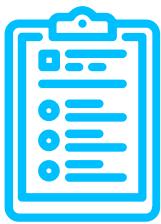


Sesión 2.0

Segmentación

U-net, R-CNN, DeepLab

1.



Segmentación *y detección*



Segmentación de imágenes



Entrada

Segmented

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

Etiquetas semánticas

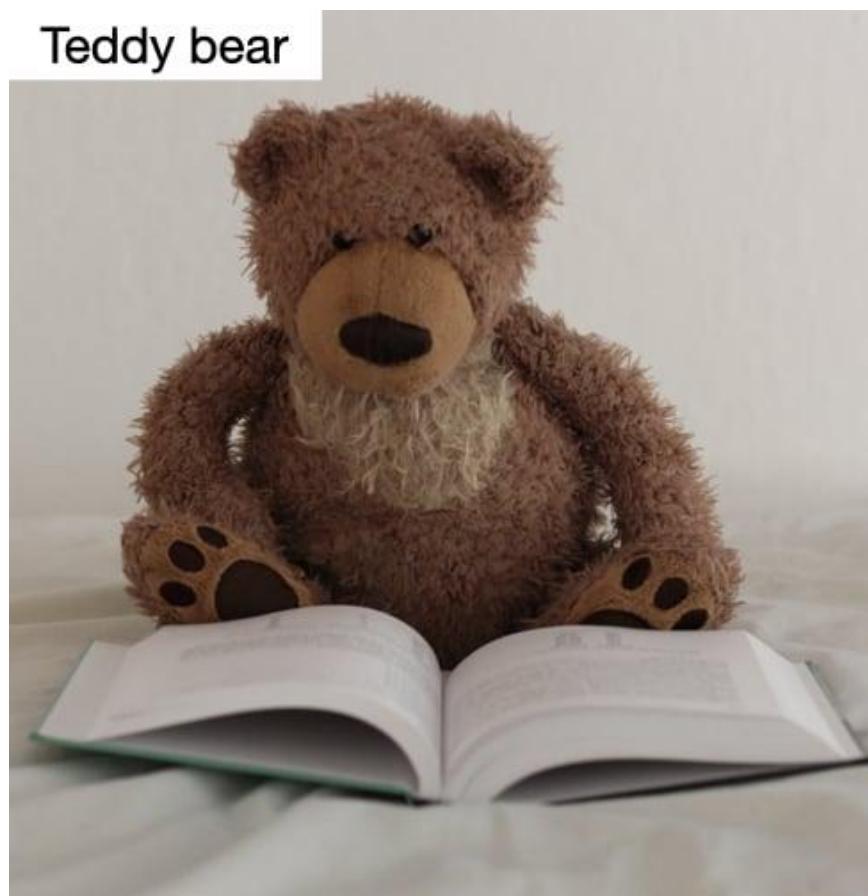
One-hot



TRANSFORMATEC

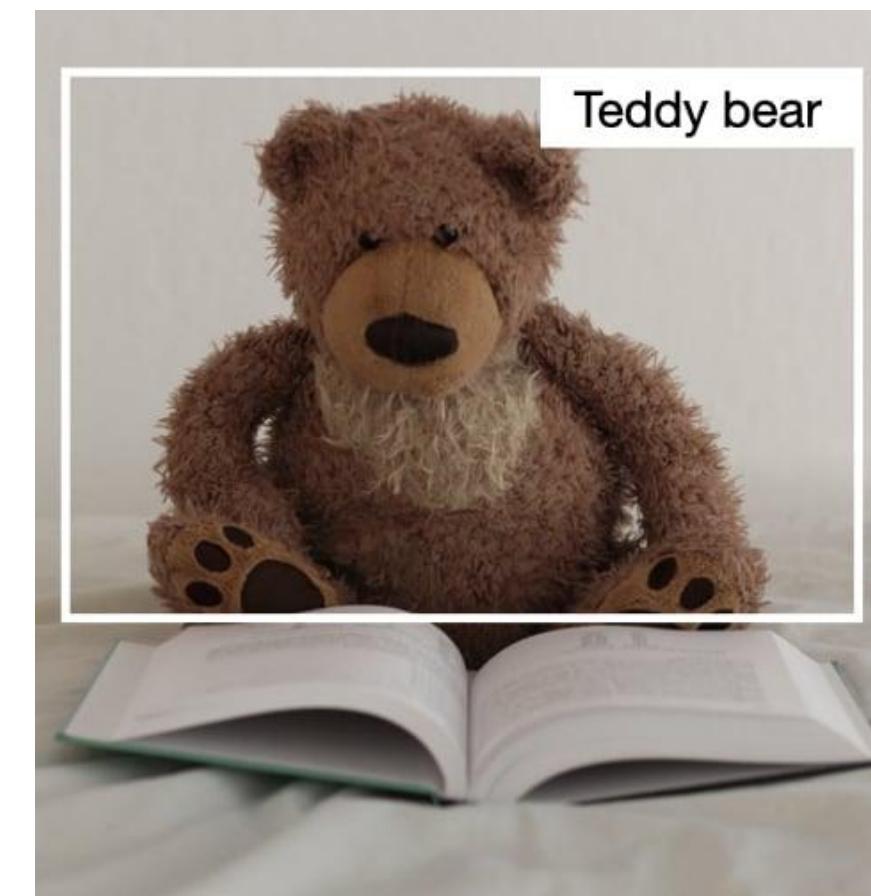
Segmentación y detección

Clasificación de imágenes



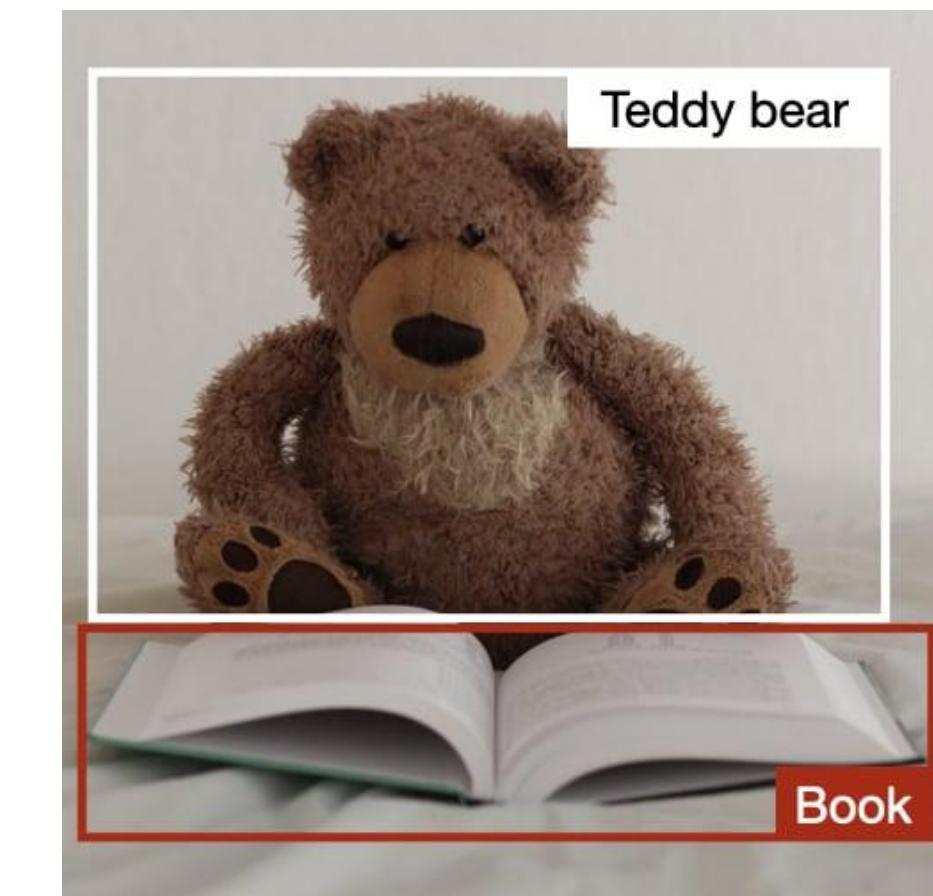
- Clasifica una imagen.
- Predice la probabilidad de que un objeto esté en una imagen.

Clasificación con localización



- Detecta un objeto en una imagen.
- Predice la probabilidad de que un objeto esté en una imagen, además de su ubicación.

Detección

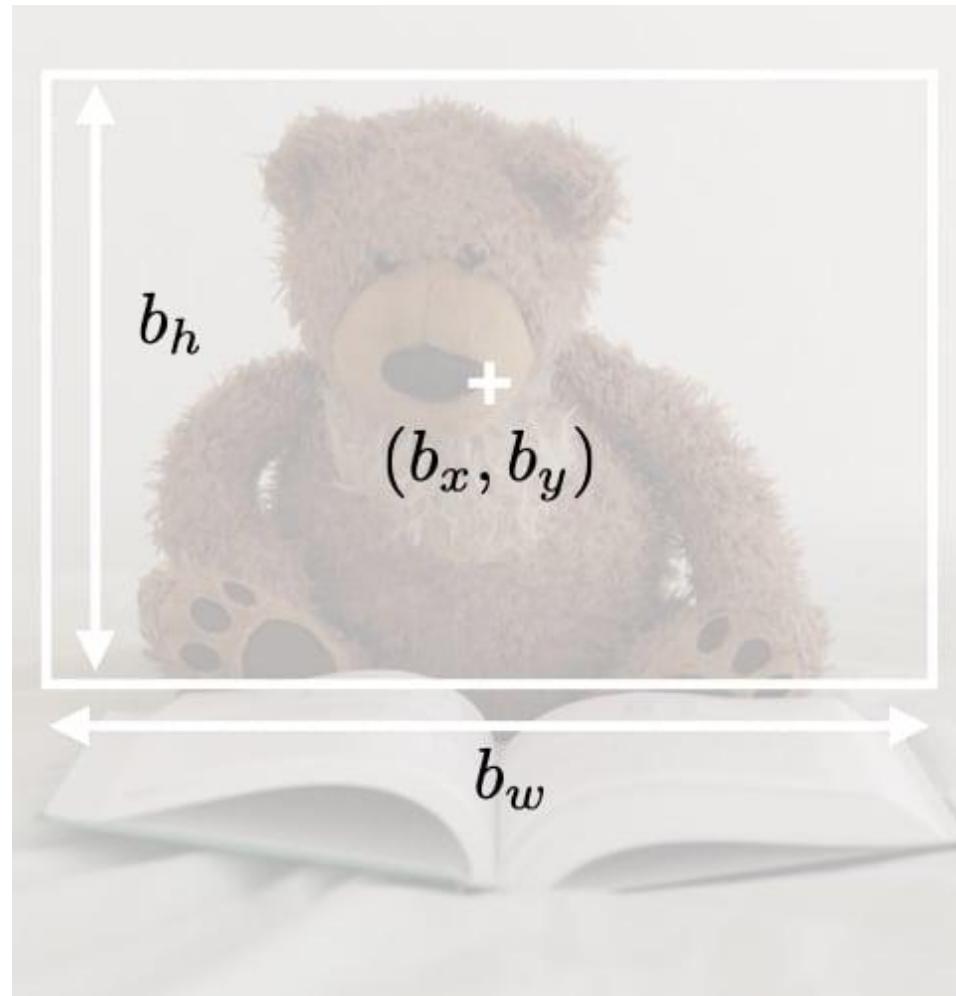


- Detecta varios objetos en una imagen.
- Predice la probabilidad de que algunos objetos estén en una imagen, además estima sus ubicaciones.



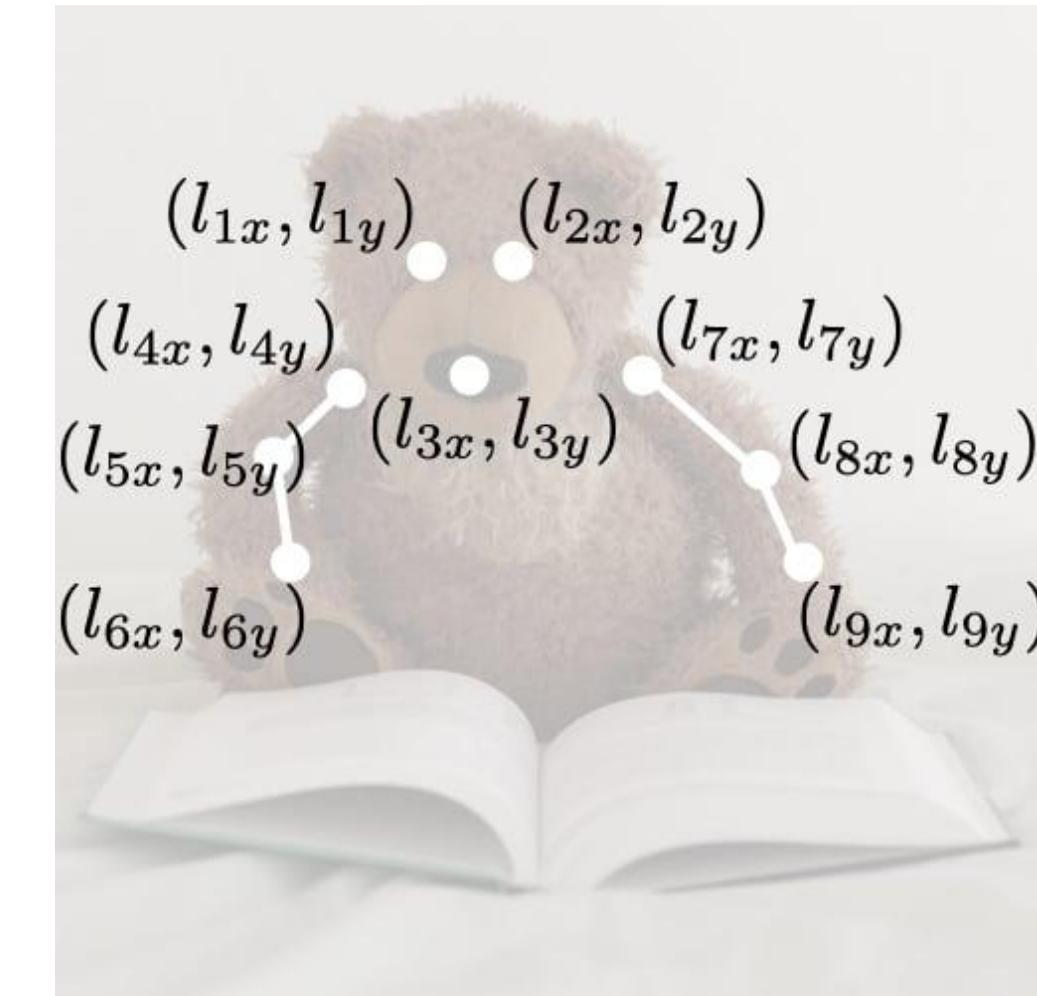
Segmentación y detección

Bounding box



- Detecta la parte de la imagen donde se encuentra el objeto.

Landmark



- Detecta una forma o características de un objeto (por ejemplo, ojos).
- Más granular.



Segmentación



Imagen de entrada



Semantic segmentation



Instance segmentation

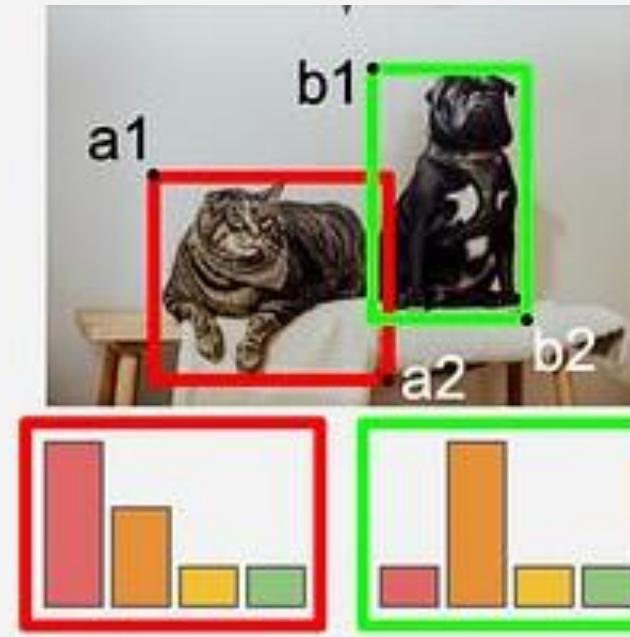


Panoptic segmentation

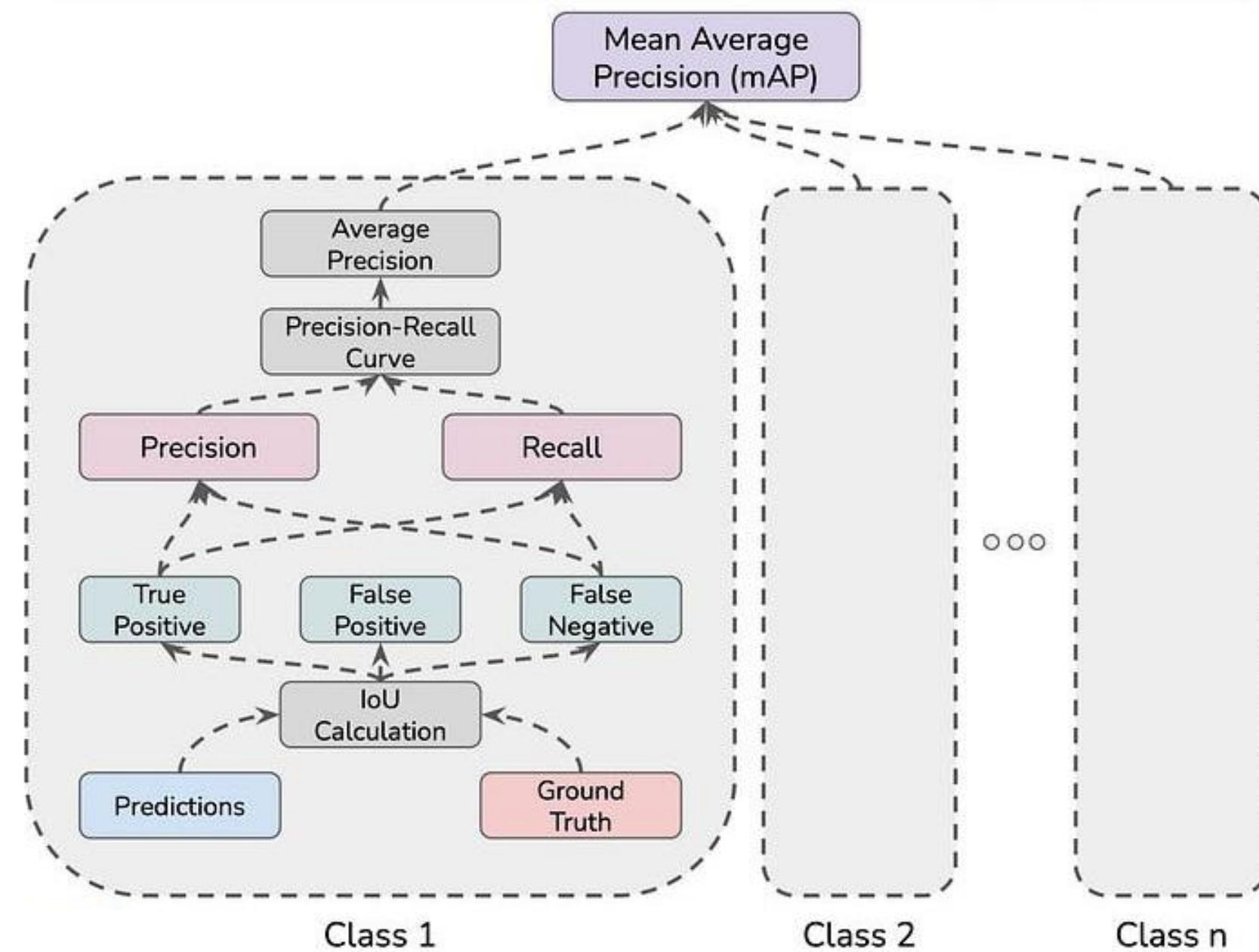


Mean average precision (mAP)

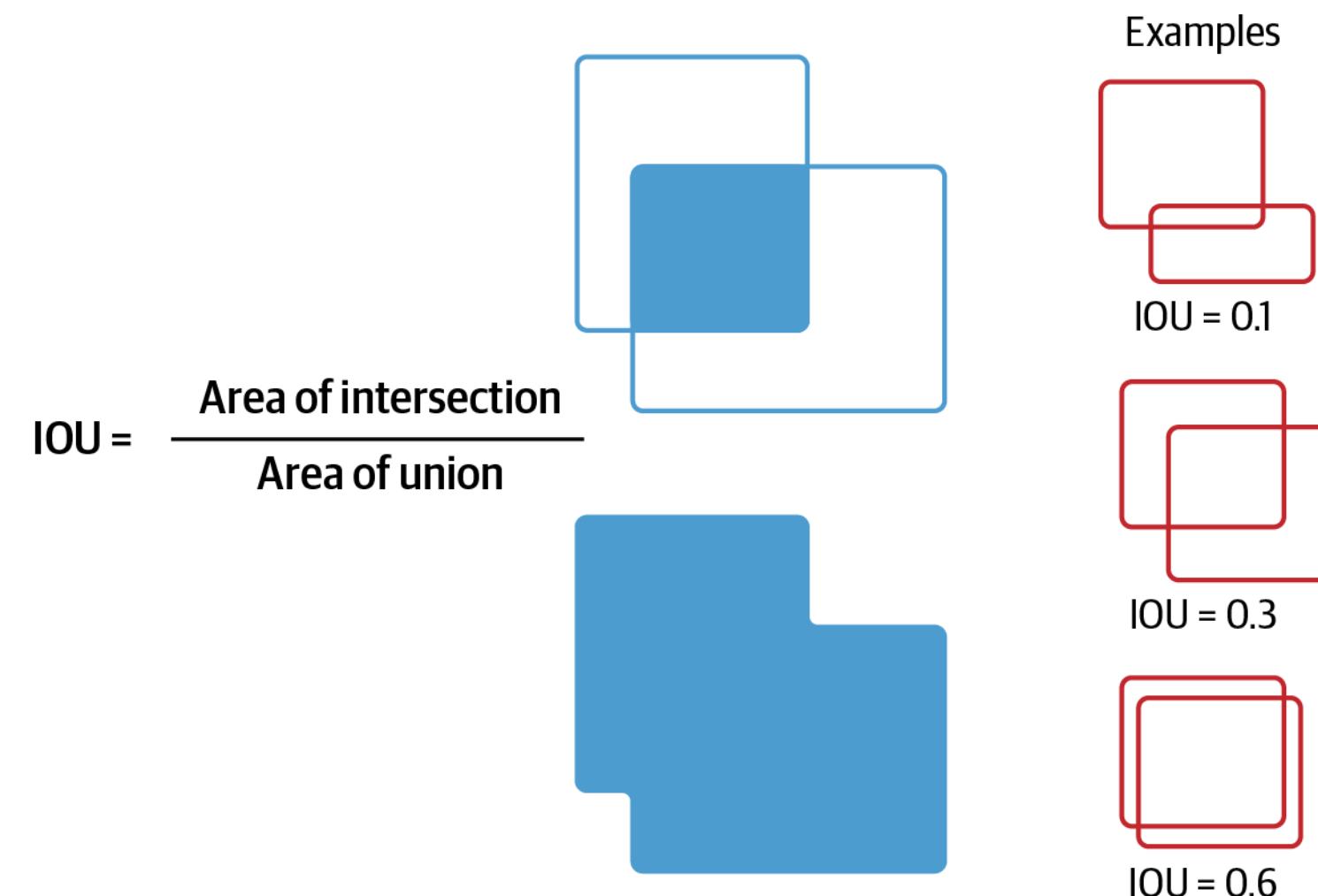
$$mAP = \frac{1}{n} \sum_{k=1}^n AP_k$$



Mean average precision (mAP)

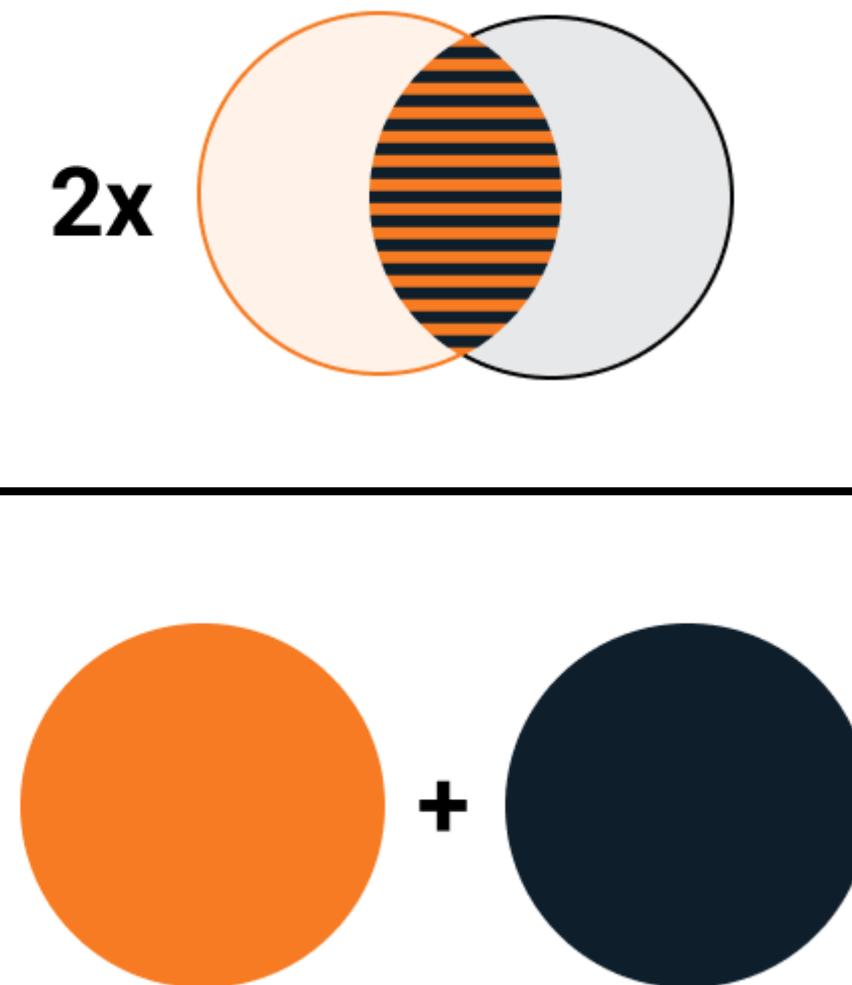


Intersection over Union (IoU)

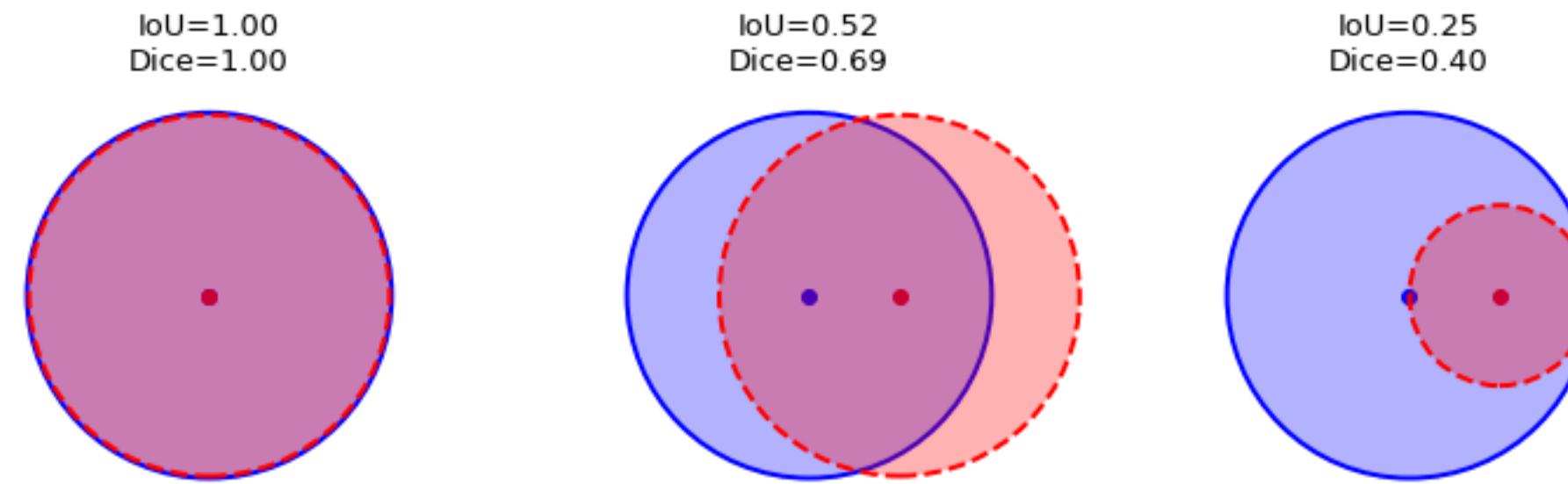


Métricas

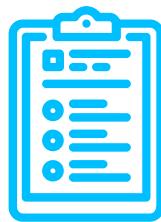
$$\text{DICE} = \frac{2|A \cap B|}{|A| + |B|} = \frac{2x}{x + y}$$



IoU vs DICE



2.



