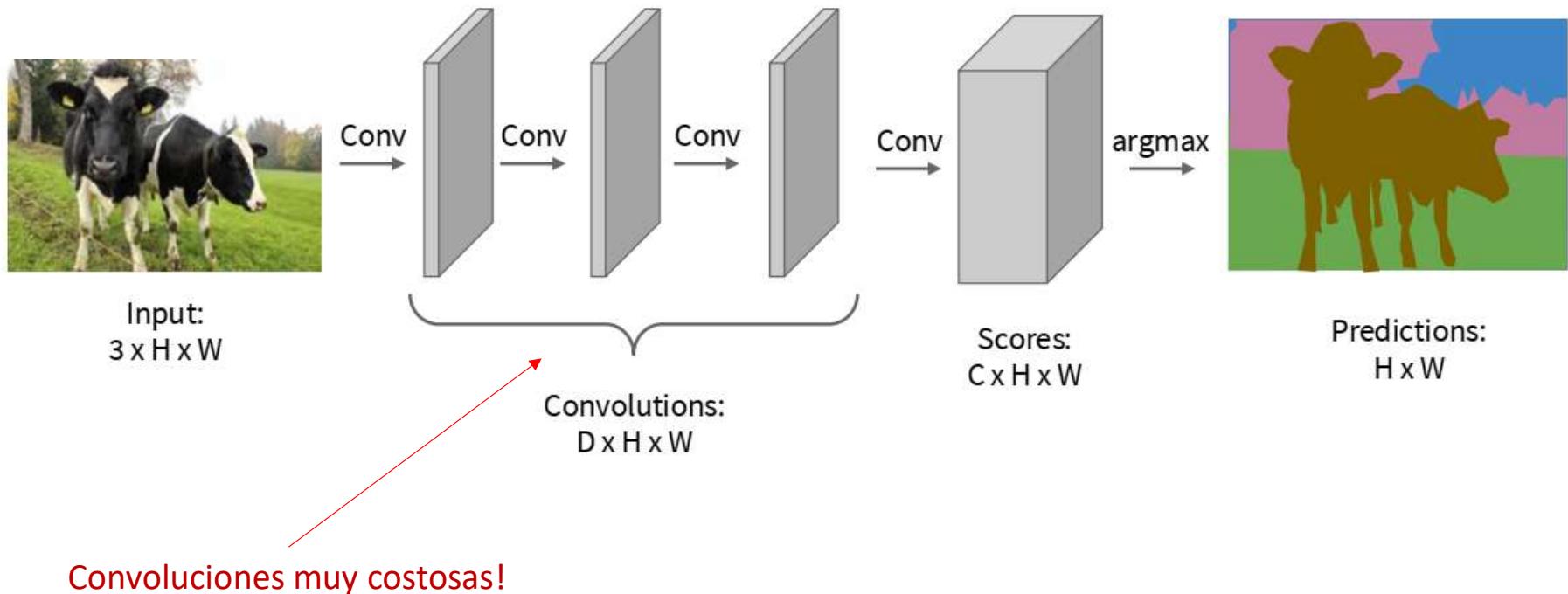
Segmentación semántica: Convolución

Diseñar una red con solo capas convolucionales sin downsampling para hacer predicciones a nivel de píxel



Segmentación semántica: Convolución

Diseñar una red con solo capas convolucionales sin downsampling para hacer predicciones a nivel de píxel

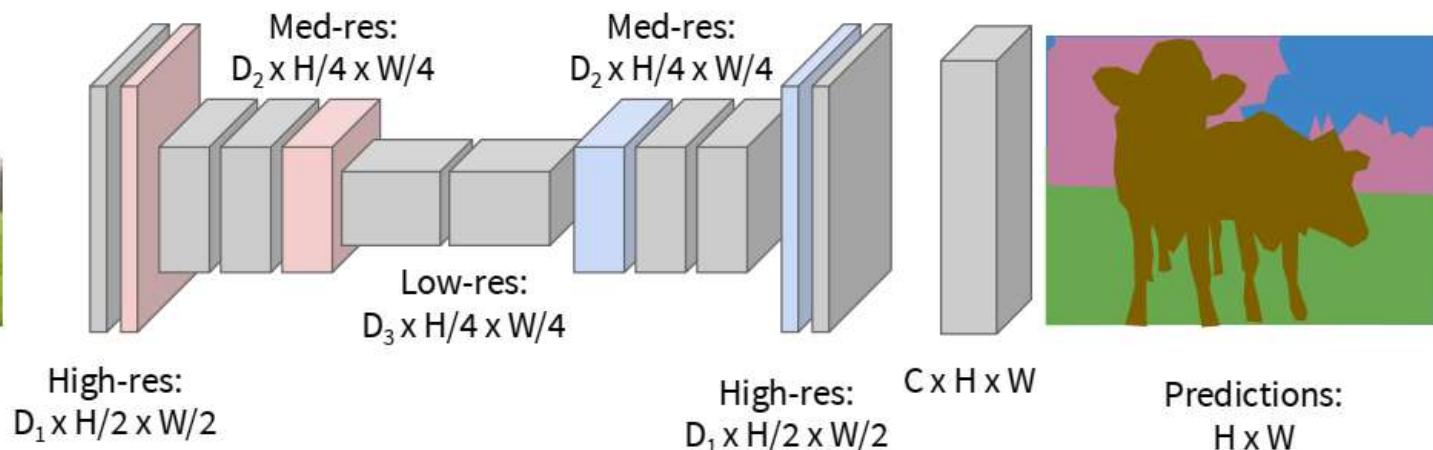


Segmentación semántica: Convolución

Idea: usar **downsampling** y **upsampling** dentro de la misma red.



Input:
 $3 \times H \times W$

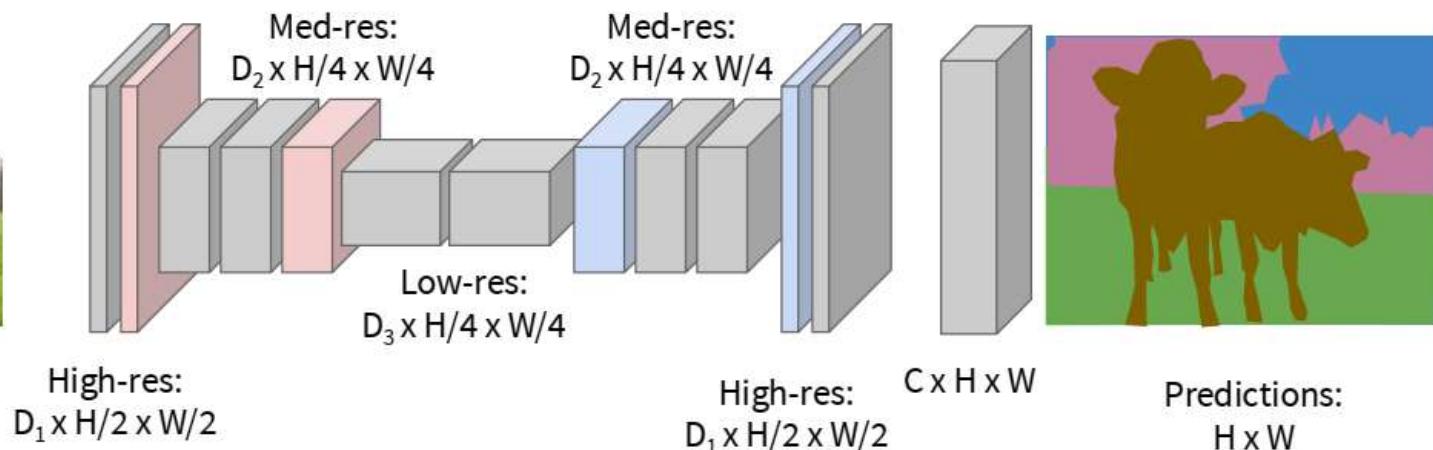


Segmentación semántica: Convolución

Idea: usar **downsampling** y **upsampling** dentro de la misma red.



Input:
 $3 \times H \times W$



Cómo se hace el upsampling?

Segmentación semántica: Convolución

Upsampling: Unpooling

Nearest Neighbor

1	2
3	4



1	1	2	2
1	1	2	2
3	3	4	4
3	3	4	4

Input: 2 x 2

Output: 4 x 4

“Bed of Nails”

1	2
3	4



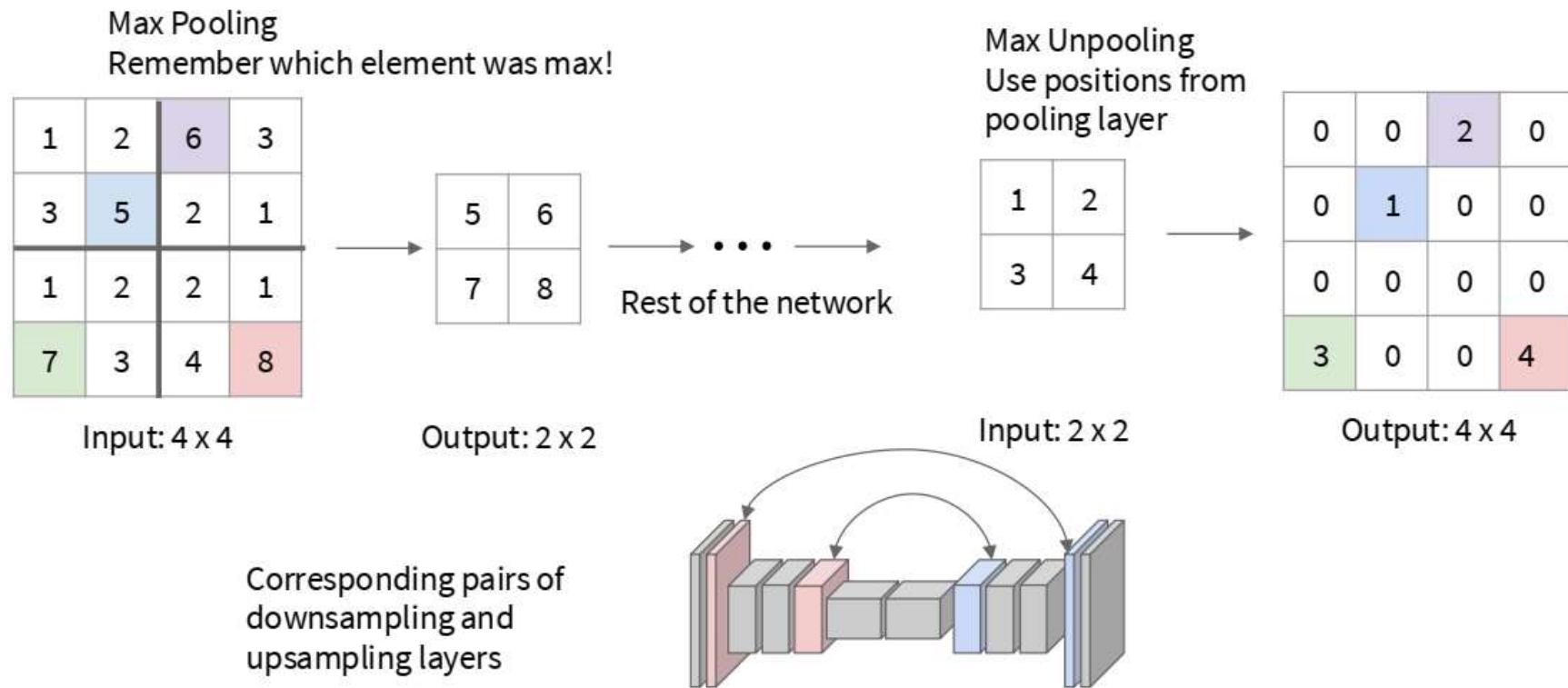
1	0	2	0
0	0	0	0
3	0	4	0
0	0	0	0

Input: 2 x 2

Output: 4 x 4

Segmentación semántica: Convolución

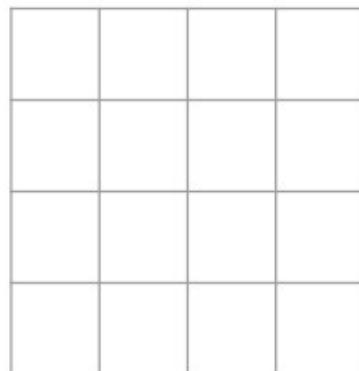
Upsampling: Max-Unpooling



Segmentación semántica: Convolución

Se puede aprender el upsampling

Recordemos: convolución normal 3×3 , stride 1, pad 1



Input: 4×4

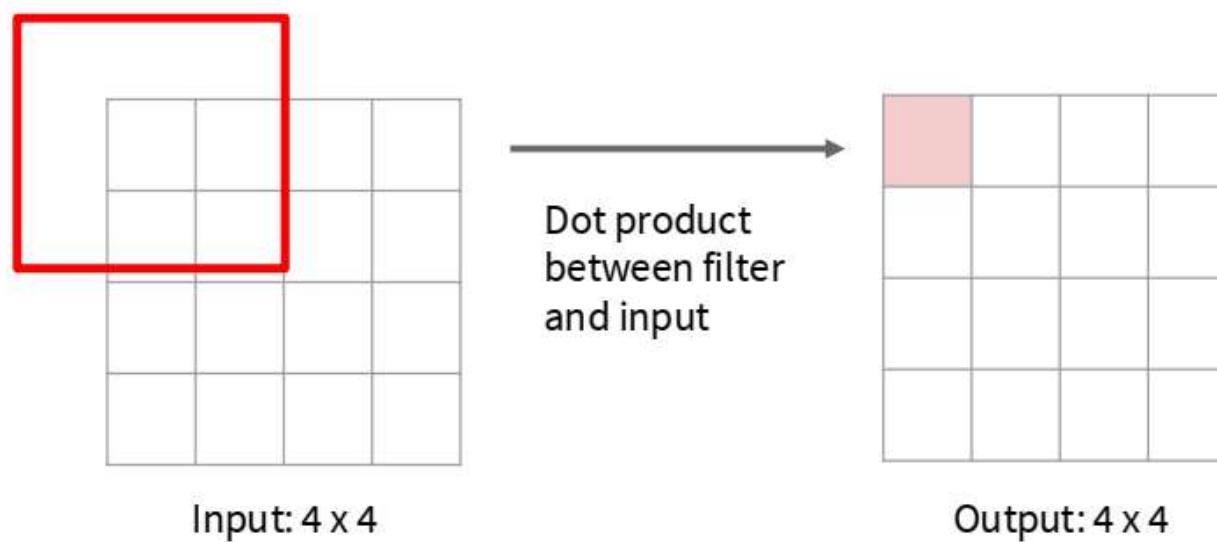


Output: 4×4

Segmentación semántica: Convolución

Se puede aprender el upsampling

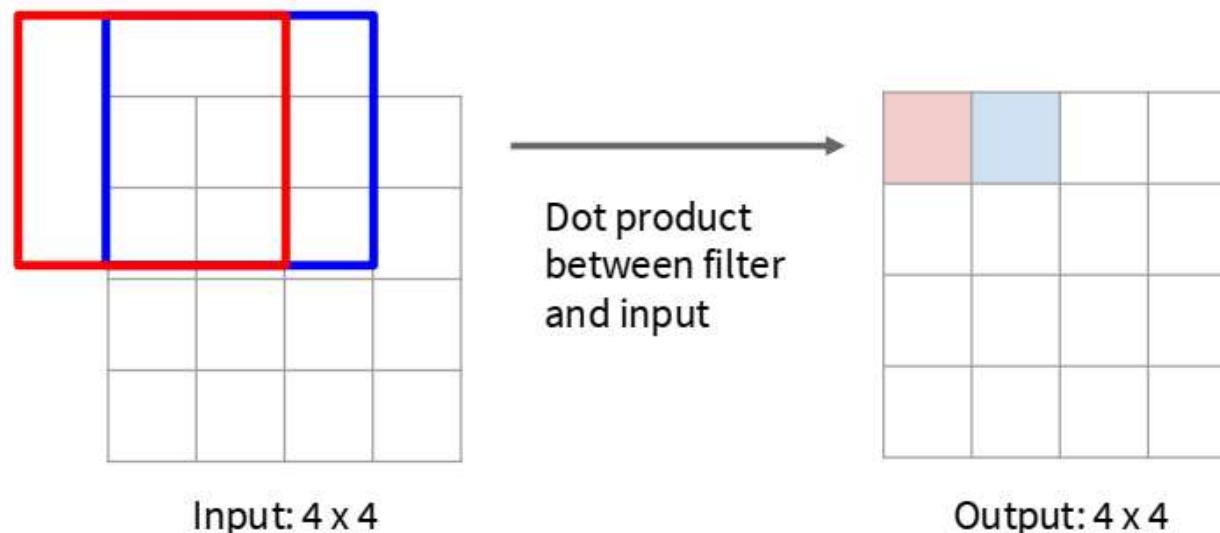
Recordemos: convolución normal 3×3 , stride 1, pad 1



Segmentación semántica: Convolución

Se puede aprender el upsampling

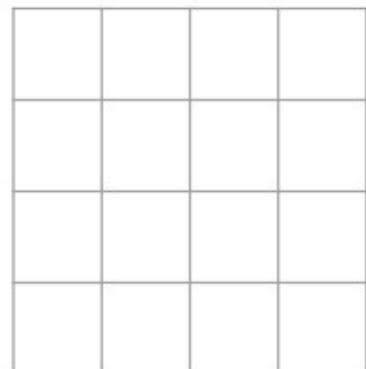
Recordemos: convolución normal 3×3 , stride 1, pad 1



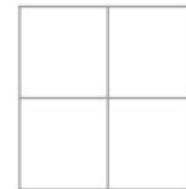
Segmentación semántica: Convolución

Se puede aprender el upsampling

Recordemos: convolución normal 3×3 , **stride 2**, pad 1



Input: 4×4

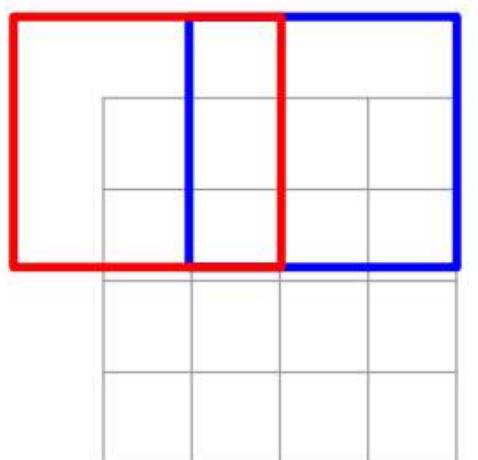


Output: 2×2

Segmentación semántica: Convolución

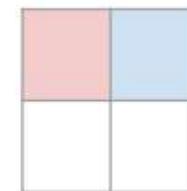
Se puede aprender el upsampling

Recordemos: convolución normal 3×3 , **stride 2**, pad 1



Input: 4×4

Dot product
between filter
and input



Output: 2×2

Filtro se mueve dos píxeles en el input por cada pixel en el output

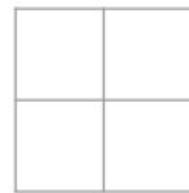
Stride: proporción entre movimiento en el input y en el output

Podemos interpretar convolución con stride como “downsampling aprendible”

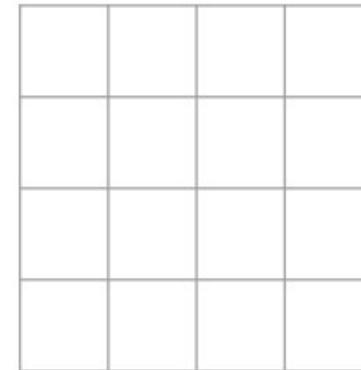
Segmentación semántica: Convolución

Se puede aprender el upsampling

Idea: convolución transpuesta 3x3 , **stride 2**, pad 1



Input: 2 x 2

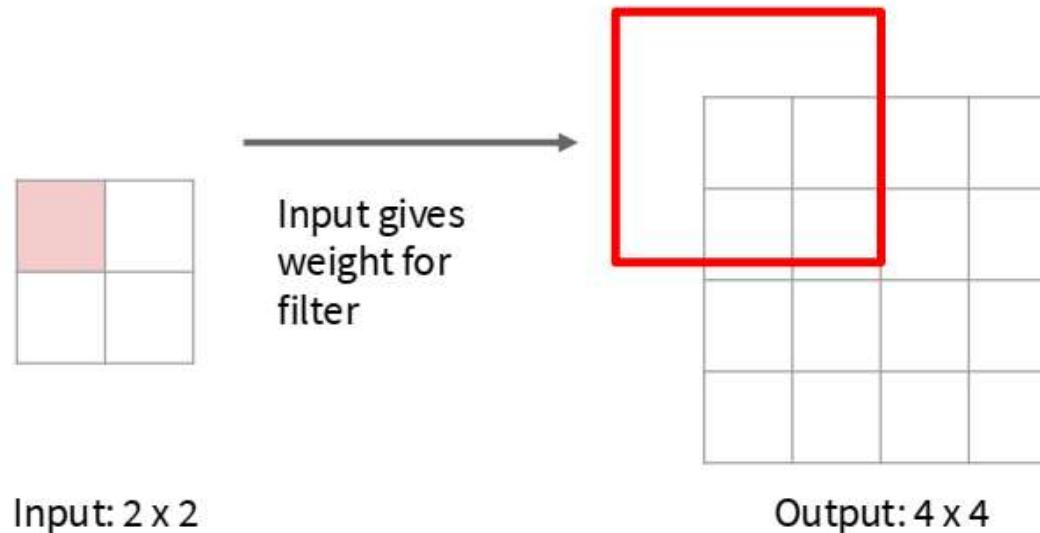


Output: 4 x 4

Segmentación semántica: Convolución

Se puede aprender el upsampling

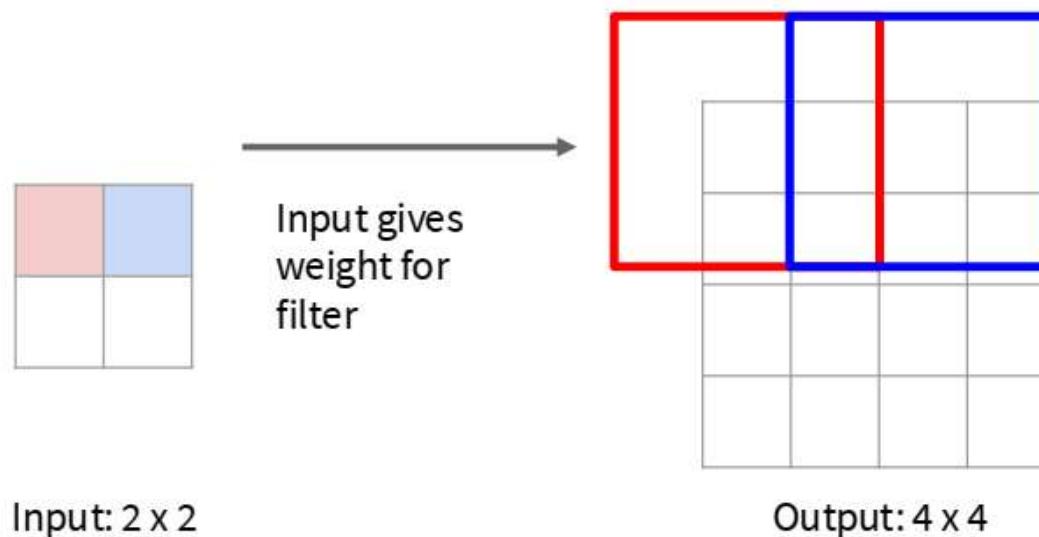
Idea: convolución transpuesta 3x3 , **stride 2**, pad 1



Segmentación semántica: Convolución

Se puede aprender el upsampling

Idea: convolución transpuesta 3x3 , **stride 2**, pad 1



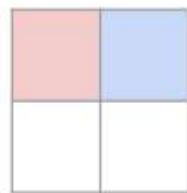
Filtro se mueve dos píxeles en el output por cada pixel en el input

Stride: proporción entre movimiento en el output y en el input

Segmentación semántica: Convolución

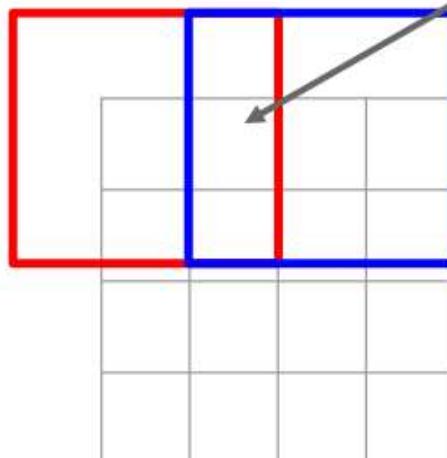
Se puede aprender el upsampling

Idea: convolución transpuesta 3x3 , **stride 2**, pad 1



Input: 2 x 2

Input gives weight for filter



Output: 4 x 4

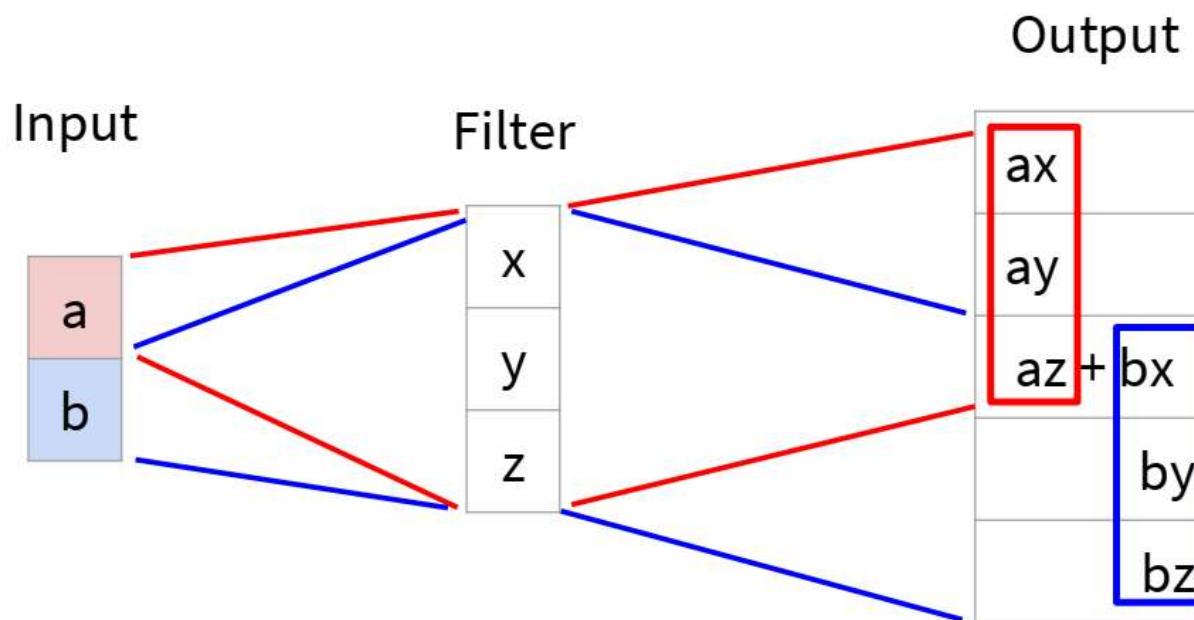
Sumas en donde
el output hace
overlap

Filtro se mueve dos píxeles en el
output por cada pixel en el input

Stride: proporción entre
movimiento en el output y en el
input

Segmentación semántica: Convolución

Se puede aprender el upsampling

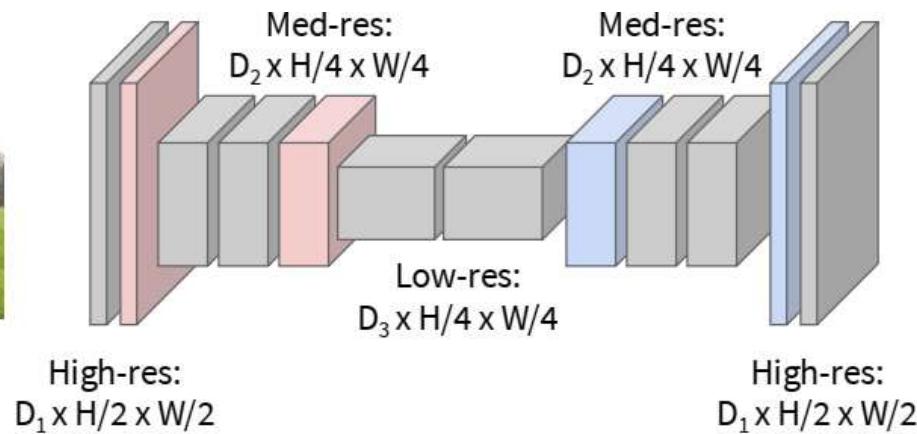


Output contiene copias del filtro ponderadas por el input, con sumas en donde hay overlap