Métodos *Tradicionales*



Threshold Segmentation

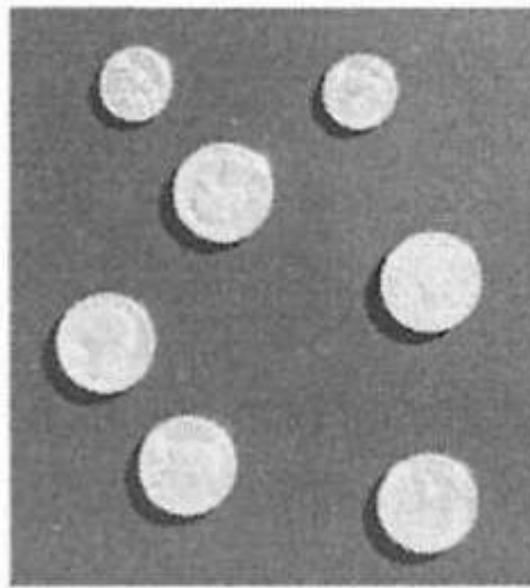
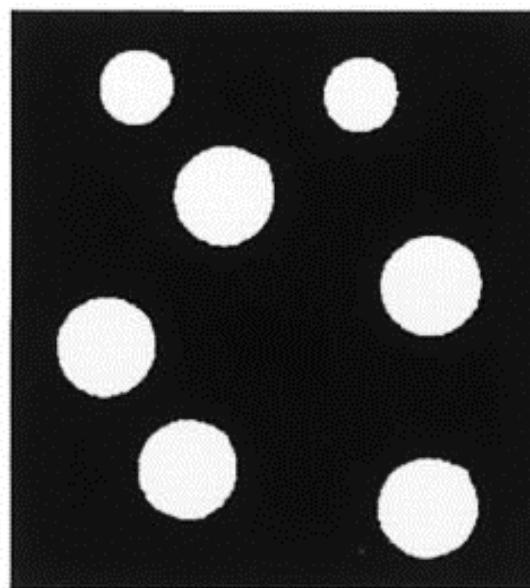
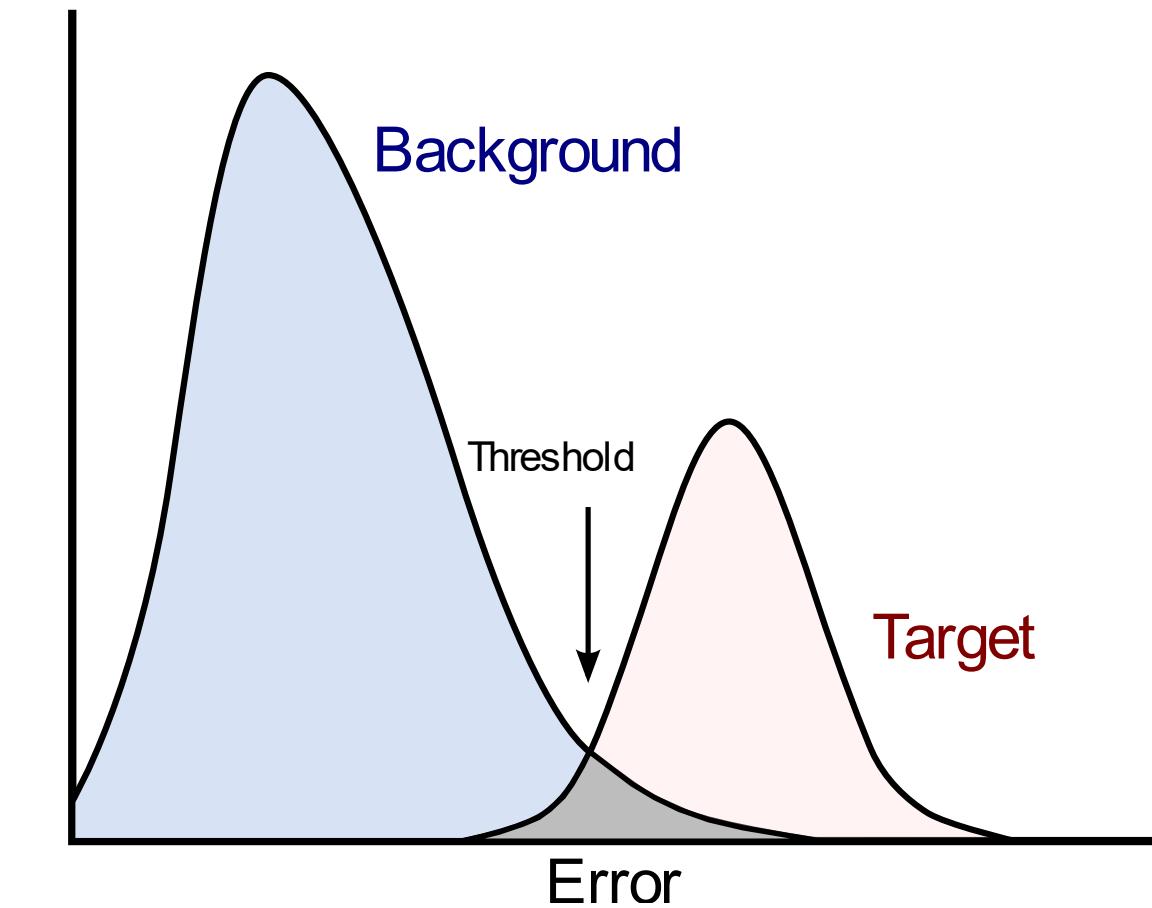


Imagen de entrada



Resultado

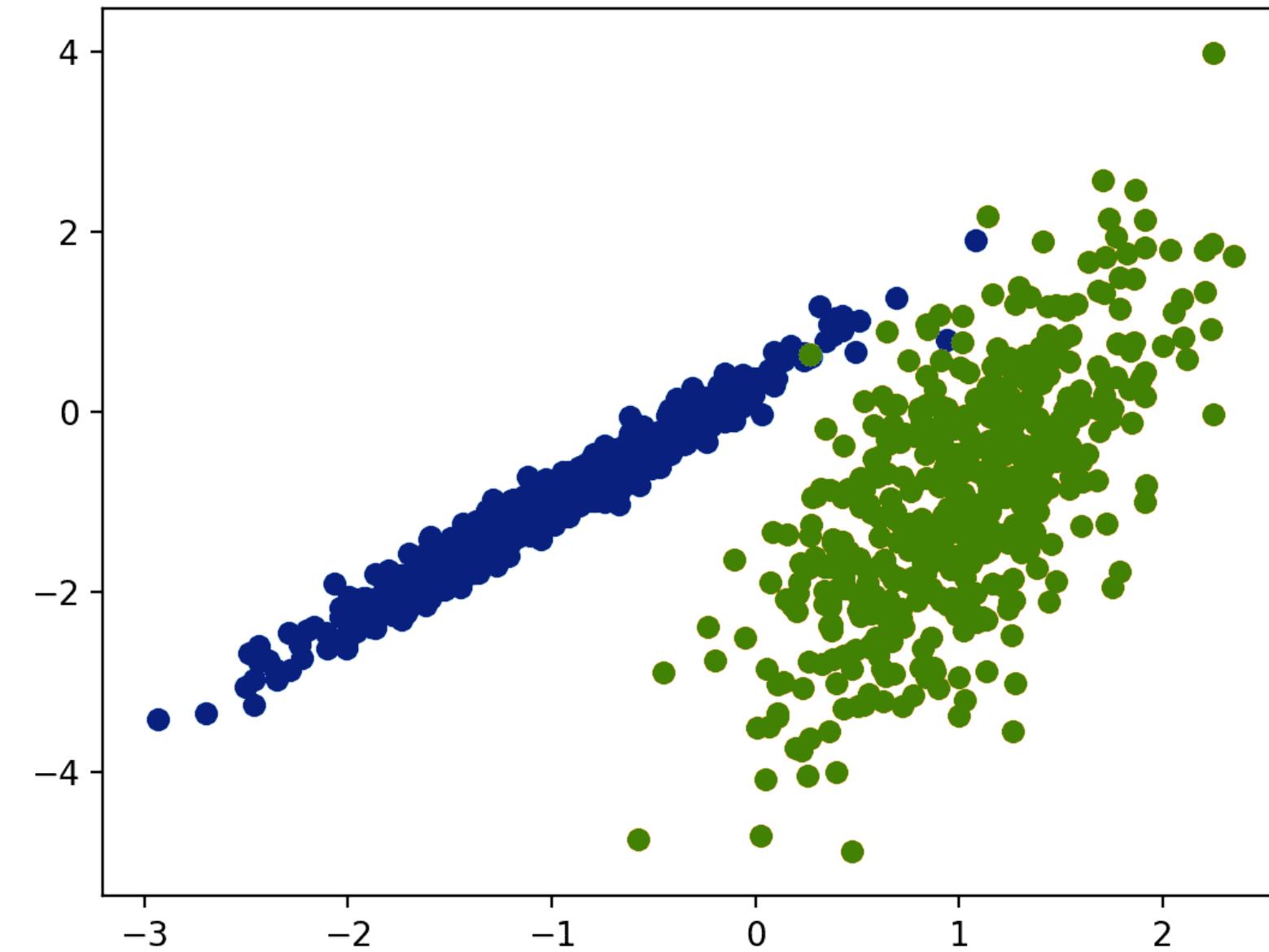
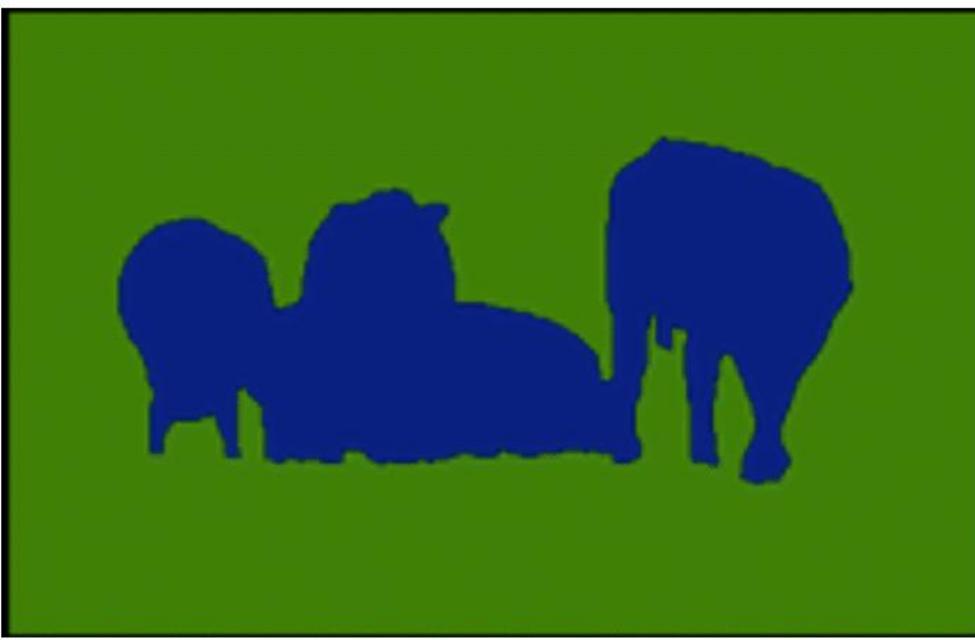


Edge-Based Segmentation

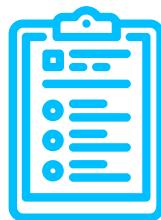


TRANSFORMATEC

Cluster-Based *Segmentation*



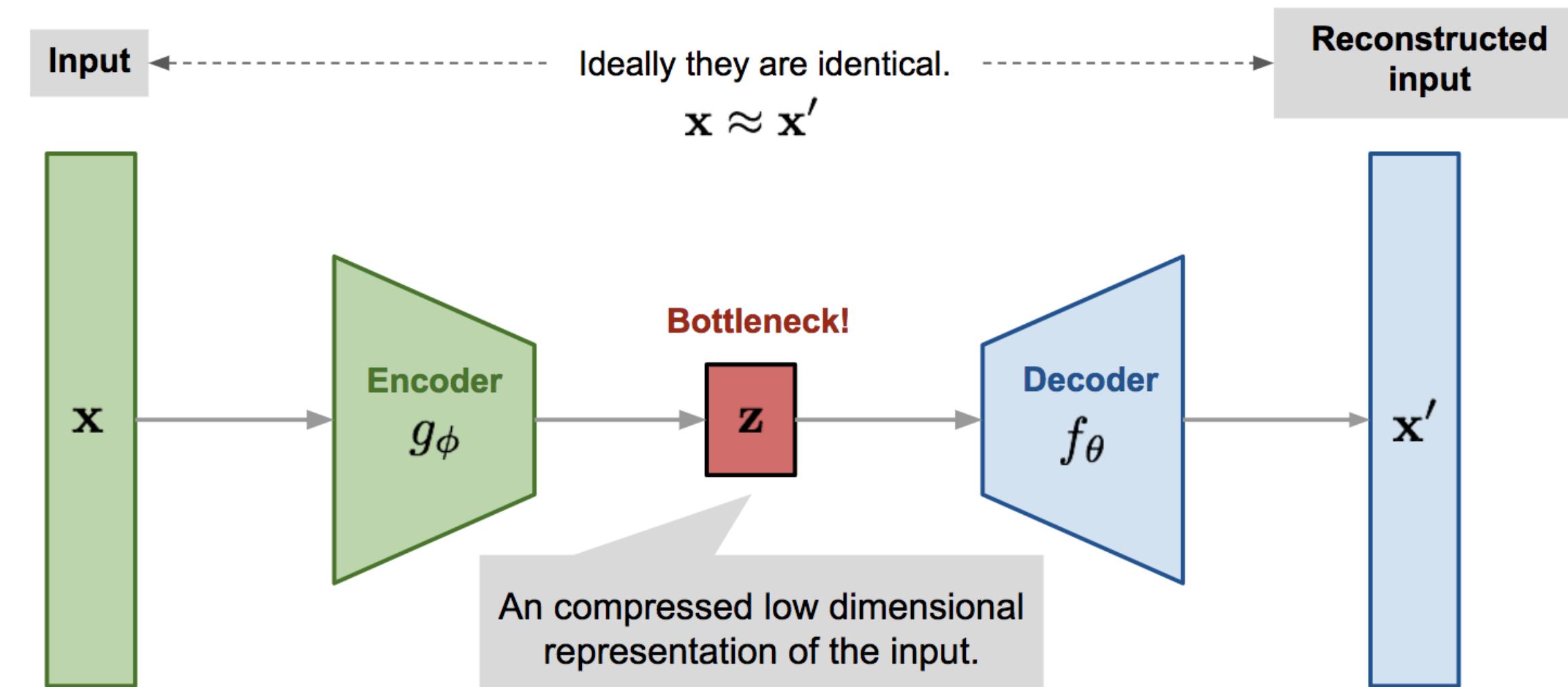
3.



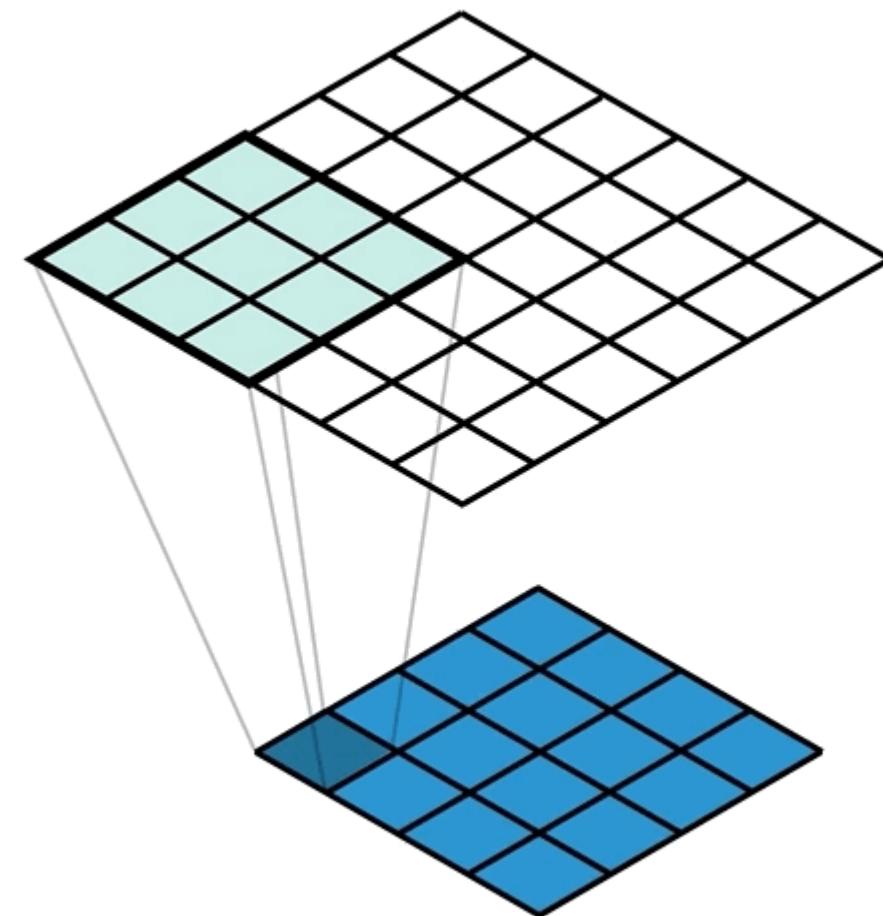
Encoder-*Decoder*



Autoencoder



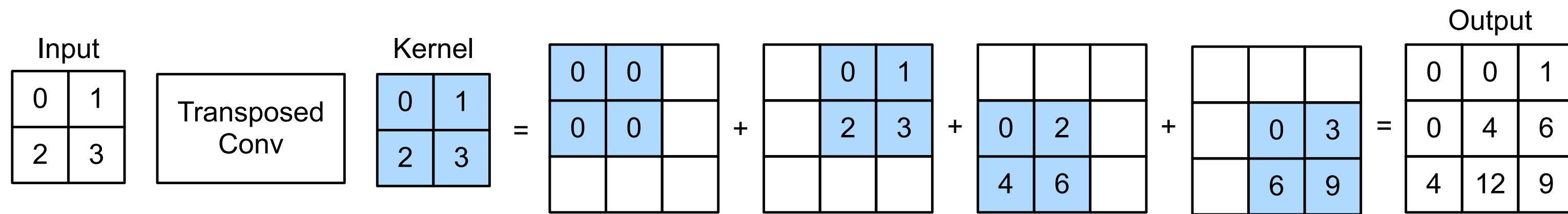
Transposed Convolution



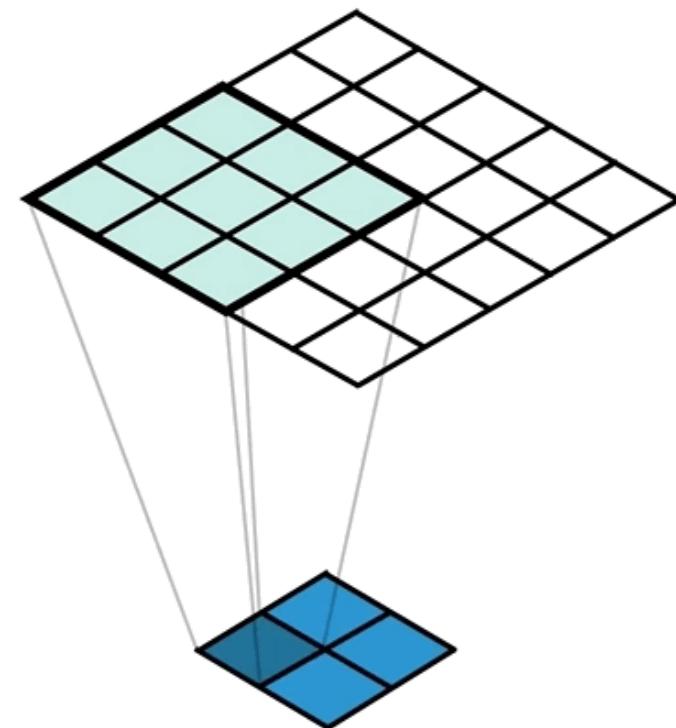
$$y(m, n) = \sum_{i,j} x(i, j) \cdot w(m - i, n - j)$$



Transposed Convolution



Transposed Convolution

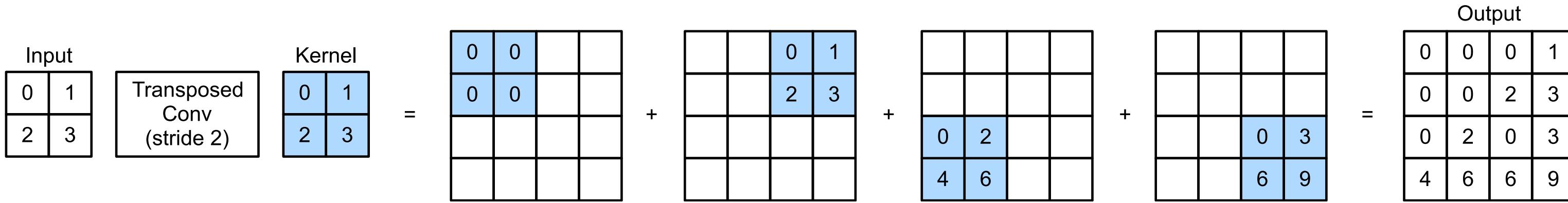
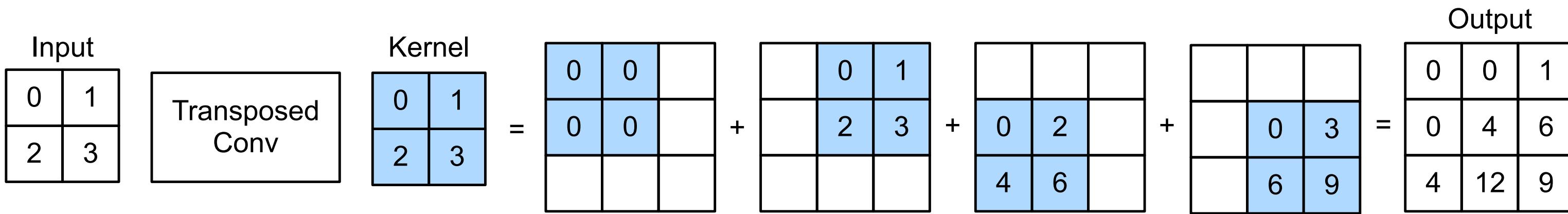