Segmentación semántica: Convolución

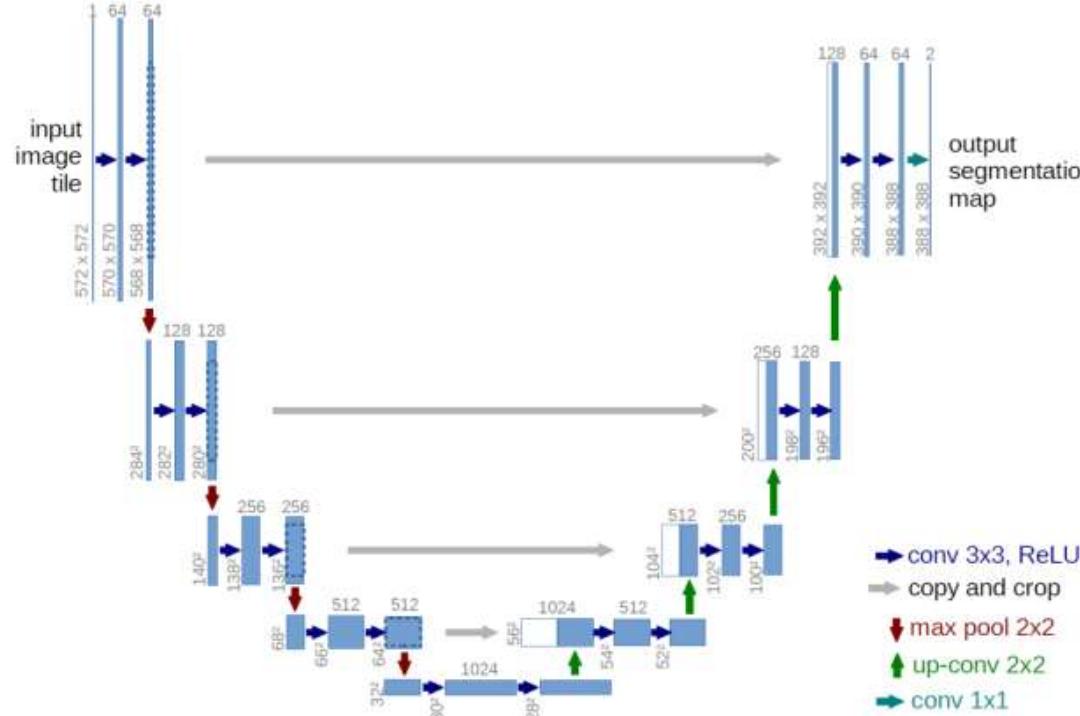


Input:
 $3 \times H \times W$



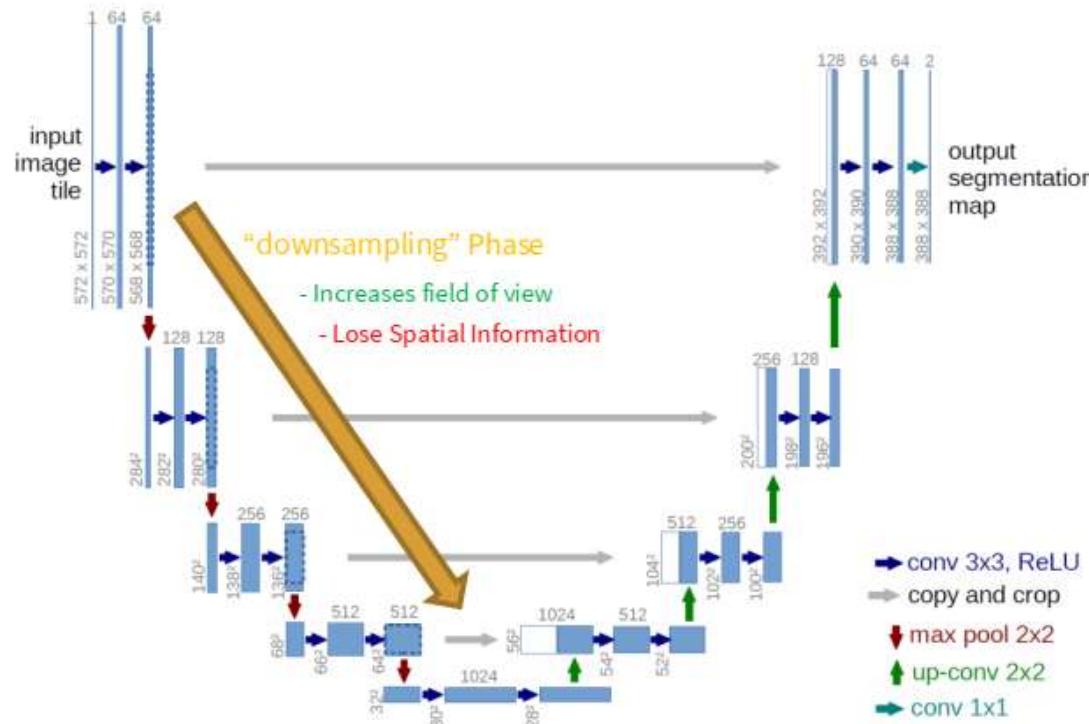
Predictions:
 $H \times W$

Segmentación semántica: Convolución UNet

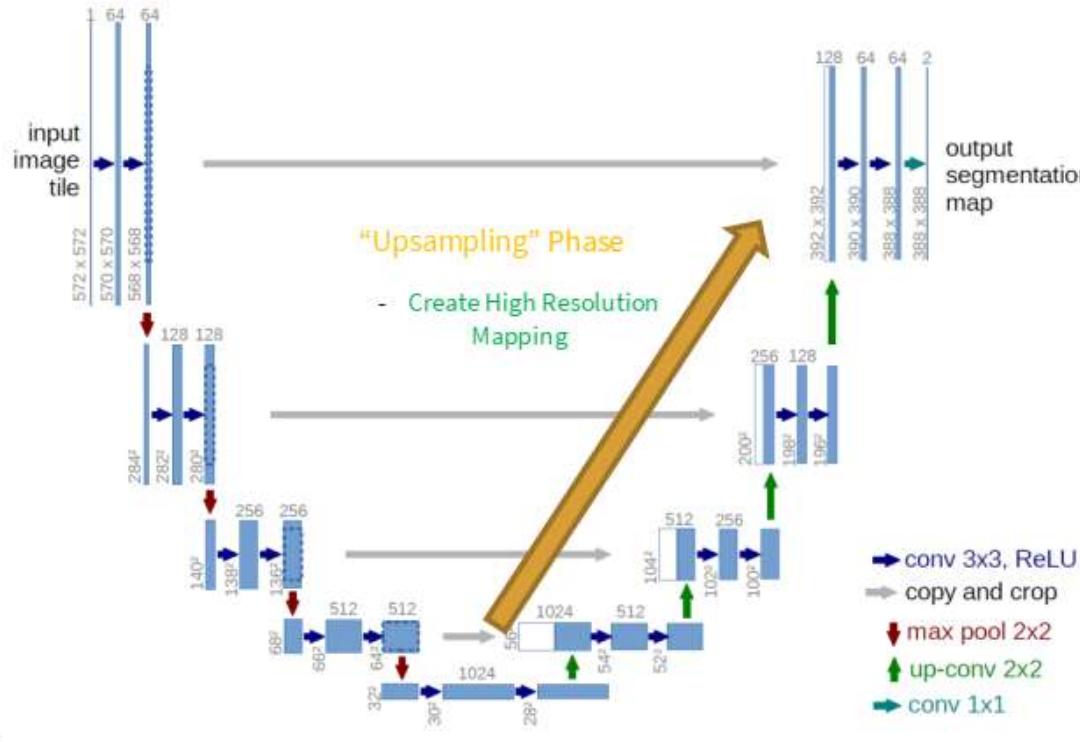


Segmentación semántica: Convolución

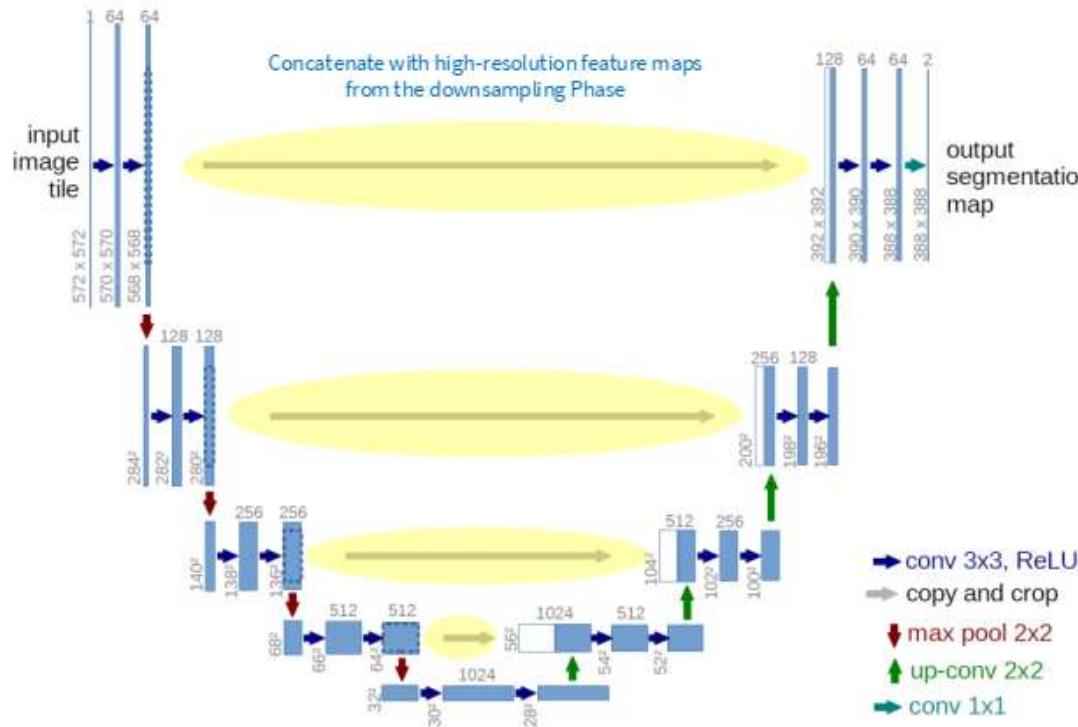
UNet



Segmentación semántica: Convolución UNet



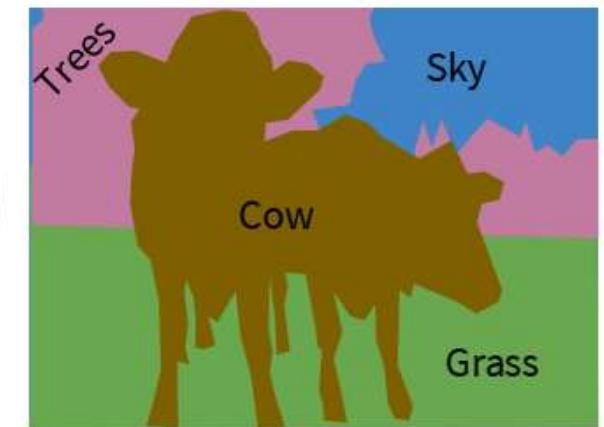
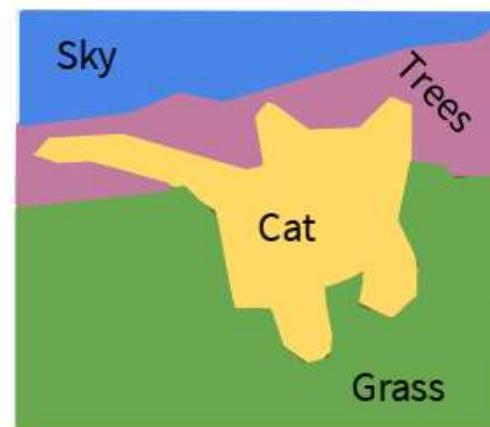
Segmentación semántica: Convolución UNet



Segmentación semántica

Etiquetar cada píxel de la imagen con una clase

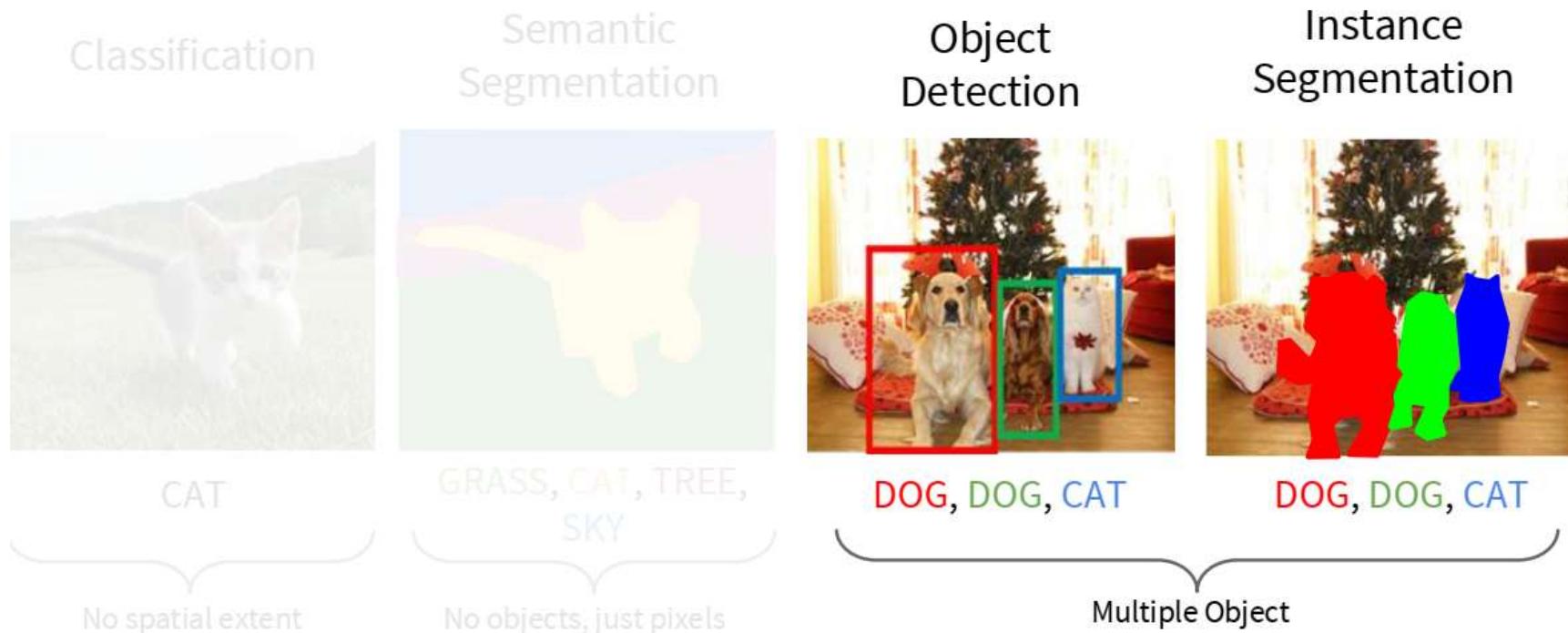
No se diferencian instancias!



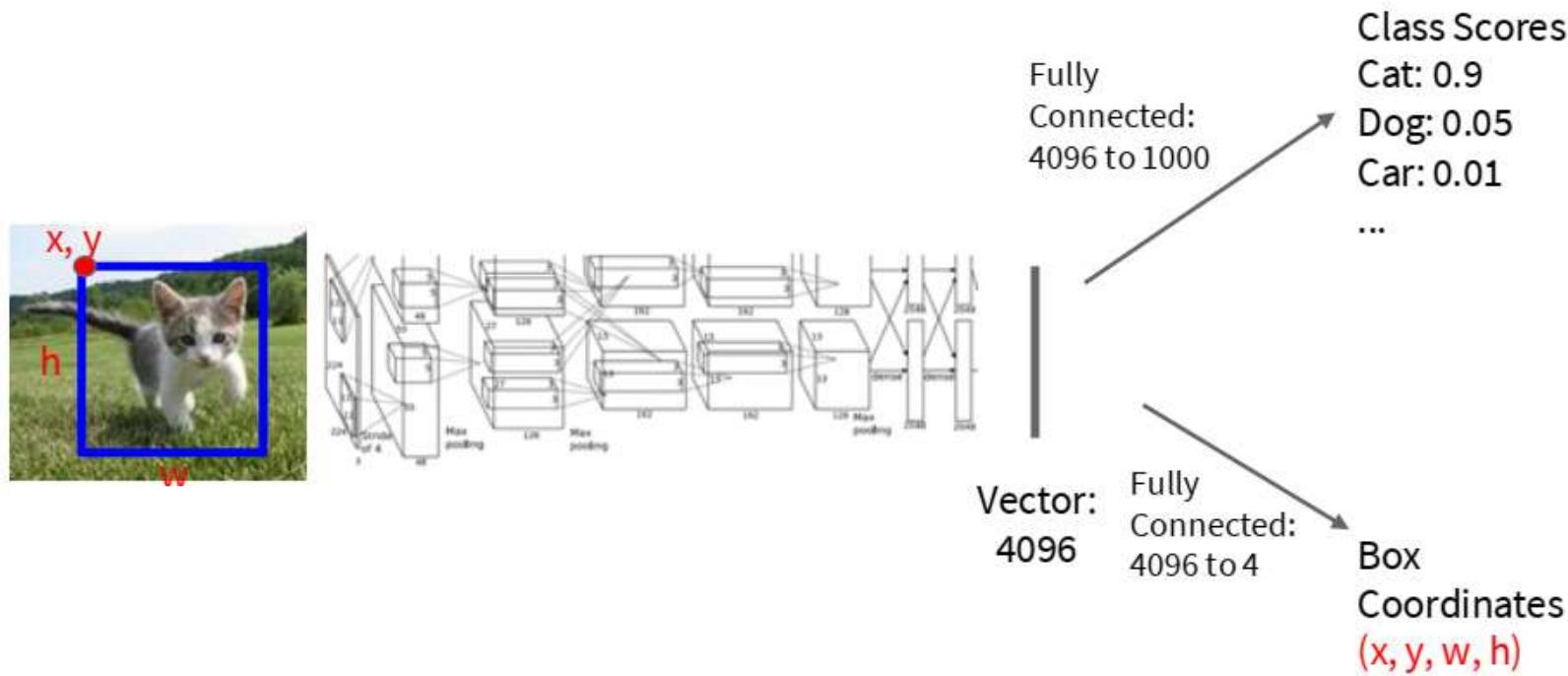
Detección de objetos



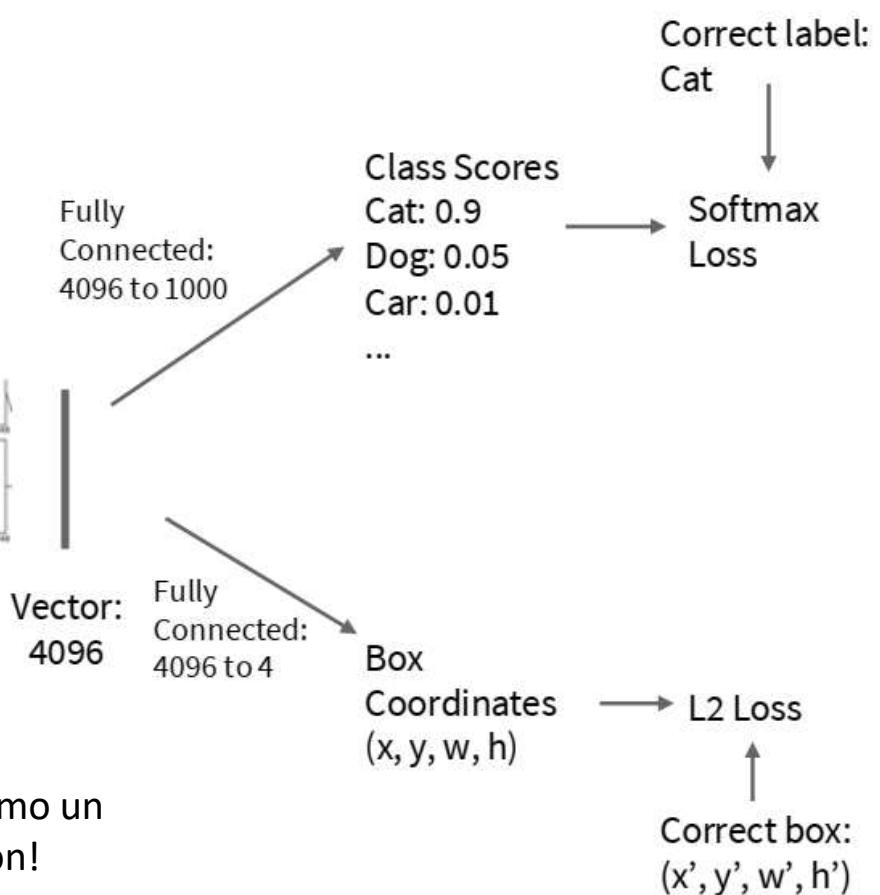
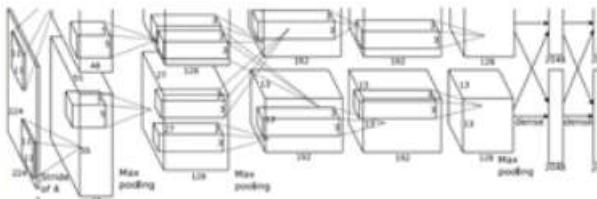
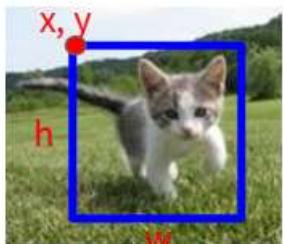
Detección de objetos



Detección de objetos: un objeto Clasificación + localización

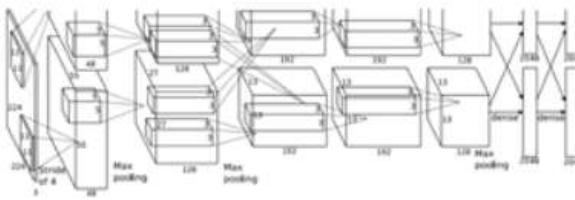
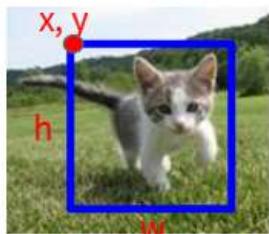


Detección de objetos: un objeto Clasificación + localización

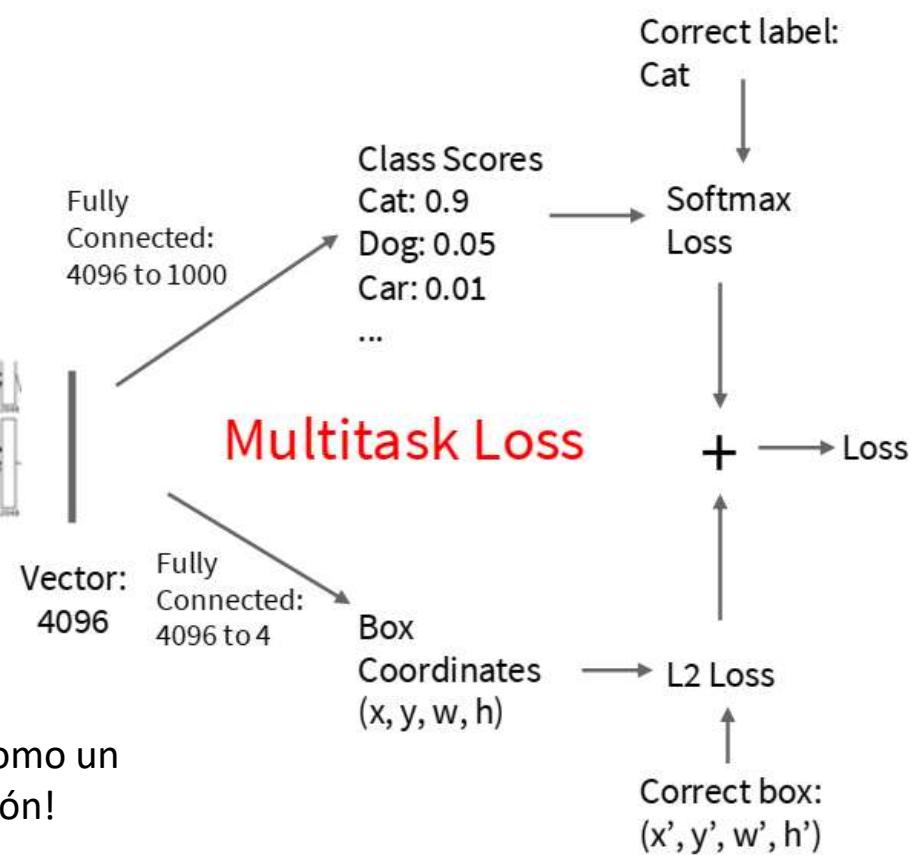


Tratar localización como un
Problema de regresión!