$$y(m, n) = \sum_{i,j} x(i, j) \cdot w(m - i \cdot s + p, n - j \cdot s + p)$$

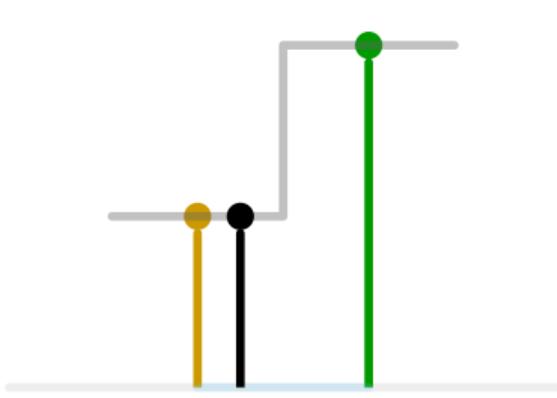
- s : stride
- p : padding



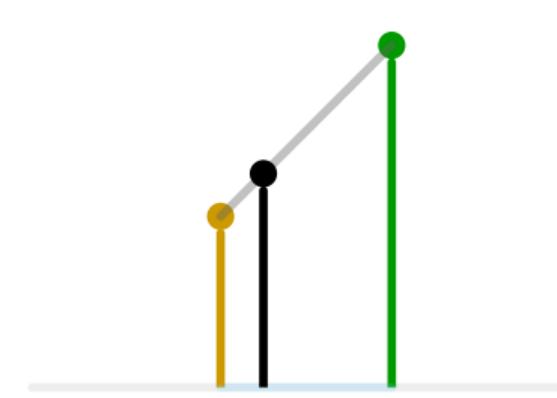
Transposed Convolution



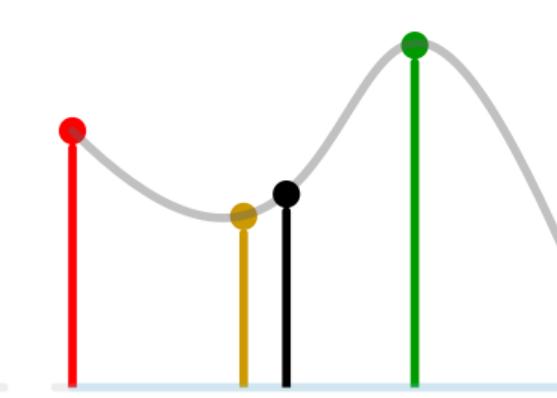
Interpolation



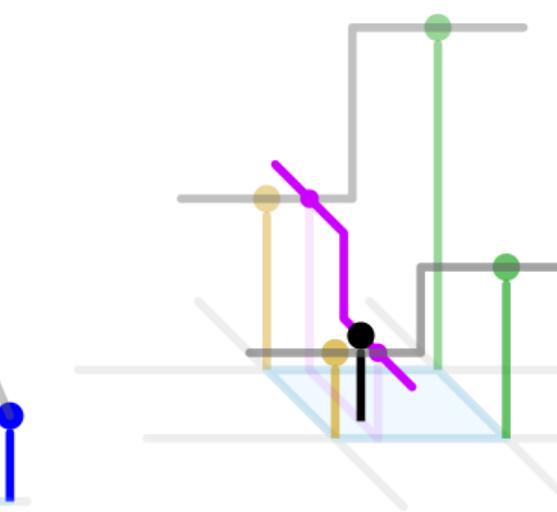
1D nearest-
neighbour



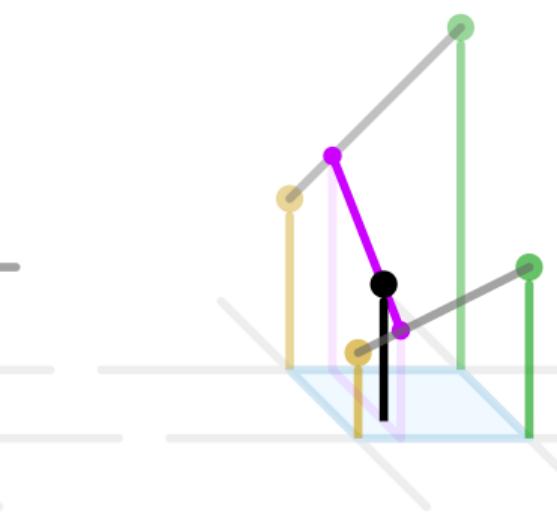
Linear



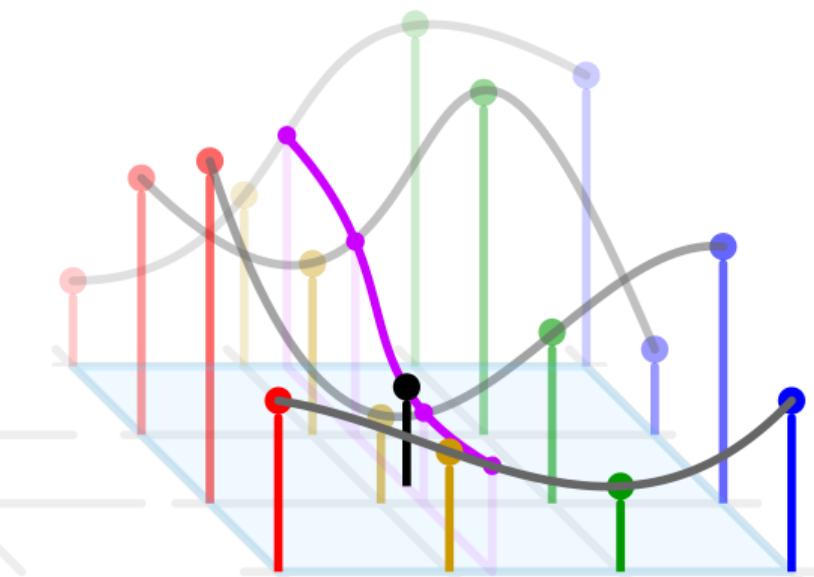
Cubic



2D nearest-
neighbour



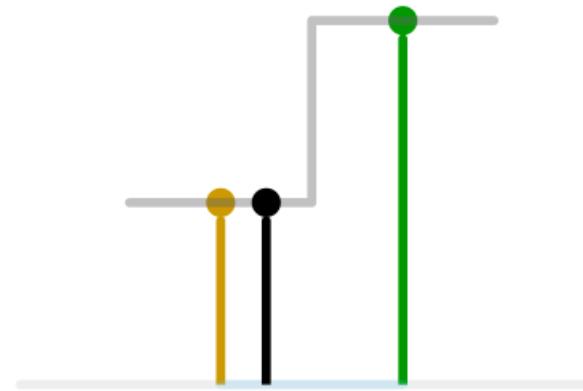
Bilinear



Bicubic

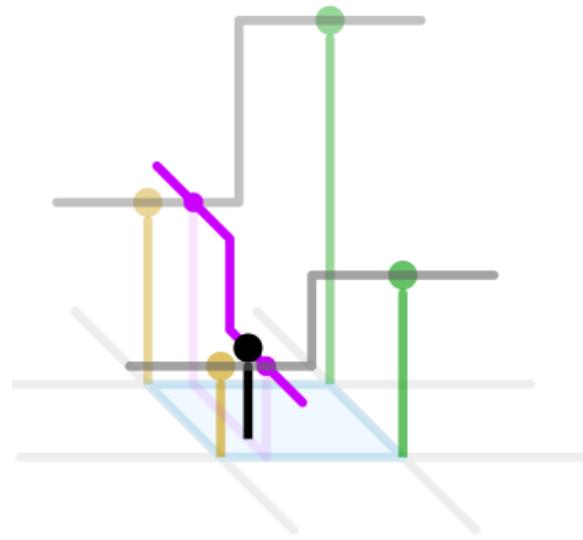


Interpolation



1D nearest-
neighbour

$$y_{\text{out}}(i) = y_{\text{in}} \left(\text{round} \left[\frac{i}{s} \right] \right)$$

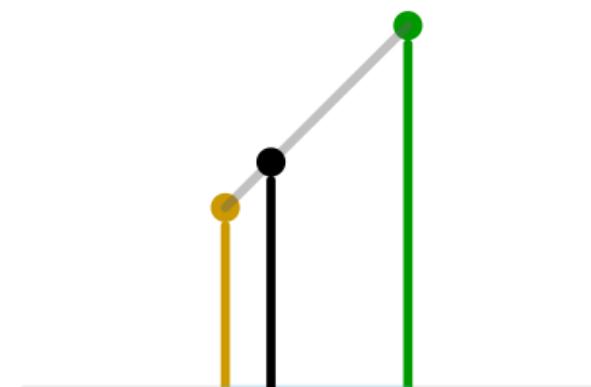


2D nearest-
neighbour

$$y_{\text{out}}(i, j) = y_{\text{in}} \left(\text{round} \left[\frac{i}{s} \right], \text{round} \left[\frac{j}{s} \right] \right)$$



Interpolation

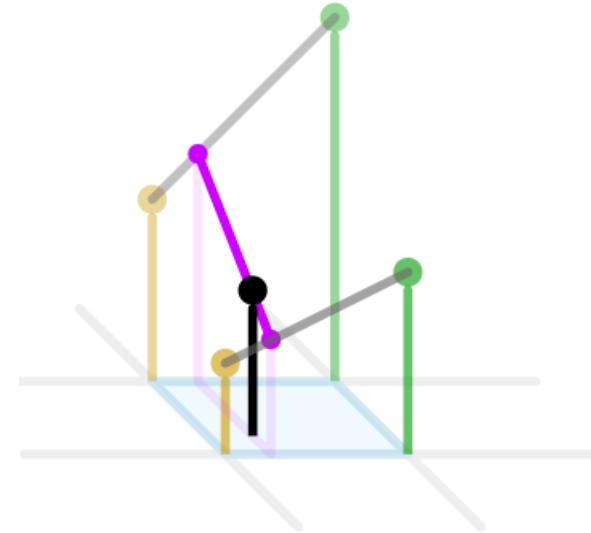


Linear

$$y_{\text{out}}(i) = (1 - a) \cdot y_{\text{in}}(x_1) + a \cdot y_{\text{in}}(x_2)$$

Donde:

$$\begin{aligned} x_1 &= \lfloor i \cdot s \rfloor & x_2 &= x_1 + 1 \\ a &= i \cdot s - x_1 \end{aligned}$$

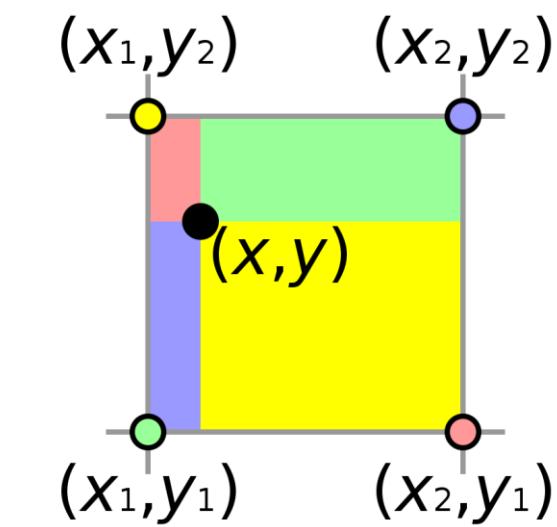


Bilinear

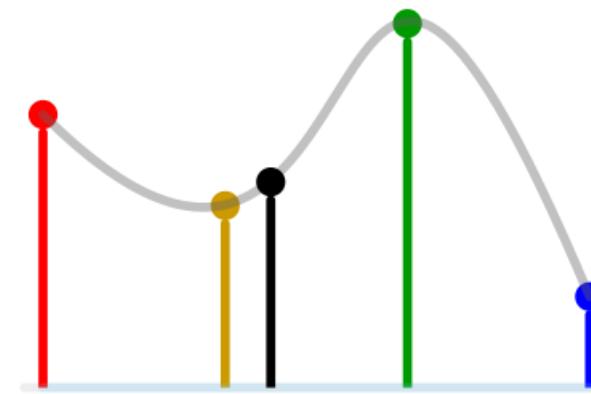
$$\begin{aligned} y_{\text{out}}(i, j) &= (1 - a)(1 - b) \cdot y_{\text{in}}(x_1, y_1) \\ &\quad + a(1 - b) \cdot y_{\text{in}}(x_2, y_1) \\ &\quad + (1 - a)b \cdot y_{\text{in}}(x_1, y_2) \\ &\quad + ab \cdot y_{\text{in}}(x_2, y_2) \end{aligned}$$

Donde:

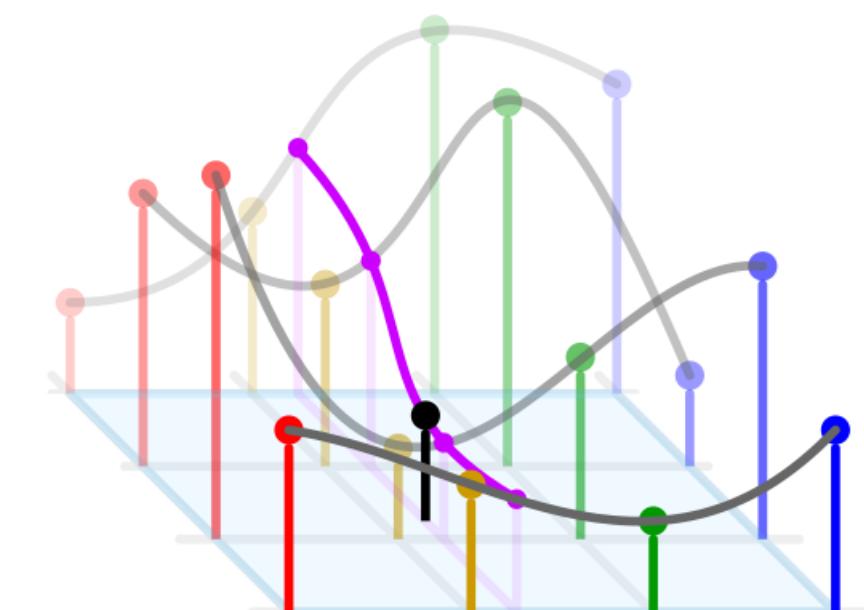
$$\begin{array}{ll} x_1 = \lfloor i \cdot s_h \rfloor & x_2 = x_1 + 1 \\ y_1 = \lfloor f \cdot s_w \rfloor & y_2 = y_1 + 1 \\ a = i \cdot s_h - x_1 & b = j \cdot s_w - y_1 \end{array}$$



Interpolation



Cubic



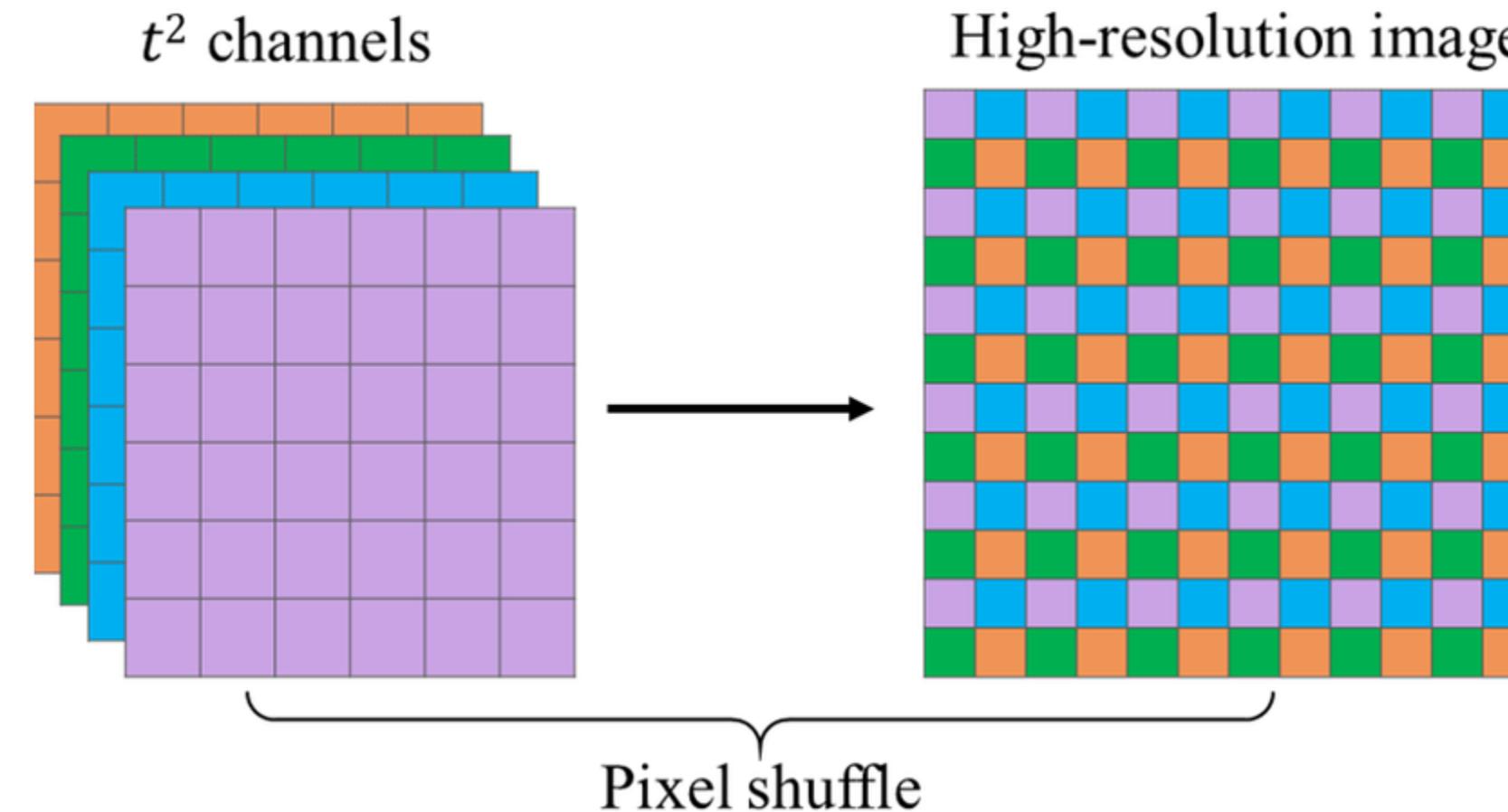
Bicubic

$$y_{\text{out}}(i) = \sum_{k=-1}^2 c_k \cdot y_{\text{in}}(x_k)$$

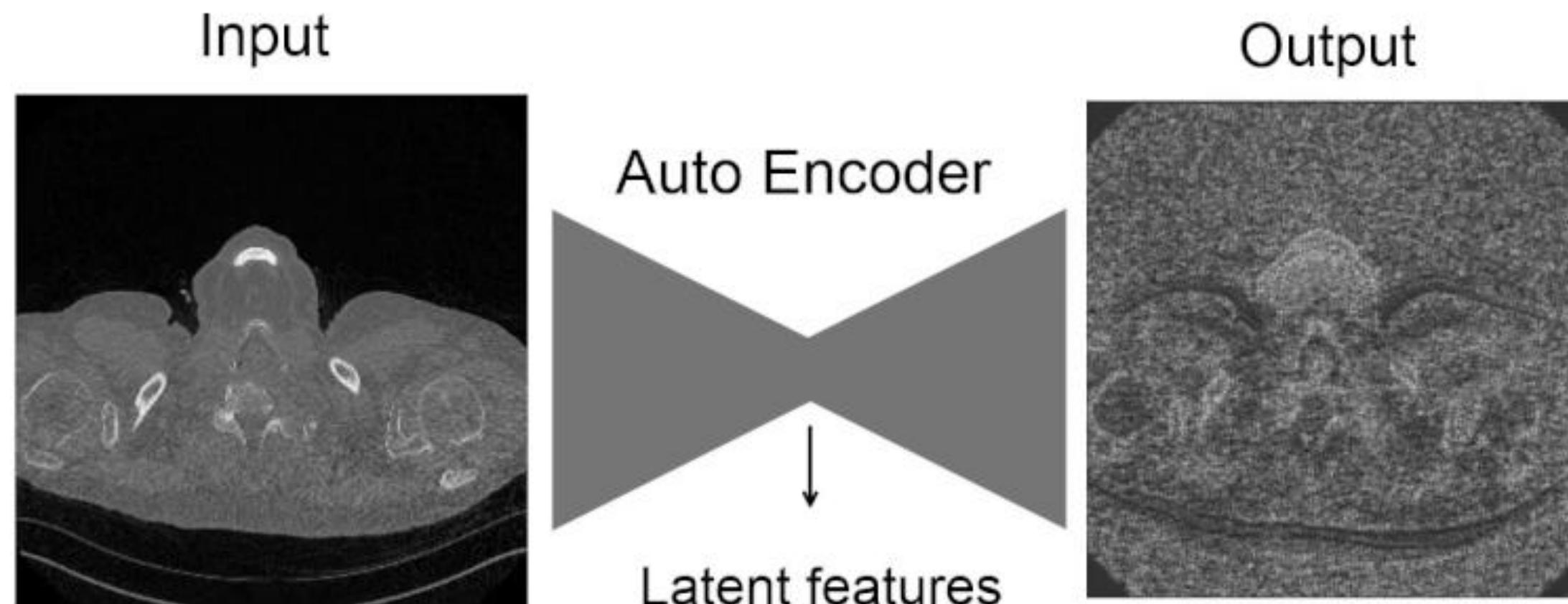
$$y_{\text{out}}(i, j) = \sum_{m=-1}^2 \sum_{n=-1}^2 c_{mn} \cdot y_{\text{in}}(x_m, y_n)$$



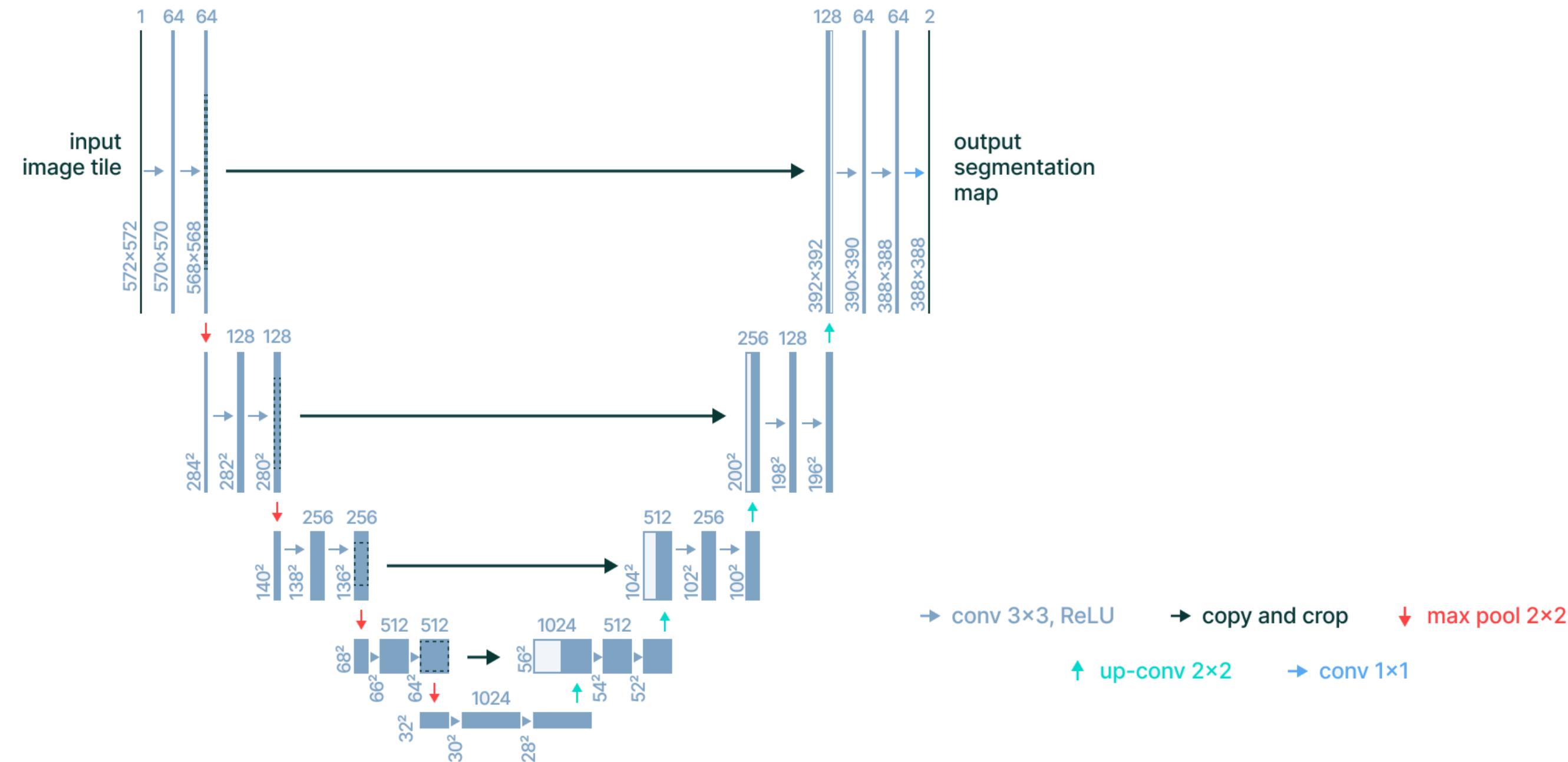
Pixel *Shuffle*



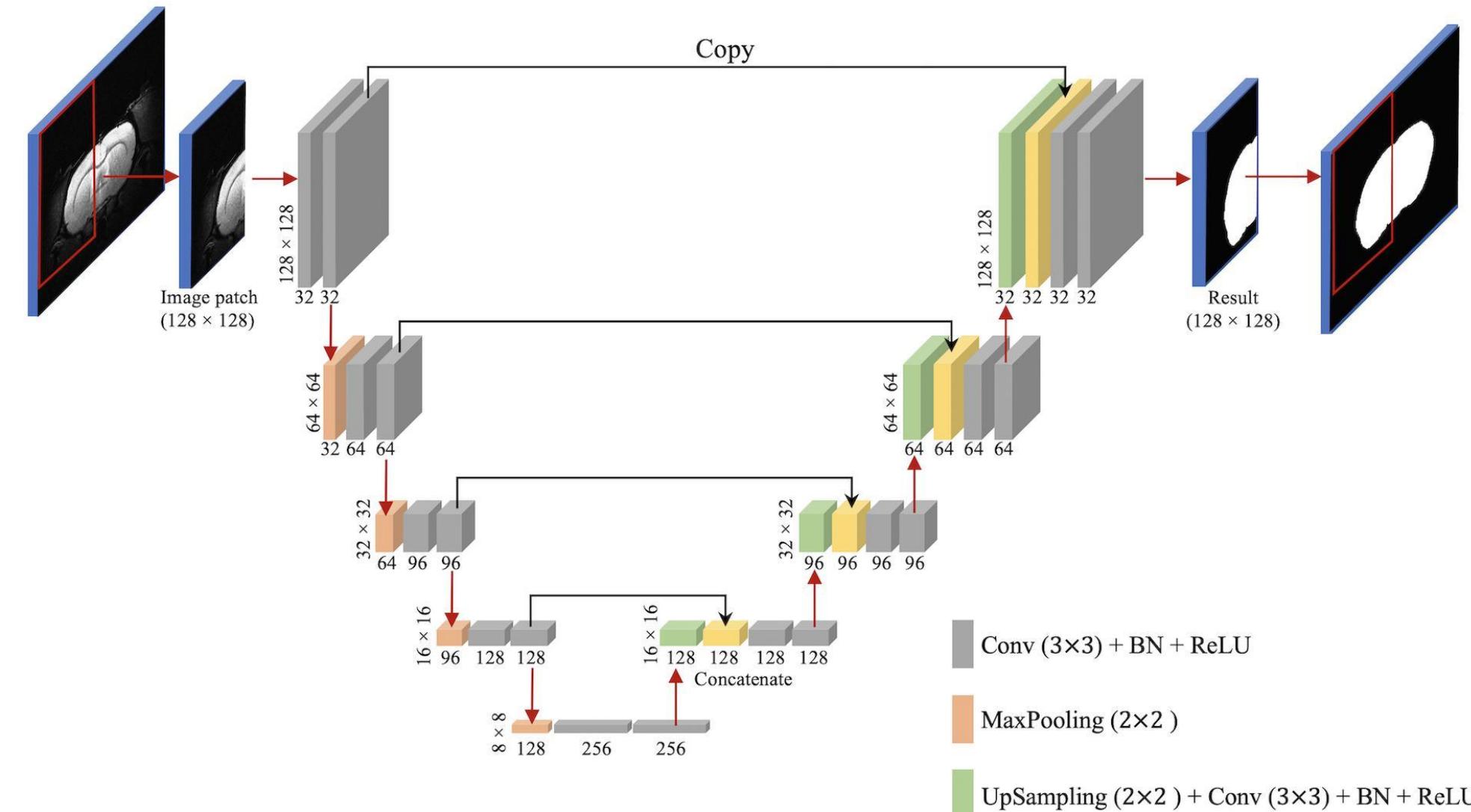
Autoencoder



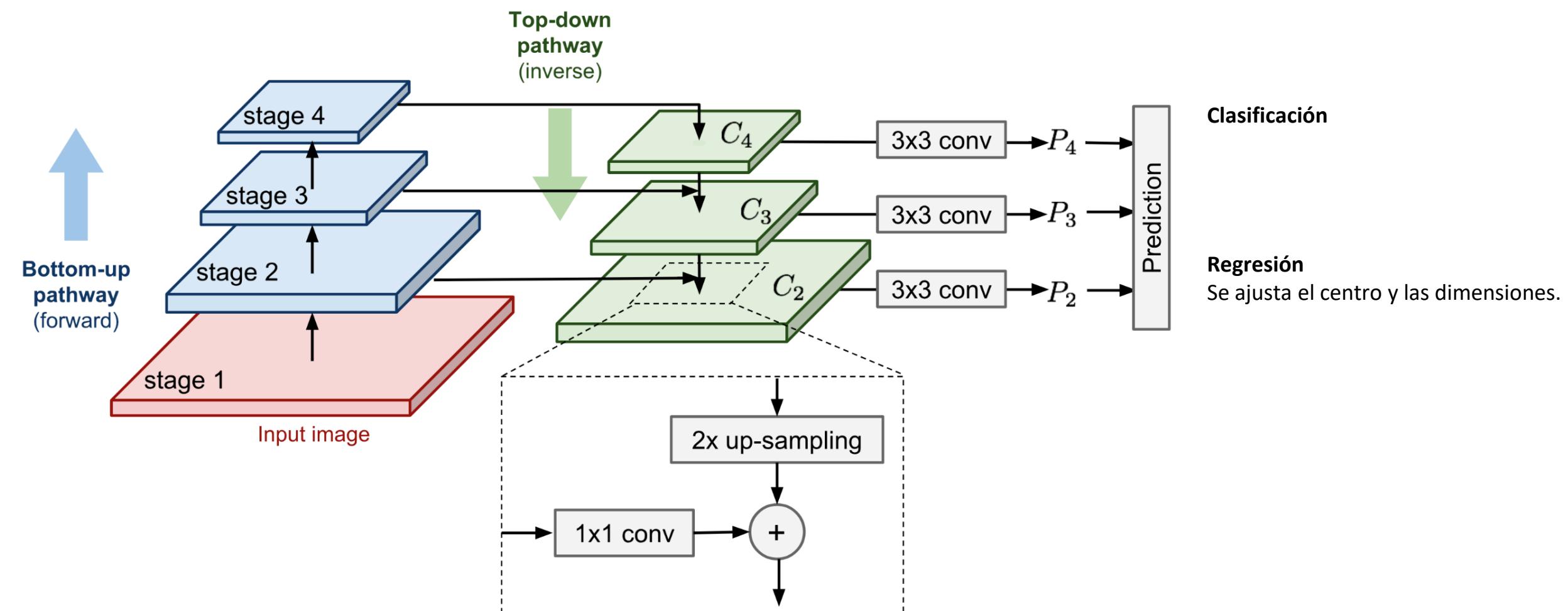
U-Net



U-Net



Feature Pyramid Network (FPN)



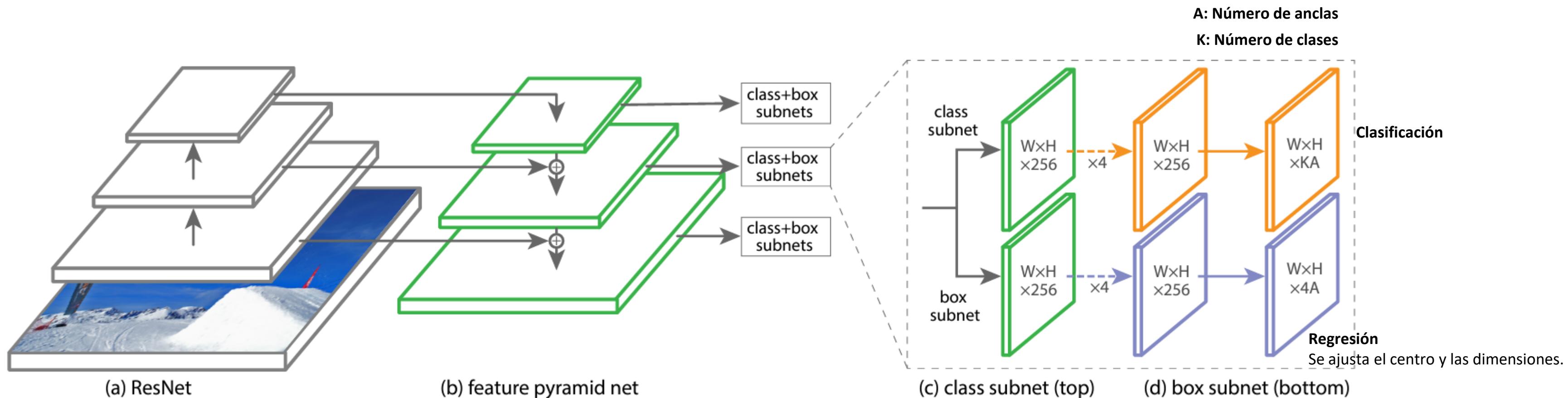
Clasificación

Regresión

Se ajusta el centro y las dimensiones.