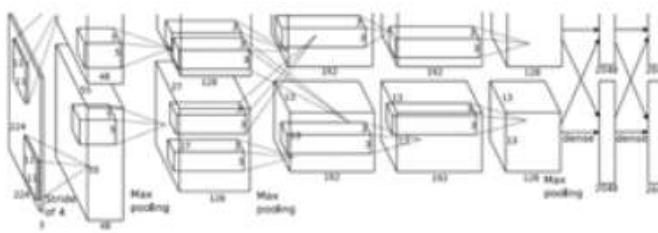
Detección de objetos: un objeto Clasificación + localización



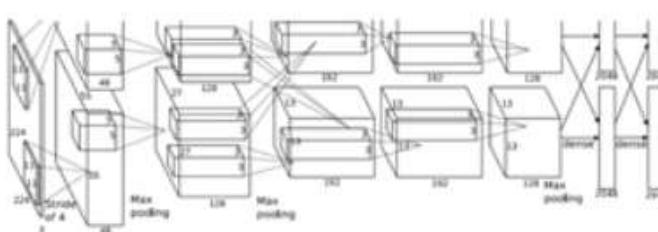
Tratar localización como un
Problema de regresión!



Detección de objetos: múltiples objetos



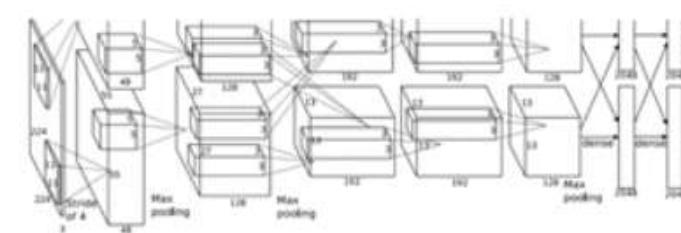
CAT: (x, y, w, h)



DOG: (x, y, w, h)

DOG: (x, y, w, h)

CAT: (x, y, w, h)



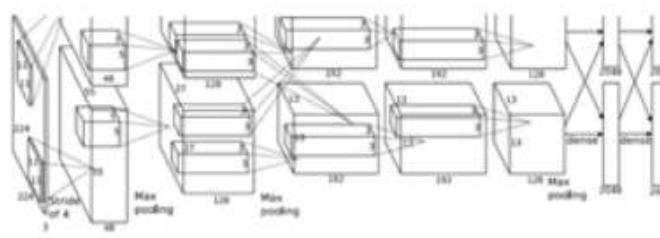
DUCK: (x, y, w, h)

DUCK: (x, y, w, h)

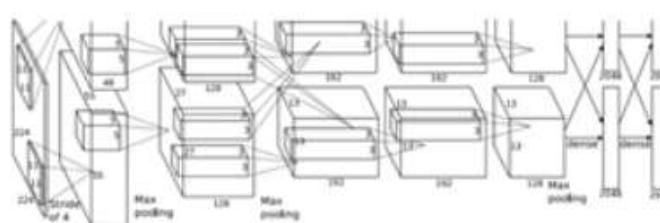
....

Detección de objetos: múltiples objetos

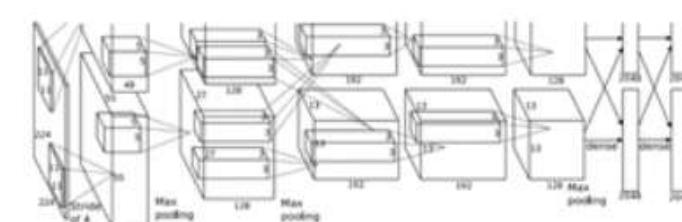
Cada imagen necesita un número diferente de salidas



CAT: (x, y, w, h) 4 números



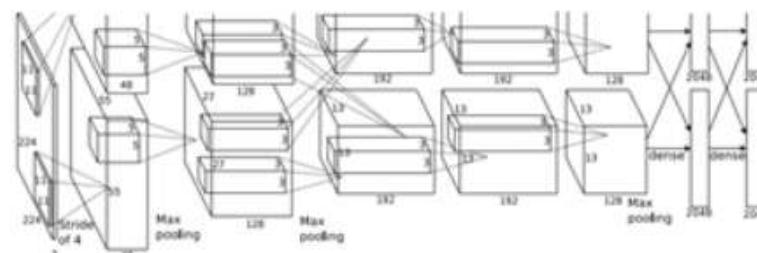
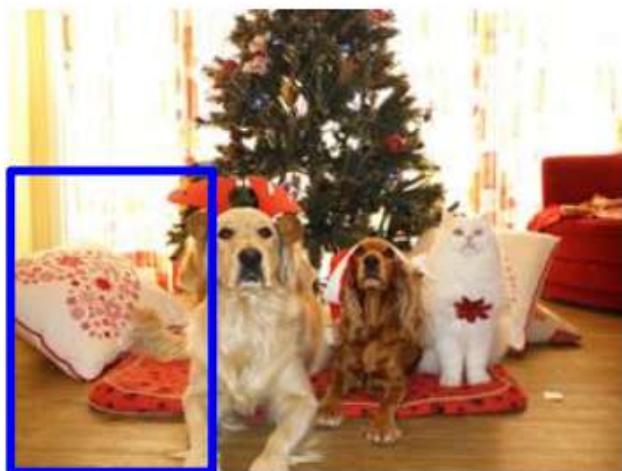
DOG: (x, y, w, h)
DOG: (x, y, w, h) 12 números
CAT: (x, y, w, h)



DUCK: (x, y, w, h)
DUCK: (x, y, w, h) Muchos números
....

Detección de objetos: múltiples objetos

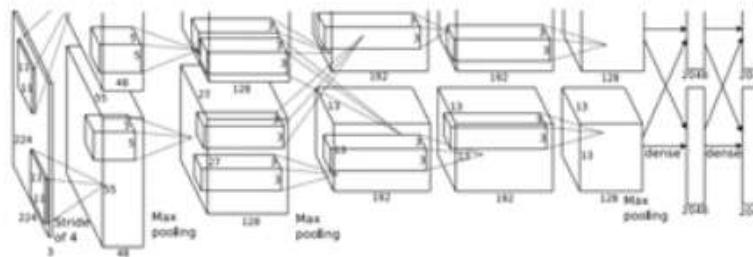
Aplicar una CNN a muchos diferentes crops de la imagen,
CNN clasifica cada crop como objeto o fondo



Dog? NO
Cat? NO
Background? YES

Detección de objetos: múltiples objetos

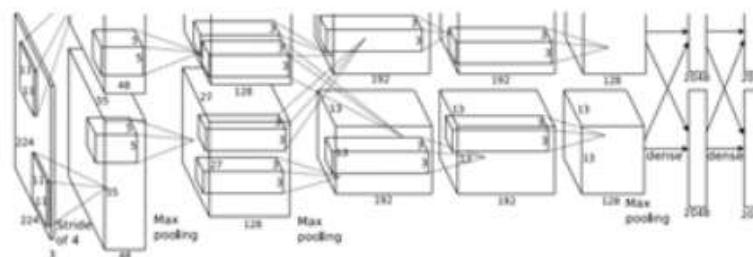
Aplicar una CNN a muchos diferentes crops de la imagen,
CNN clasifica cada crop como objeto o fondo



Dog? YES
Cat? NO
Background? NO

Detección de objetos: múltiples objetos

Aplicar una CNN a muchos diferentes crops de la imagen,
CNN clasifica cada crop como objeto o fondo

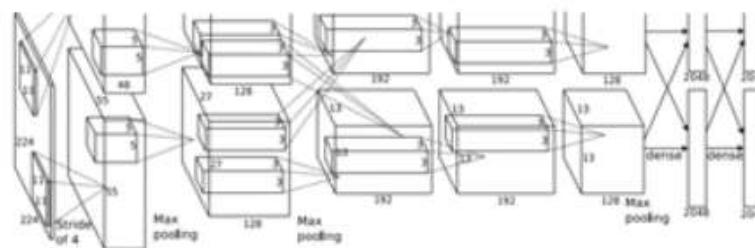


Dog? YES
Cat? NO
Background? NO

Detección de objetos: múltiples objetos

Cuál es el problema con este enfoque?

Aplicar una CNN a muchos diferentes crops de la imagen,
CNN clasifica cada crop como objeto o fondo

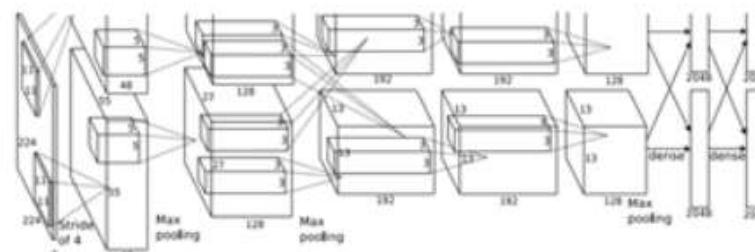
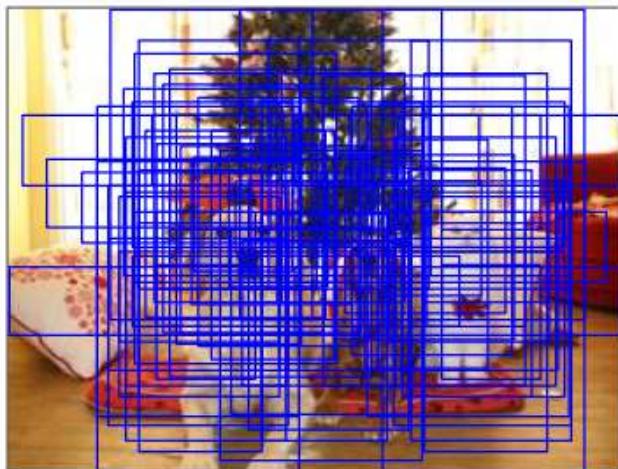


Dog? NO
Cat? YES
Background? NO

Detección de objetos: múltiples objetos

Necesitas aplicar CNN a un gran número de Lugares, escalas y aspect ratios. Es muy caro!

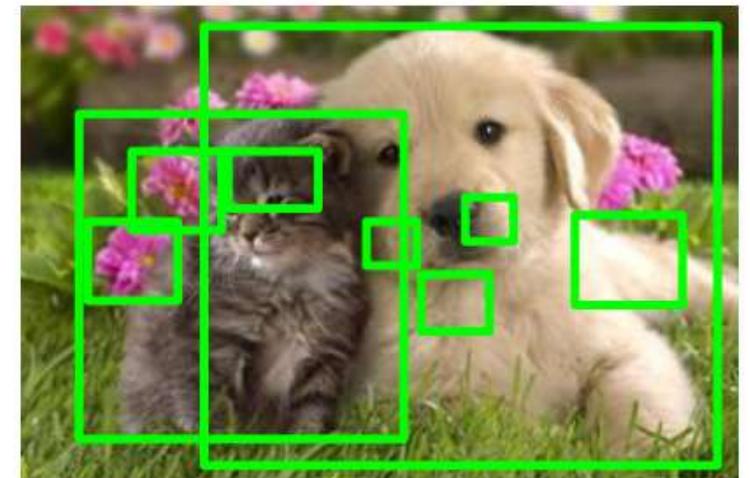
Aplicar una CNN a muchos diferentes crops de la imagen, CNN clasifica cada crop como objeto o fondo



Dog? NO
Cat? YES
Background? NO

Region proposals: Selective Search

- Encontrar regiones con potenciales objetos
- Selective Search entrega 2000 regiones en pocos segundos usando CPU

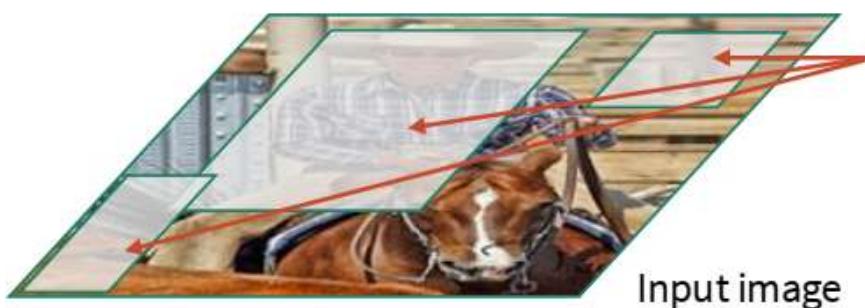


R-CNN



Input image

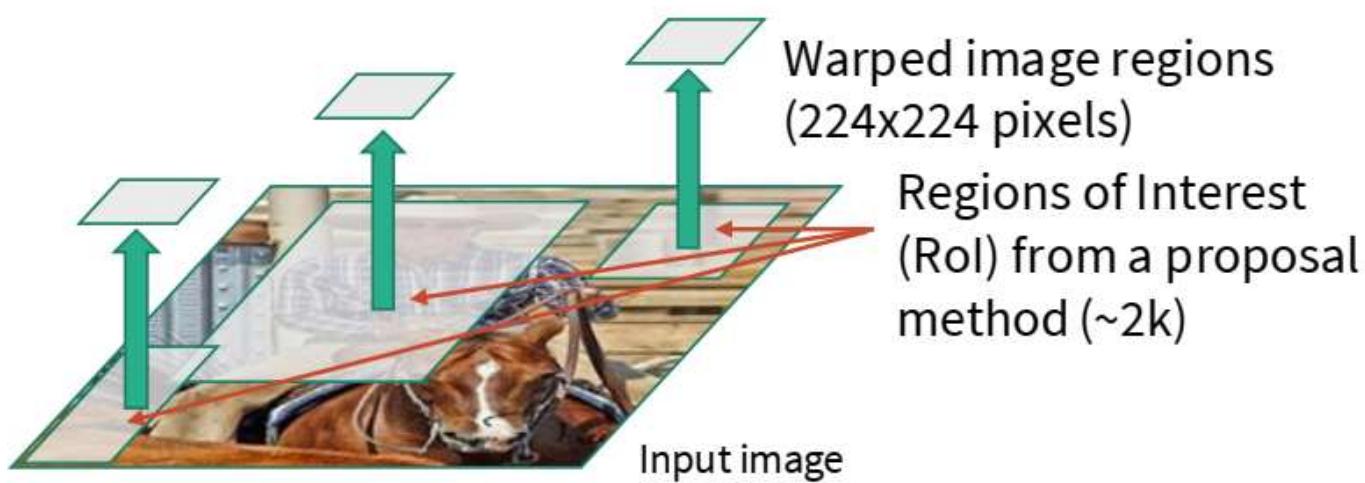
R-CNN



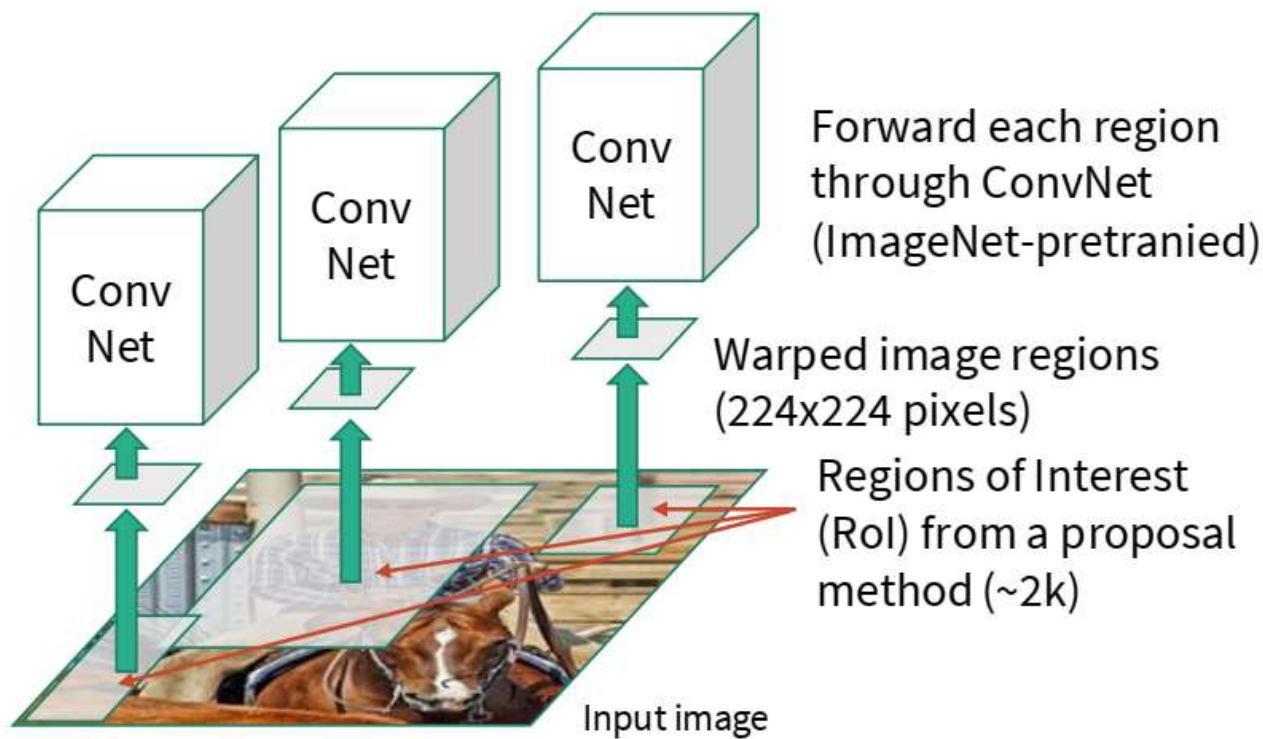
Input image

Regions of Interest
(RoI) from a proposal
method (~2k)

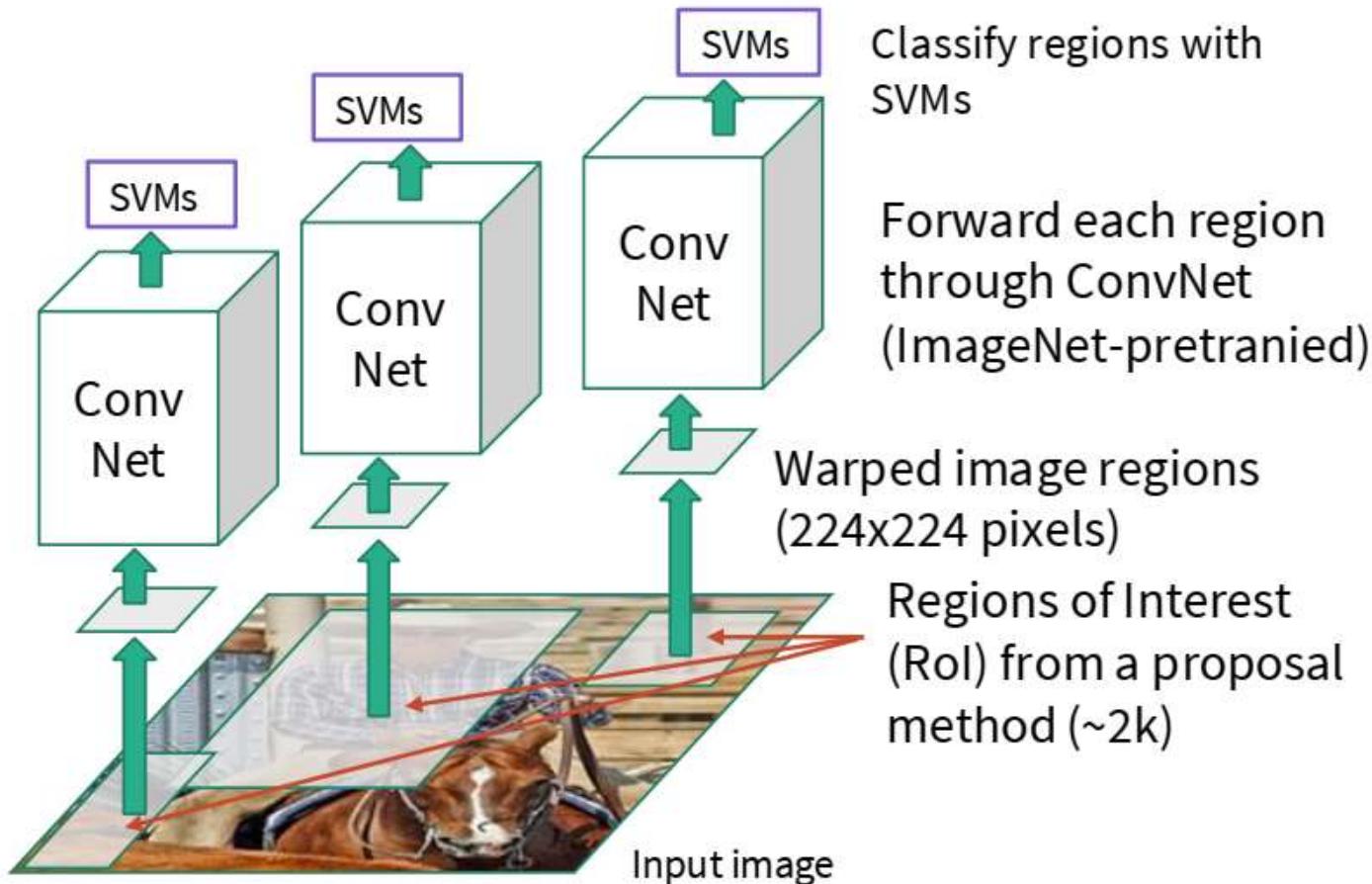
R-CNN



R-CNN

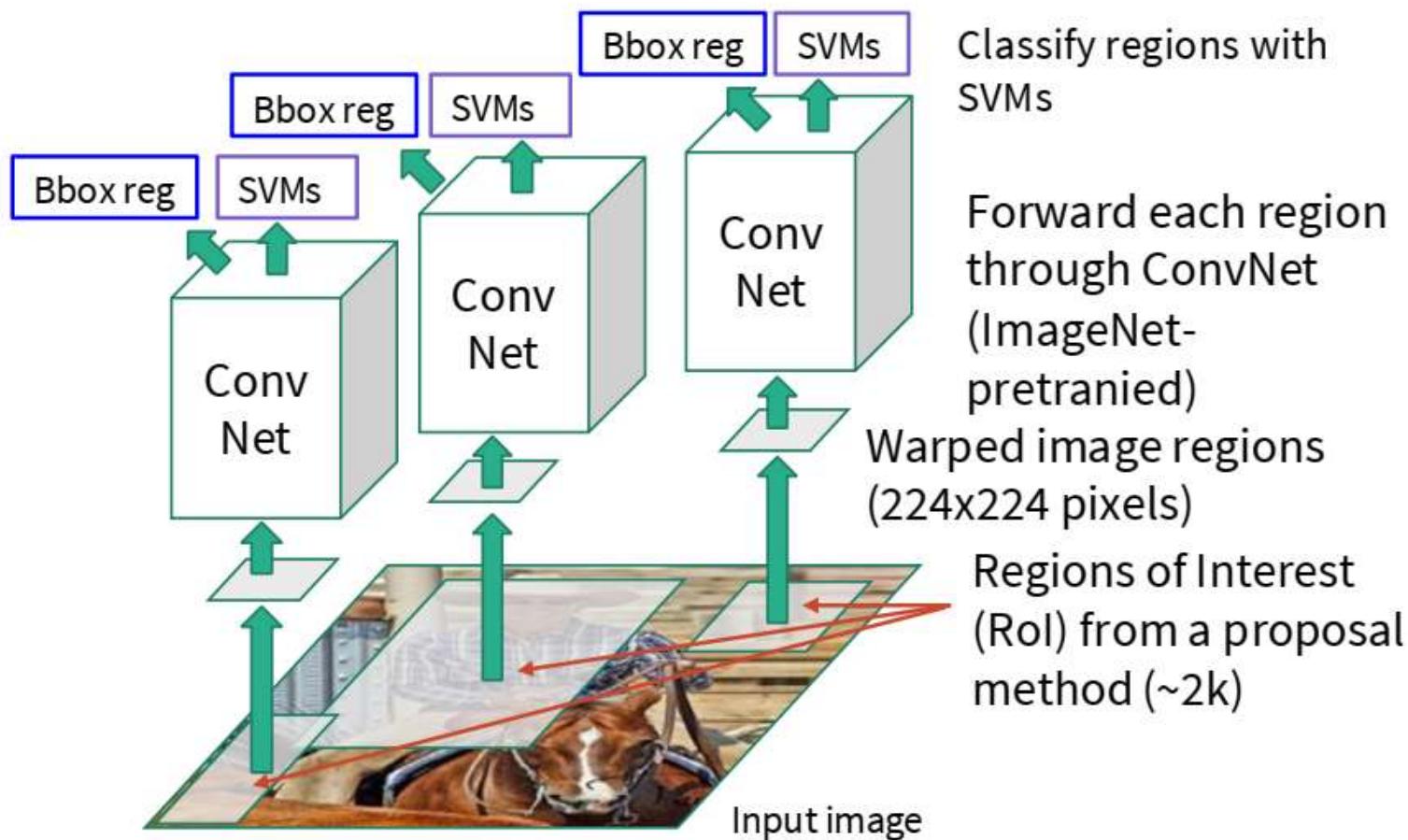


R-CNN



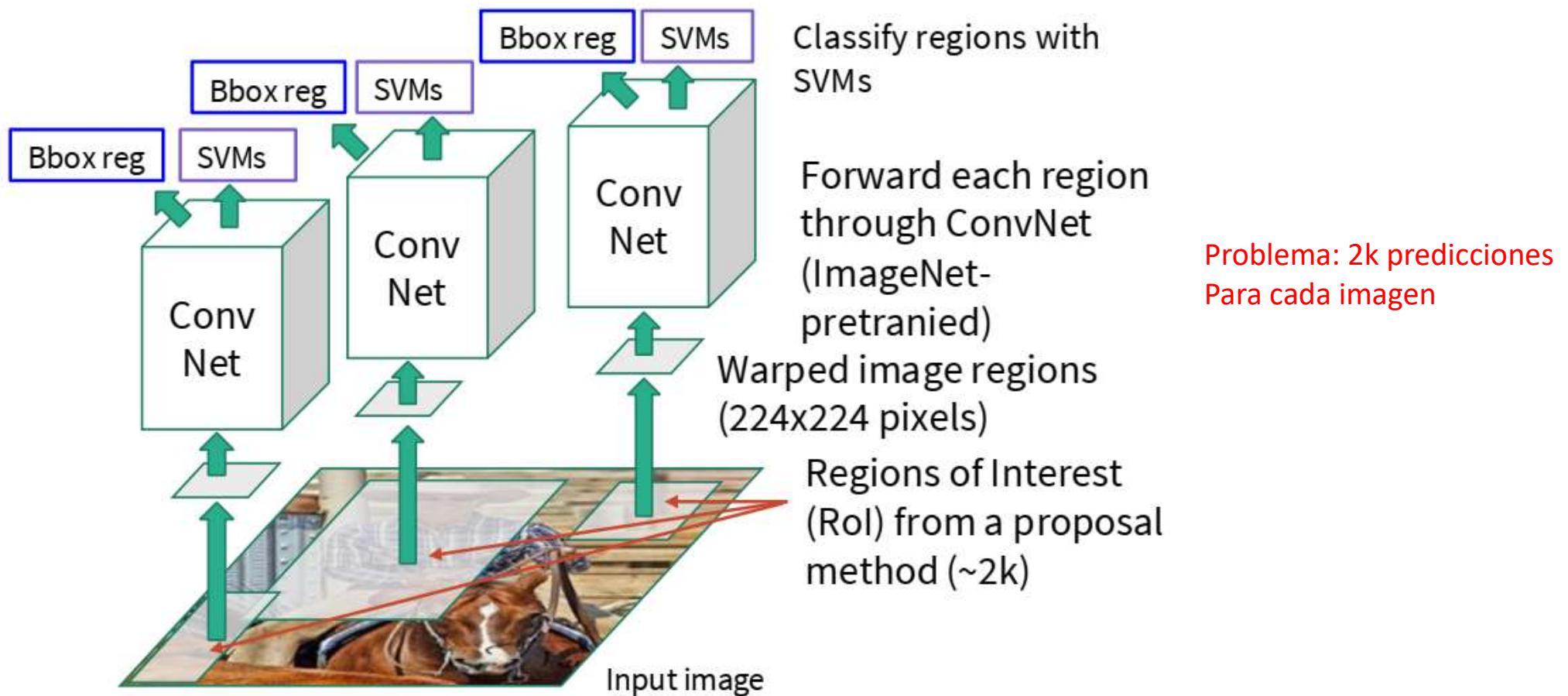
R-CNN

Predict “corrections” to the RoI: 4 numbers: (dx, dy, dw, dh)