RetinaNet



TRANSFORMATEC

Tsung-Yi Lin et al. (2018) "Focal Loss for Dense Object Detection".
IEEE International Conference on Computer Vision ICCV.

RetinaNet

Cross-entropy

$$\text{CE} = \begin{cases} -\log p, & y = 1 \\ -\log(1 - p), & \text{otro caso} \end{cases}$$



TRANSFORMATEC

Tsung-Yi Lin et al. (2018) "Focal Loss for Dense Object Detection".
IEEE International Conference on Computer Vision ICCV.

RetinaNet

Cross-entropy

$$\text{CE} = \begin{cases} -\log p, & y = 1 \\ -\log(1 - p), & \text{otro caso} \end{cases}$$

$$p_t = \begin{cases} p, & y = 1 \\ 1 - p, & \text{otro caso} \end{cases}$$



RetinaNet

Cross-entropy

$$\text{CE} = \begin{cases} -\log p, & y = 1 \\ -\log(1 - p), & \text{otro caso} \end{cases}$$

$$p_t = \begin{cases} p, & y = 1 \\ 1 - p, & \text{otro caso} \end{cases}$$

$p_t \rightarrow 1$
La clasificación es correcta



RetinaNet

Cross-entropy

$$\text{CE} = \begin{cases} -\log p, & y = 1 \\ -\log(1 - p), & \text{otro caso} \end{cases}$$

$$p_t = \begin{cases} p, & y = 1 \\ 1 - p, & \text{otro caso} \end{cases}$$

$p_t \rightarrow 1$
La clasificación es correcta

Resumiendo: $\text{CE} = -\log p_t$



RetinaNet

Cross-entropy

$$\text{CE} = \begin{cases} -\log p, & y = 1 \\ -\log(1 - p), & \text{otro caso} \end{cases}$$

$$p_t = \begin{cases} p, & y = 1 \\ 1 - p, & \text{otro caso} \end{cases}$$

$p_t \rightarrow 1$
La clasificación es correcta

Resumiendo: $\text{CE} = -\log p_t$

Asignamos pesos basados en p_t : $\text{FL} = -(1 - p_t)^\gamma \log p_t$



RetinaNet

Cross-entropy

$$\text{CE} = \begin{cases} -\log p, & y = 1 \\ -\log(1-p), & \text{otro caso} \end{cases}$$

$$p_t = \begin{cases} p, & y = 1 \\ 1-p, & \text{otro caso} \end{cases}$$

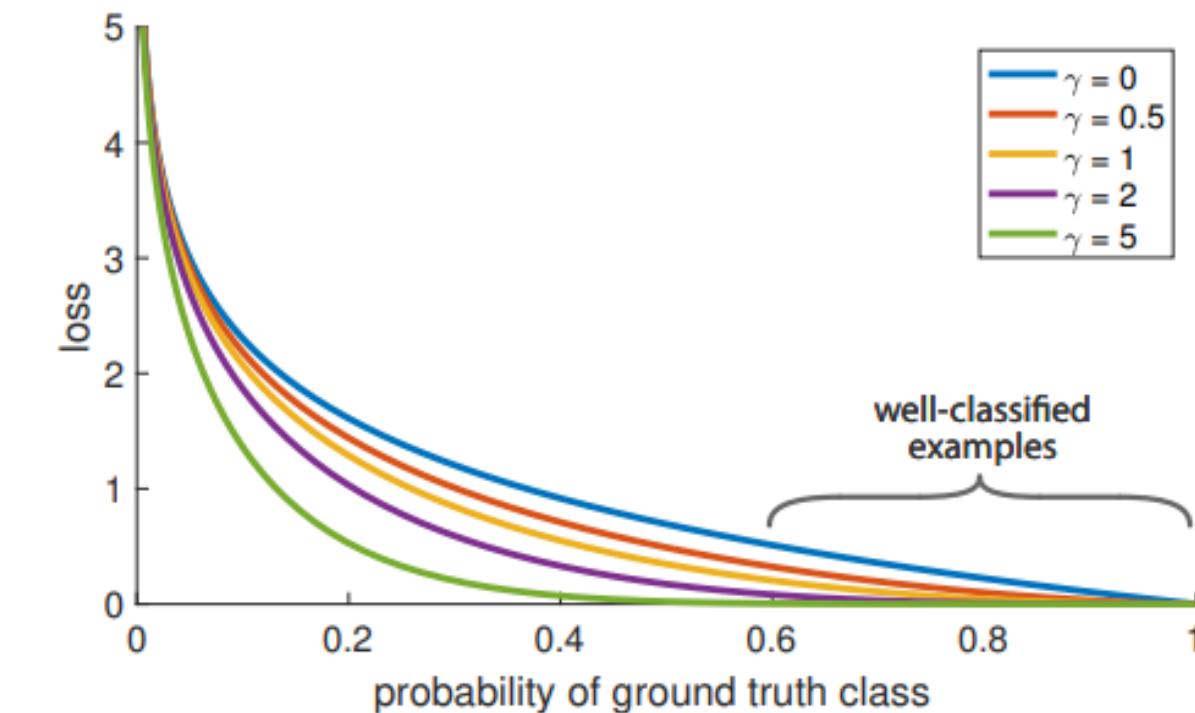
$p_t \rightarrow 1$
La clasificación es correcta

Resumiendo: $\text{CE} = -\log p_t$

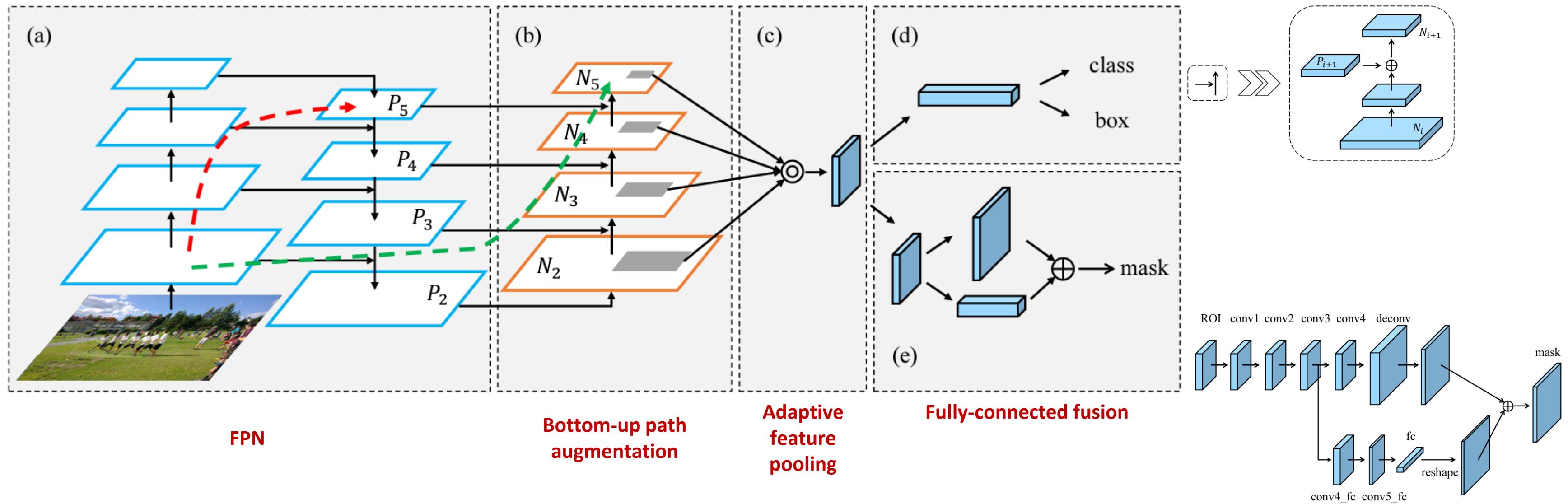
Asignamos pesos basados en p_t : $\text{FL} = -(1-p_t)^\gamma \log p_t$

$(1-p_t) \rightarrow 0$ Caso conocido

$(1-p_t) \rightarrow 1$ Caso desconocido



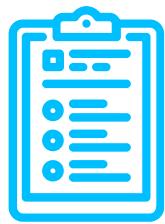
PA*Net*



TRANSFORMATEC

Shu Liu et al. (2018) "Path Aggregation Network for Instance Segmentation".
Proceedings of the IEEE conference on computer vision and pattern recognition (pp. 8759-8768).

4.



Mask R-CNN

TRANSFORMATEC

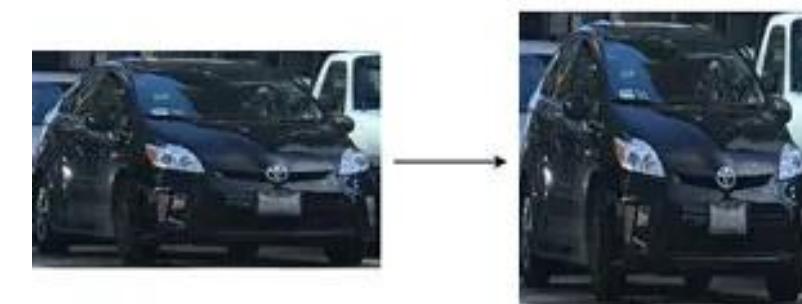
> Reinventa el mundo <



Sliding-window detectors



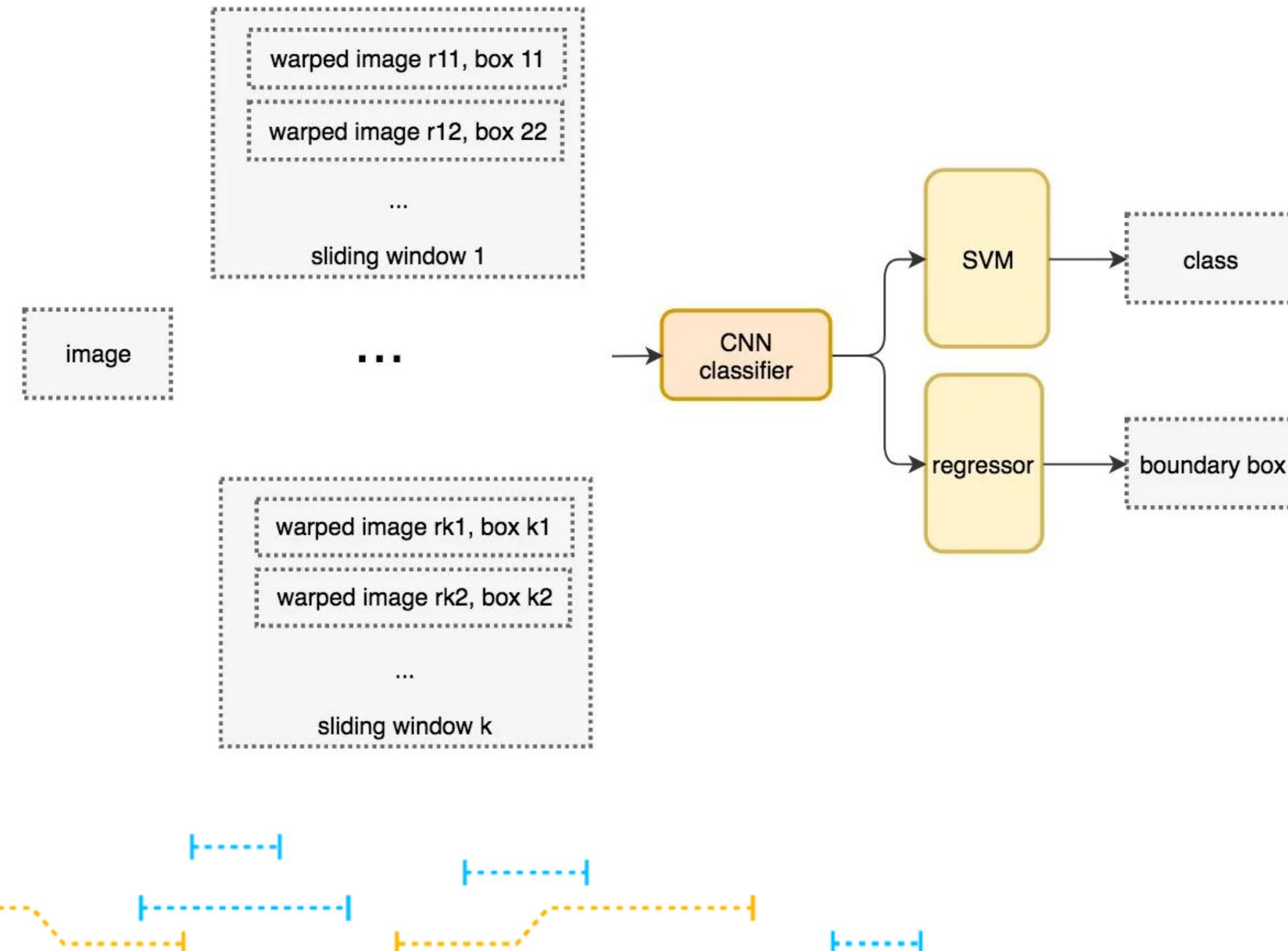
Sliding window



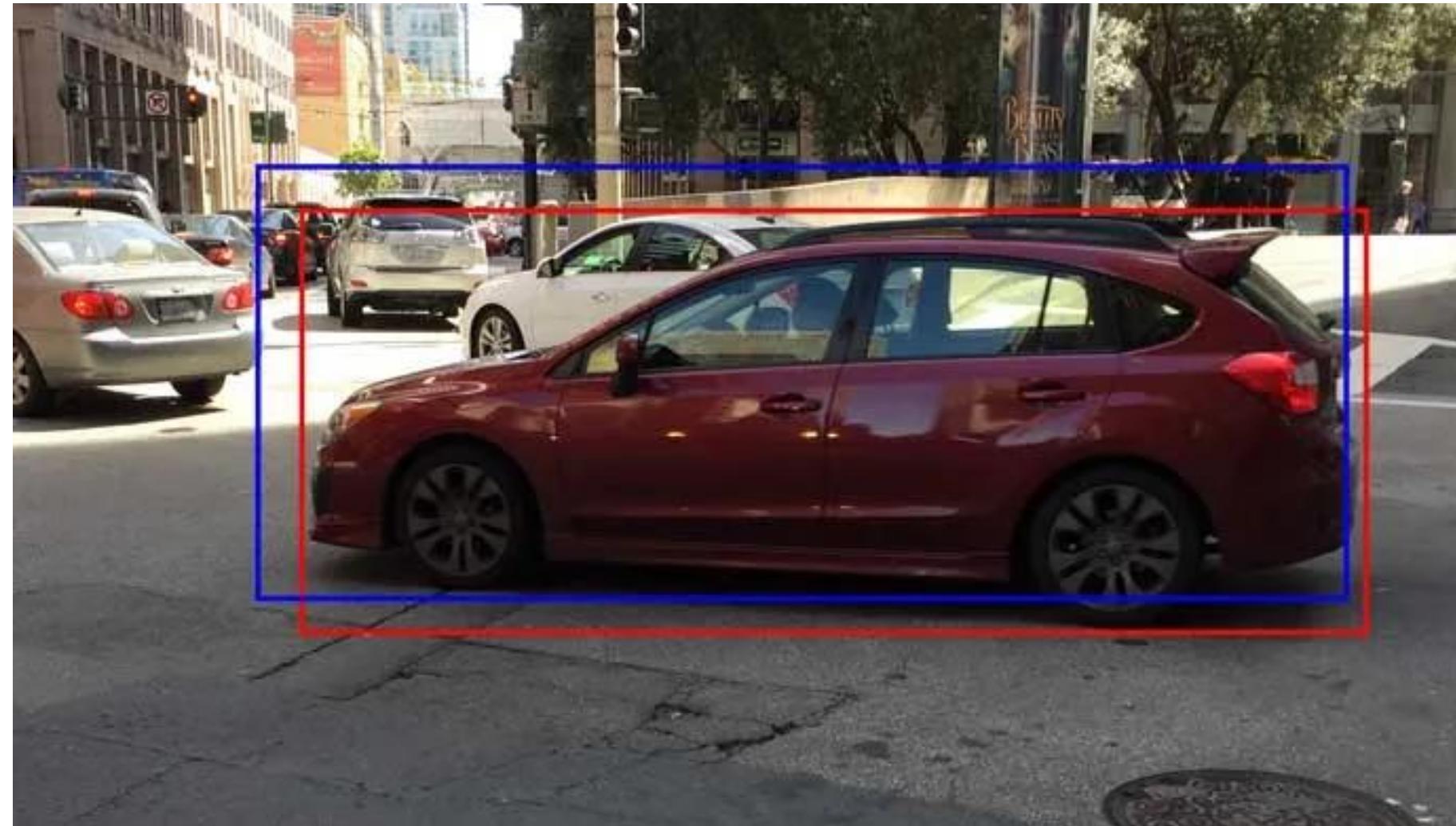
Warp



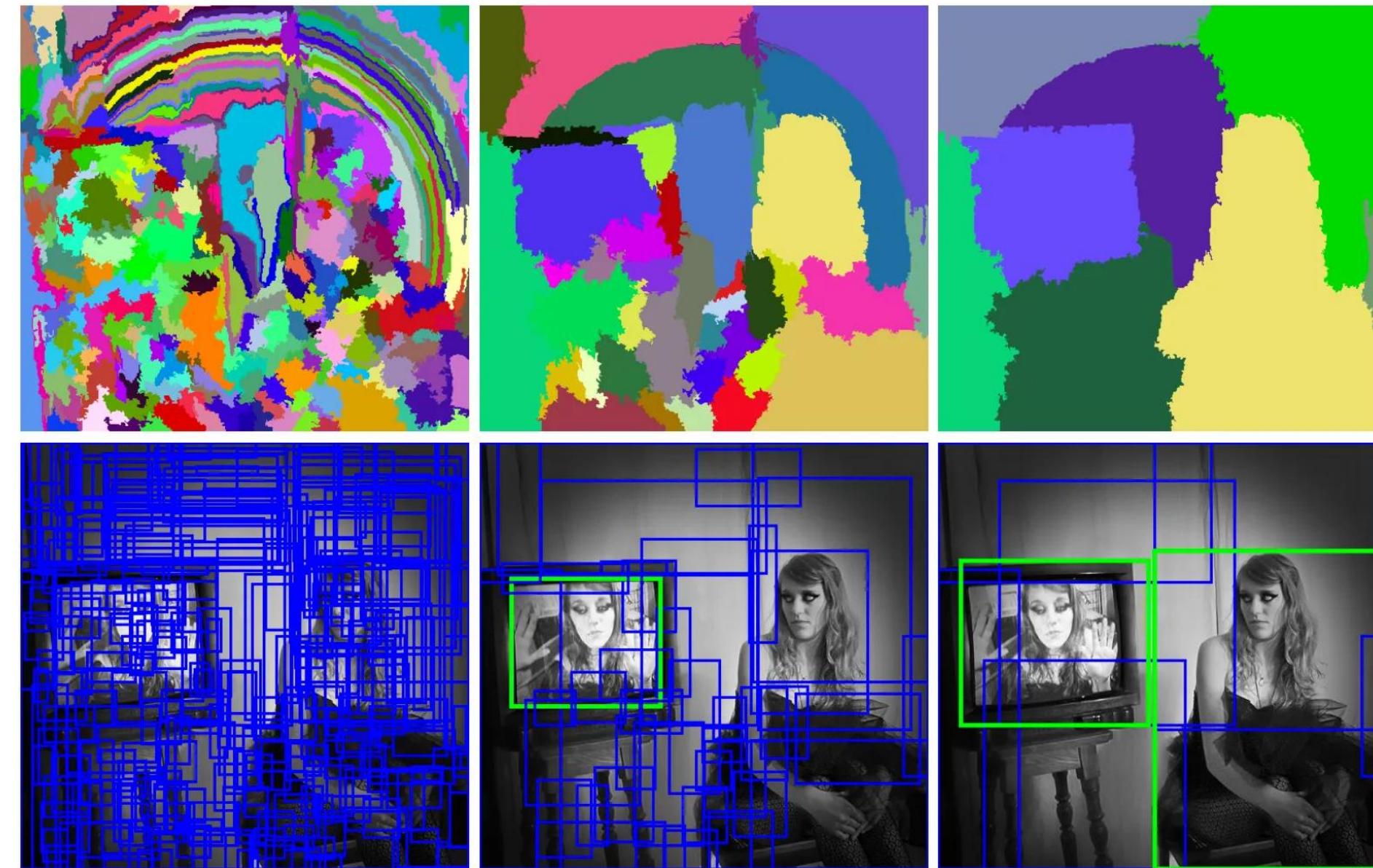
Sliding-window detectors



Sliding-window detectors



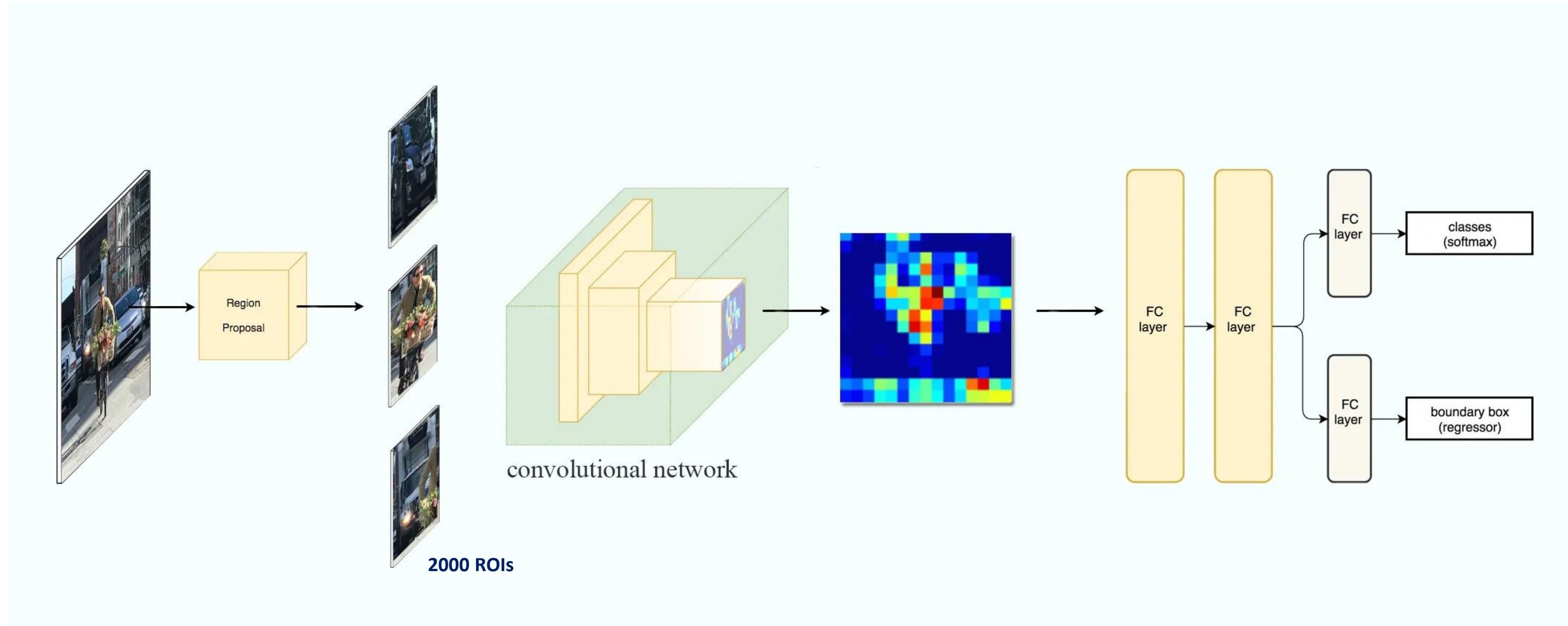
Selective Search



TRANSFORMATEC

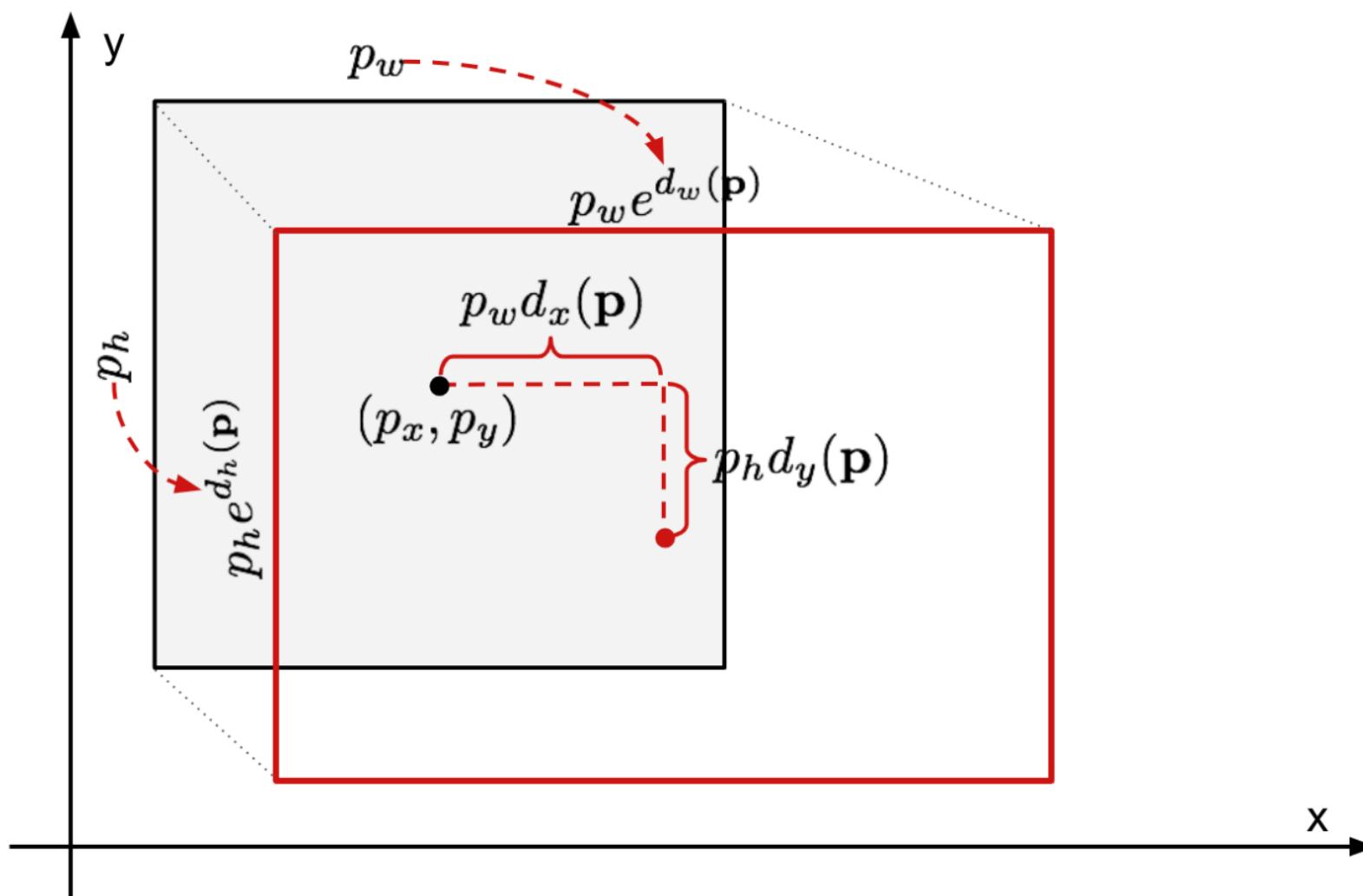
Koen E. A. van de Sande, et al. (2011) "Segmentation as Selective Search for Object Recognition".
2011 International Conference on Computer Vision (pp. 1879-1886). IEEE.