R-CNN

Predict “corrections” to the RoI: 4 numbers: (dx, dy, dw, dh)

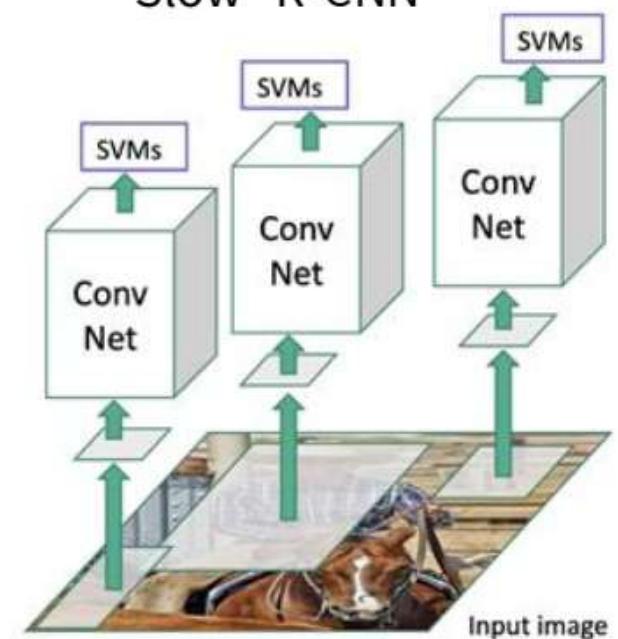


Fast R-CNN

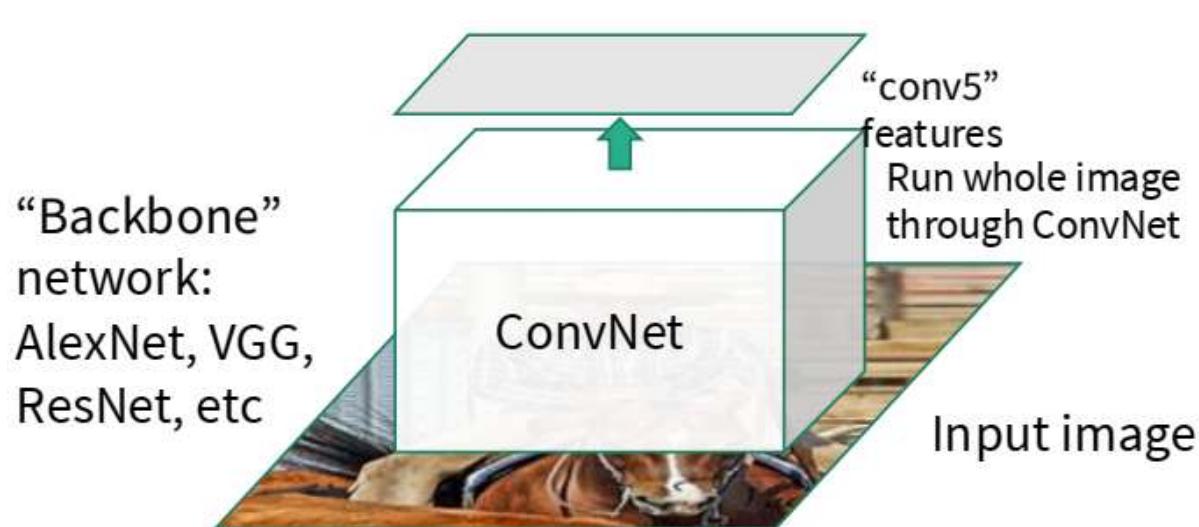


Input image

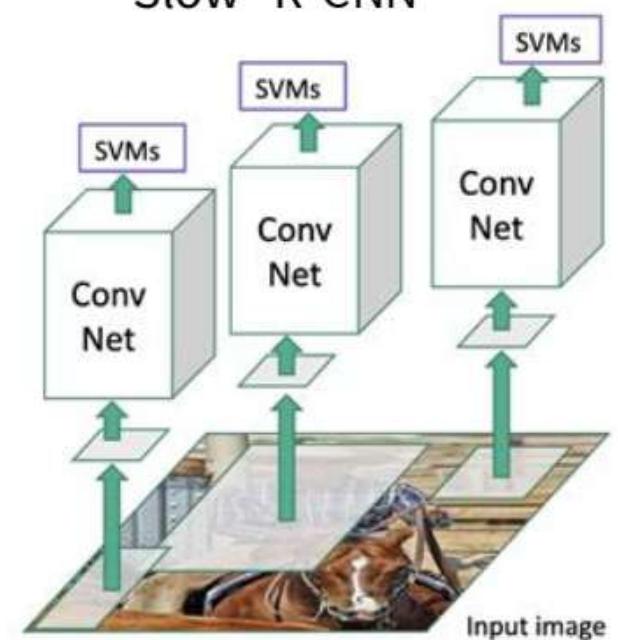
“Slow” R-CNN



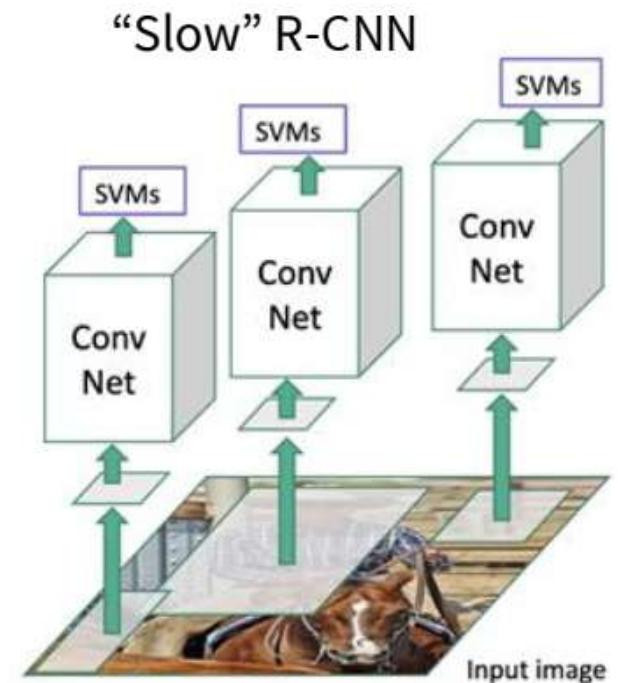
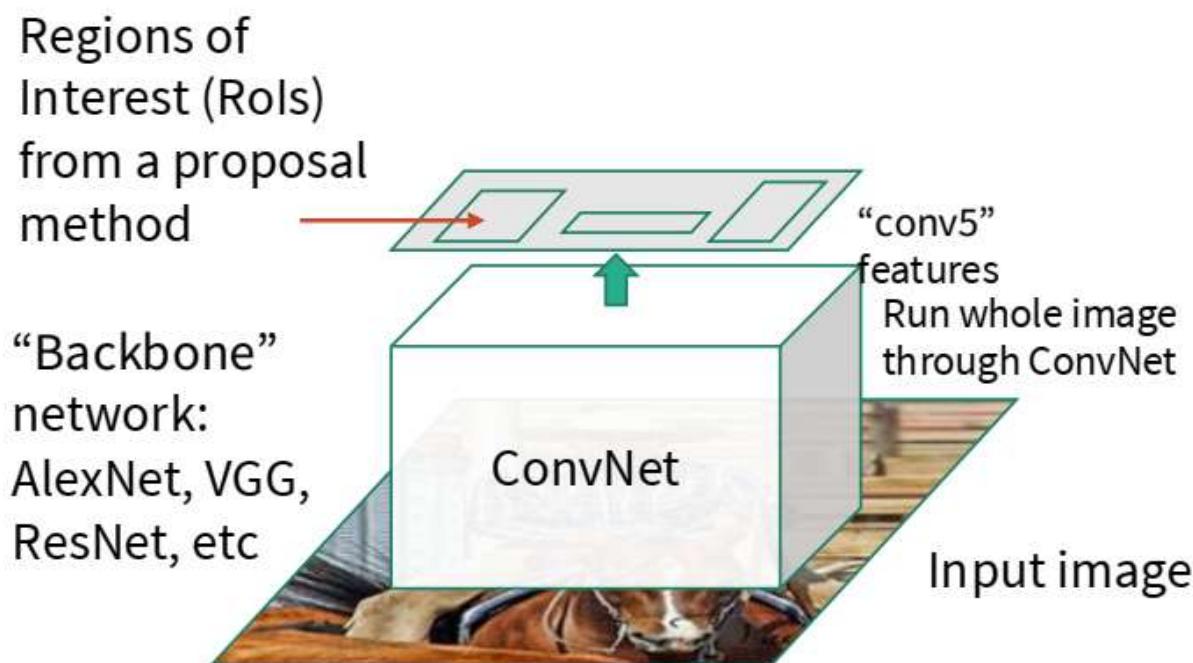
Fast R-CNN



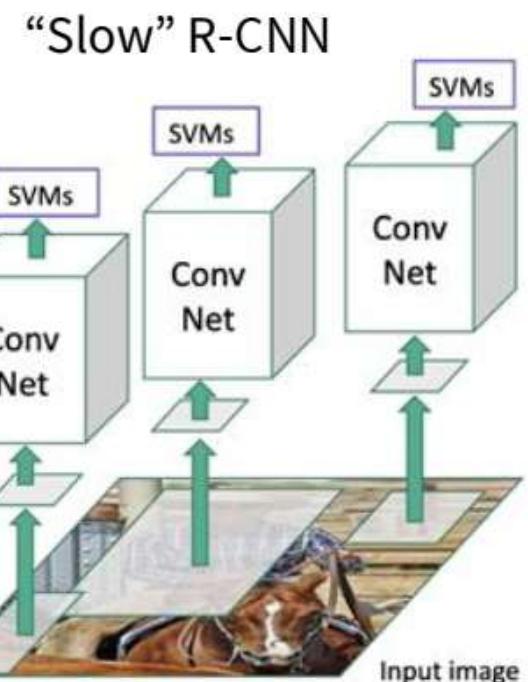
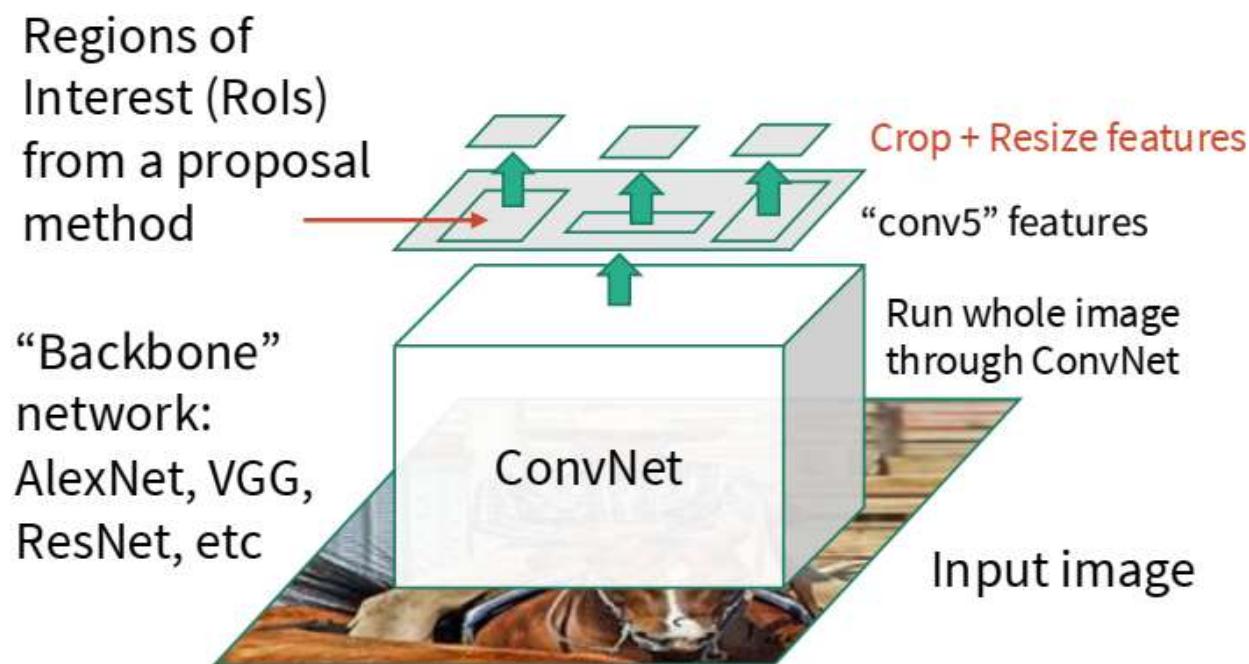
“Slow” R-CNN



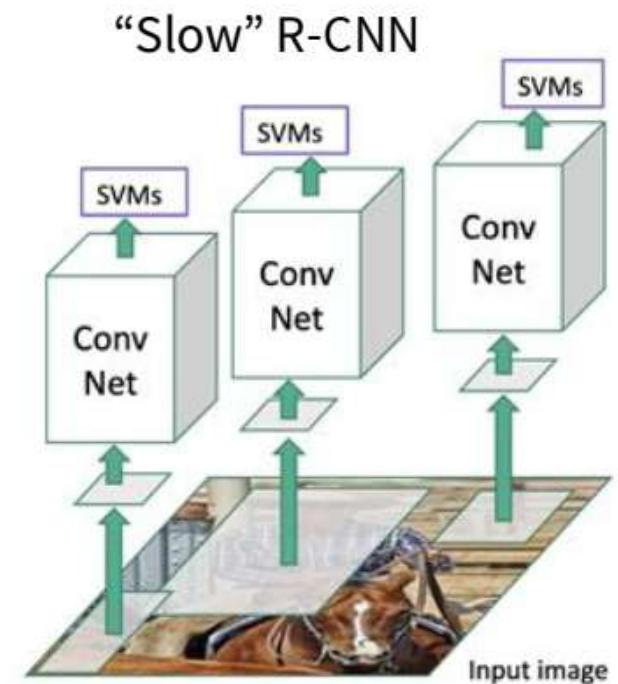
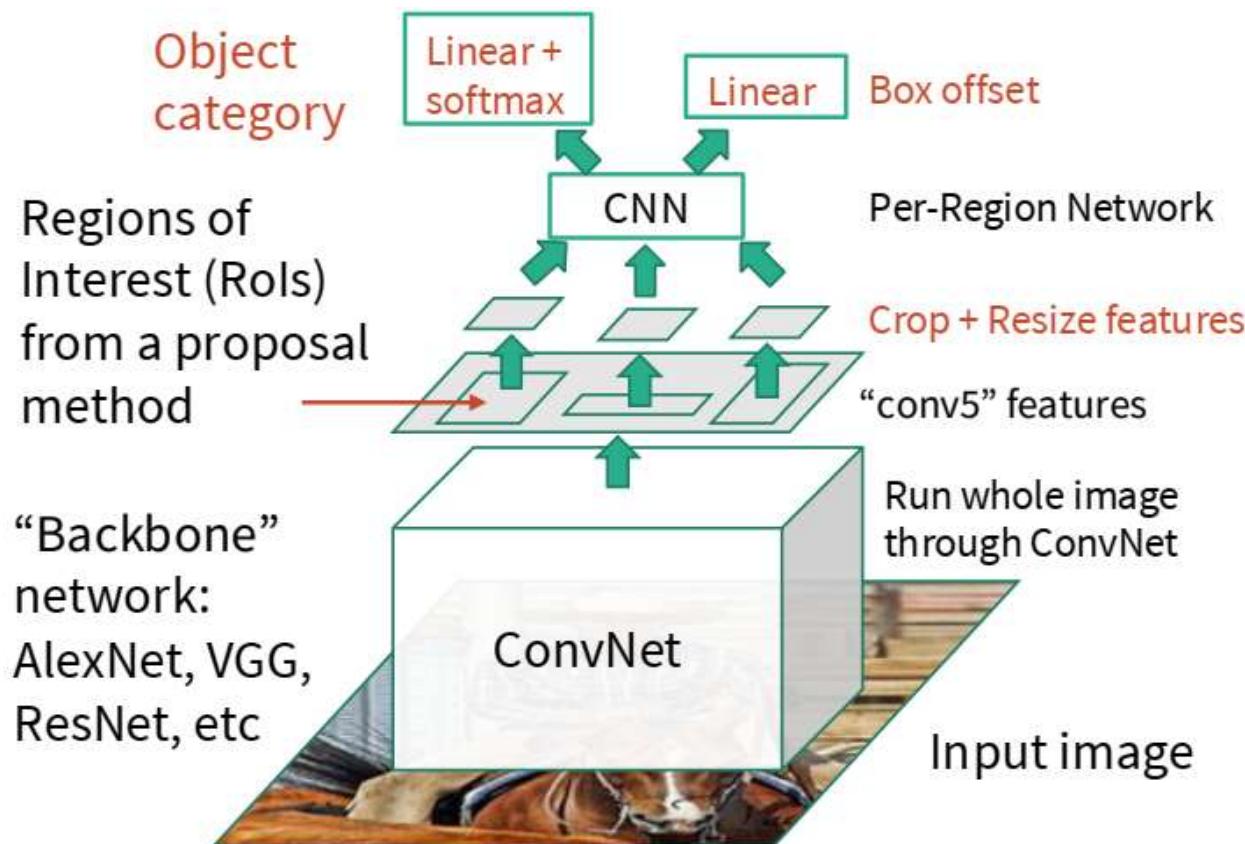
Fast R-CNN



Fast R-CNN



Fast R-CNN



Region Proposal Network



Input Image
(e.g. $3 \times 640 \times 480$)

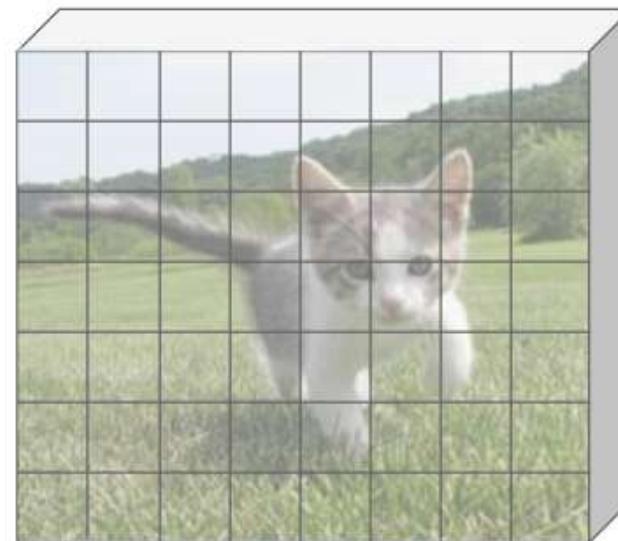
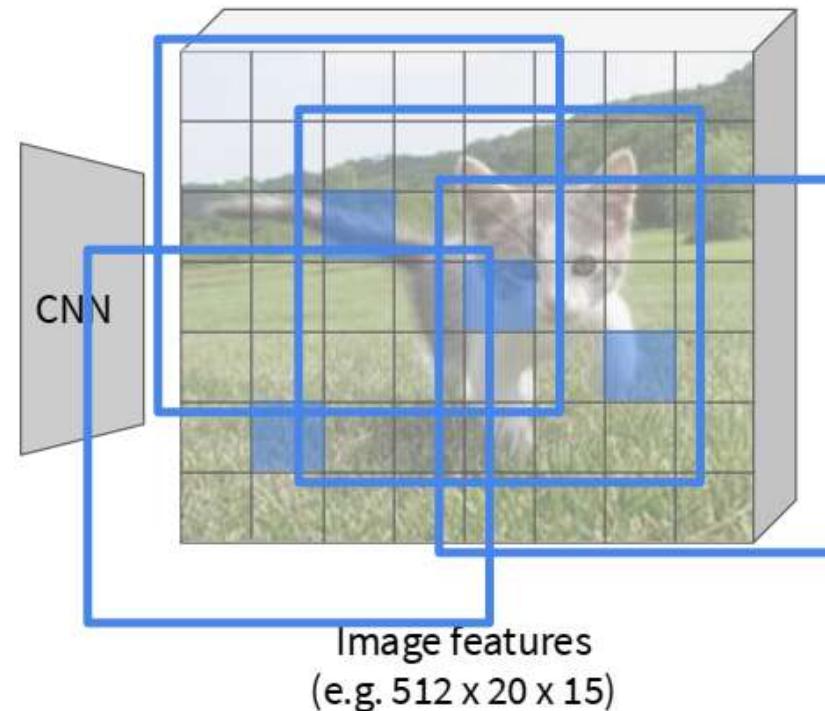


Image features
(e.g. $512 \times 20 \times 15$)

Region Proposal Network



Input Image
(e.g. 3 x 640 x 480)

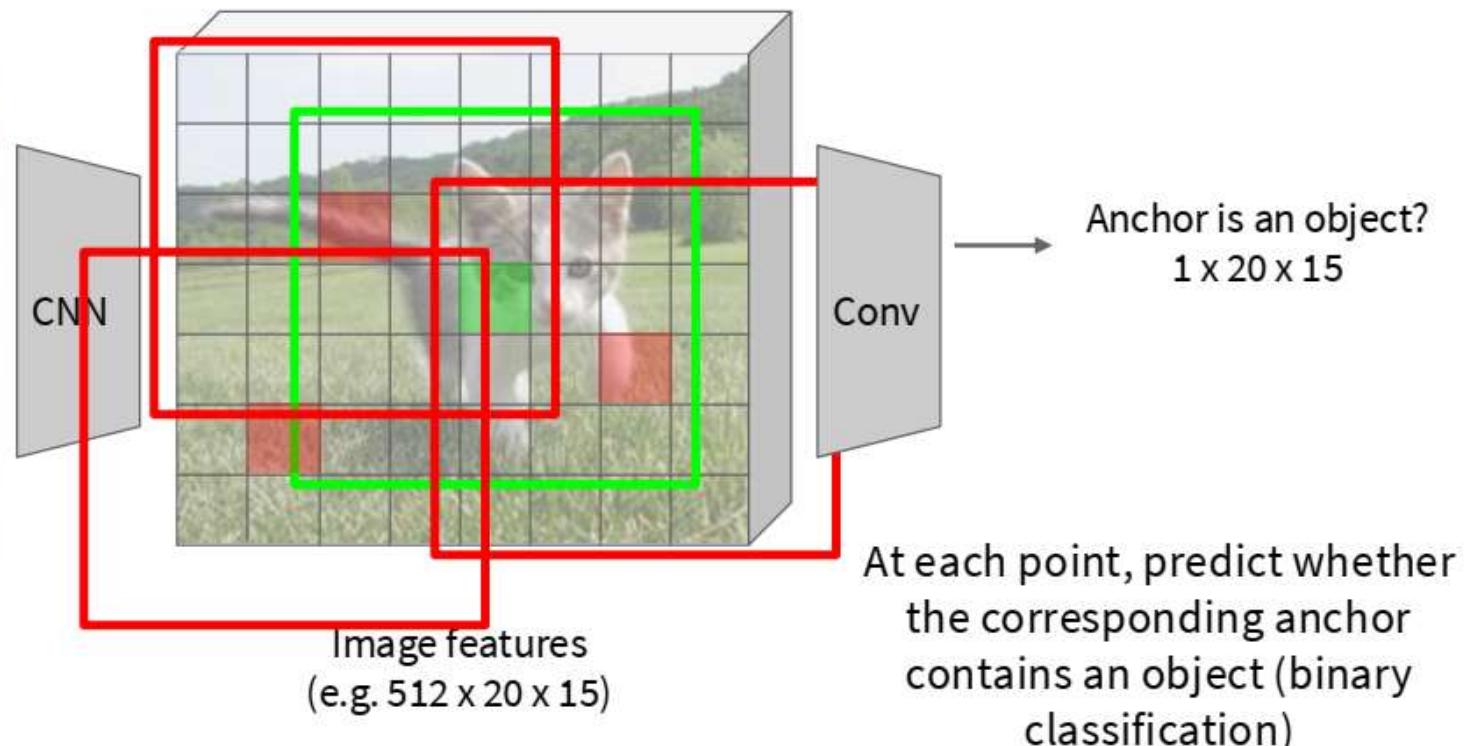


Imagina un anchor box de
Tamaño fijo en cada punto del
Feature map

Region Proposal Network



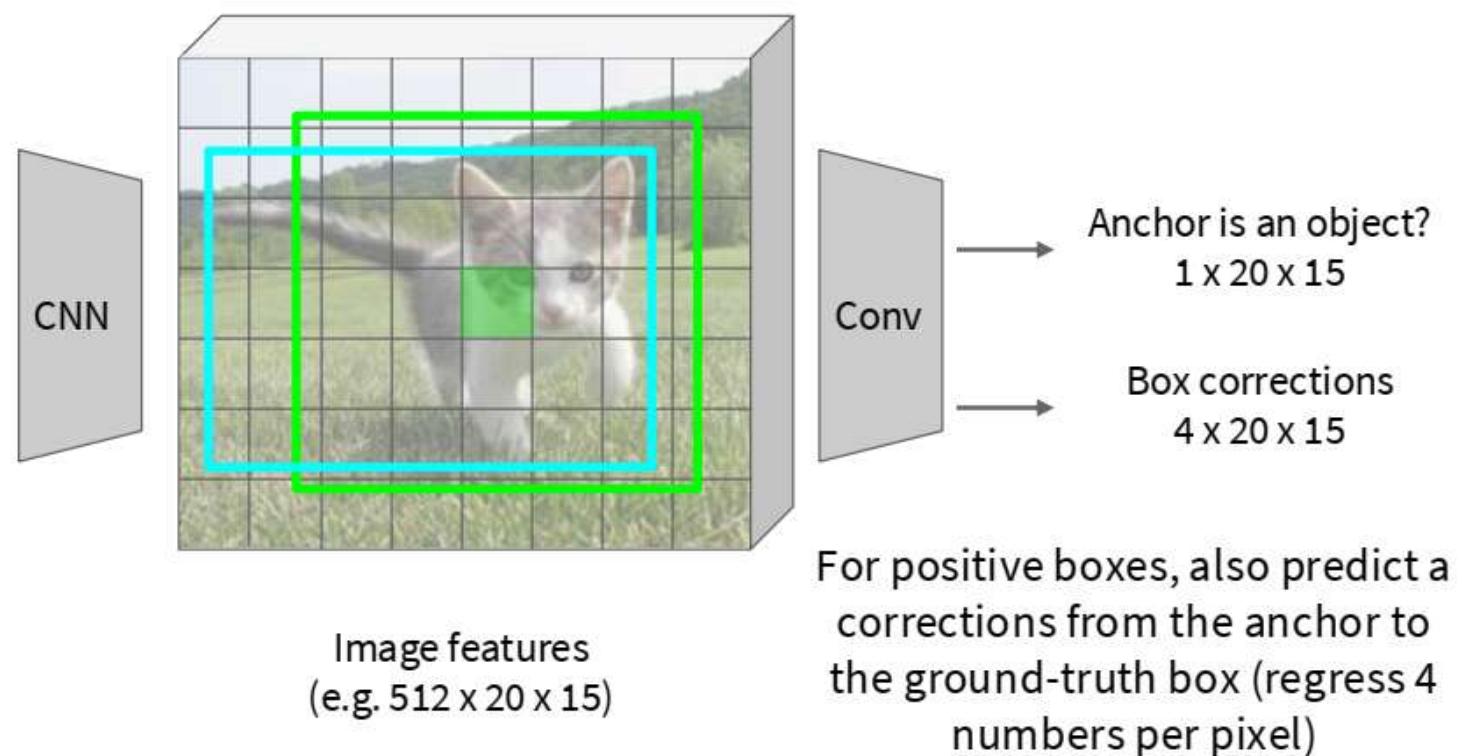
Input Image
(e.g. $3 \times 640 \times 480$)