R-CNN



R-CNN

Bounding Box Regression

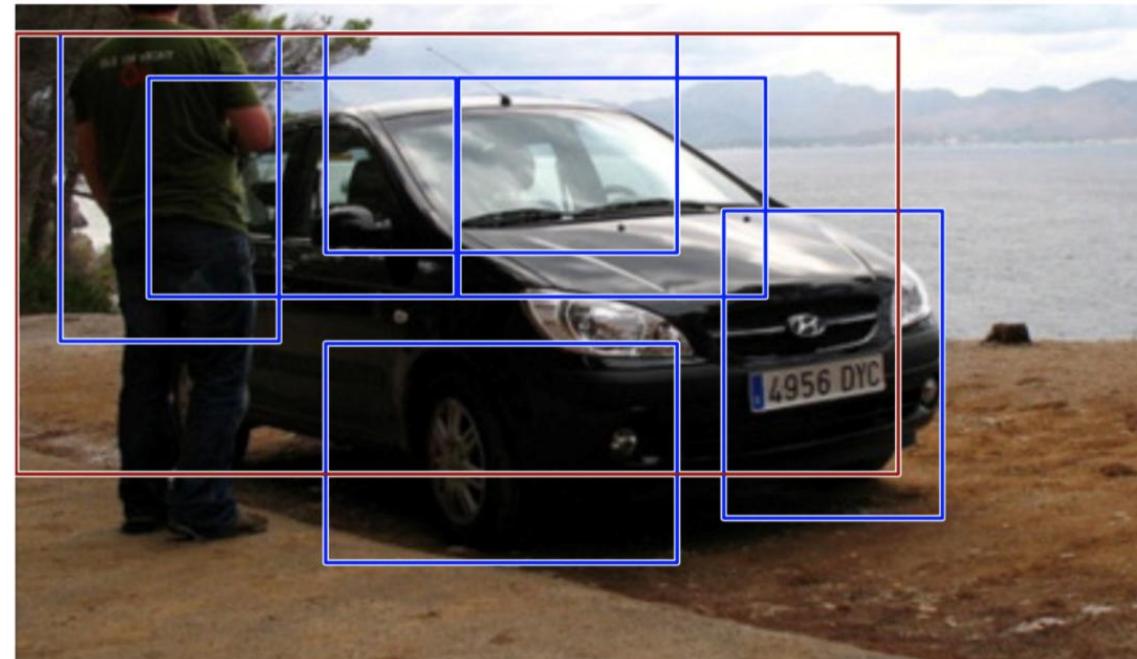


Ross Girshick, et al. (2014) "Rich feature hierarchies for accurate object detection and semantic segmentation". Proceedings of the IEEE conference on computer vision and pattern recognition. 2014. p. 580-587.

R-CNN

Non-Maximum Suppression

1. Los cuadros delimitadores se ordenan de acuerdo a sus puntuaciones, de mayor a menor.
2. Se selecciona el ROI con la puntuación más alta y se designa como una detección real.
3. Se eliminan todos los ROIs restantes que tienen IoU altos con los ROIs previamente seleccionados.



Before non-max suppression



After non-max suppression



Ross Girshick, et al. (2014) "Rich feature hierarchies for accurate object detection and semantic segmentation". Proceedings of the IEEE conference on computer vision and pattern recognition. 2014. p. 580-587.

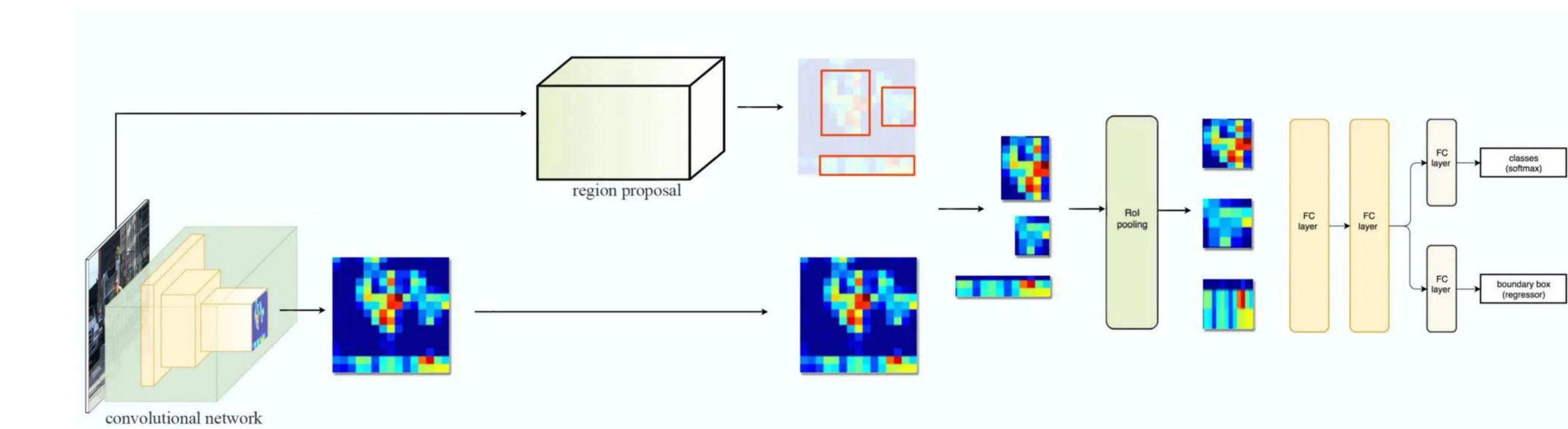
R-CNN

R-CNN es lento en entrenamiento e inferencia



Ross Girshick, et al. (2014) "Rich feature hierarchies for accurate object detection and semantic segmentation". Proceedings of the IEEE conference on computer vision and pattern recognition. 2014. p. 580-587.

Fast R-CNN

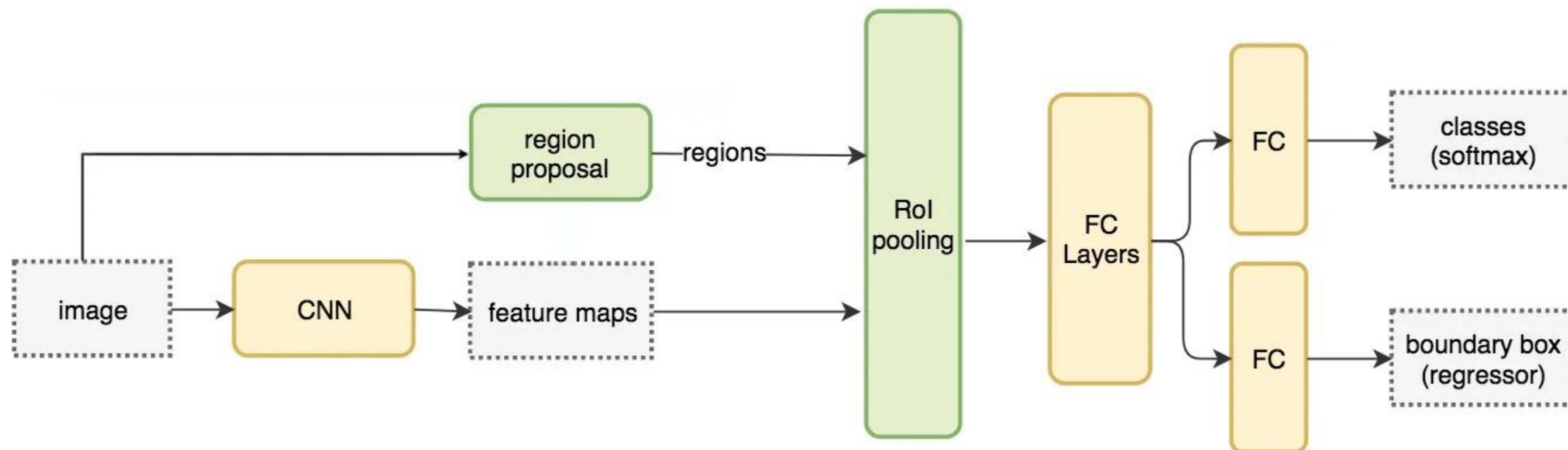


TRANSFORMATEC

Ross Girshick (2015) "Fast R-CNN".

Proceedings of the IEEE international conference on computer vision. 2015. p. 1440-1448.

Fast R-CNN



Fast R-CNN

ROI pooling

0.1	0.3	0.2	0.3	0.2	0.6	0.8	0.9
0.4	0.5	0.1	0.4	0.7	0.1	0.4	0.3
0.2	0.1	0.3	0.8	0.6	0.2	0.1	0.1
0.4	0.6	0.2	0.1	0.3	0.6	0.1	0.2
0.1	0.8	0.3	0.3	0.5	0.3	0.3	0.3
0.2	0.9	0.4	0.5	0.1	0.1	0.1	0.2
0.3	0.1	0.8	0.6	0.3	0.3	0.6	0.5
0.5	0.5	0.2	0.1	0.1	0.2	0.1	0.2

0.1	0.3	0.2	0.3	0.2	0.6	0.8	0.9
0.4	0.5	0.1	0.4	0.7	0.1	0.4	0.3
0.2	0.1	0.3	0.8	0.6	0.2	0.1	0.1
0.4	0.6	0.2	0.1	0.3	0.6	0.1	0.2
0.1	0.8	0.3	0.3	0.5	0.3	0.3	0.3
0.2	0.9	0.4	0.5	0.1	0.1	0.1	0.2
0.3	0.1	0.8	0.6	0.3	0.3	0.6	0.5
0.5	0.5	0.2	0.1	0.1	0.2	0.1	0.2

Buscamos el máximo valor

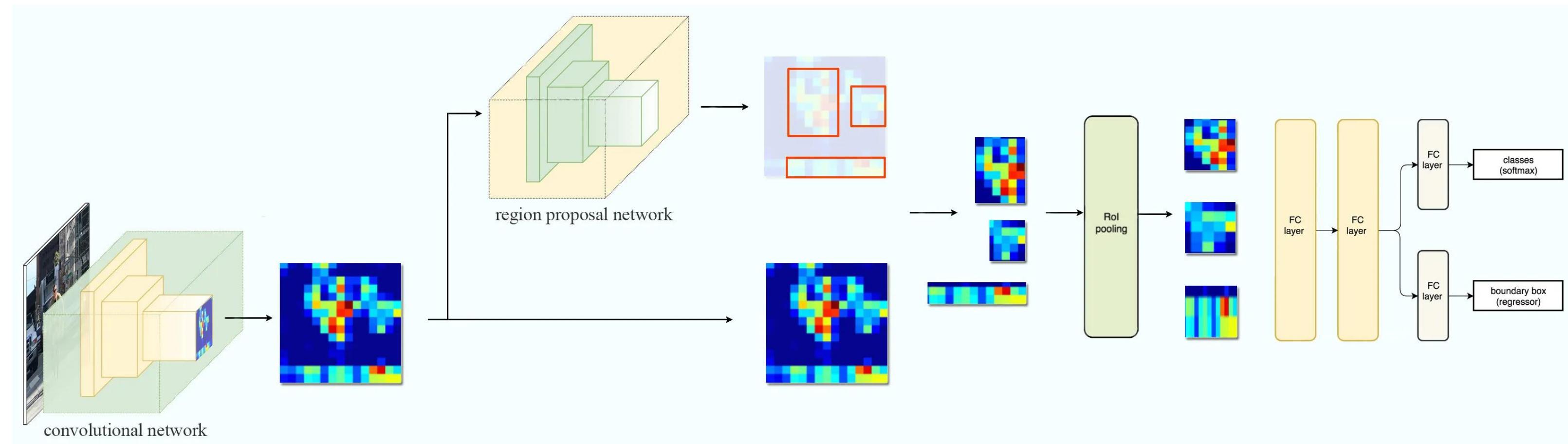


0.8	0.6
0.9	0.6

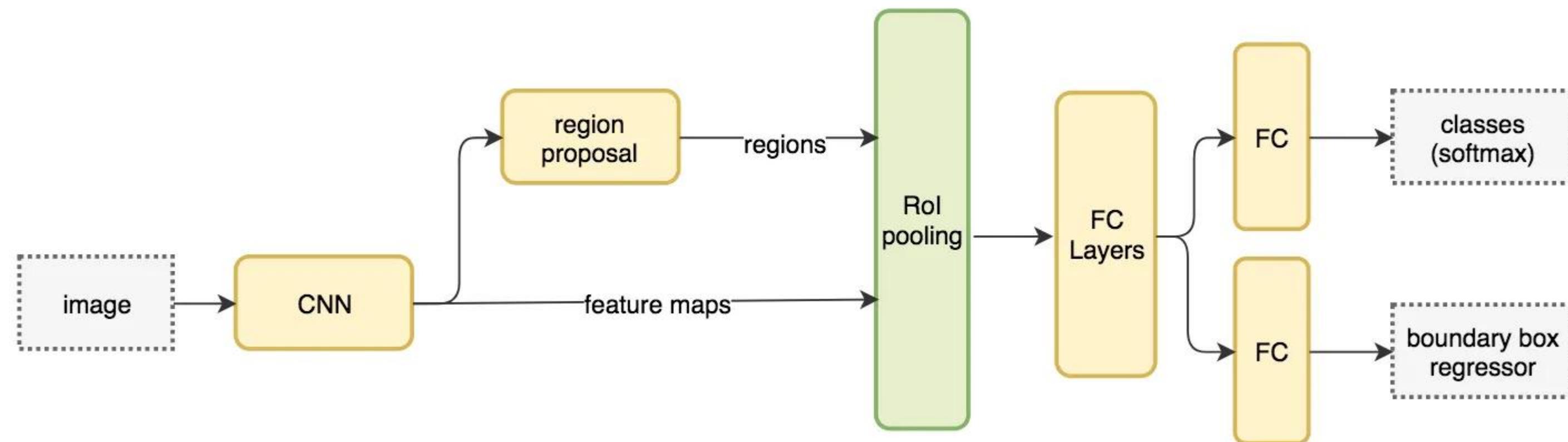
TRANSFORMATEC

Ross Girshick (2015) "Fast R-CNN".
Proceedings of the IEEE international conference on computer vision. 2015. p. 1440-1448.

Faster R-CNN

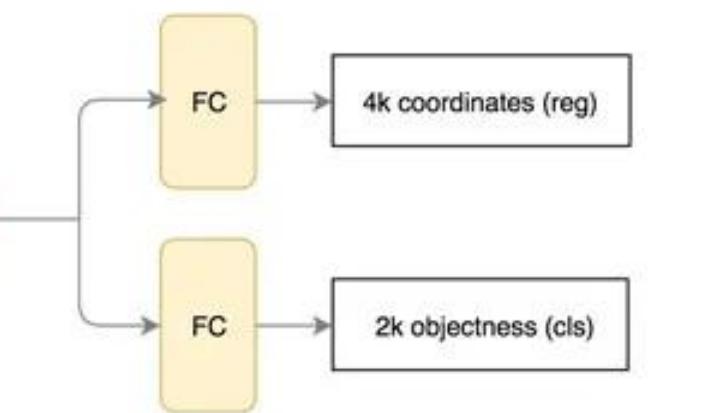
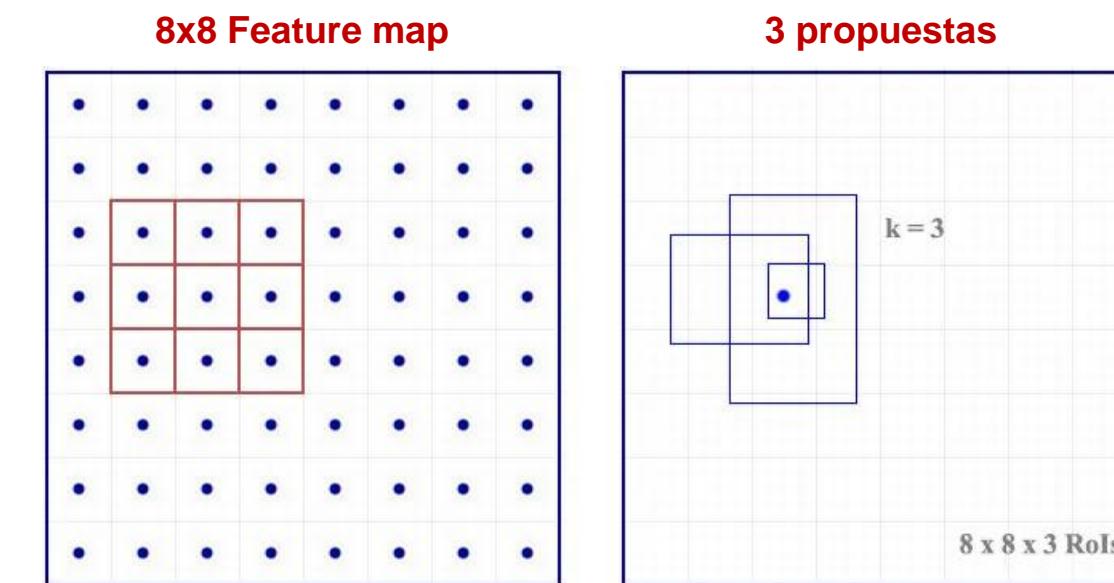
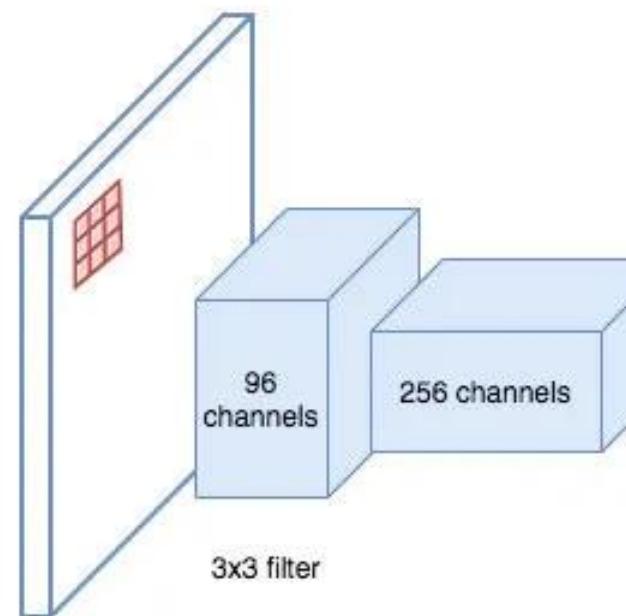


Faster *R-CNN*



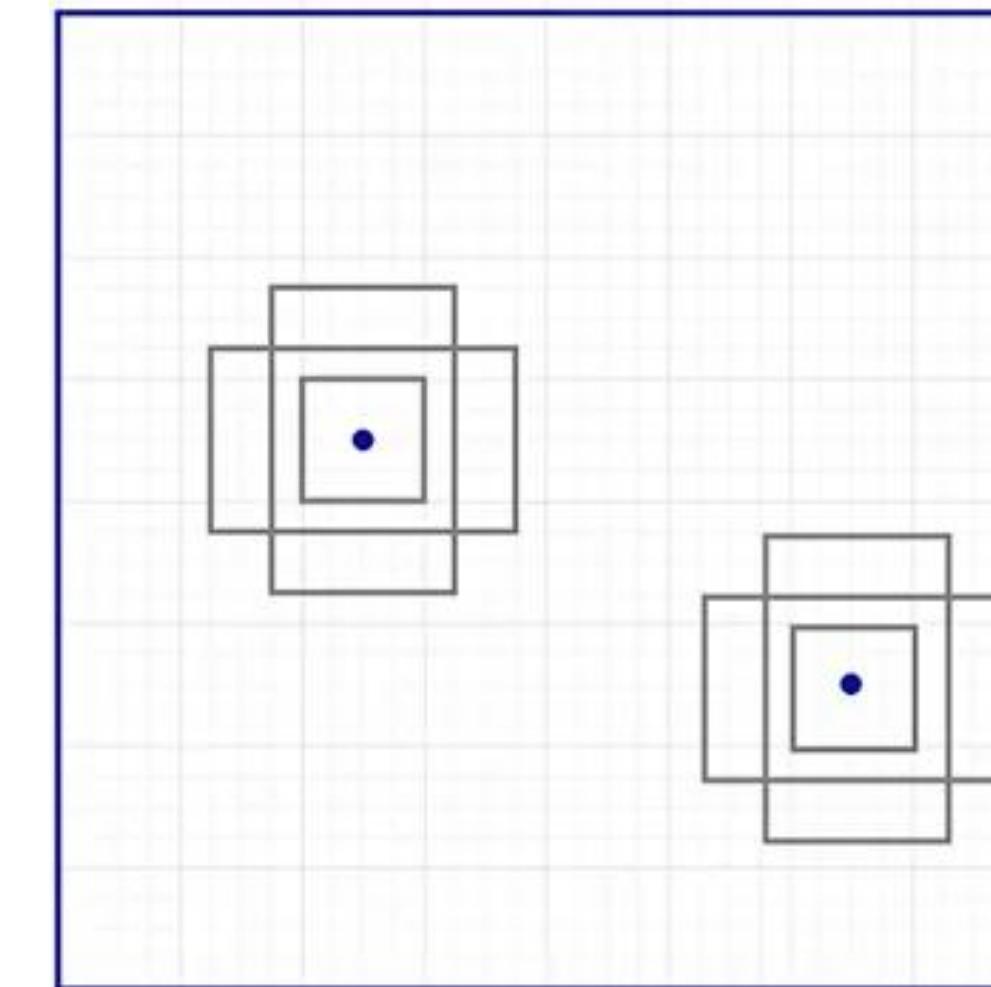
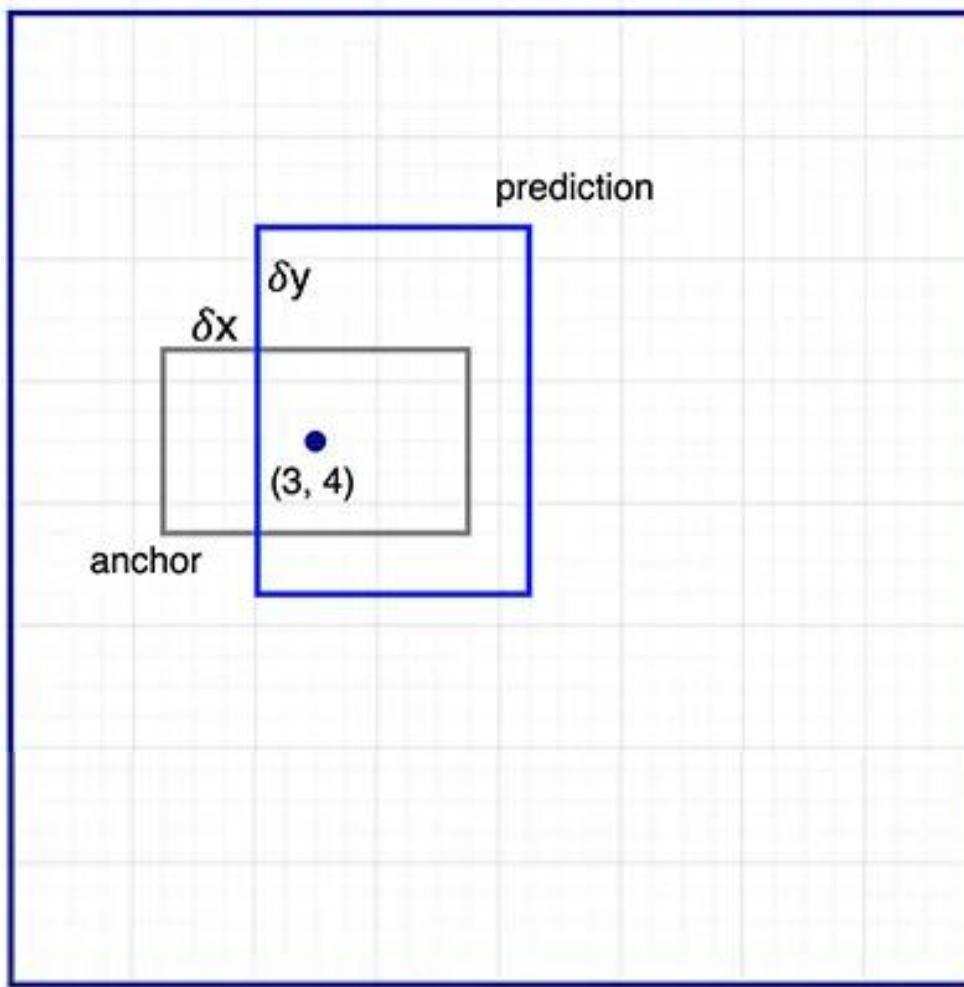
Faster R-CNN