Region Proposal Network



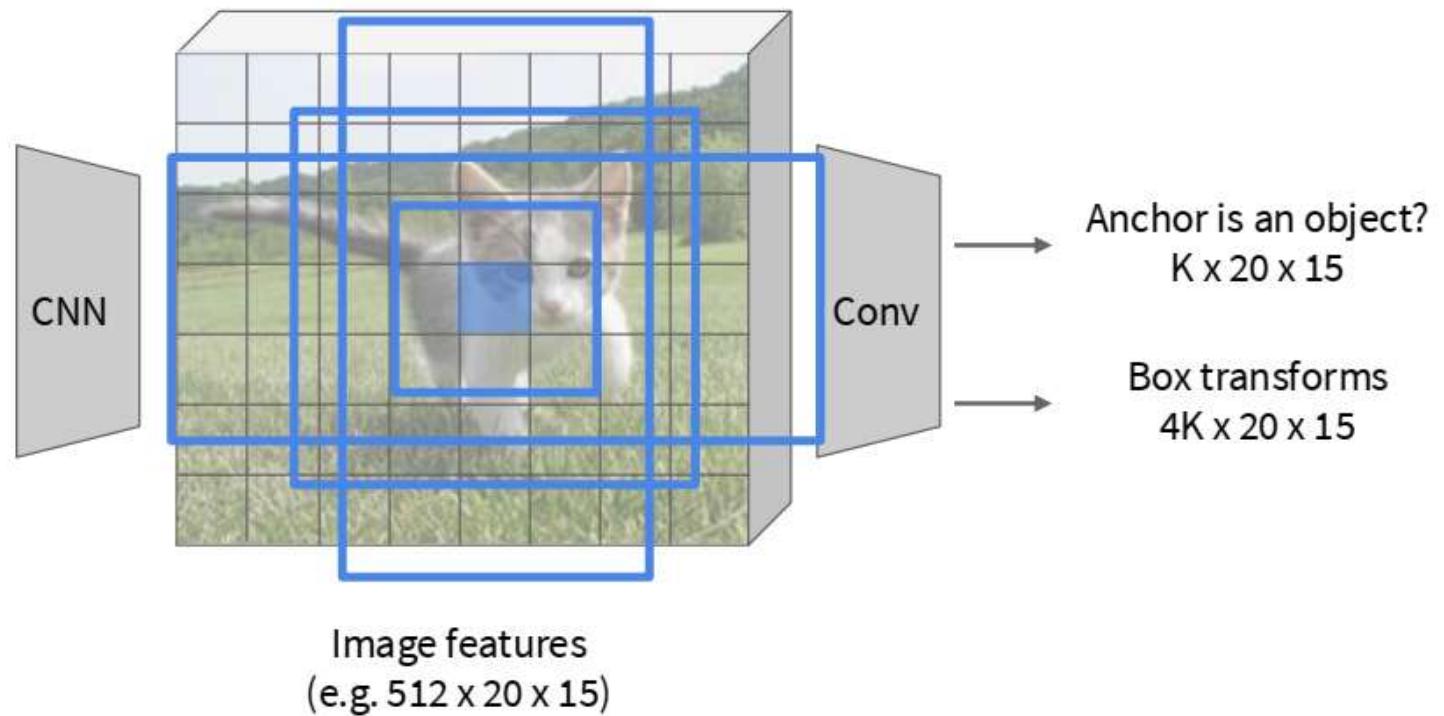
Input Image
(e.g. $3 \times 640 \times 480$)



Region Proposal Network



Input Image
(e.g. $3 \times 640 \times 480$)



Region Proposal Network



Input Image
(e.g. 3 x 640 x 480)

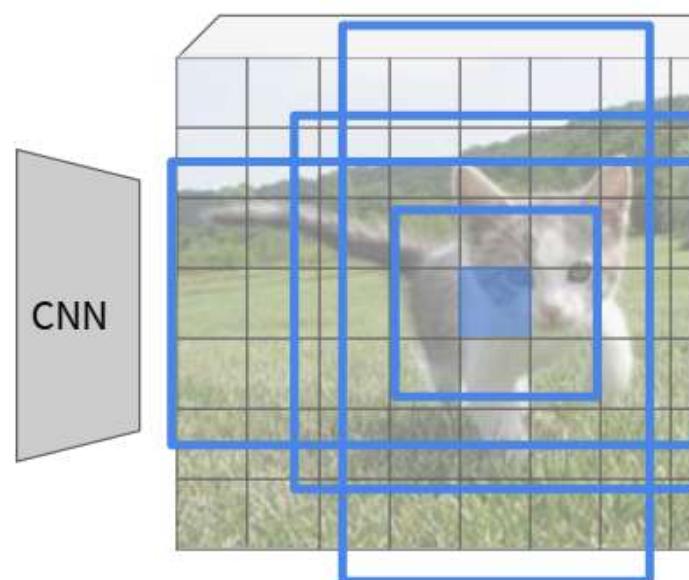


Image features
(e.g. 512 x 20 x 15)

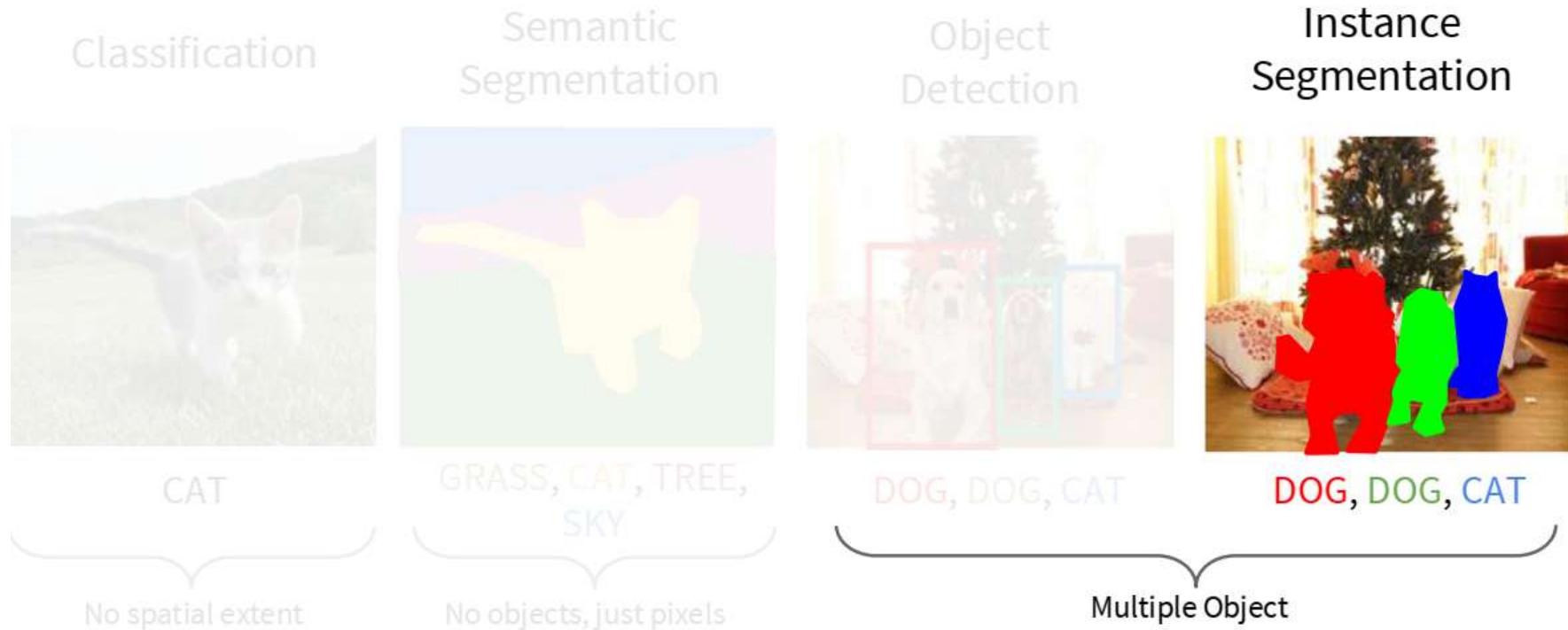
En la práctica se usan K diferentes Anchor boxes de diferentes tamaño Y escalas en cada punto

Anchor is an object?
 $K \times 20 \times 15$

Box transforms
 $4K \times 20 \times 15$

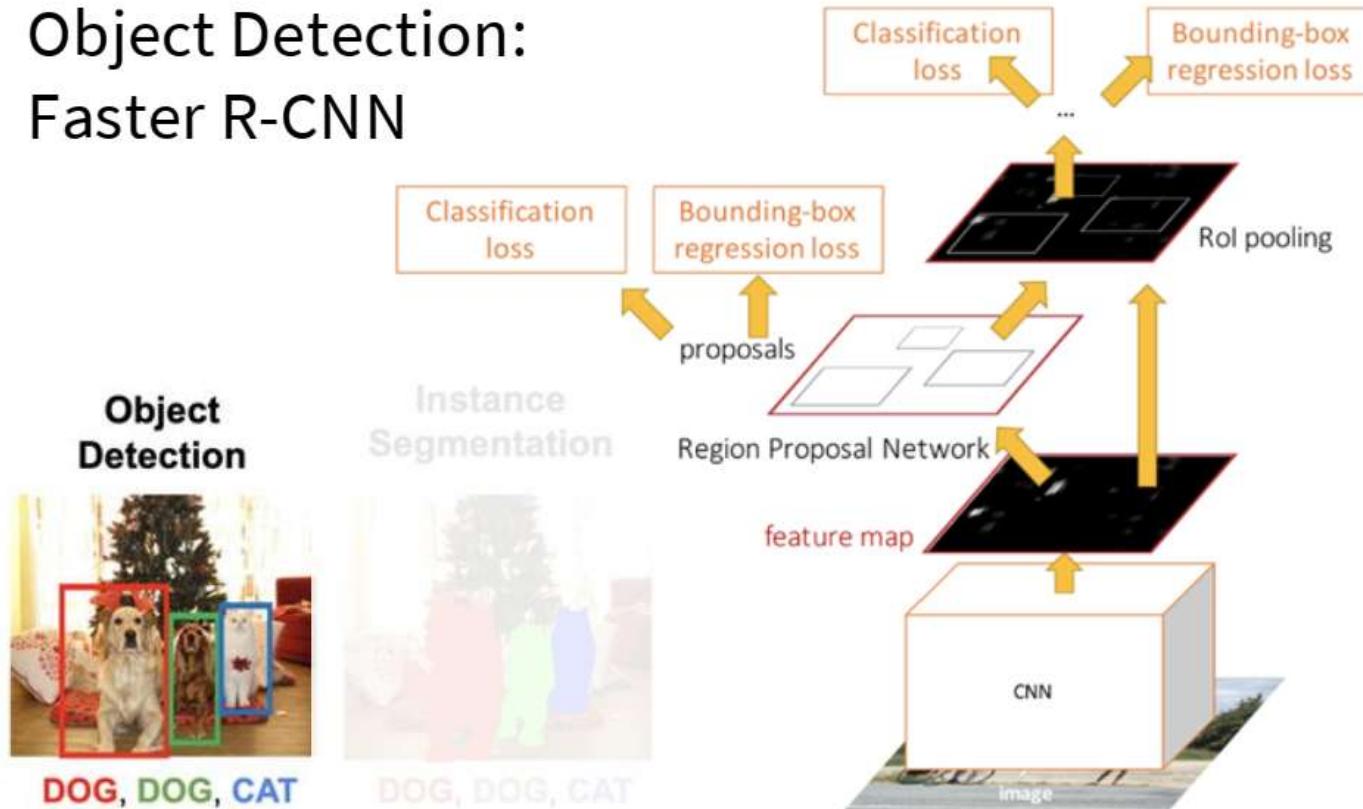
Sort the $K \times 20 \times 15$ boxes by their “objectness” score, take top ~ 300 as our proposals

Instance Segmentation



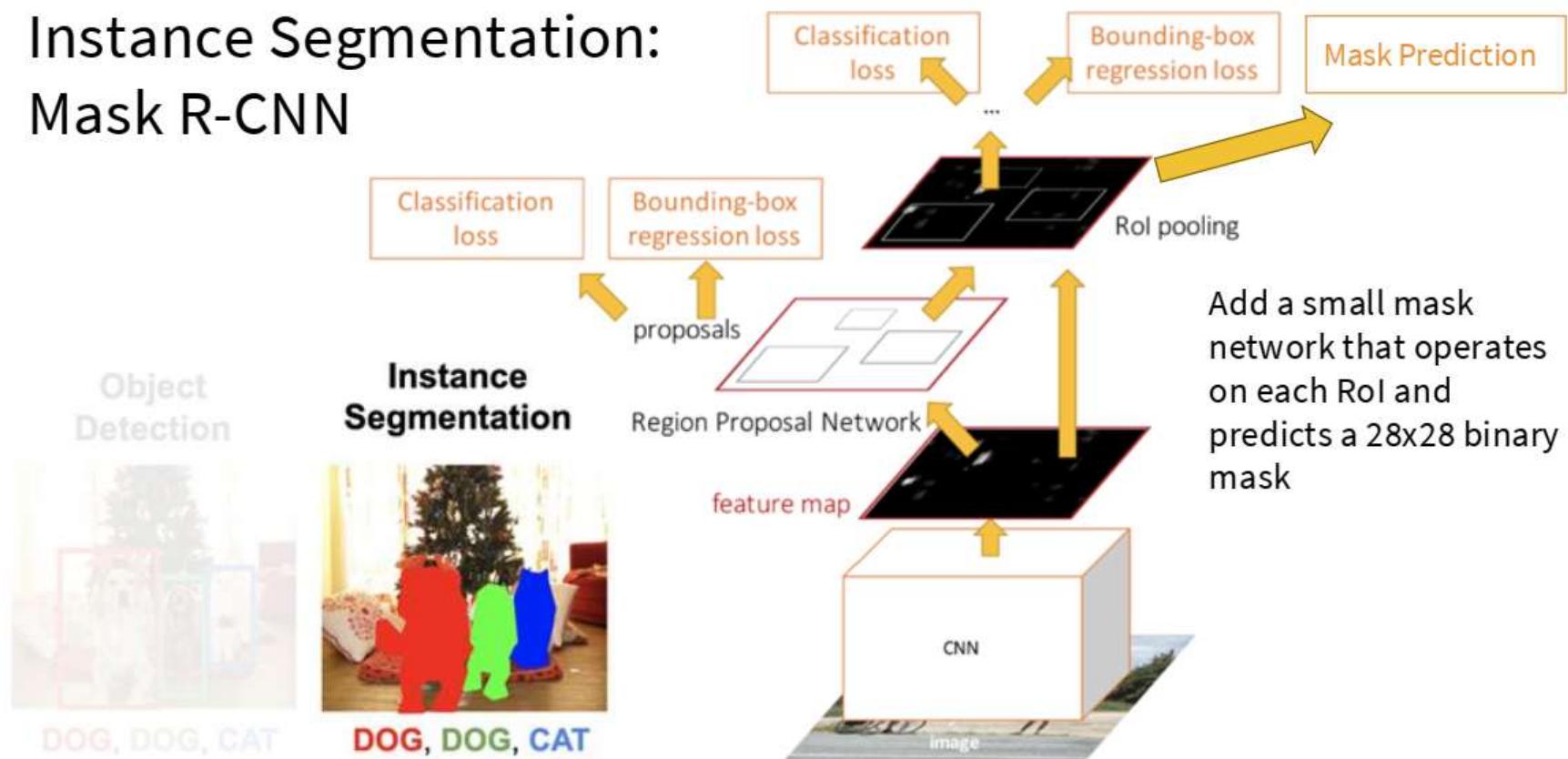
Instance Segmentation

Object Detection: Faster R-CNN



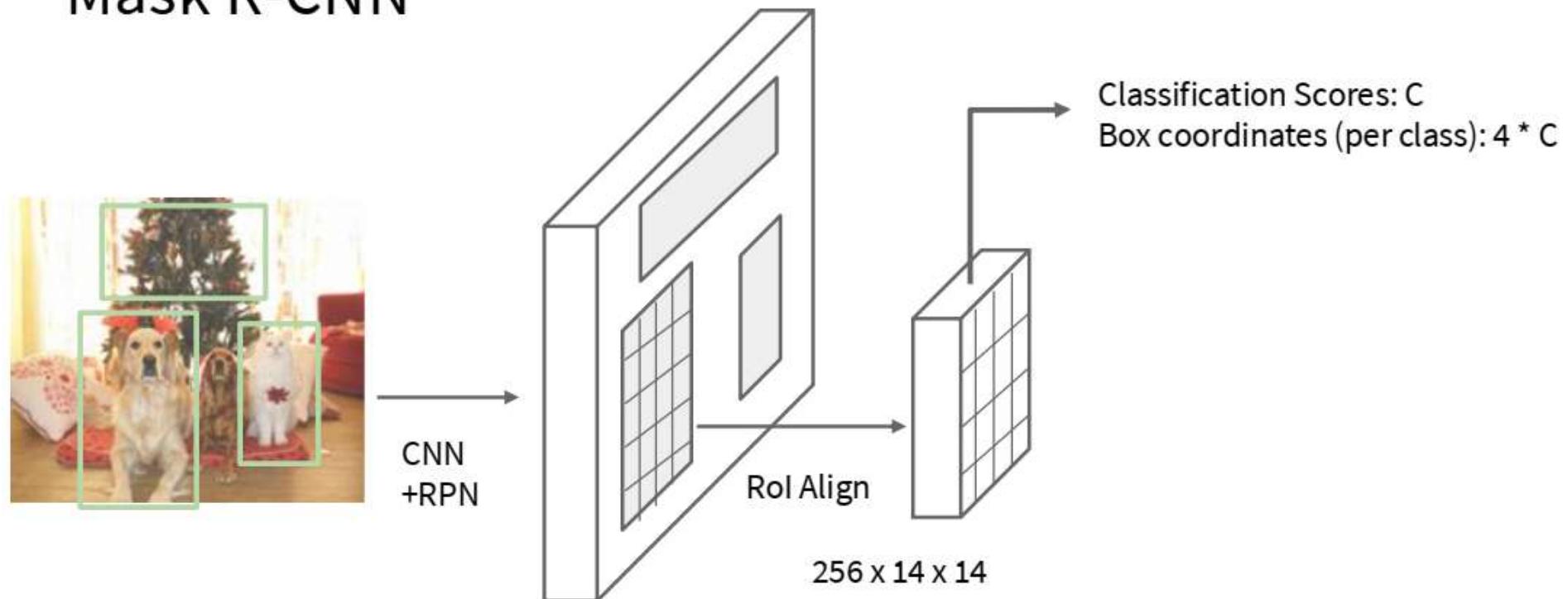
Instance Segmentation

Instance Segmentation: Mask R-CNN



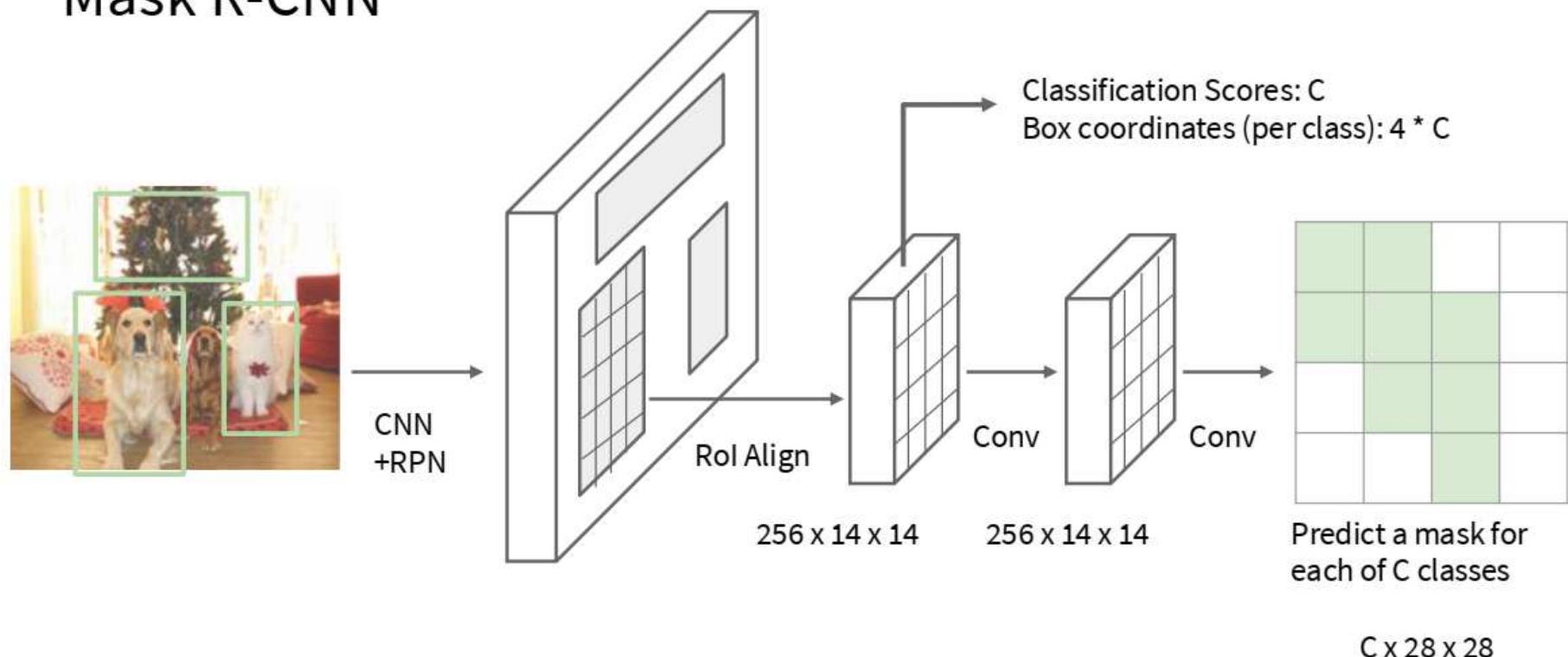
Instance Segmentation

Mask R-CNN



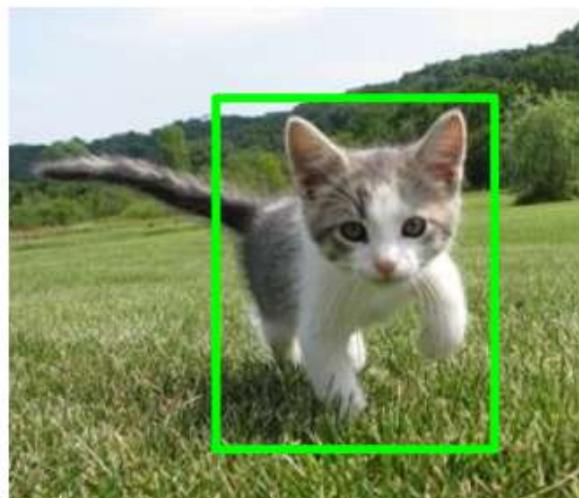
Instance Segmentation

Mask R-CNN



Instance Segmentation

ROI Pool



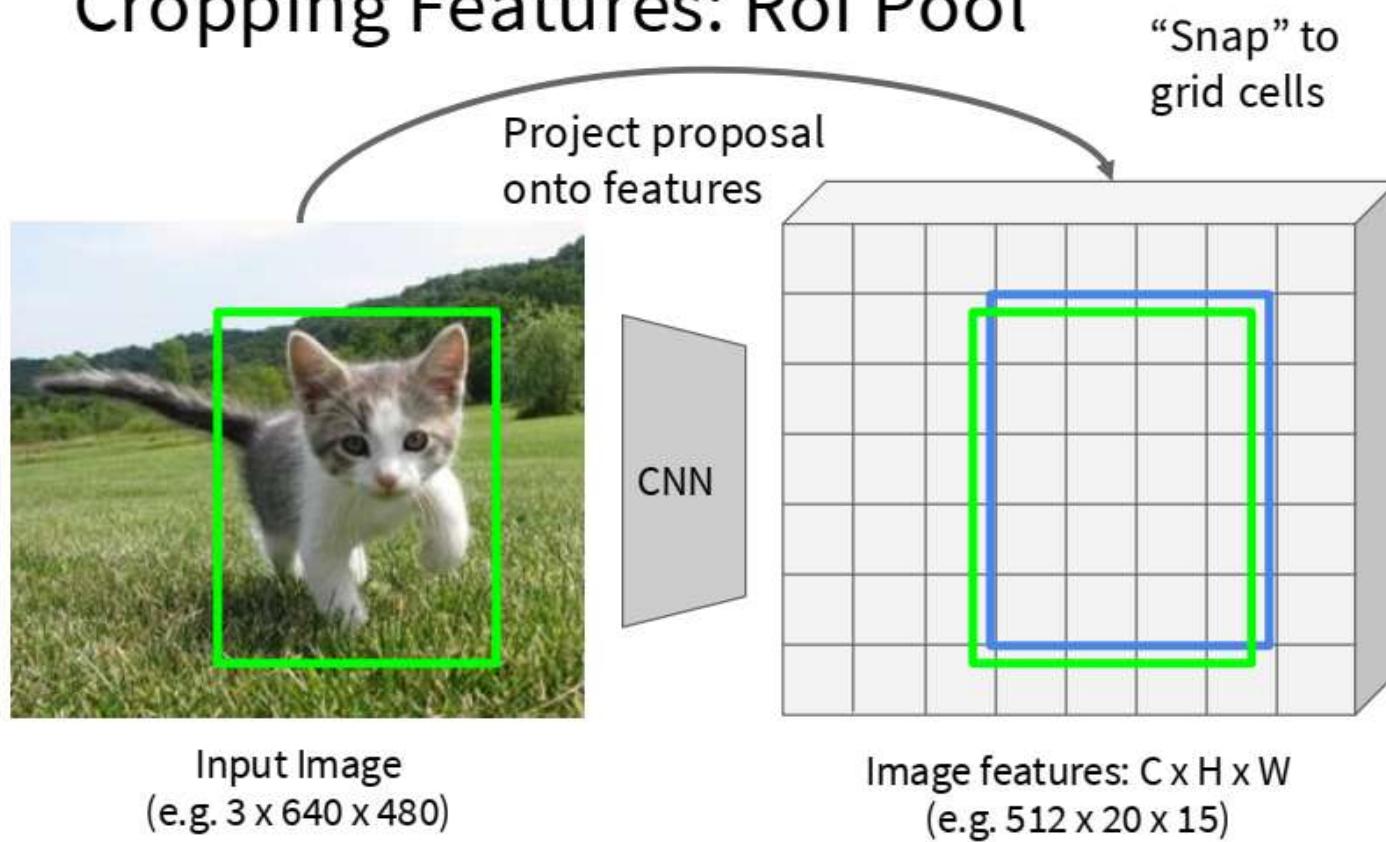
Input Image
(e.g. $3 \times 640 \times 480$)