Region proposal network

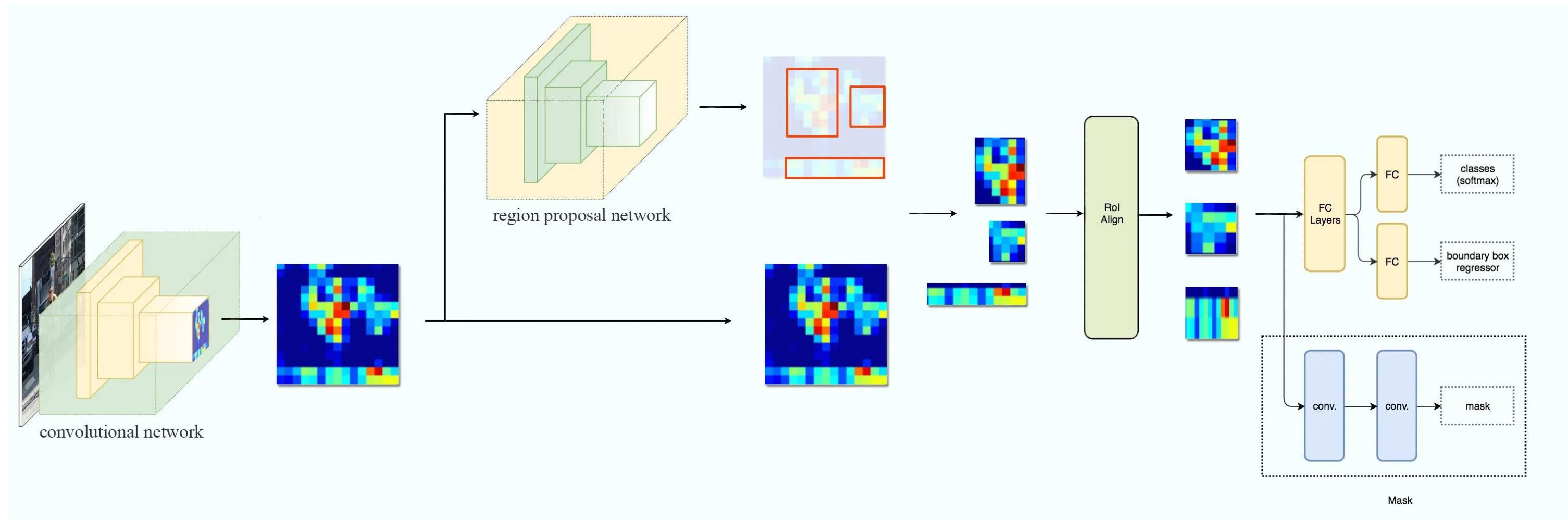


Faster *R-CNN*

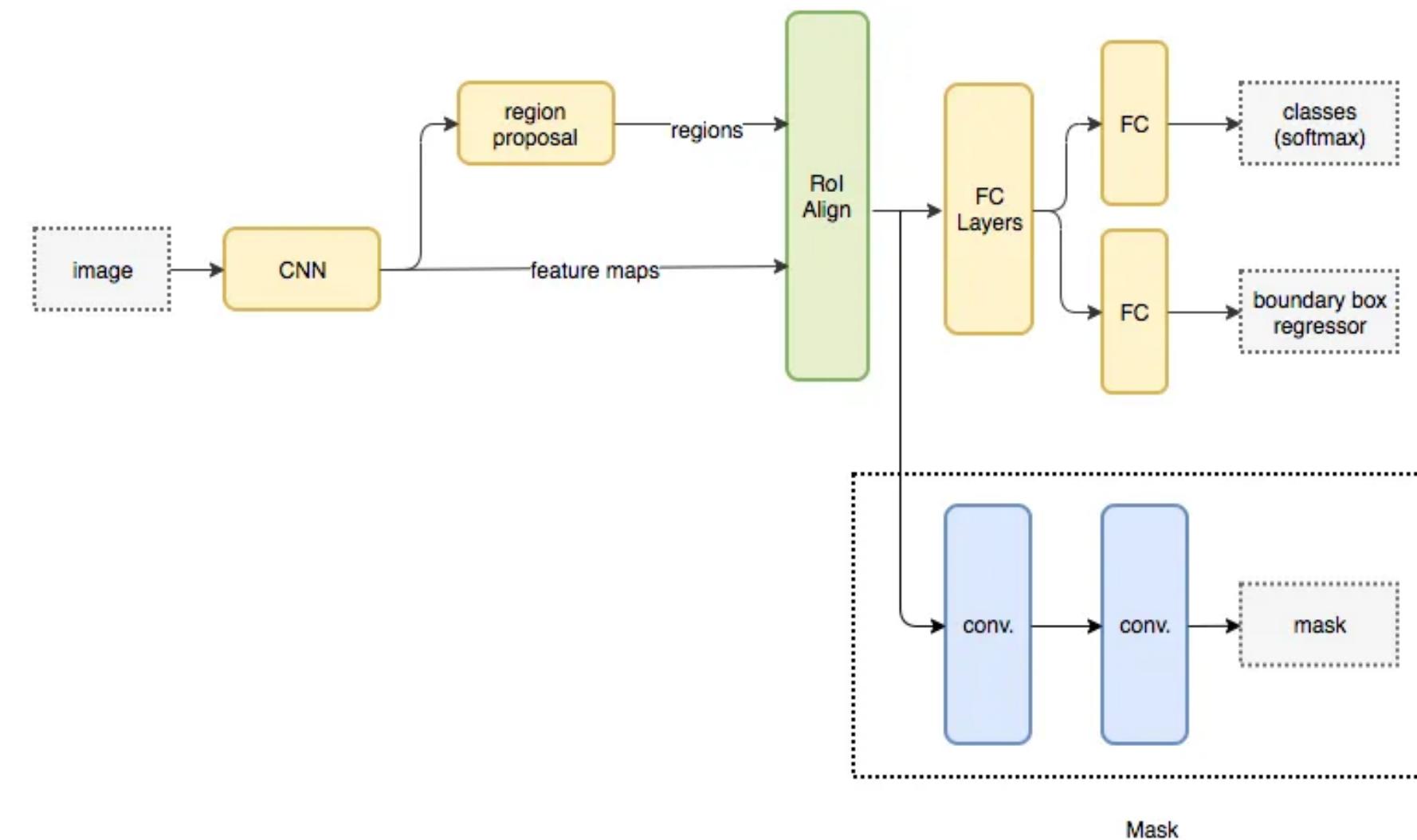
Region proposal network



Mask R-CNN



Mask R-CNN



Mask *R-CNN*

ROI pooling

0.1	0.3	0.2	0.3	0.2	0.6	0.8	0.9
0.4	0.5	0.1	0.4	0.7	0.1	0.4	0.3
0.2	0.1	0.3	0.8	0.6	0.2	0.1	0.1
0.4	0.6	0.2	0.1	0.3	0.6	0.1	0.2
0.1	0.8	0.3	0.3	0.5	0.3	0.3	0.3
0.2	0.9	0.4	0.5	0.1	0.1	0.1	0.2
0.3	0.1	0.8	0.6	0.3	0.3	0.6	0.5
0.5	0.5	0.2	0.1	0.1	0.2	0.1	0.2

0.8	0.6
0.9	0.6

ROI Align

0.1	0.3	0.2	0.3	0.2	0.6	0.8	0.9
0.4	0.5	0.1	0.4	0.7	0.1	0.4	0.3
0.2	0.1	0.3	0.8	0.6	0.2	0.1	0.1
0.4	0.6	0.2	0.	0.3	0.6	0.1	0.2
0.1	0.8	0.3	0.3	0.5	0.3	0.3	0.3
0.2	0.9	0.4	0.5	0.1	0.1	0.1	0.2
0.3	0.1	0.8	0.6	0.3	0.3	0.6	0.5
0.5	0.5	0.2	0.1	0.1	0.2	0.1	0.2

0.88	0.6
0.9	0.6



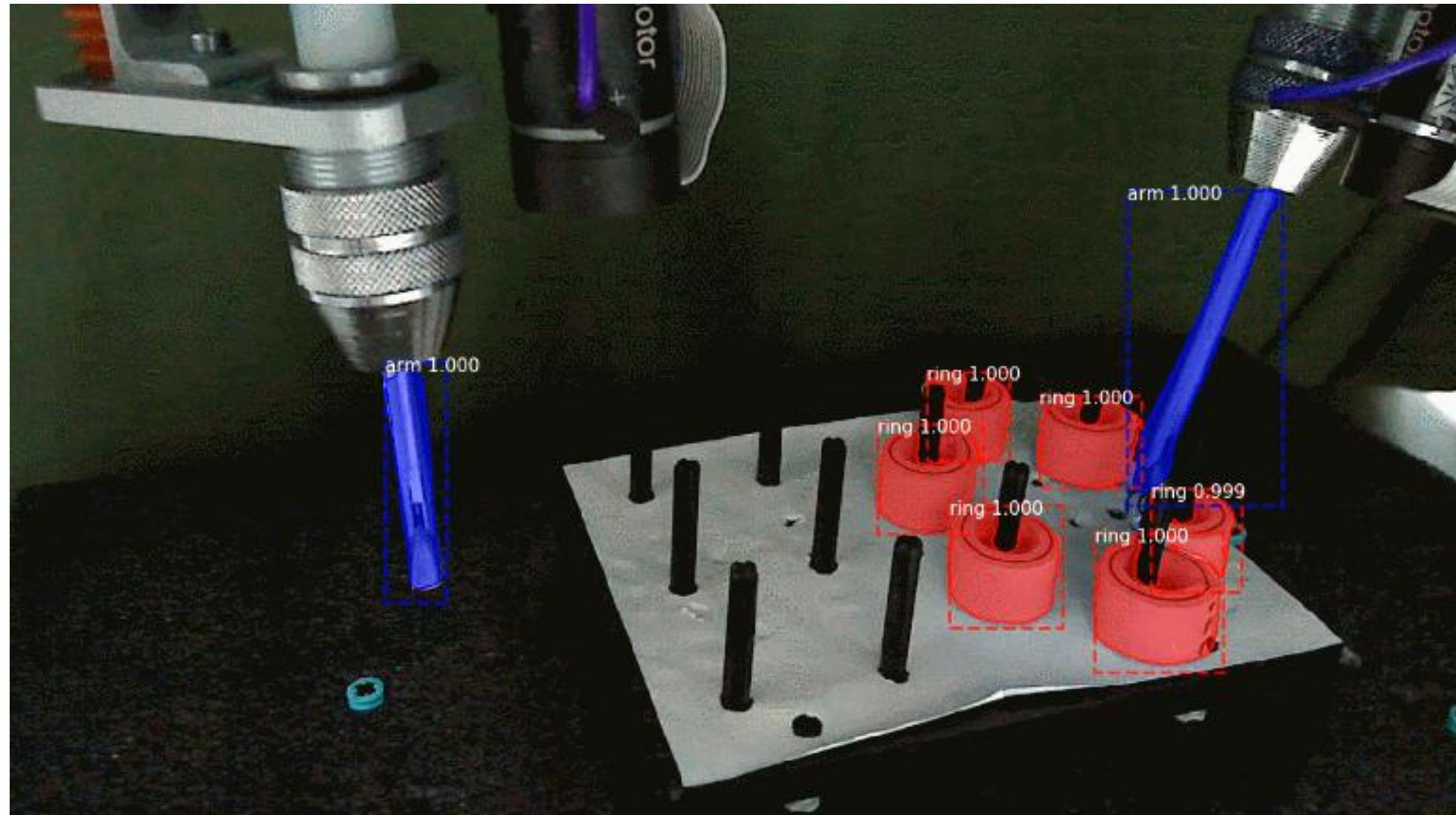
TRANSFORMATEC

Kaiming He et al. (2017) "Mask R-CNN".
Proceedings of the IEEE international conference on computer vision (pp. 2961-2969).

Mask R-CNN



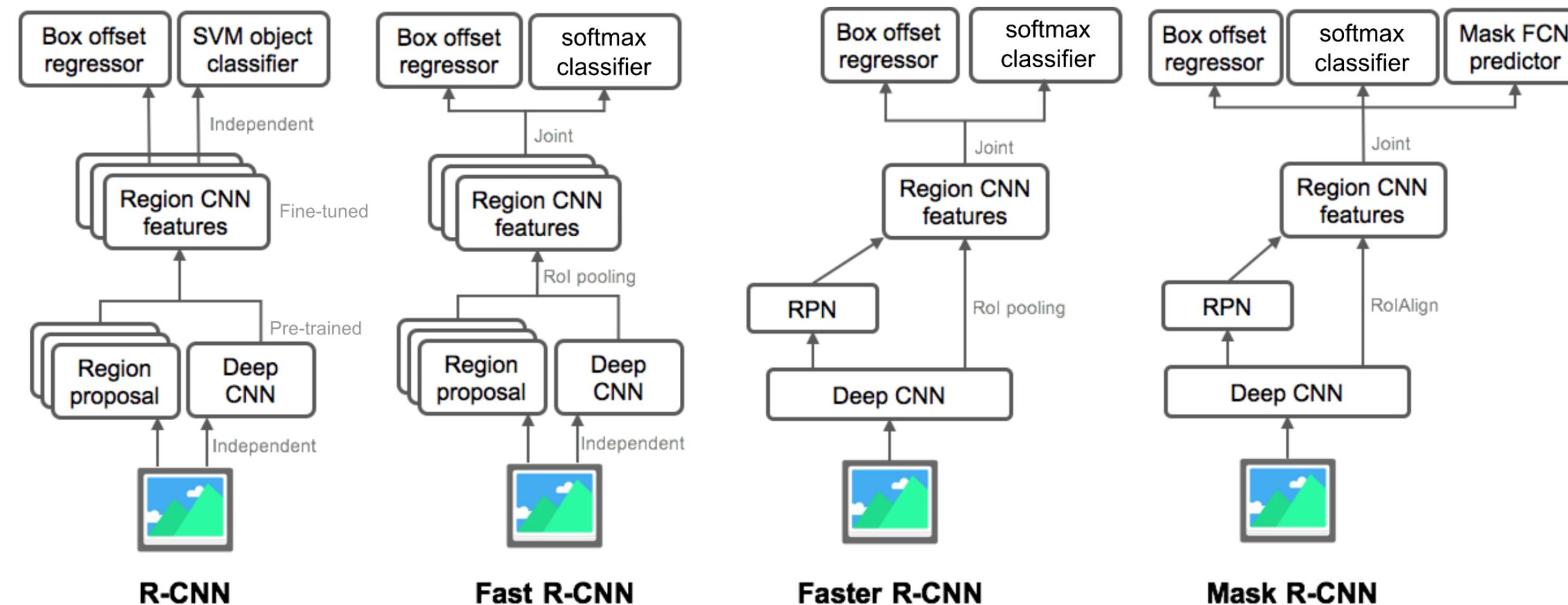
Mask R-CNN



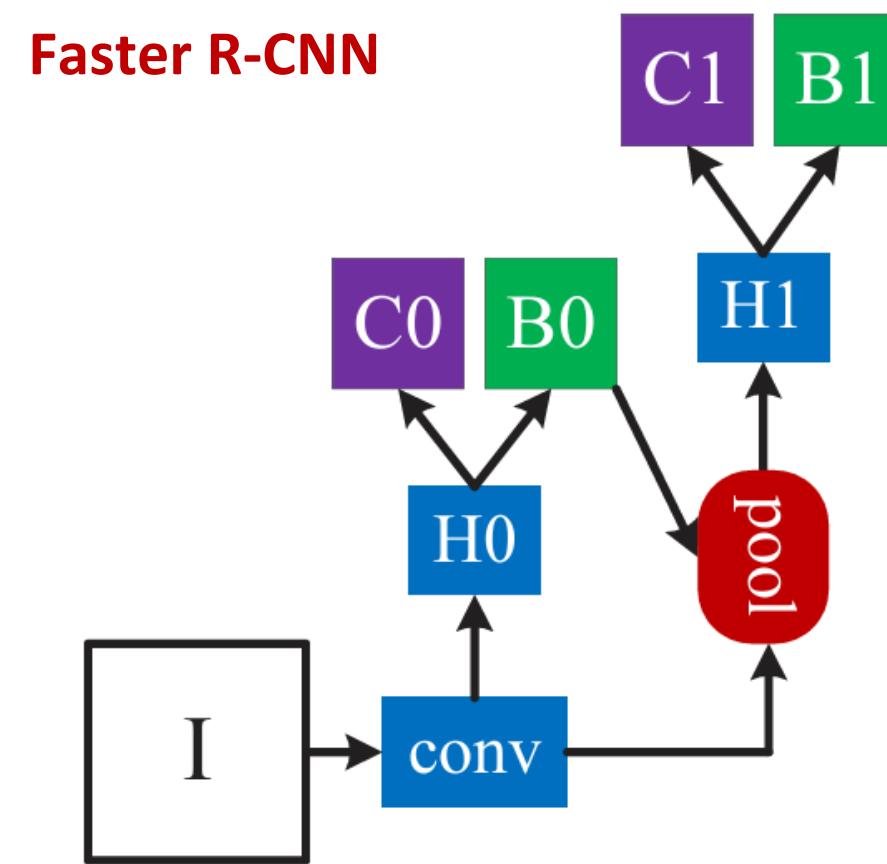
TRANSFORMATEC

Kaiming He et al. (2017) "Mask R-CNN".
Proceedings of the IEEE international conference on computer vision (pp. 2961-2969).

R-CNN family

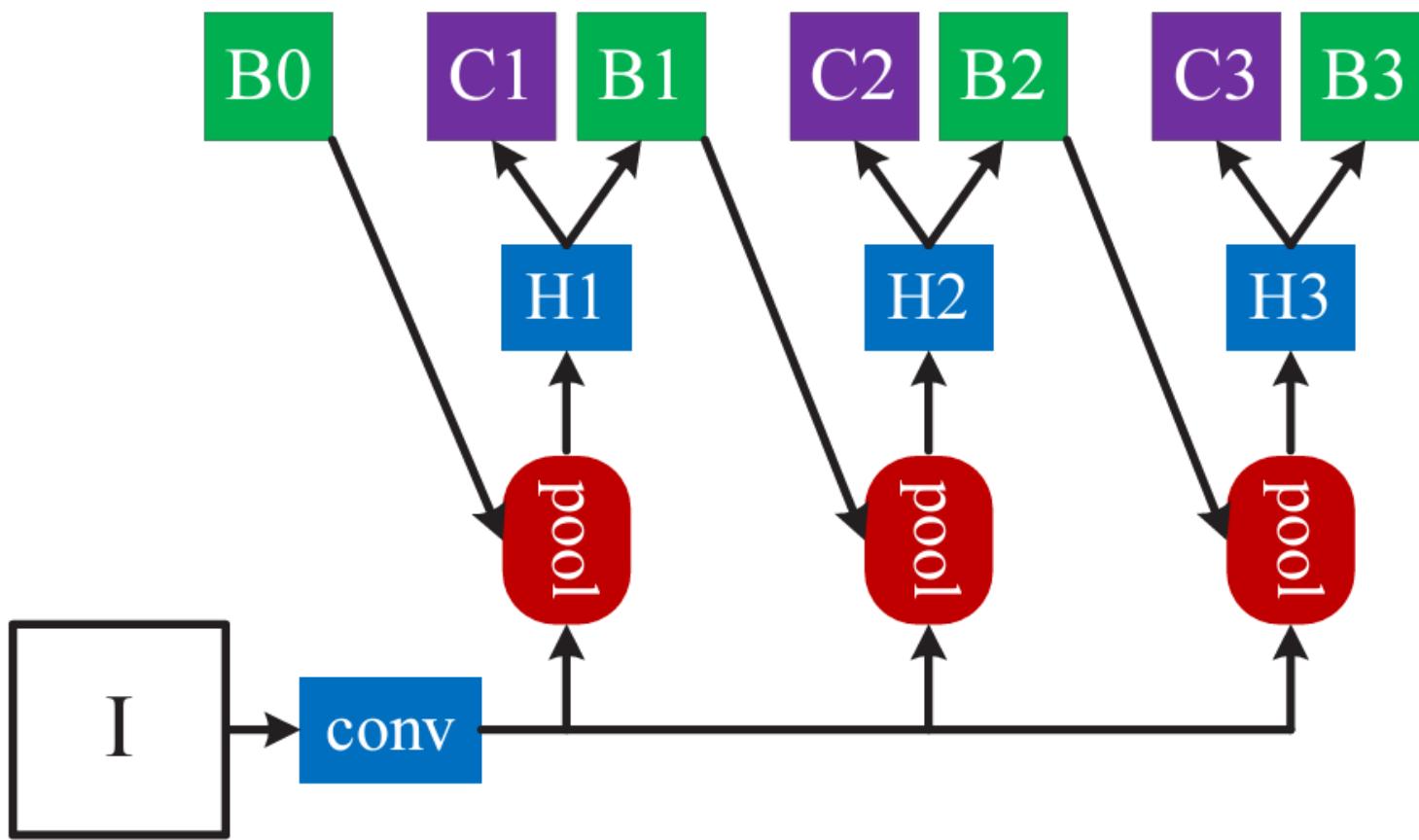
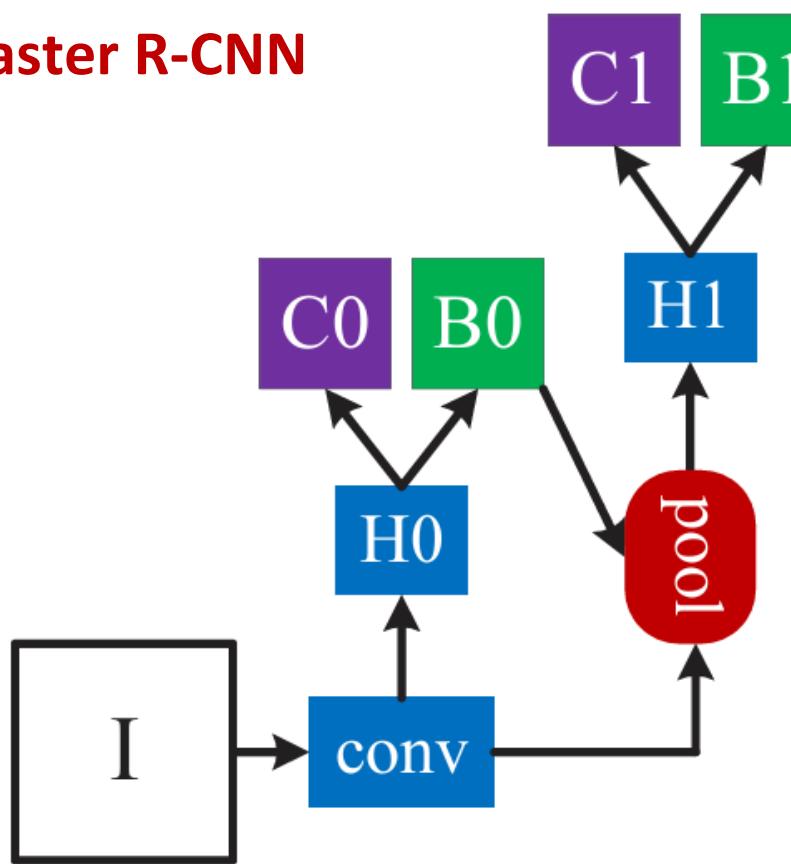


Cascade *R-CNN*



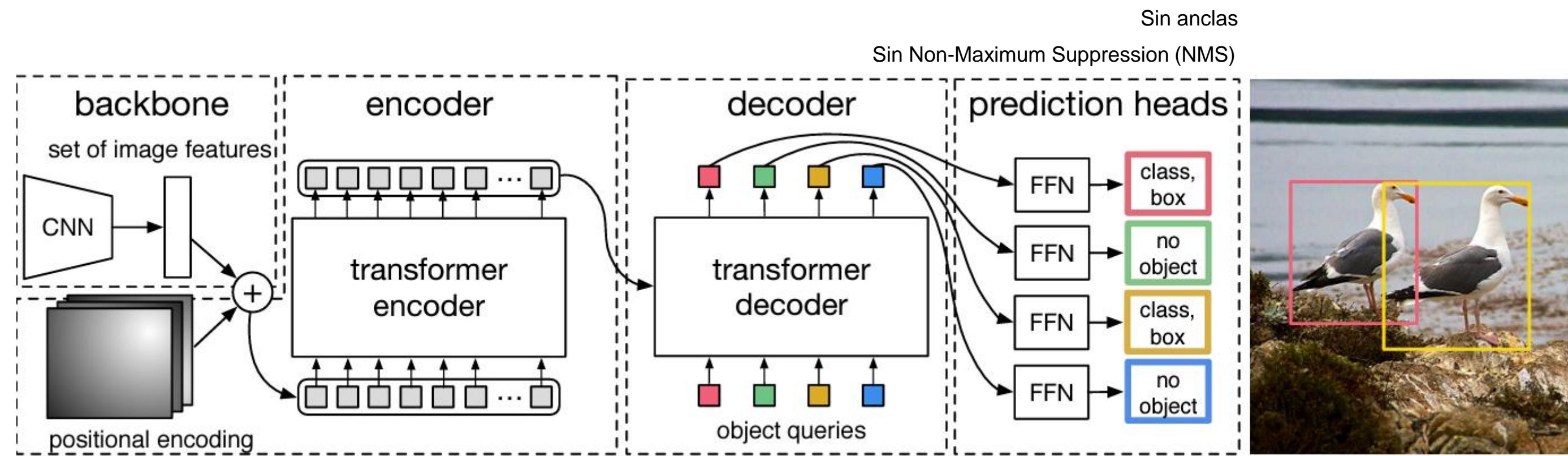
Cascade *R-CNN*

Faster R-CNN



DETR

(DEtection TRansformer)

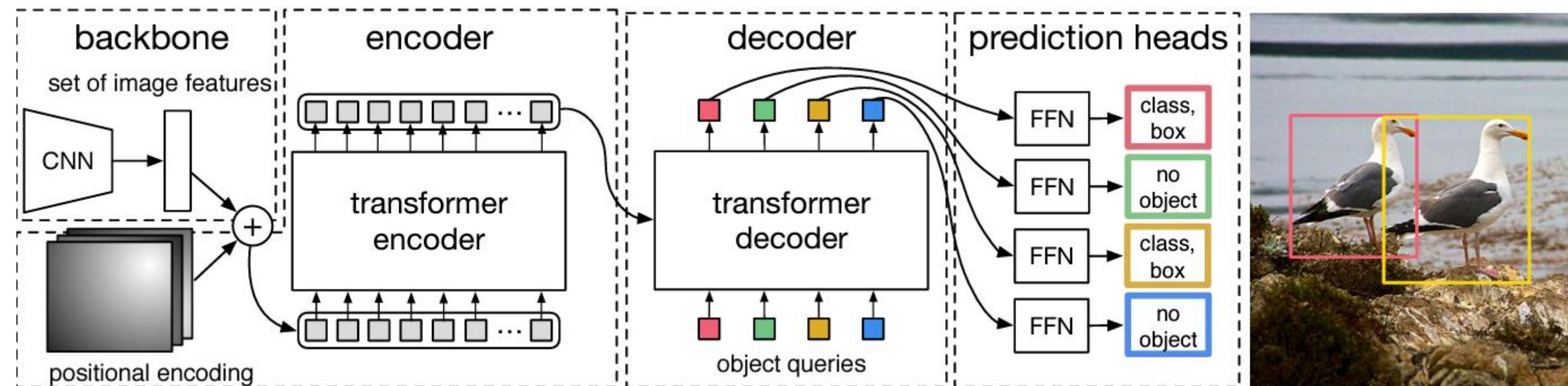


TRANSFORMATEC

Nicolas Carion et al. (2020) "End-to-End Object Detection with Transformers". European conference on computer vision. Cham: Springer International Publishing, 2020.