Image features: $C \times H \times W$
(e.g. $512 \times 20 \times 15$)

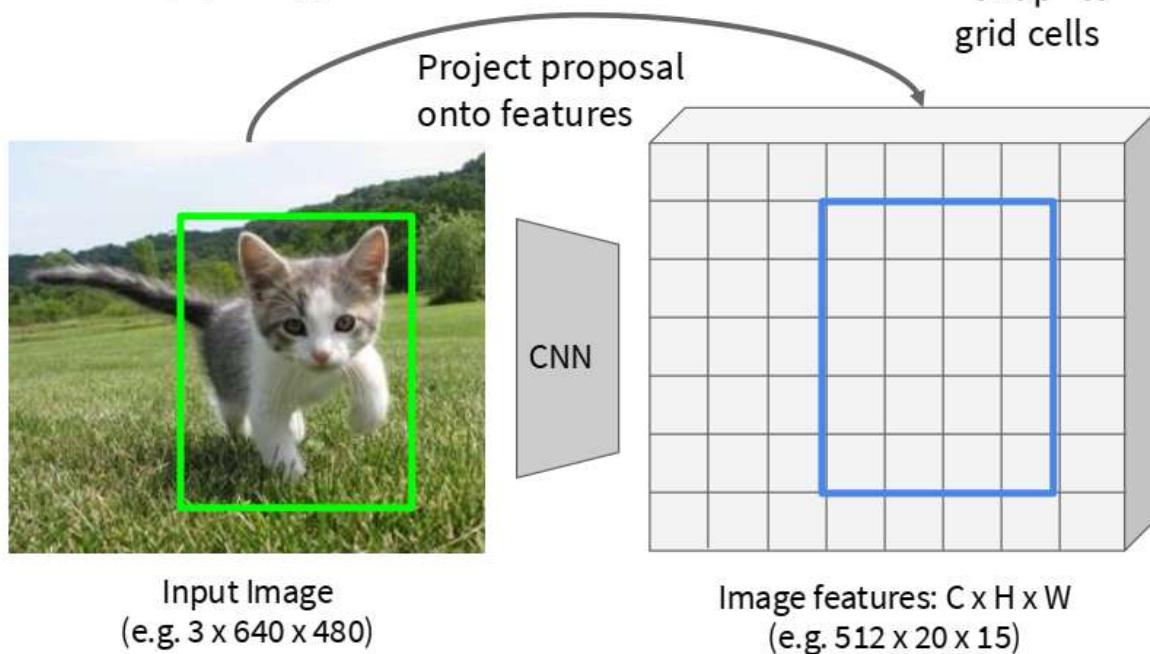
Instance Segmentation

Cropping Features: RoI Pool



Instance Segmentation

Cropping Features: RoI Pool

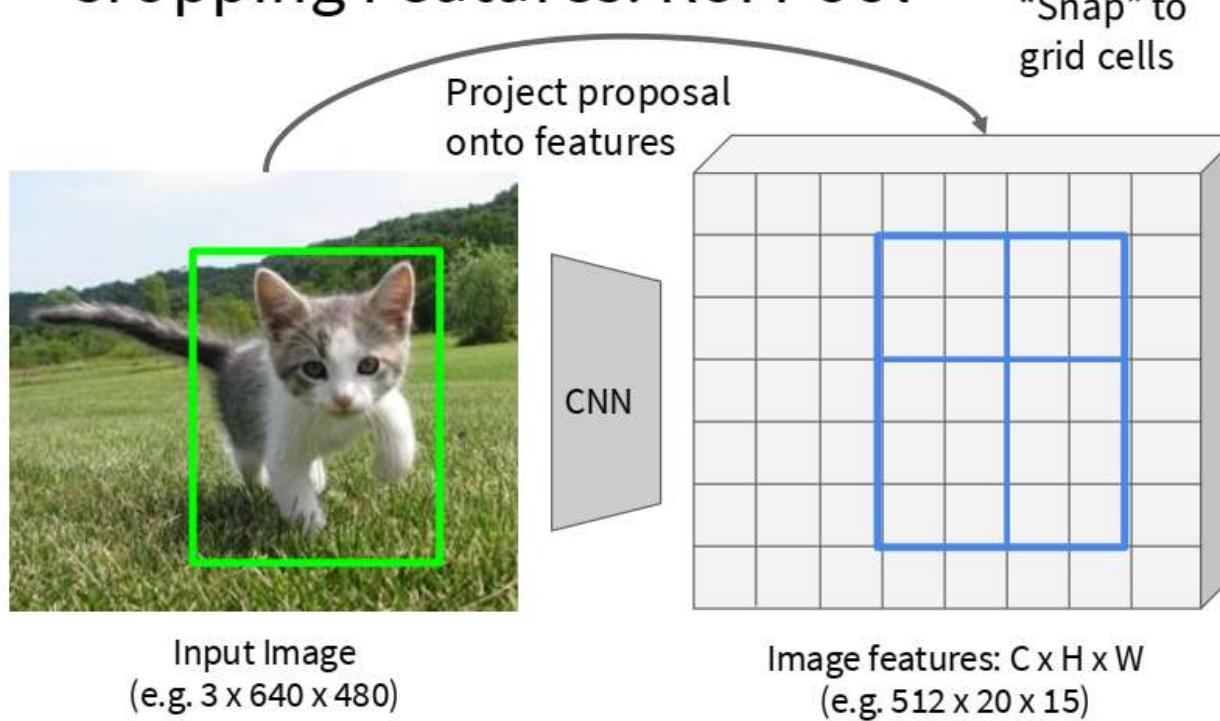


Cómo transformamos la región 512x5x4 a 512x2x2, por ejemplo?

Instance Segmentation

Dividir en grid de 2x2 de regiones más o menos iguales

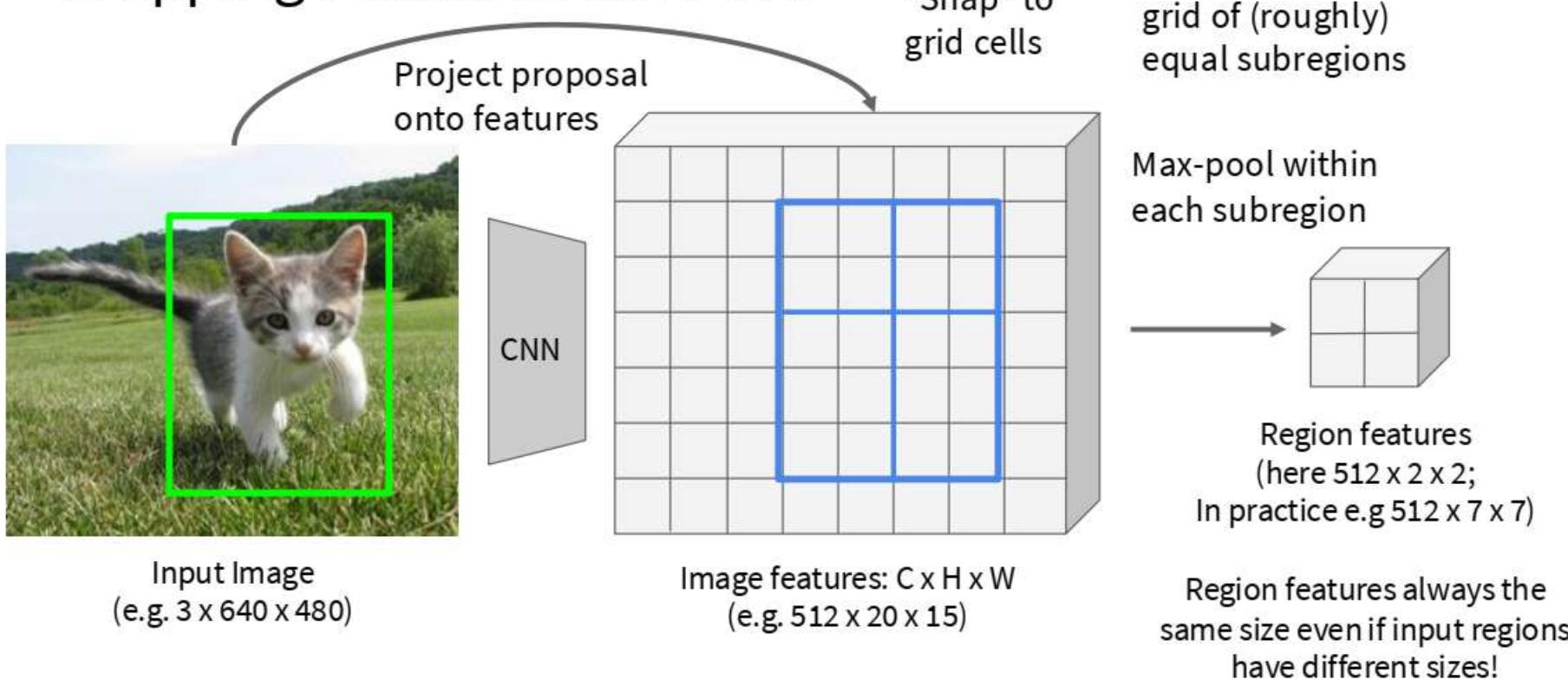
Cropping Features: RoI Pool



Cómo transformamos la región 512x5x4 a 512x2x2, por ejemplo?