DETR

(DEtection TRansformer)



$$\mathcal{L}_{\text{Hungarian}}(y, \hat{y}) = \sum_{i=1}^N \left[-\log \hat{p}_{\hat{\sigma}(i)}(c_i) + \mathbb{1}_{\{c_i \neq \emptyset\}} \mathcal{L}_{\text{box}}(b_i, \hat{b}_{\hat{\sigma}(i)}) \right]$$

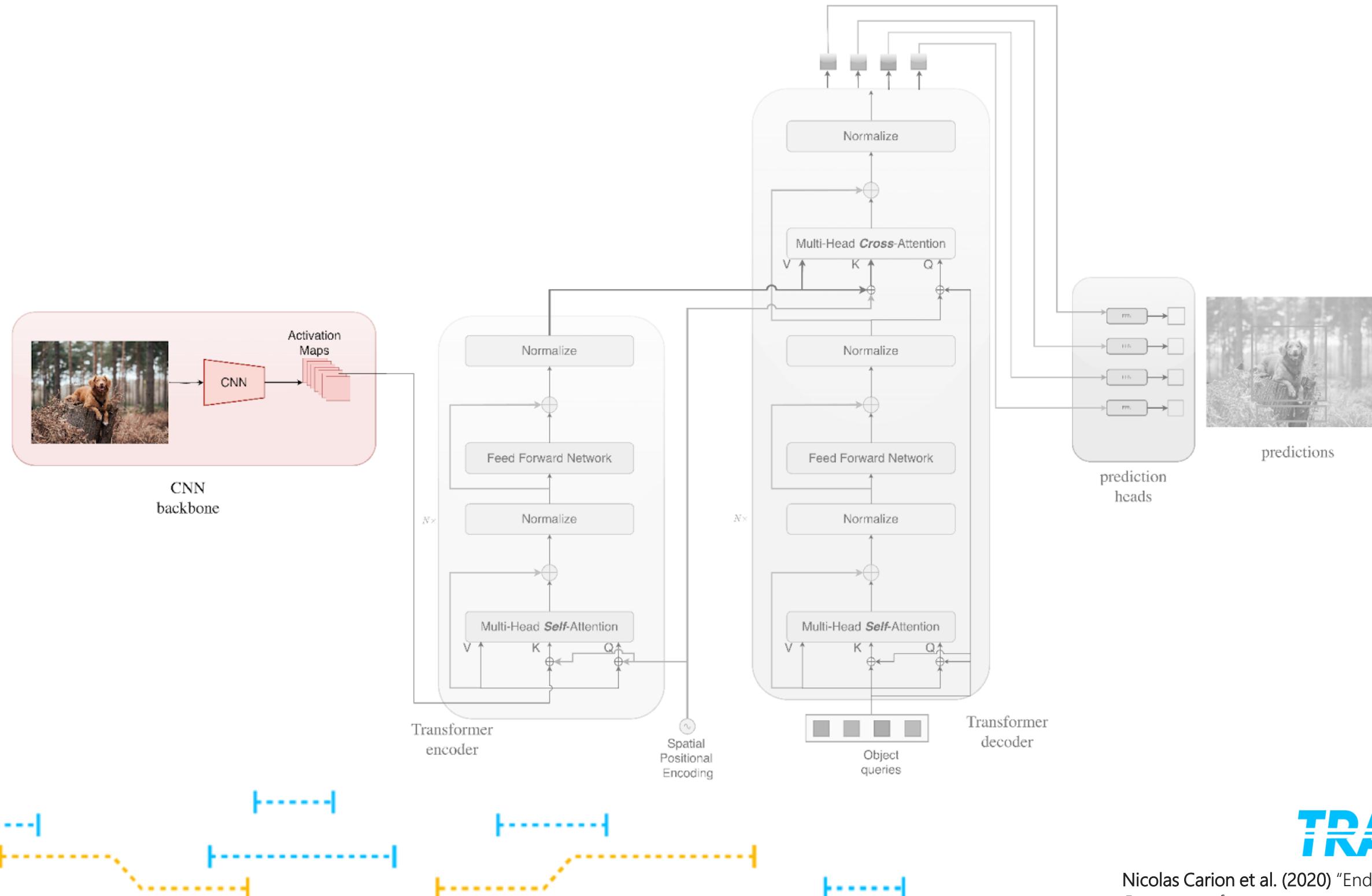


TRANSFORMATEC

Nicolas Carion et al. (2020) "End-to-End Object Detection with Transformers". European conference on computer vision. Cham: Springer International Publishing, 2020.

DETR

(DEtection TRansformer)

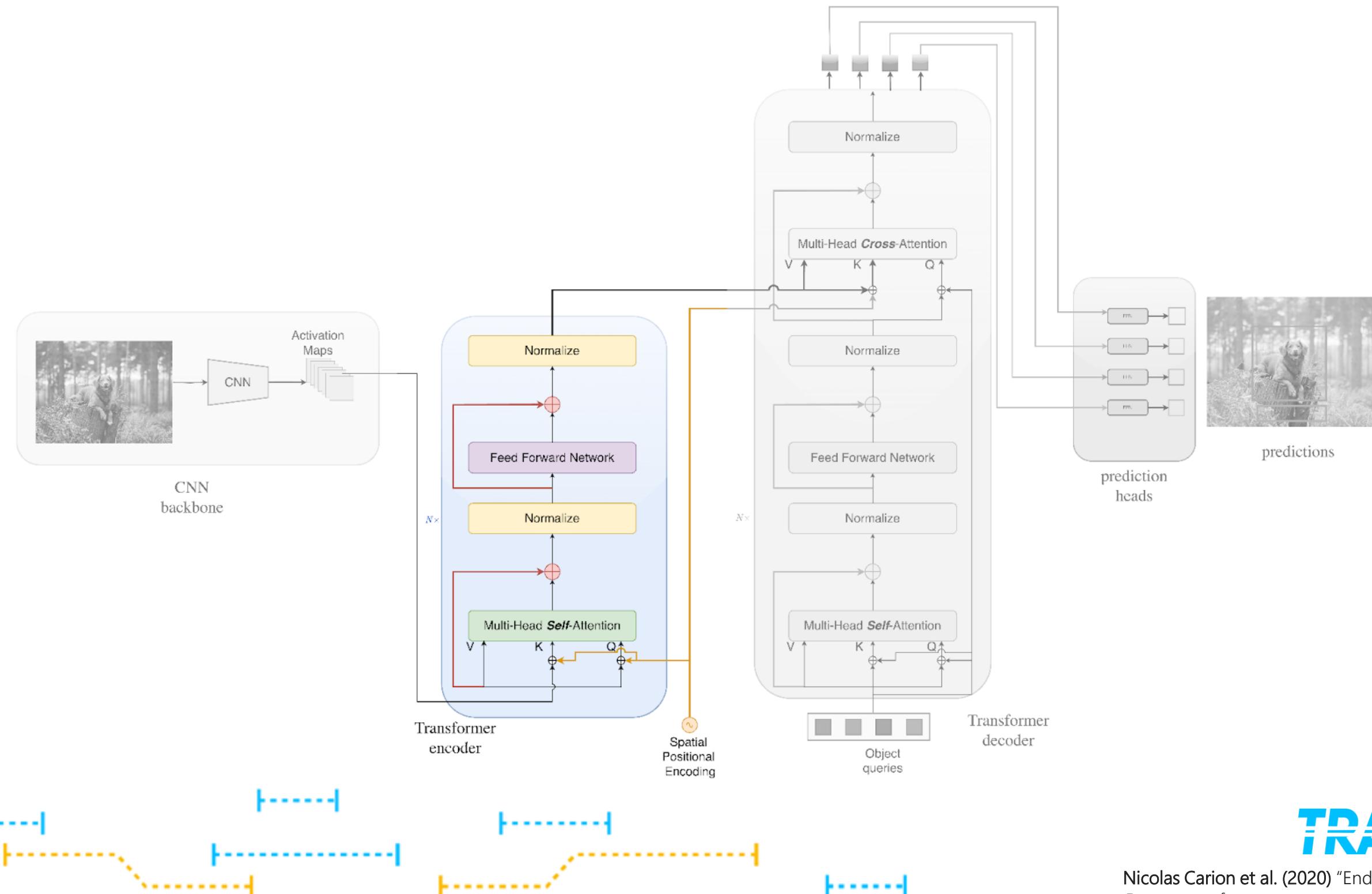


TRANSFORMATEC

Nicolas Carion et al. (2020) "End-to-End Object Detection with Transformers".
European conference on computer vision. Cham: Springer International Publishing, 2020.

DETR

(DEtection TRansformer)

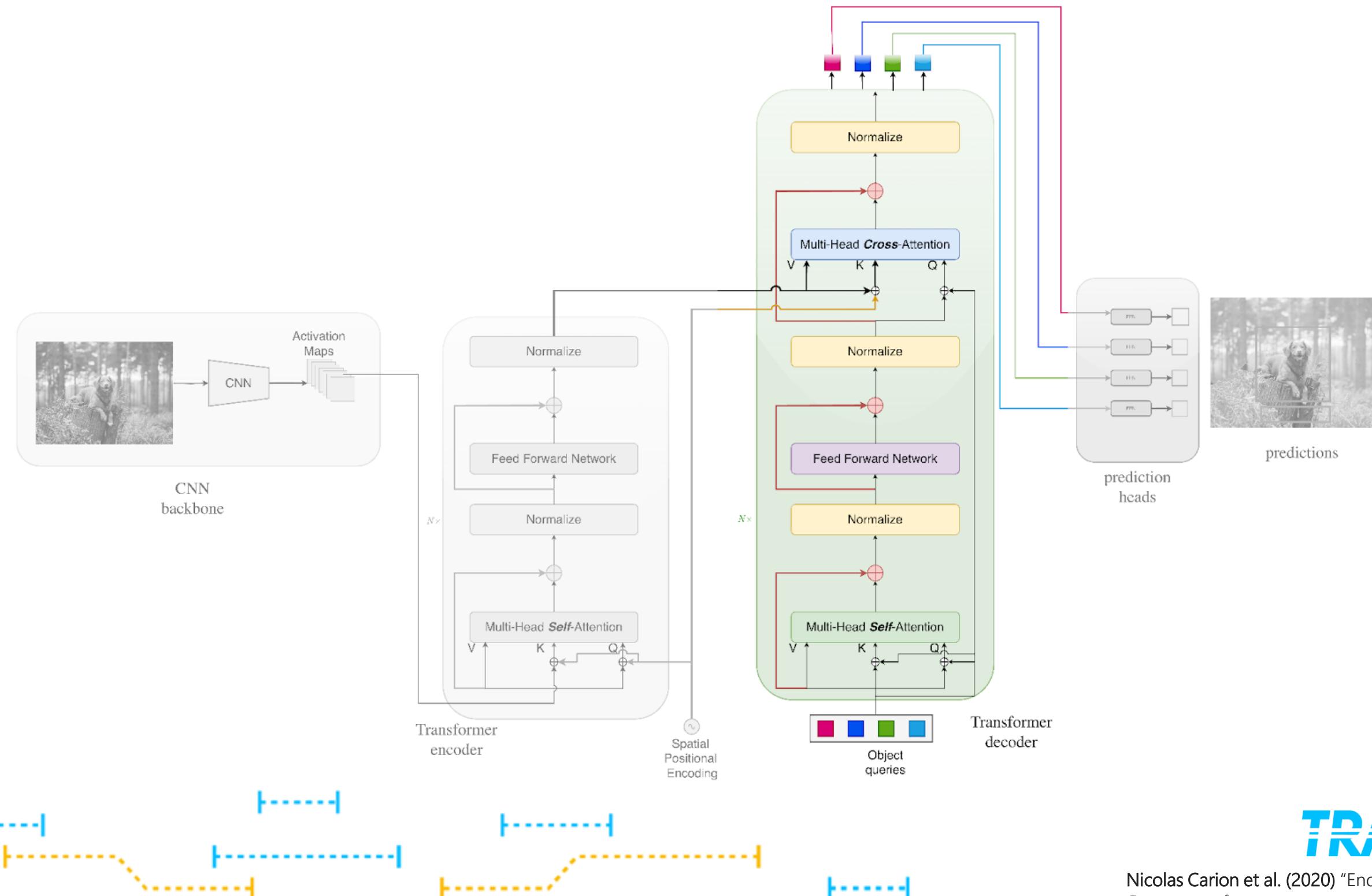


TRANSFORMATEC

Nicolas Carion et al. (2020) "End-to-End Object Detection with Transformers". European conference on computer vision. Cham: Springer International Publishing, 2020.

DETR

(DEtection TRansformer)

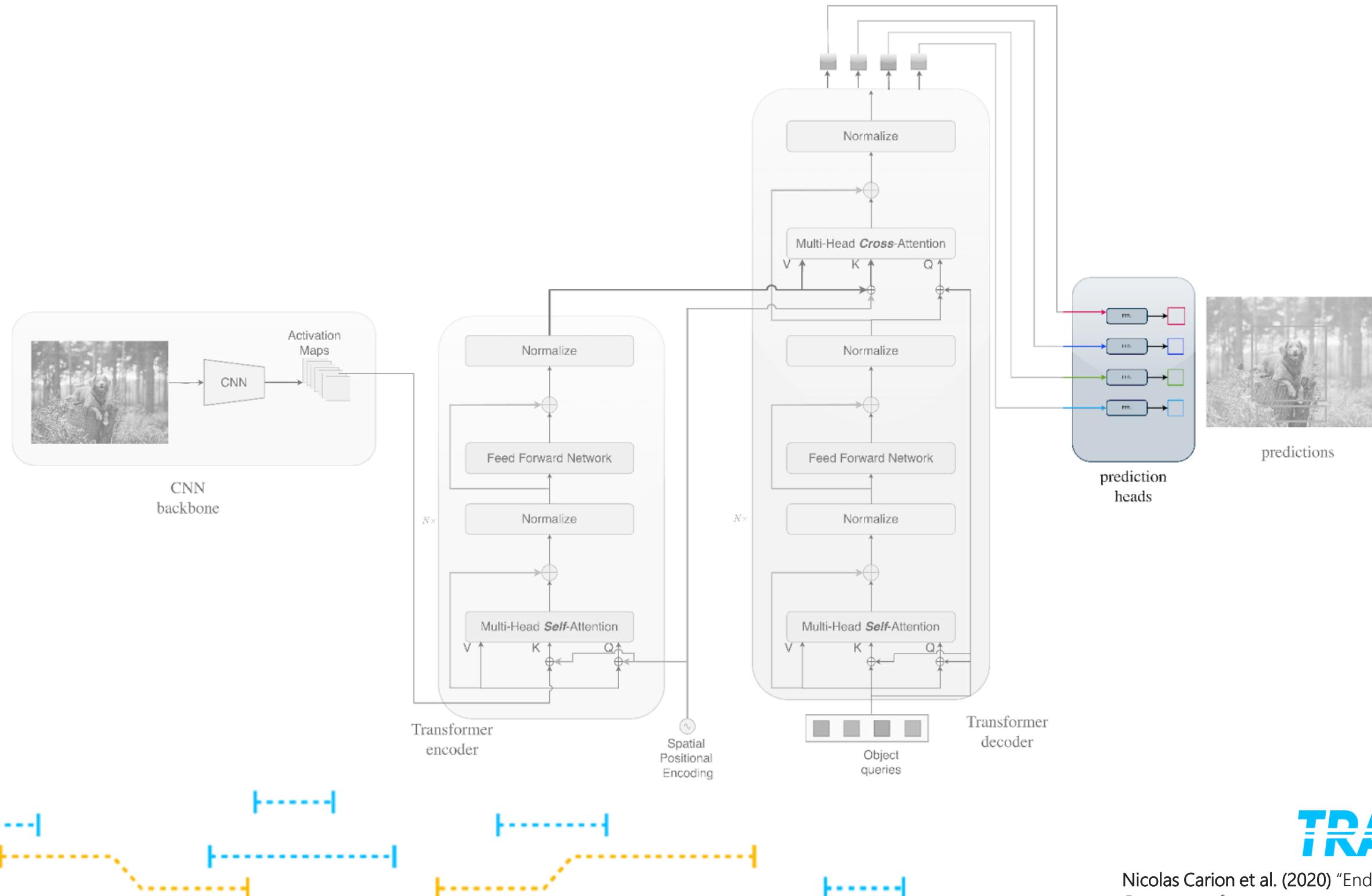


TRANSFORMATEC

Nicolas Carion et al. (2020) "End-to-End Object Detection with Transformers". European conference on computer vision. Cham: Springer International Publishing, 2020.

DETR

(DEtection TRansformer)

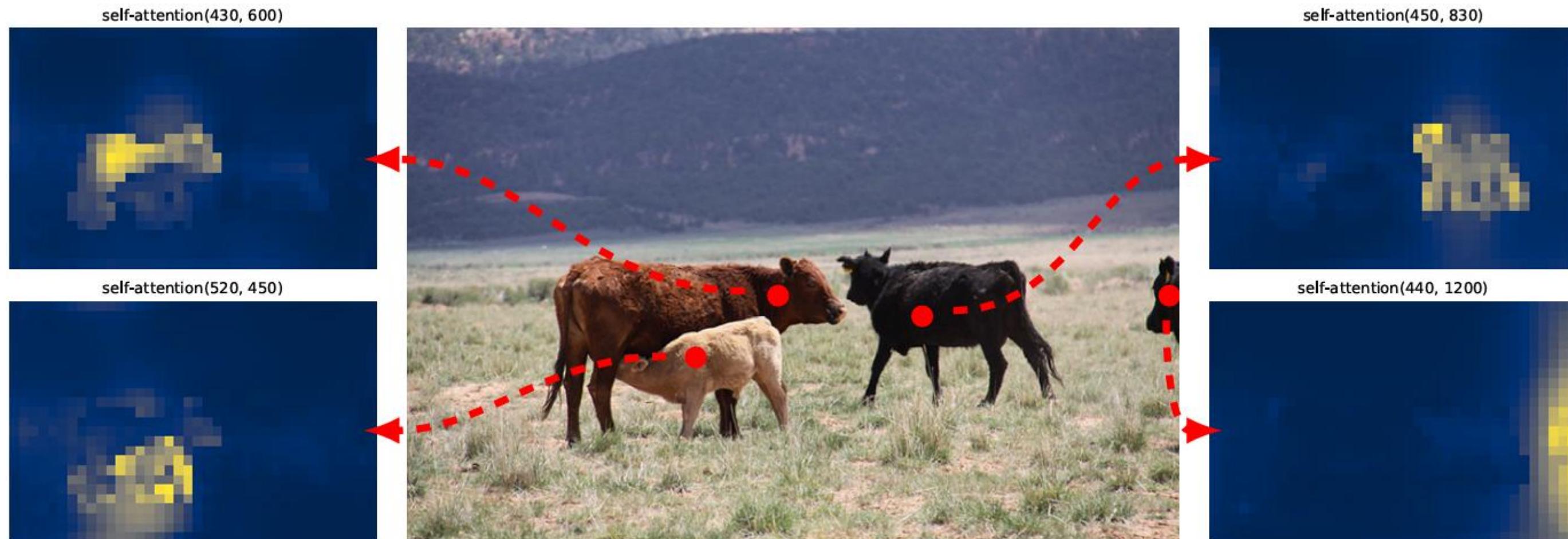


TRANSFORMATEC

Nicolas Carion et al. (2020) "End-to-End Object Detection with Transformers".
European conference on computer vision. Cham: Springer International Publishing, 2020.

DETR

(DEtection TRansformer)

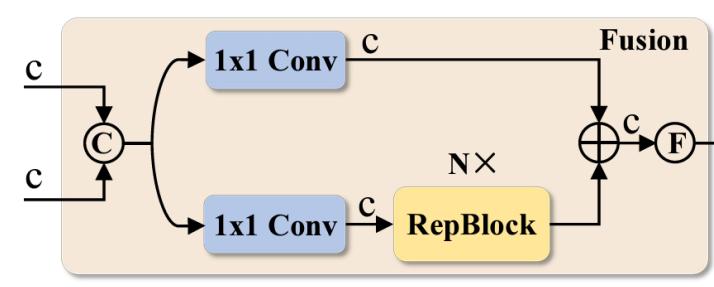
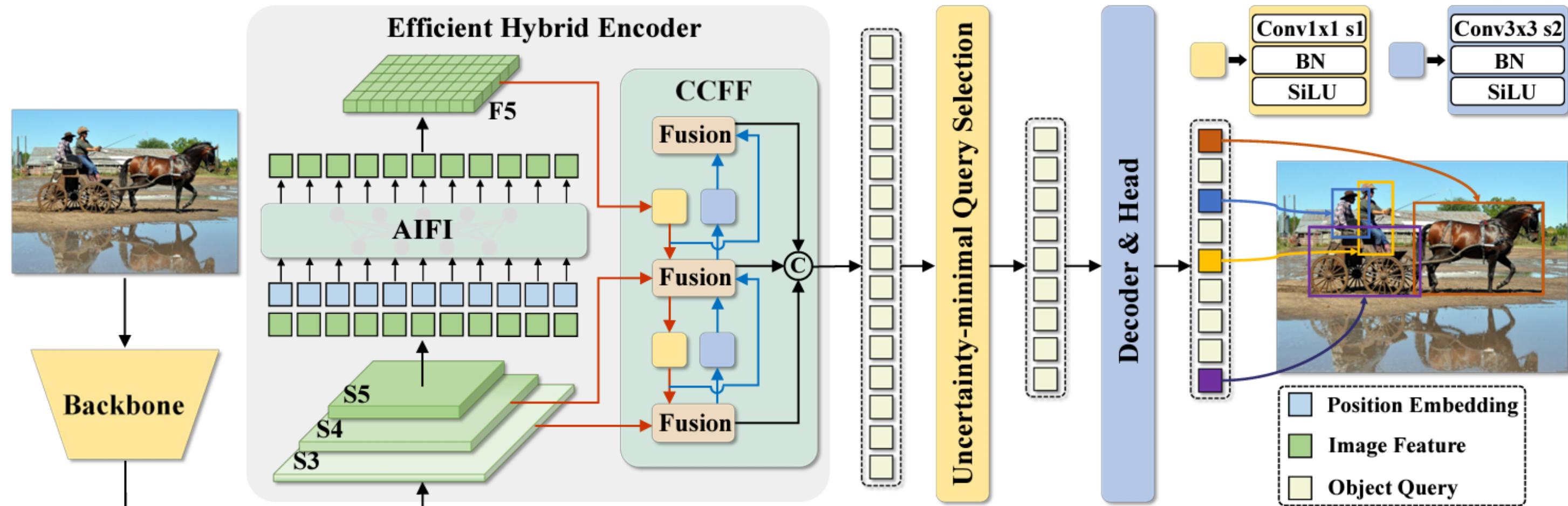


TRANSFORMATEC

Nicolas Carion et al. (2020) "End-to-End Object Detection with Transformers".
European conference on computer vision. Cham: Springer International Publishing, 2020.