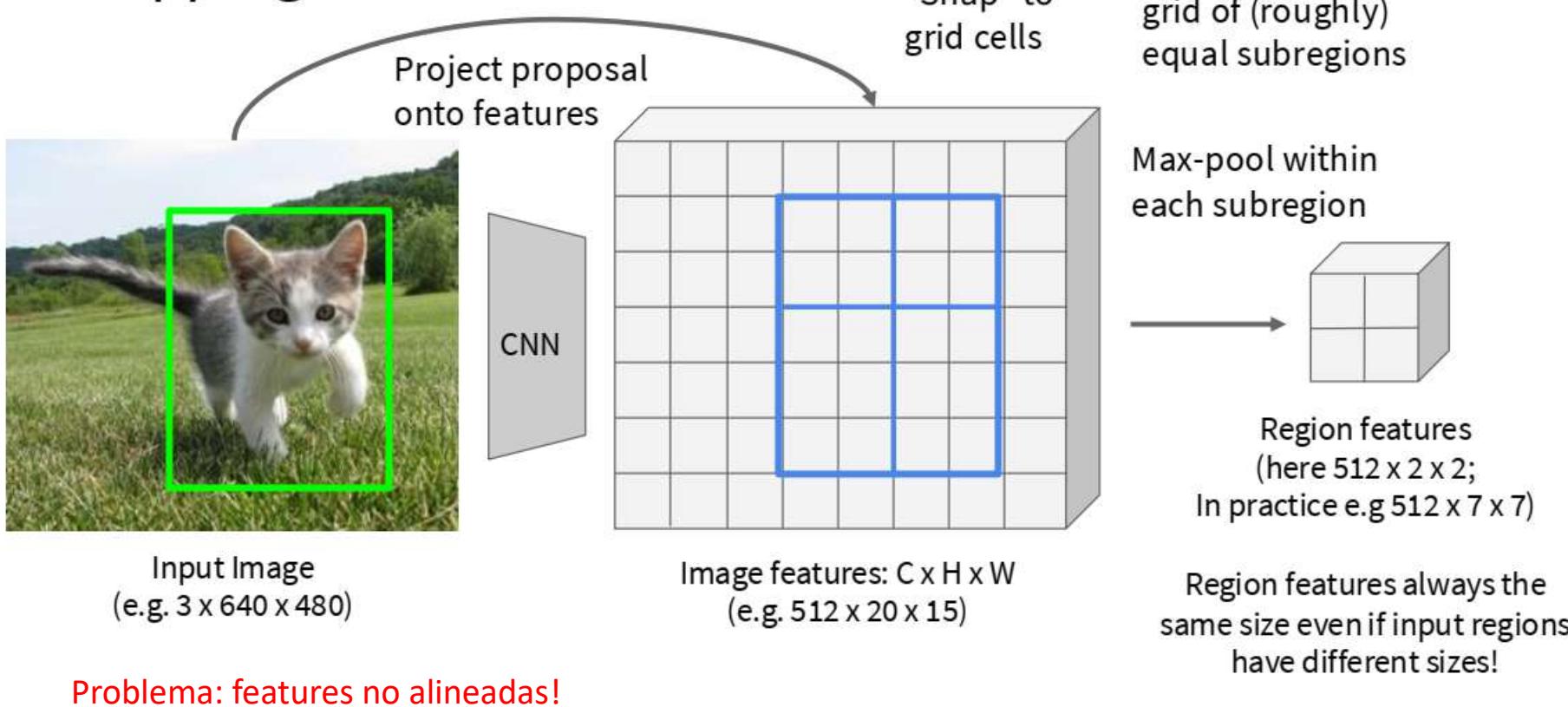
Instance Segmentation

Cropping Features: RoI Pool



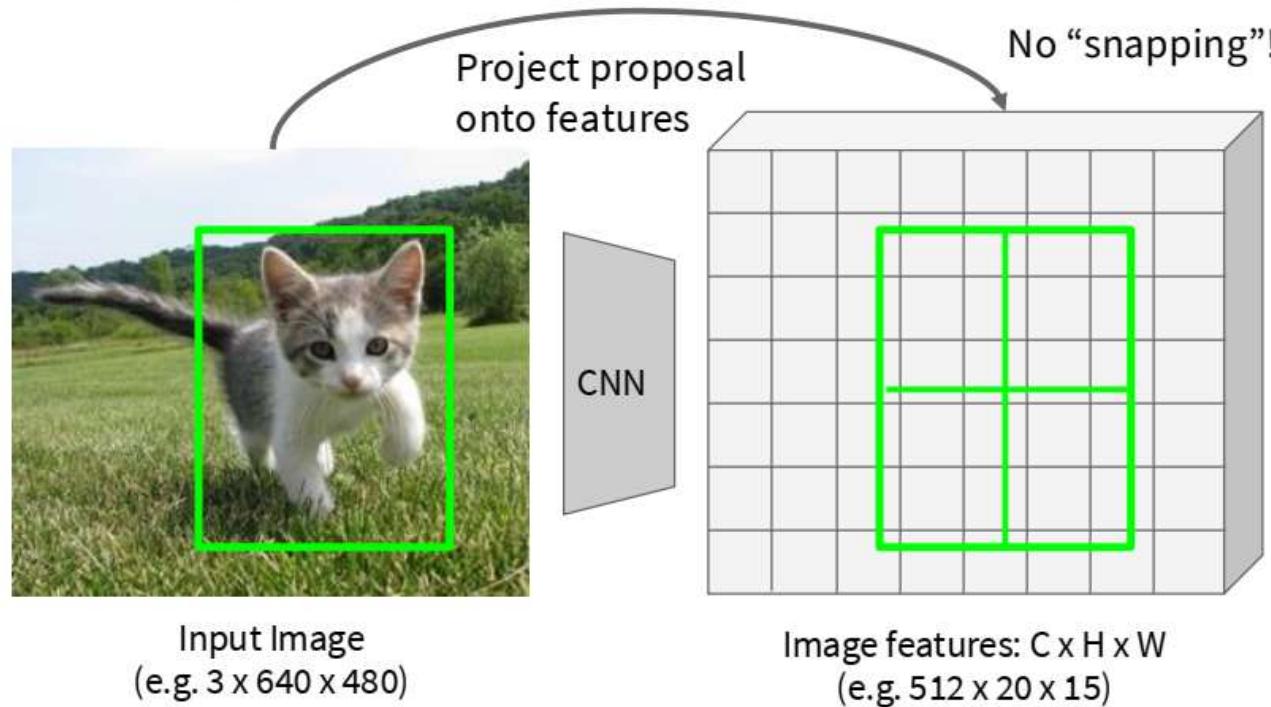
Instance Segmentation

Cropping Features: RoI Pool



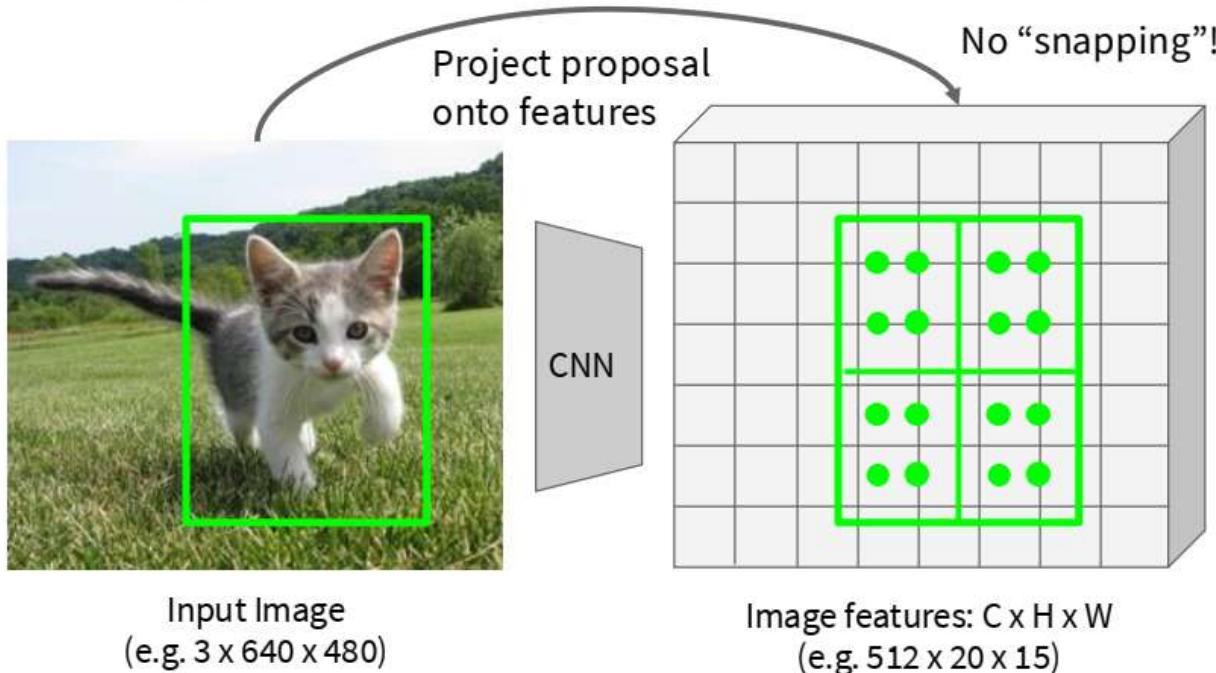
Instance Segmentation

Cropping Features: RoI Align



Instance Segmentation

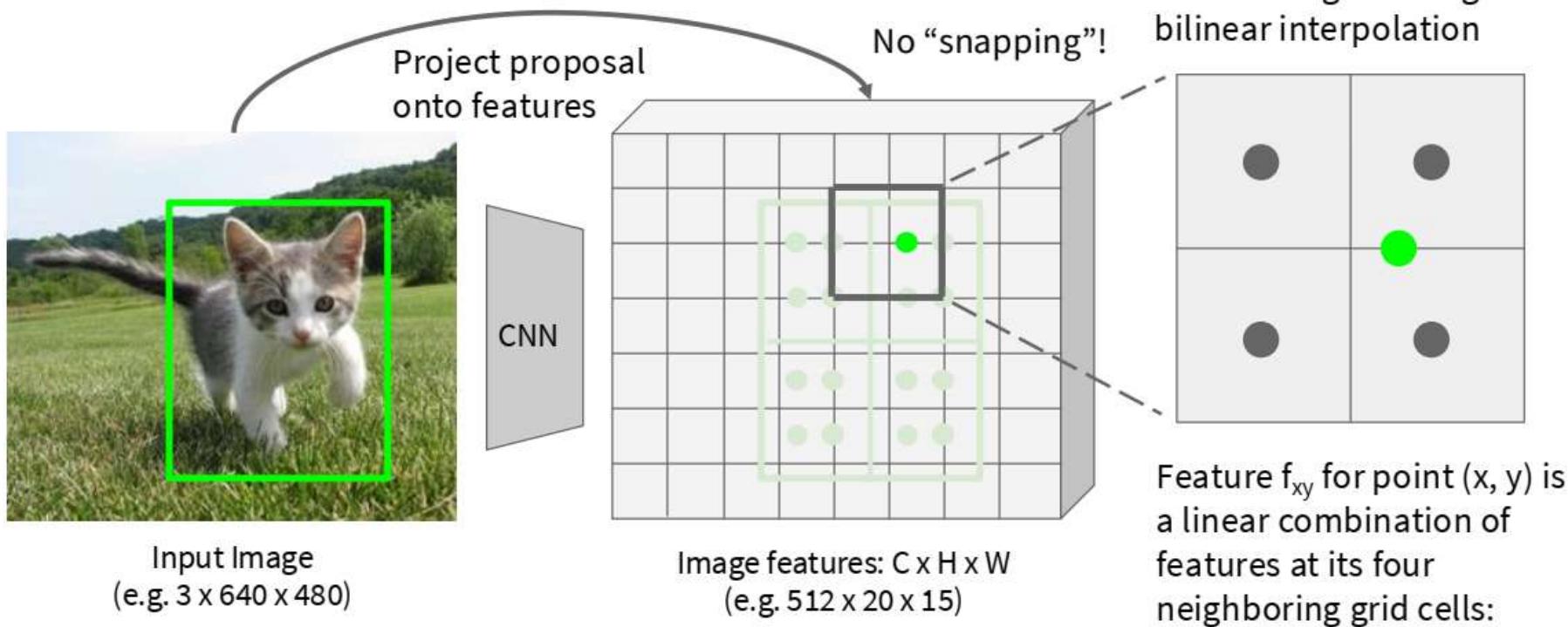
Cropping Features: RoI Align



Sample at regular points in
each subregion using
bilinear interpolation

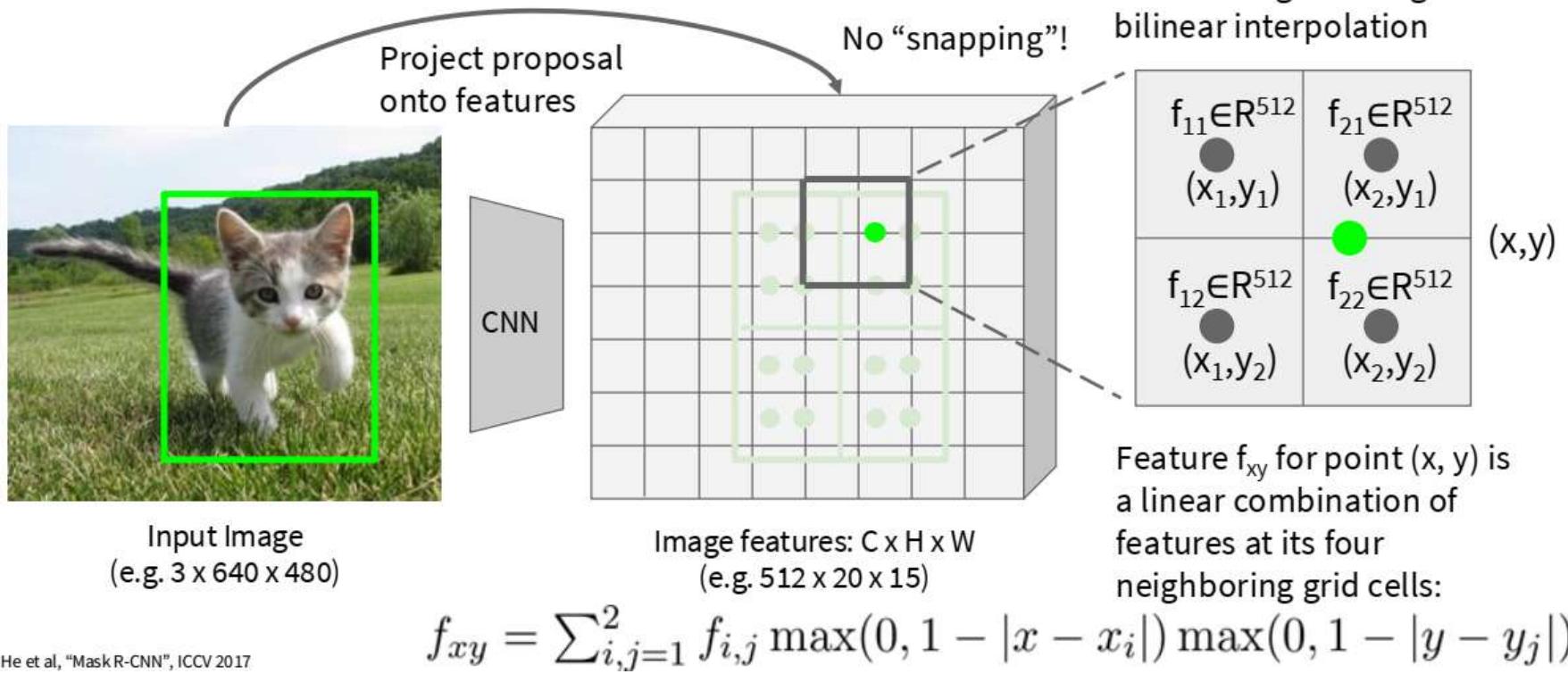
Instance Segmentation

Cropping Features: RoI Align



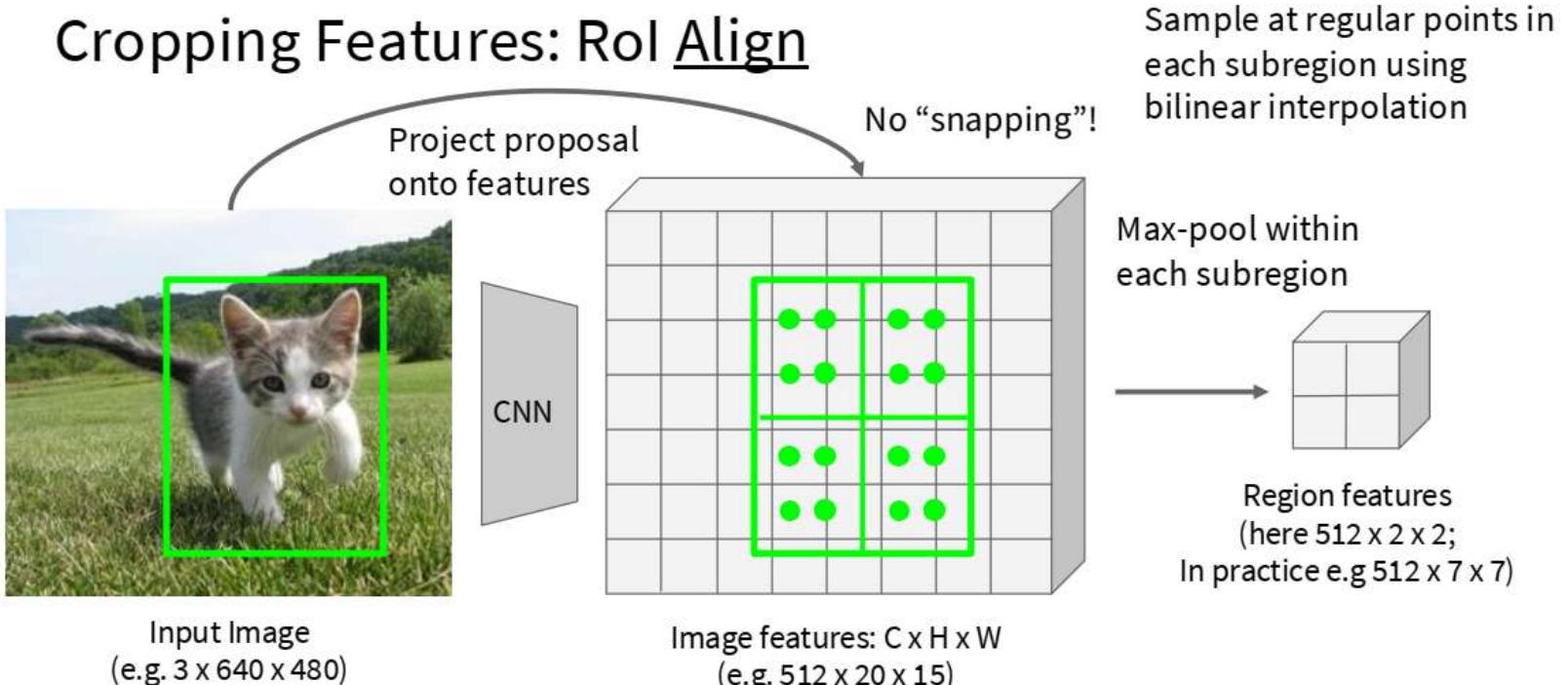
Instance Segmentation

Cropping Features: RoI Align



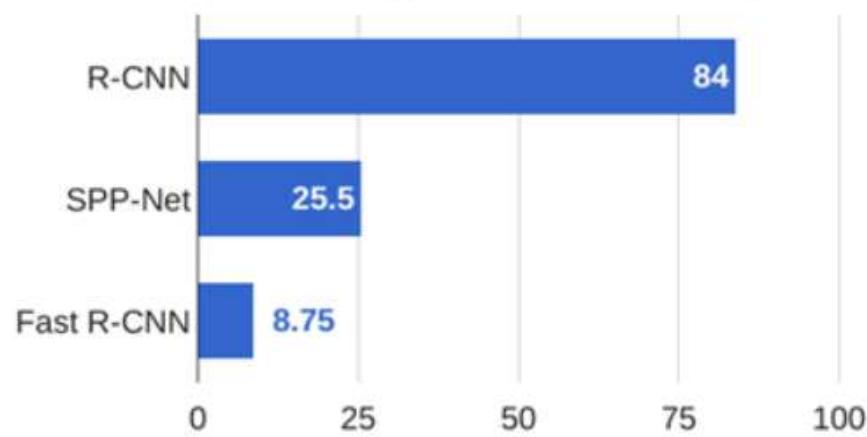
He et al, "Mask R-CNN", ICCV 2017

Instance Segmentation

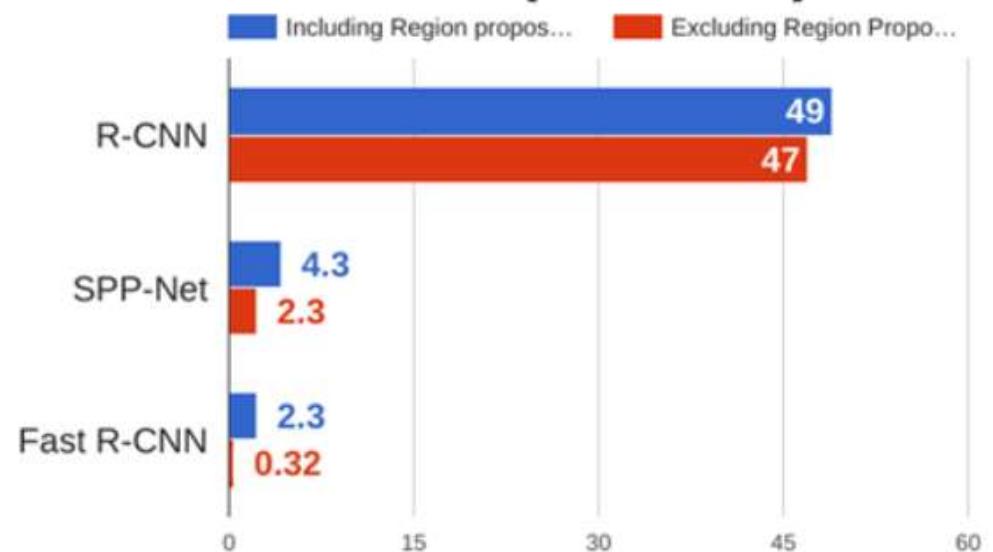


Timing

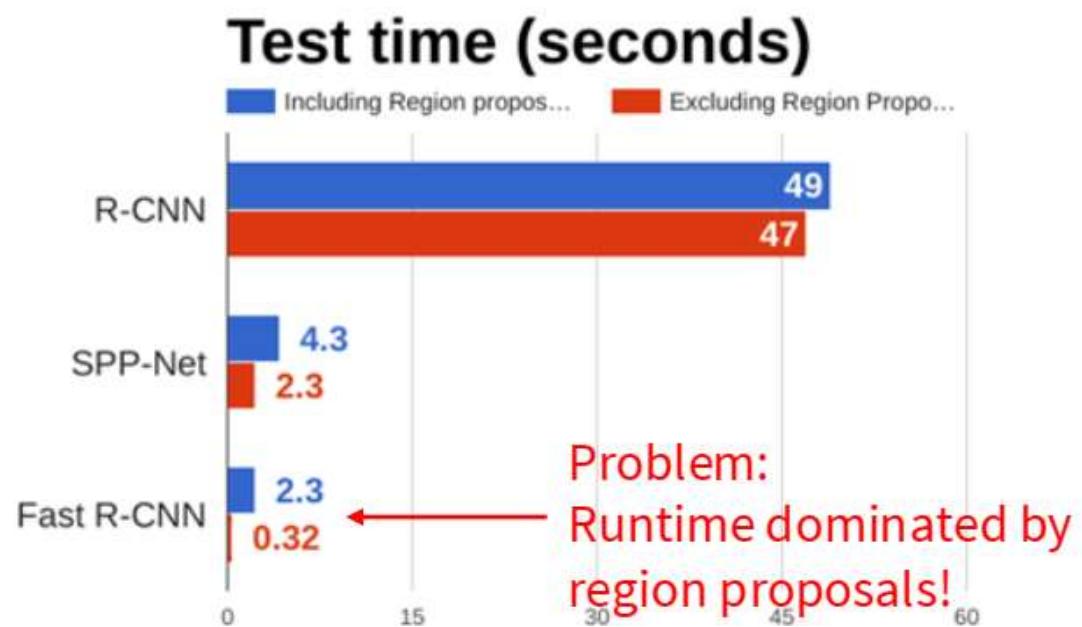
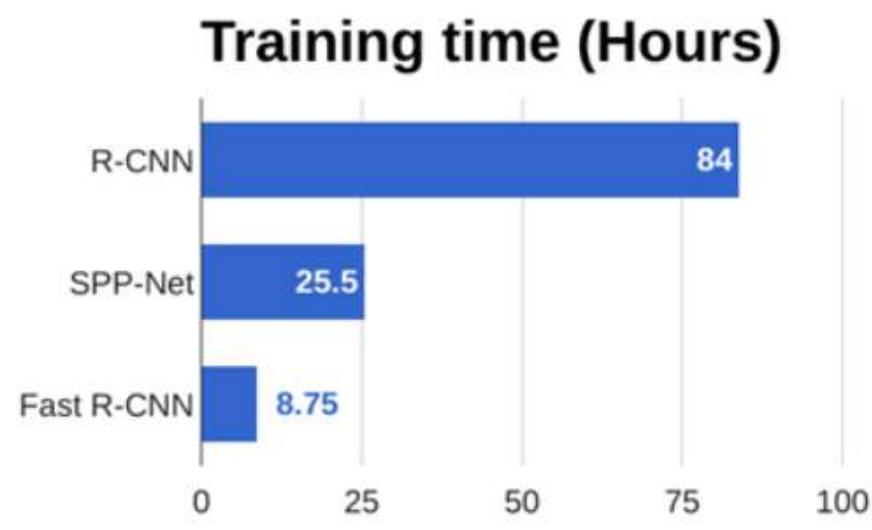
Training time (Hours)



Test time (seconds)



Timing

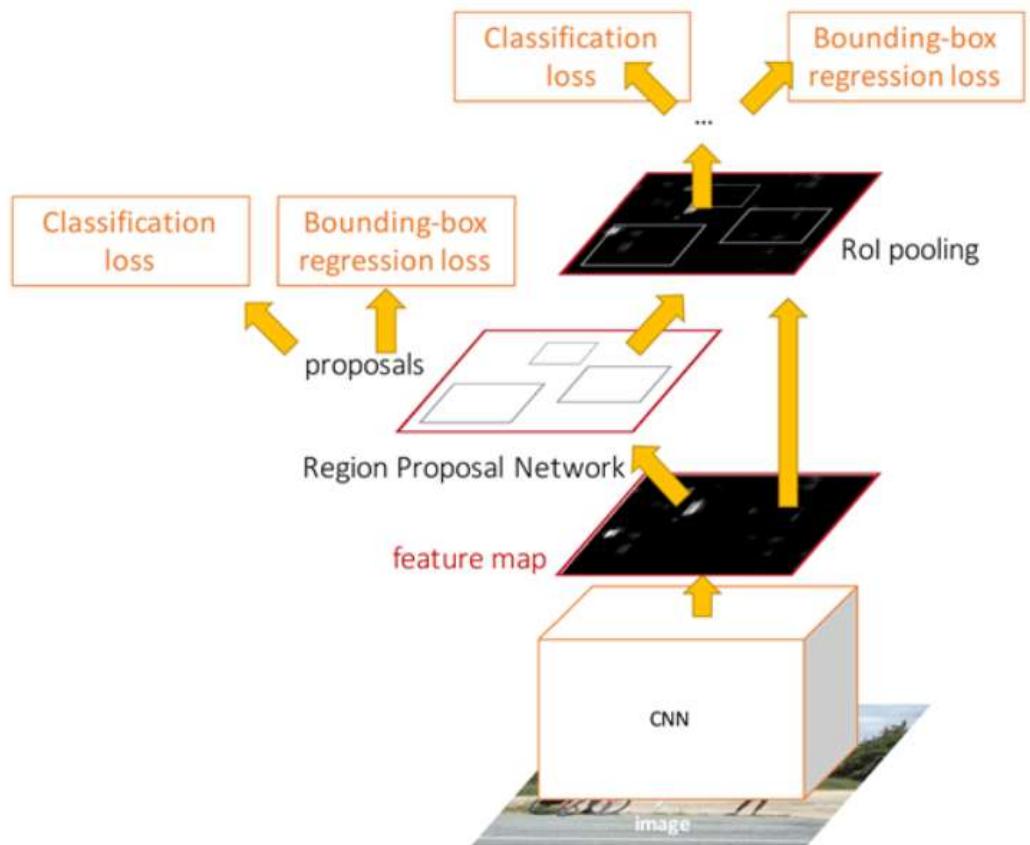


Faster R-CNN

Que la CNN se encargue de los proposals!

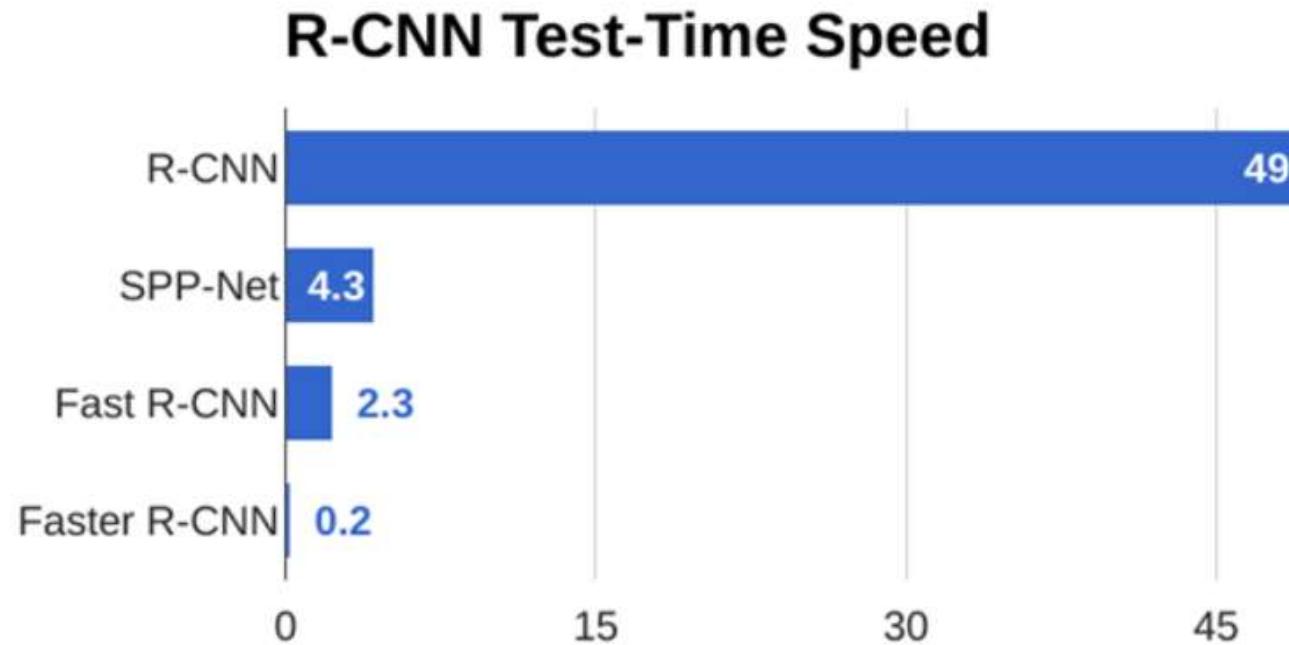
Se entrena en conjunto:

- RPN clasifica objeto/no objeto
- RPN regresión coordenadas box
- Score de clasificación final
- Coordenadas de box final



Faster R-CNN

Que la CNN se encargue de los proposals!



Faster R-CNN:

Make CNN do proposals!

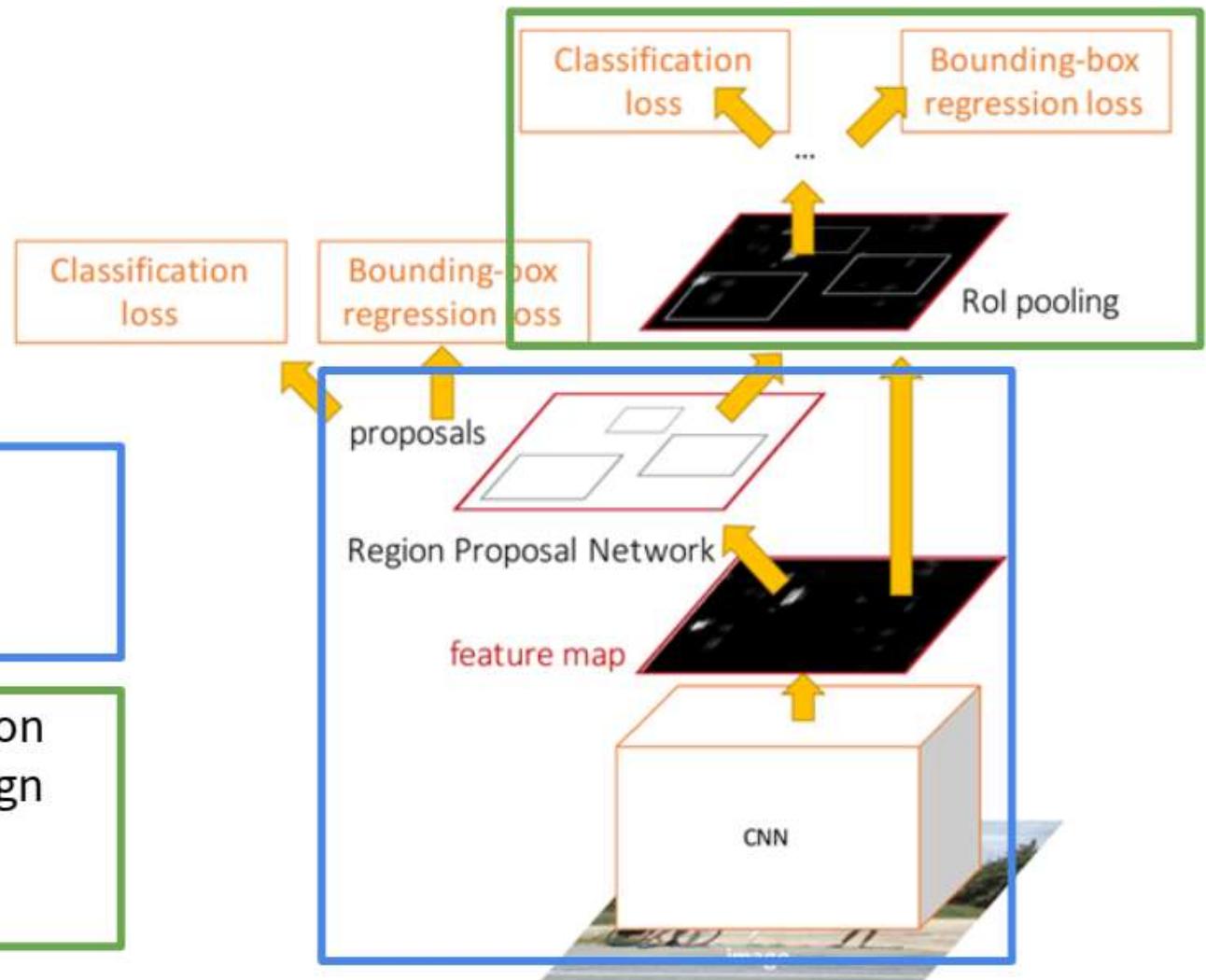
Faster R-CNN is a
Two-stage object detector

First stage: Run once per image

- Backbone network
- Region proposal network

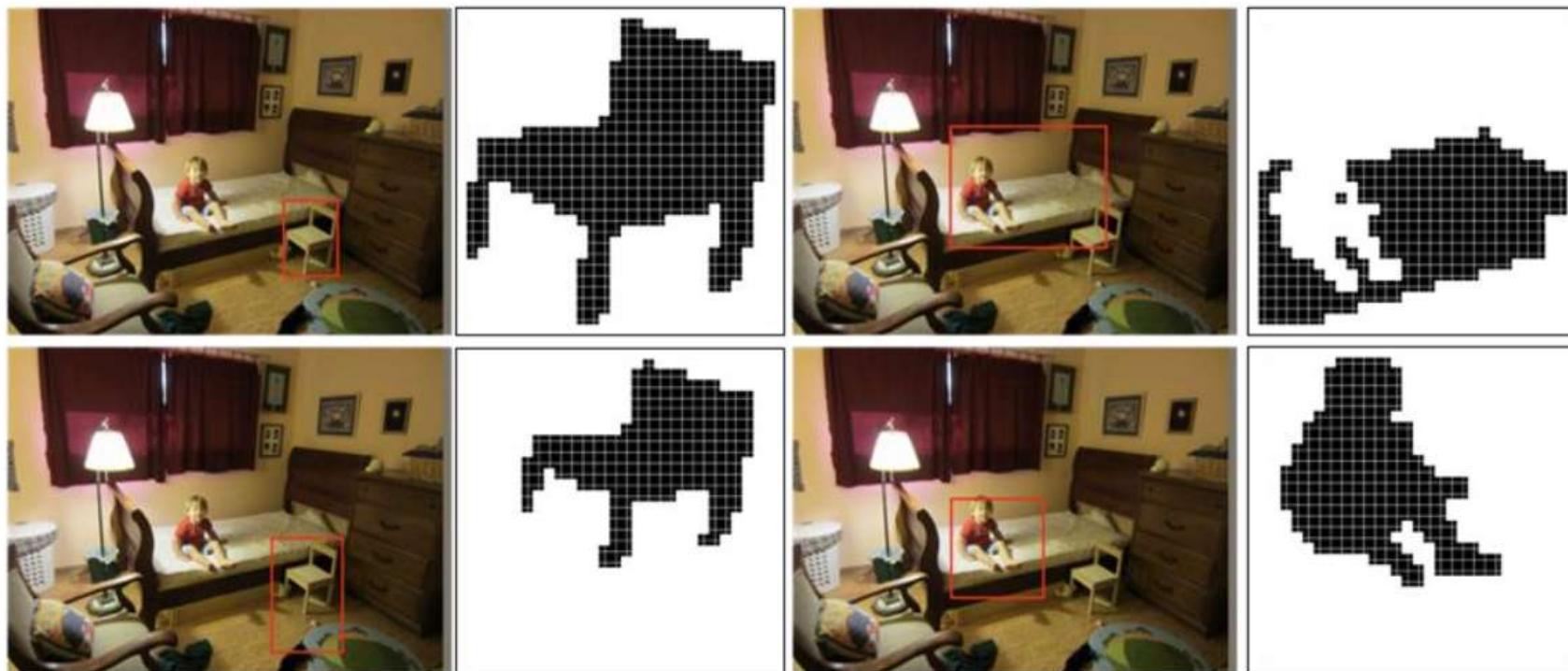
Second stage: Run once per region

- Crop features: RoI pool / align
- Predict object class
- Prediction bbox offset

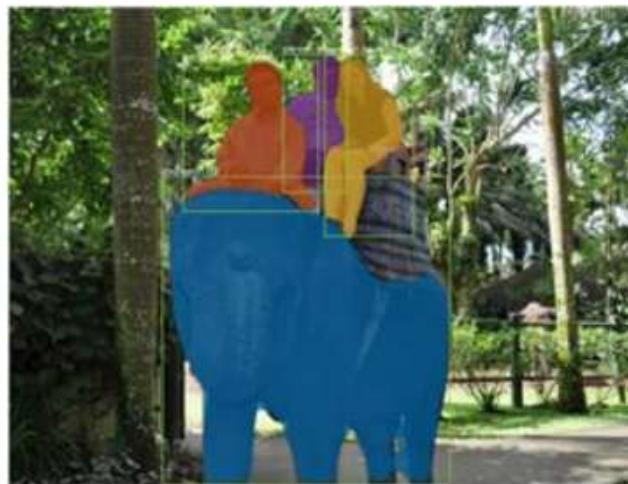


Instance Segmentation

Mask R-CNN: Example Mask Training Targets



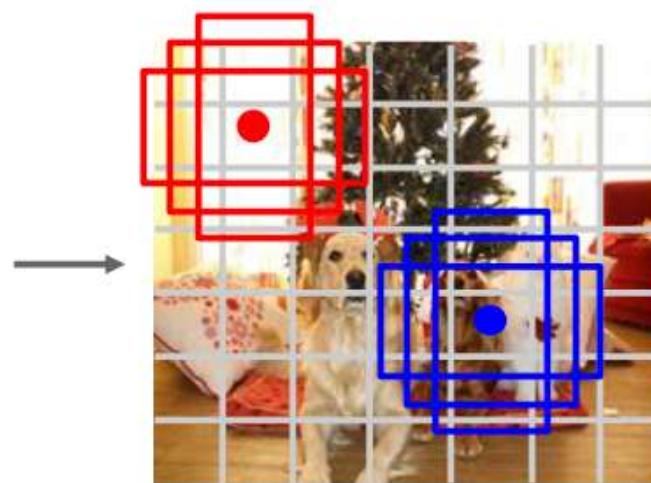
Instance Segmentation



Detección de Objetos Single-Stage



Input image
 $3 \times H \times W$



Divide image into grid
 7×7

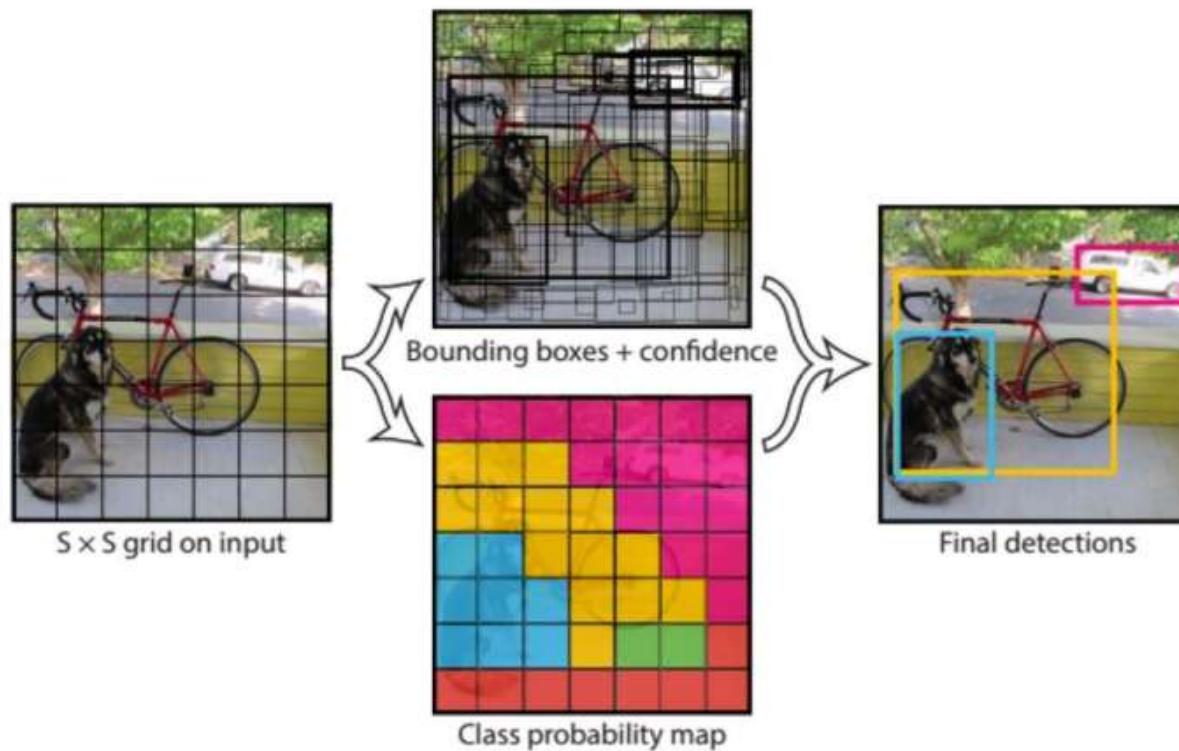
Image a set of base boxes
centered at each grid cell
Here $B = 3$

Dentro de cada celda:

- Regresión de cada B base box a un box final con 5 números ($dx, dy, dh, dw, confidence$)
- Predecir scores para cada una de las C clases (background es una clase).
- Parece una RPN, pero específico para las clases

Output:
 $7 \times 7 \times (5 * B + C)$

YOLO (You Only Look Once)



YOLO (You Only Look Once)



SxS Grid

YOLO (You Only Look Once)

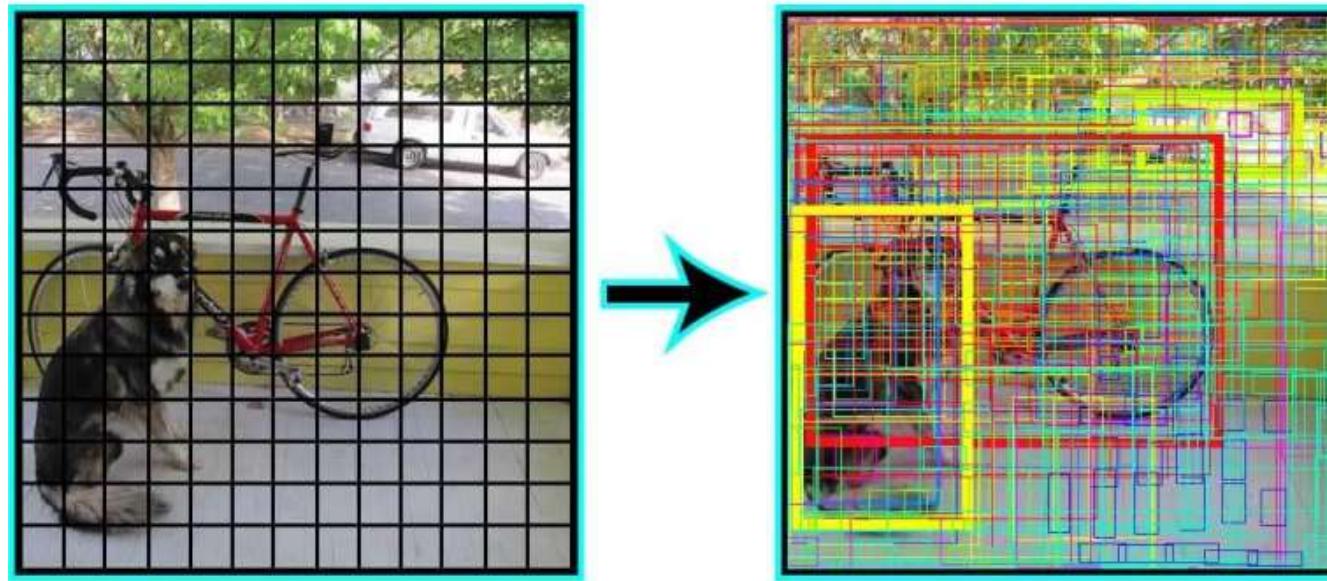


SxS Grid

Para cada box generar:

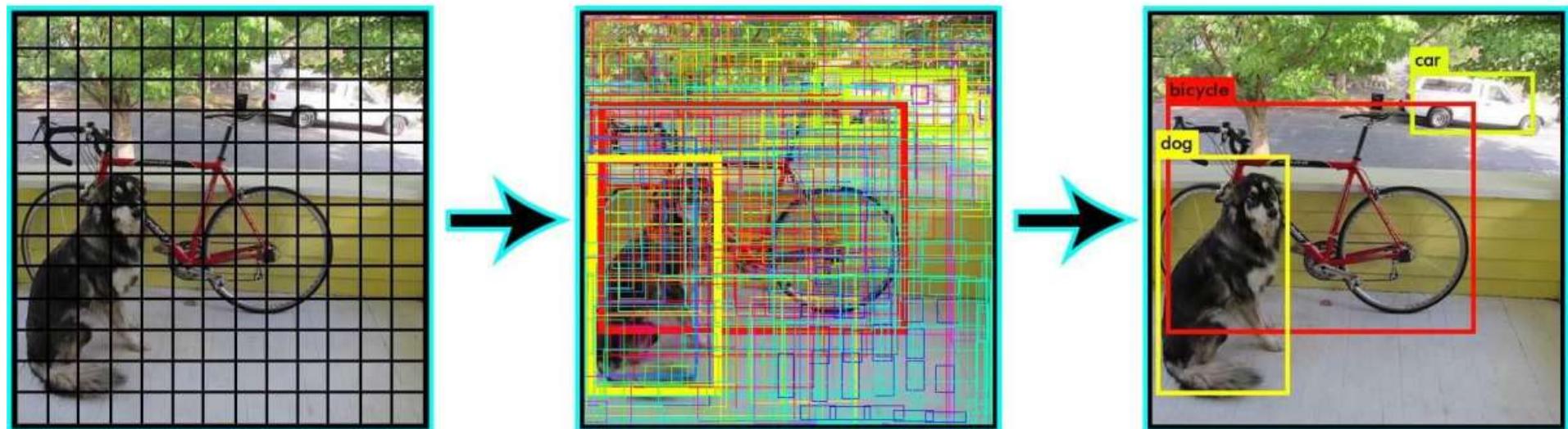
- $P(\text{object})$: probabilidad de que el box contenga un objeto
- B bounding boxes (x, y, d, w)
- $P(\text{class})$: probabilidad de pertenecer a una clase

YOLO (You Only Look Once)



Muchos bounding boxes con
Diferentes probabilidades

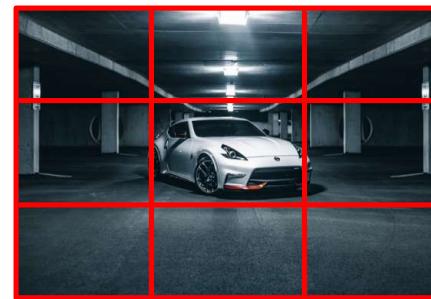
YOLO (You Only Look Once)



Yolo (You Only Look Once)

- Imagen se partitiona en bloques
- Cada bloque es responsable de sus detecciones

Un anchor es un marco rectangular
en donde puede ocurrir una detección

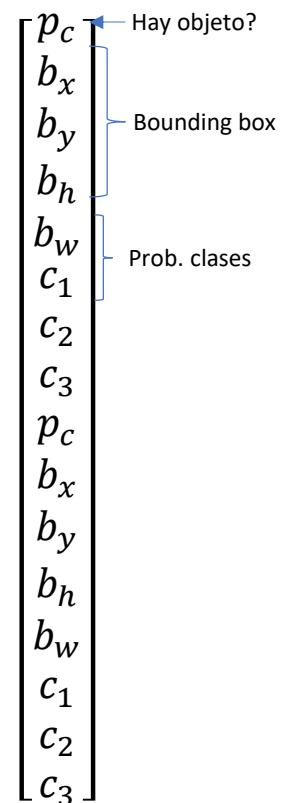


y is $3 \times 3 \times 2 \times 8$

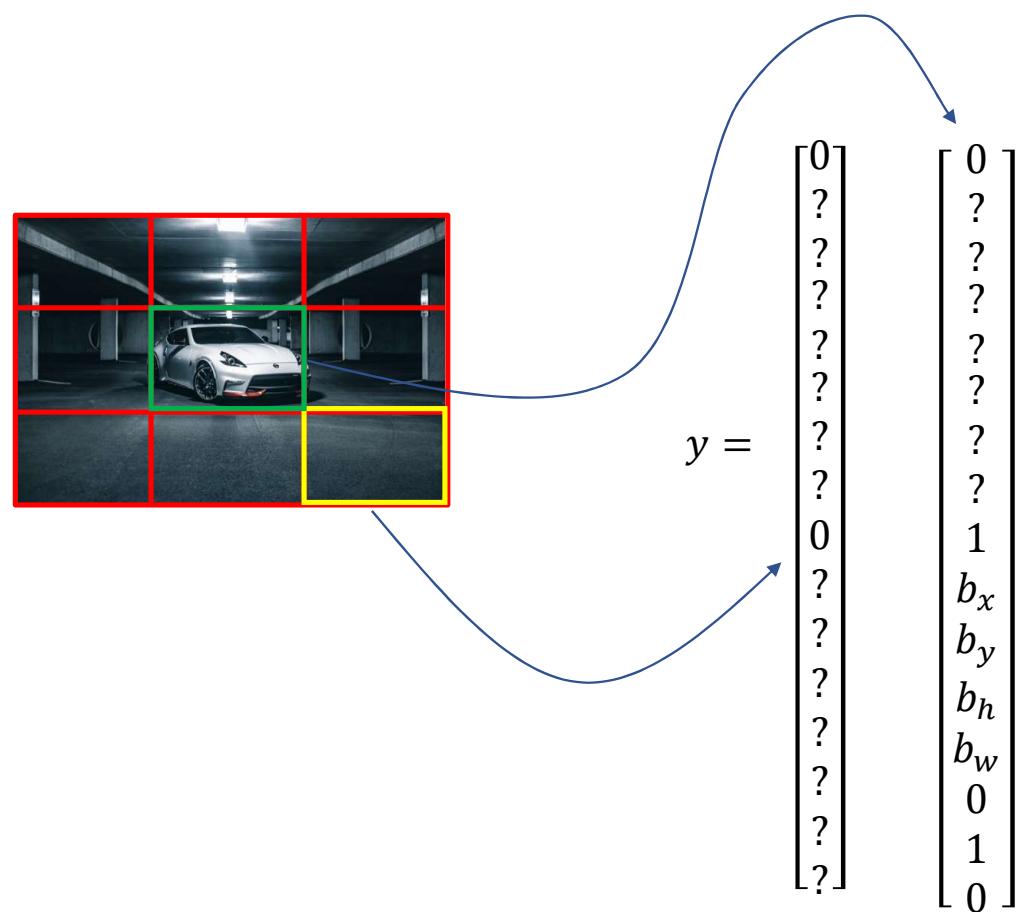
Cada celda tiene B anchors ->

- Clases
1. Person
 2. Car
 3. Motorbike

$y =$



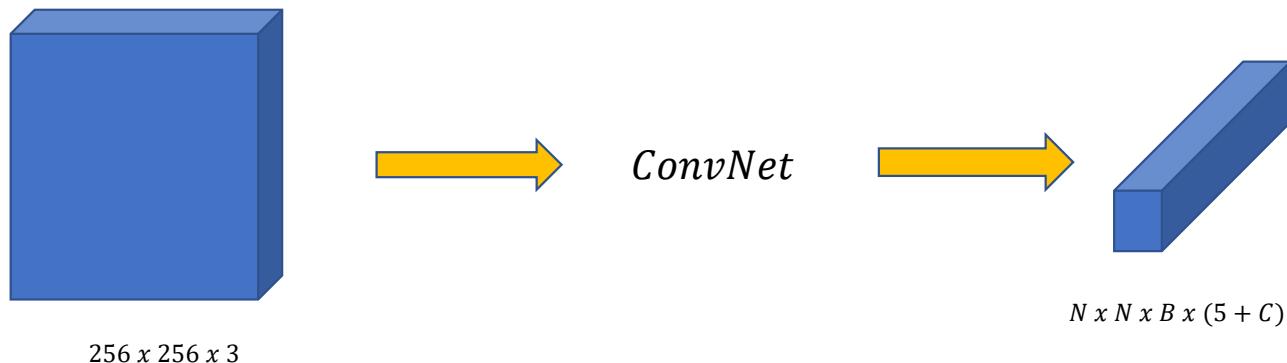
Yolo (You Only Look Once)



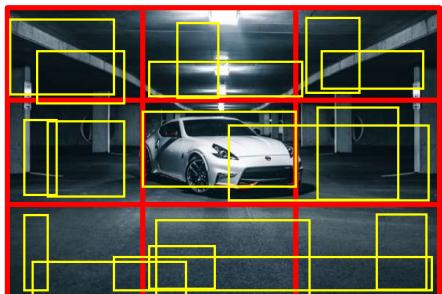
Clases

1. Person
2. Car
3. Motorbike

Yolo (You Only Look Once)

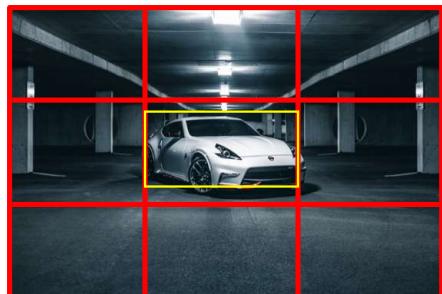


Yolo (You Only Look Once)



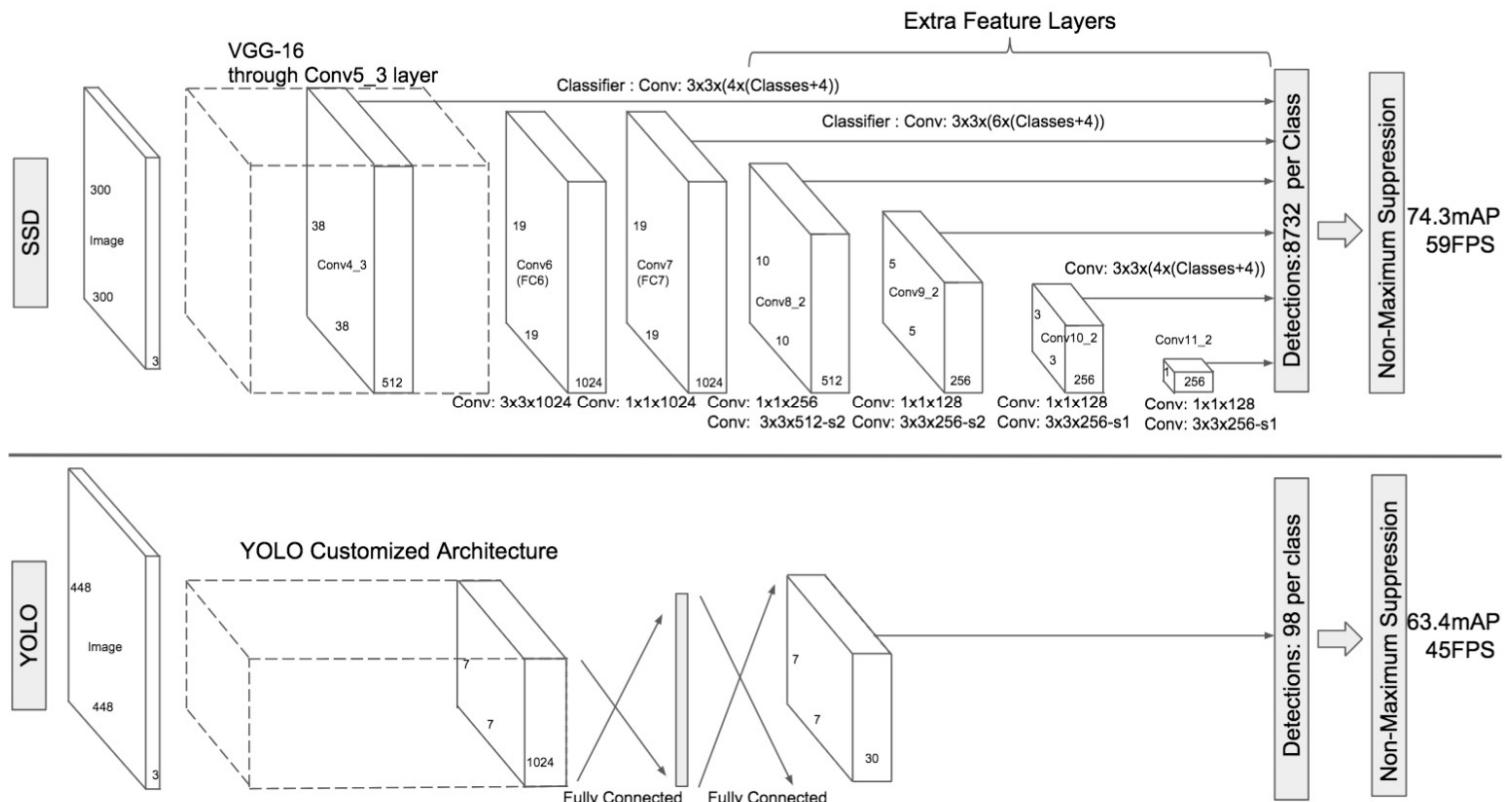
Para cada celda hay B anchor boxes

Se remueven las bajas confidencias



Para cada clase, usamos supresión de no máximos para generar predicciones finales

Single Shot Detection



Comparison of methods

Model	PASCAL VOC 2007	PASCAL VOC 2010	PASCAL VOC 2012	COCO 2015 (IoU=0.5)	COCO 2015 (IoU=0.75)	COCO 2015 (Official Metric)	COCO 2016 (IoU=0.5)	COCO 2016 (IoU=0.75)	COCO 2016 (Official Metric)	Real Time Speed
R-CNN	x	62.4%	x	x	x	x	x	x	x	No
Fast R-CNN	70.0%	68.8%	68.4%	x	x	x	x	x	x	No
Faster R-CNN	78.8%	x	75.9%	x	x	x	x	x	x	No
R-FCN	82.0%	x	x	53.2%	x	31.5%	x	x	x	No
YOLO	63.7%	x	57.9%	x	x	x	x	x	x	Yes
SSD	83.2%	x	82.2%	48.5%	30.3%	31.5%	x	x	x	No
YOLOv2	78.6%	x	x	44.0%	19.2%	21.6%	x	x	x	Yes

Evaluación

Datos

Datos para localizar objetos.

- Big datasets
- Tagged datasets



Pascal VOC

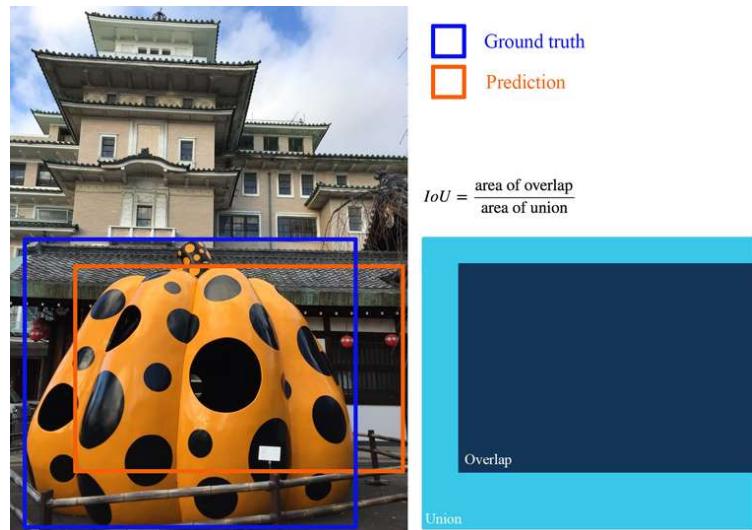


COCO dataset

Evaluación de detección

Dado un groundtruth y las detecciones de un algoritmo, cómo sabemos si el algoritmo lo hace bien?

- Una detección es exitosa si existe un alto overlap entre bounding boxes

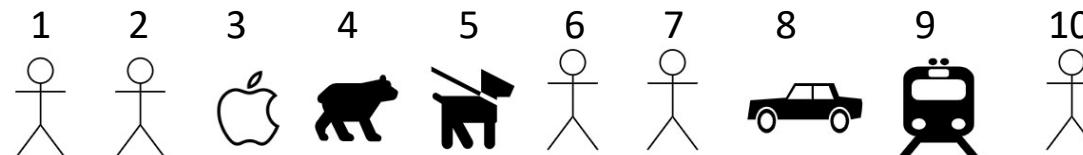


Evaluación de detección

La métrica más común es el Average Precision (AP). Esta métrica es computada para cada Clase en el dataset. Veamos un ejemplo:

- Tenemos un dataset con una clase *persona*

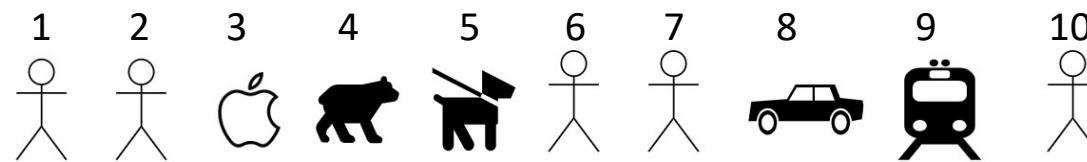
Colectamos todas las ocurrencias en donde el algoritmo predice una persona y Ordenamos de acuerdo a la confidencia



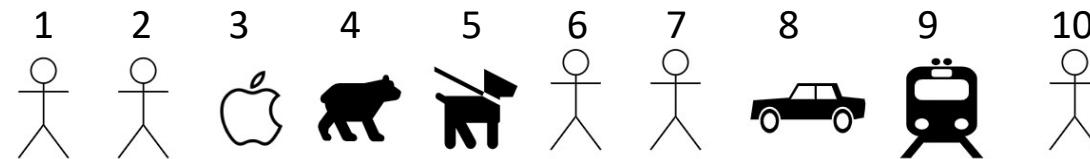
Evaluación de detección

Medimos la precisión y recall en cada posición de la lista

- **Precision:** número de personas/elementos en lista
- **Recall:** número de personas/número total de personas

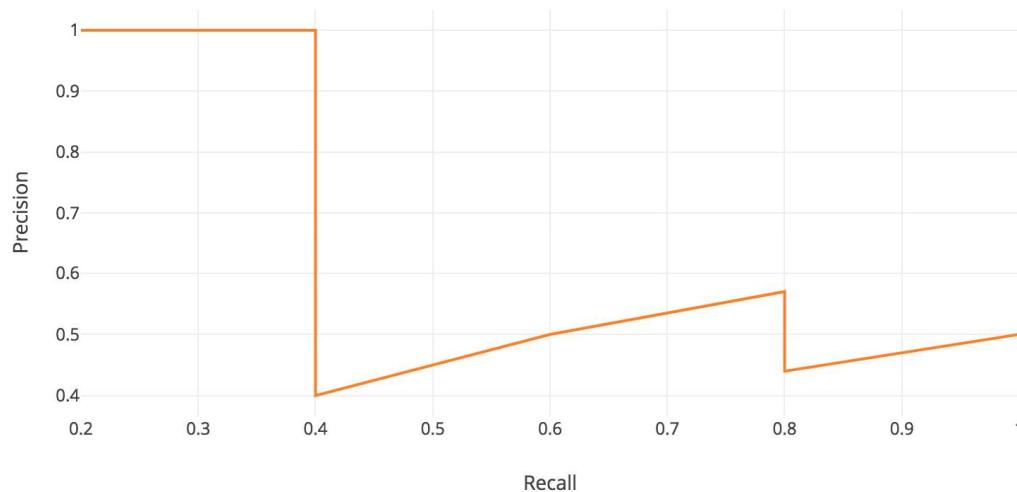


P.	1	1	0.67	0.5	0.4	0.5	0.57	0.5	0.44	0.5
R.	0.2	0.4	0.4	0.4	0.4	0.6	0.8	0.8	0.8	1.0



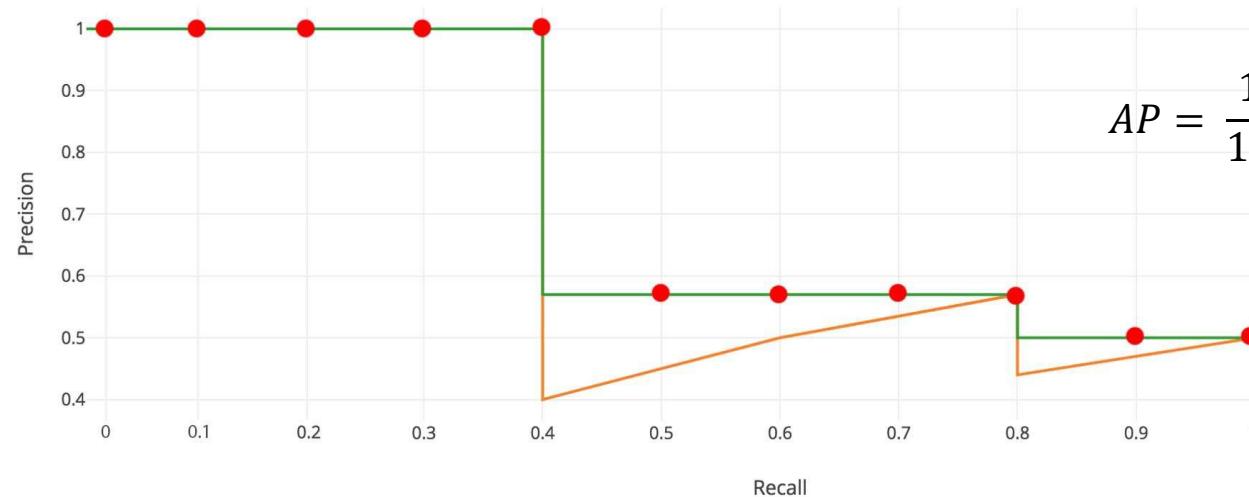
P.	1	1	0.67	0.5	0.4	0.5	0.57	0.5	0.44	0.5
R.	0.2	0.4	0.4	0.4	0.4	0.6	0.8	0.8	0.8	1.0

Ploteamos la curva recall vs. precision: precisión promedio es el área bajo la curva



1	2	3	4	5	6	7	8	9	10	
P.	1	1	0.67	0.5	0.4	0.5	0.57	0.5	0.44	0.5
R.	0.2	0.4	0.4	0.4	0.4	0.6	0.8	0.8	0.8	1.0

En la práctica, se toman valores interpolados de recall y valores monotónicamente decreciente de precisión



$$AP = \frac{1}{11}(5 * 1 + 4 * 0.57 + 2 * 0.5) = 0.75$$

Para clase *persona*

Mean Average Precision (mAP) es el Promedio de AP para todas las clases