RT-DETR

(Real-Time DEtection TRansformer)



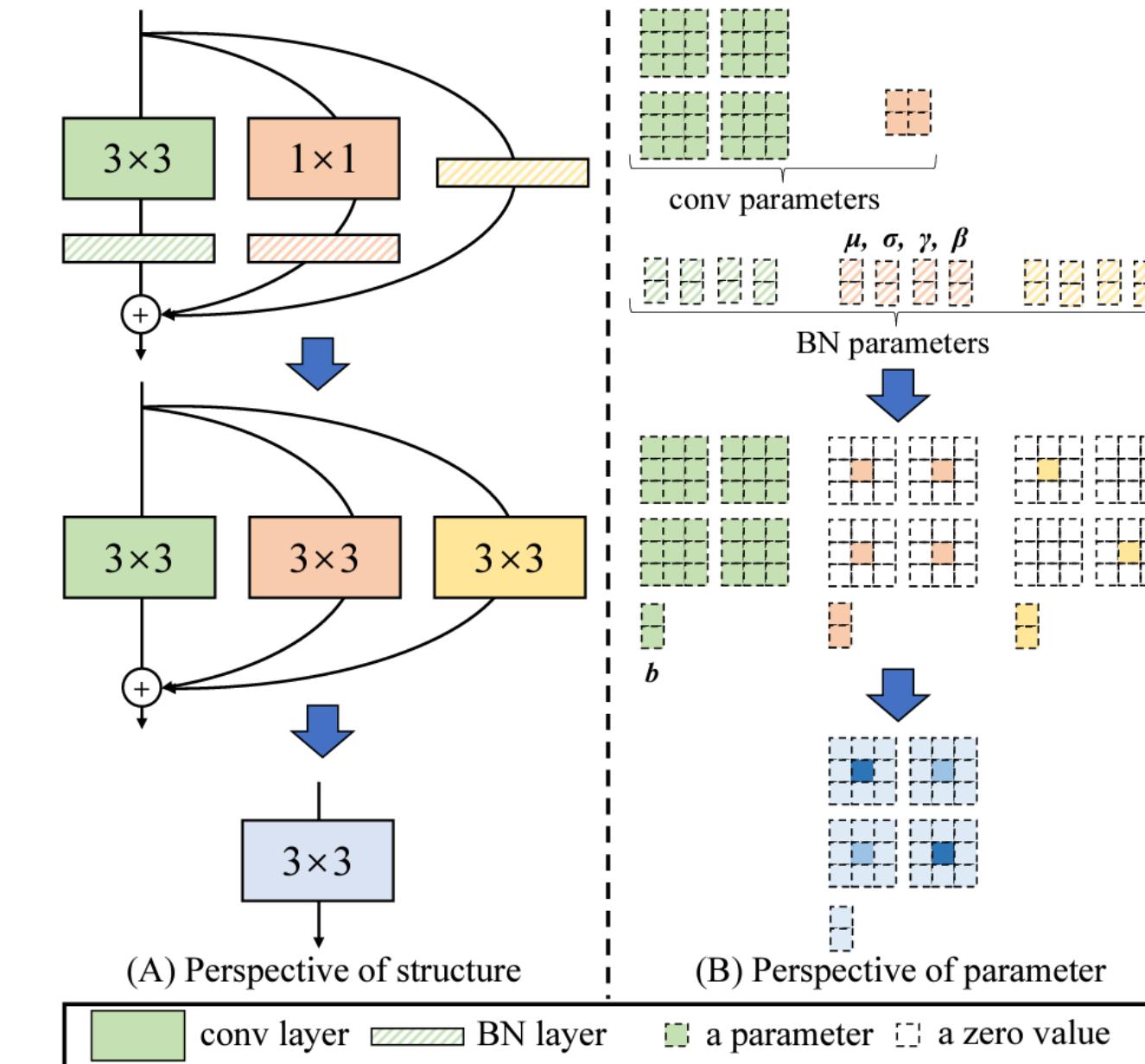
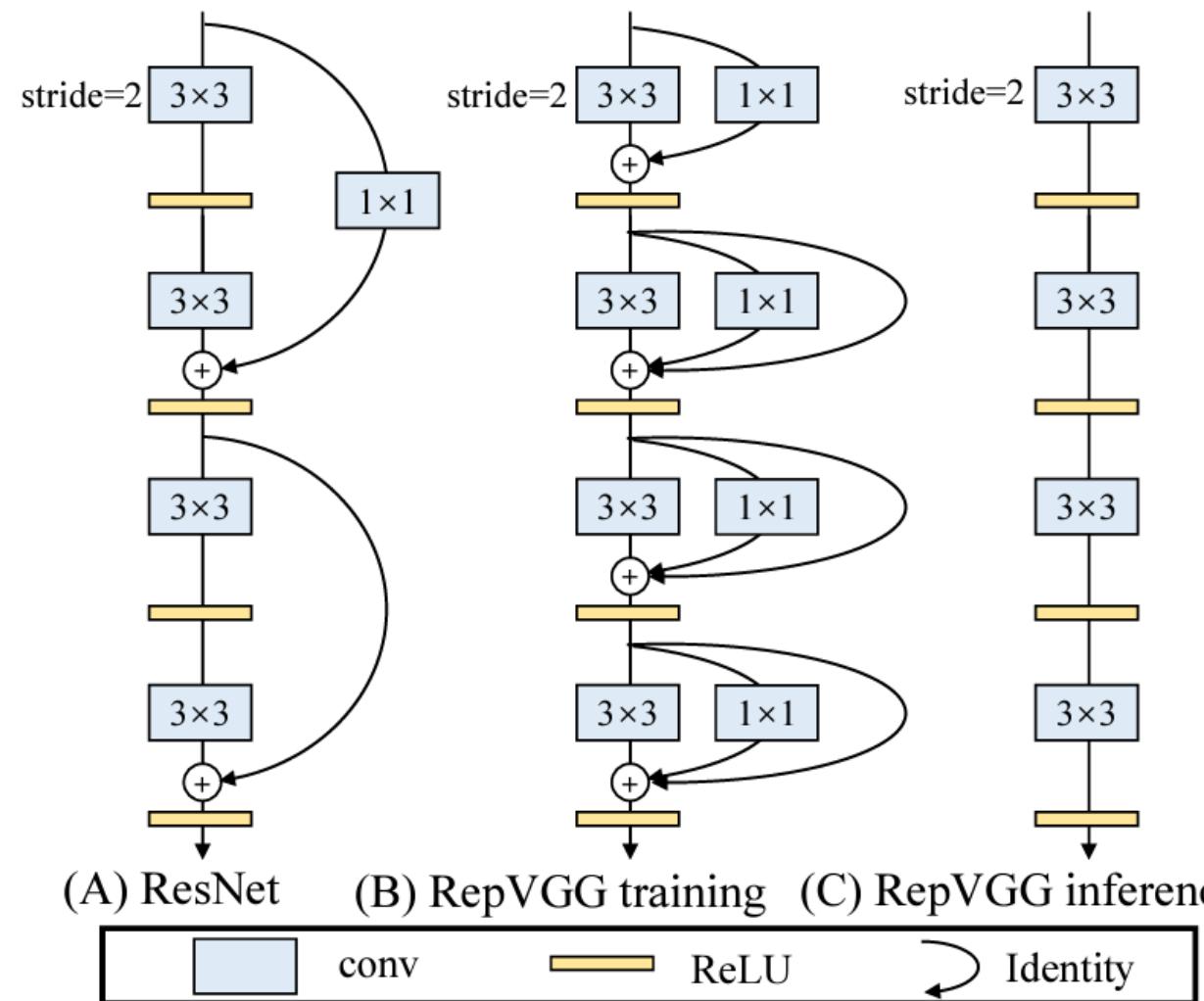
(C) Concatenate (⊕) Element-wise add (F) Flatten

AIFI: Attention-based Intra-scale Feature Interaction.

CCFF: CNN-based Cross-scale Feature Fusion

TRANSFORMATEC

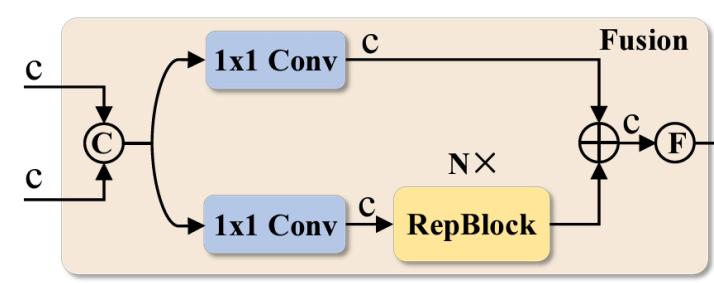
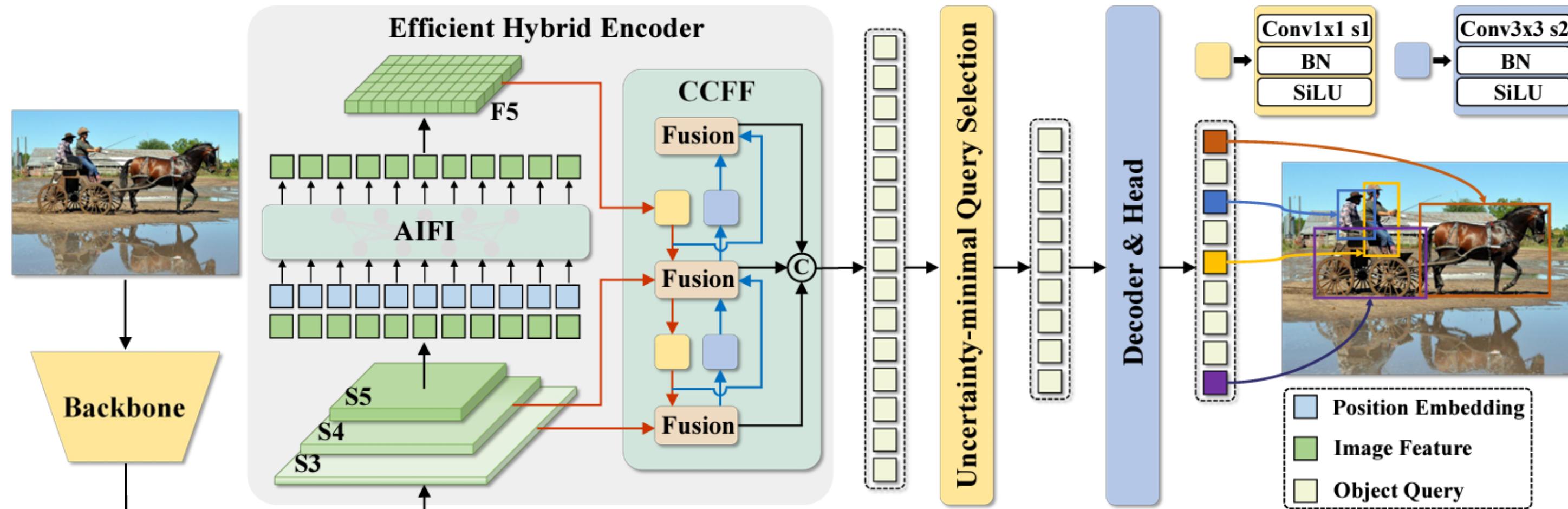
RepVGG



TRANSFORMATEC

RT-DETR

(Real-Time DEtection TRansformer)



(C) Concatenate (⊕) Element-wise add (F) Flatten

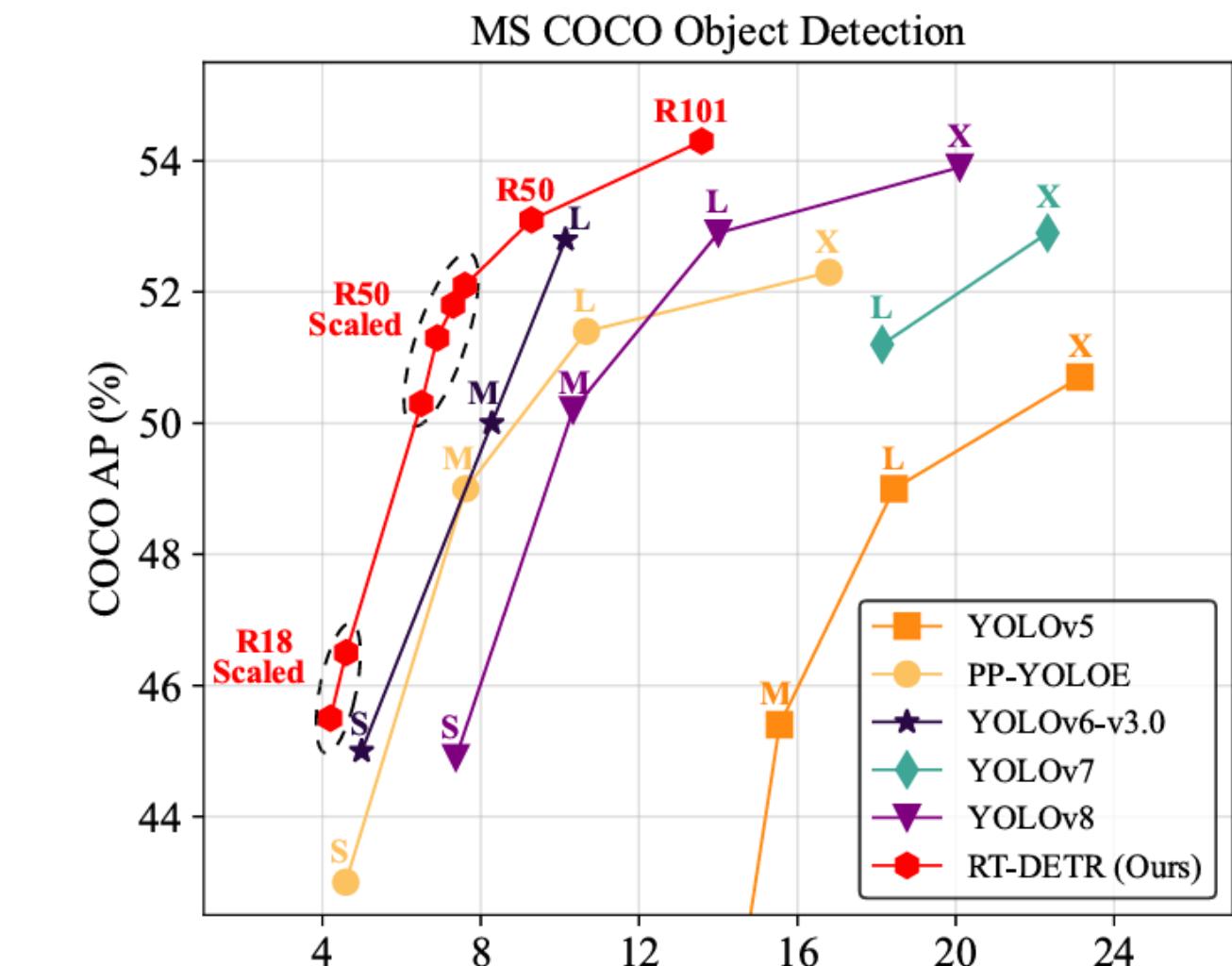
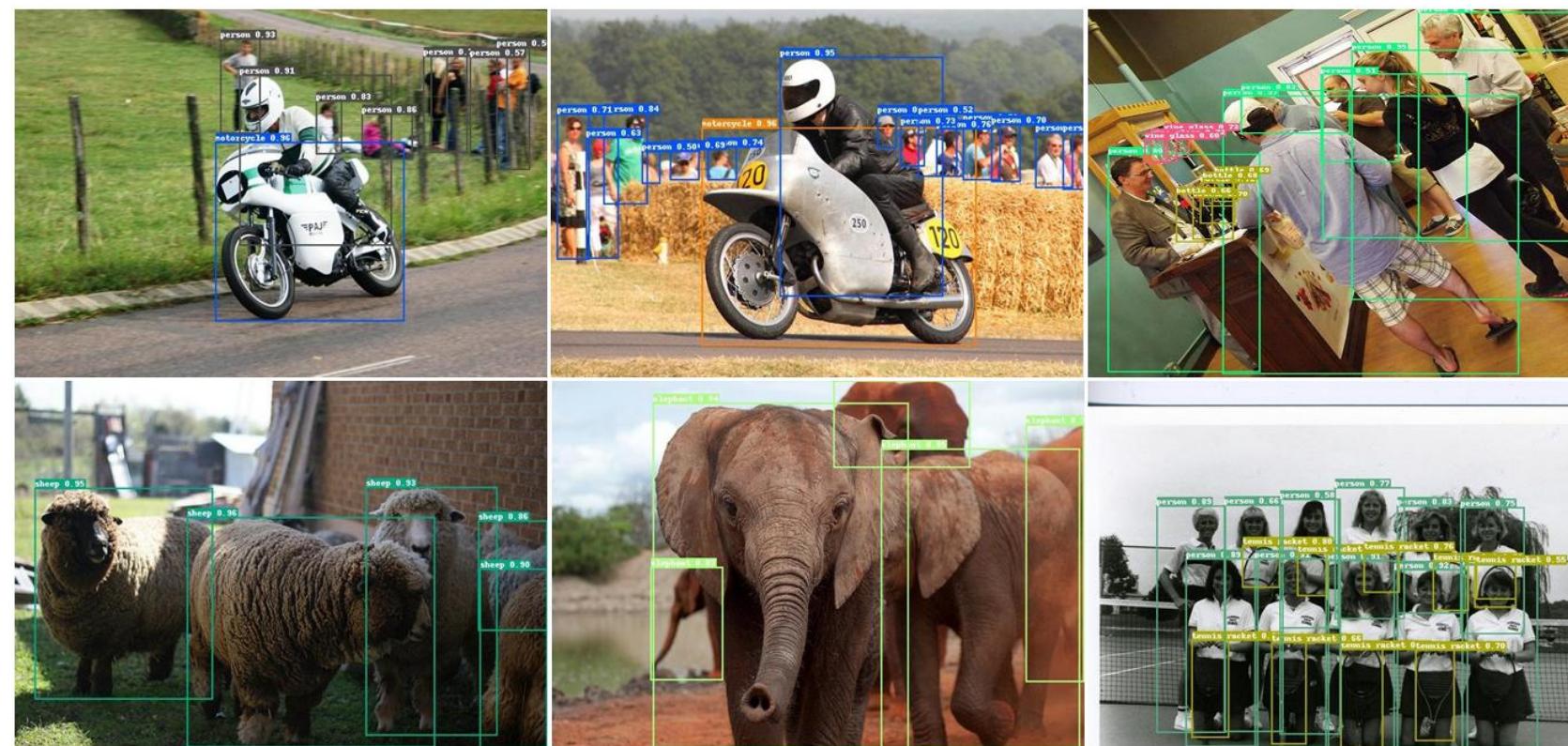
AIFI: Attention-based Intra-scale Feature Interaction.

CCFF: CNN-based Cross-scale Feature Fusion

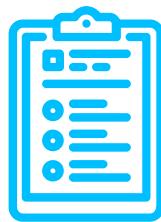
TRANSFORMATEC

RT-DETR

(Real-Time DEtection TRansformer)



5.



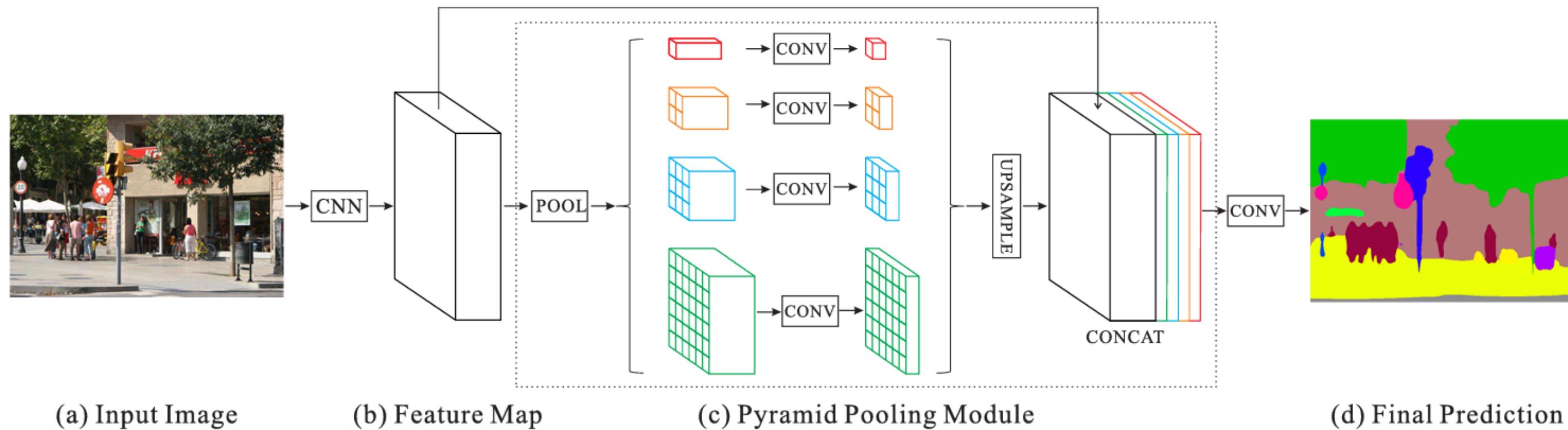
Deep*Lab*

TRANSFORMATEC

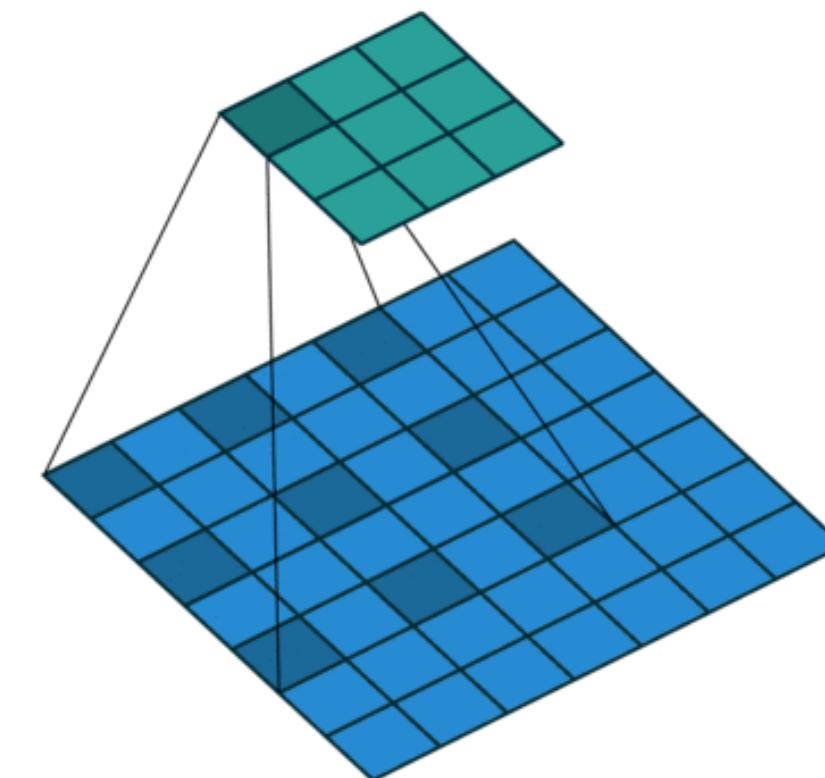
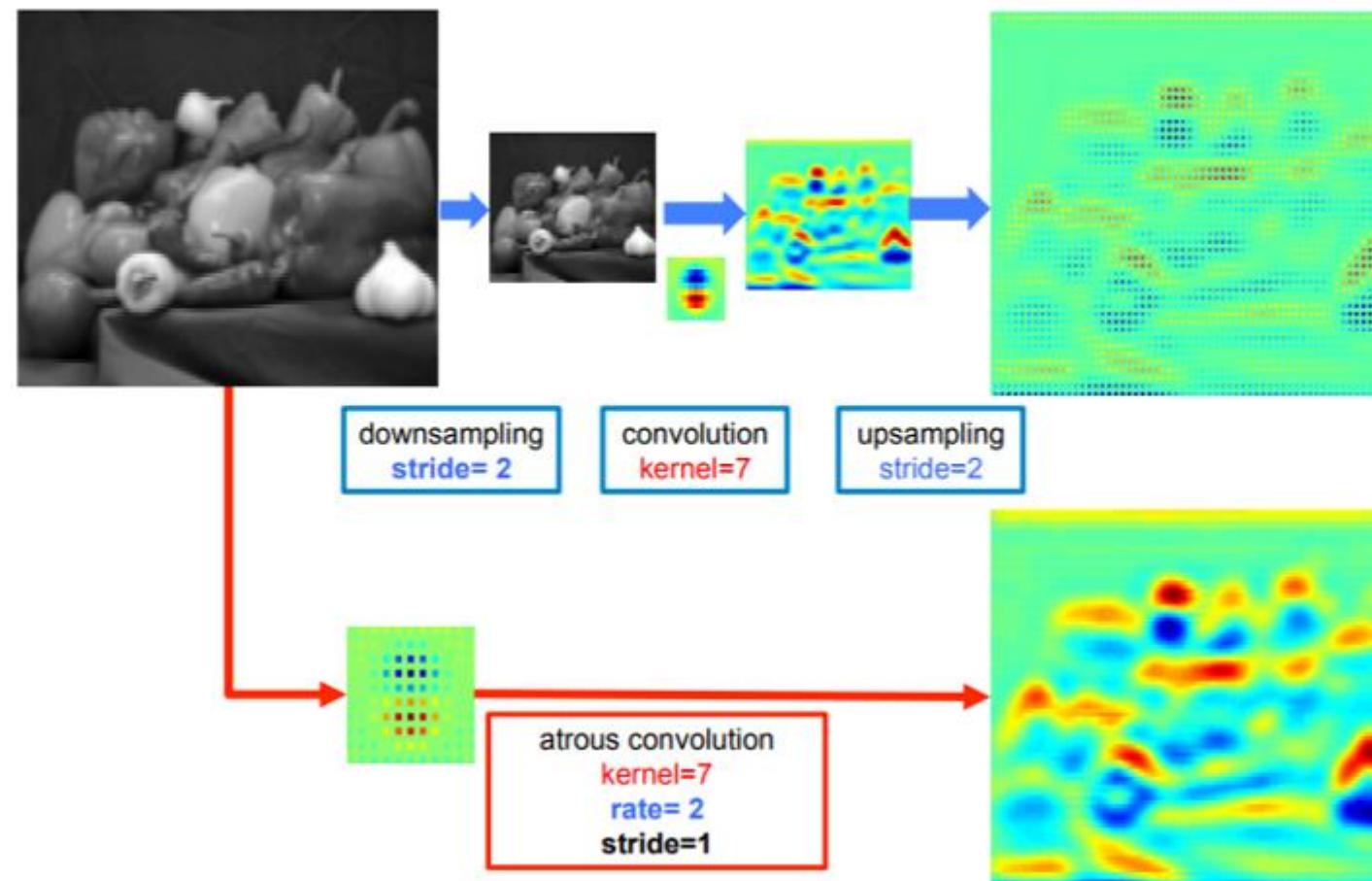
> Reinventa el mundo <



Pyramid Scene Parsing Network (*PSPNet*)



DeepLab

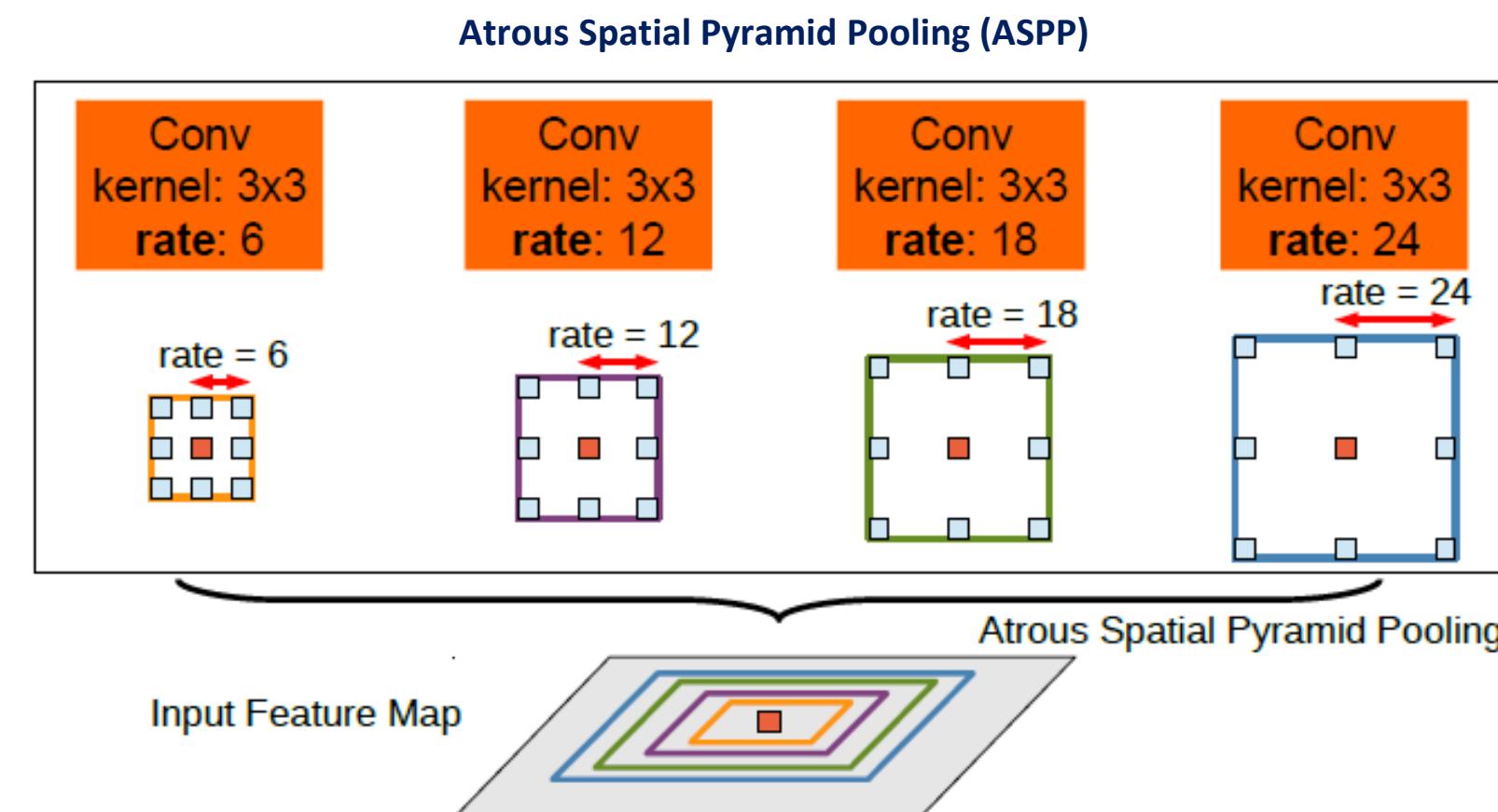


Atrous convolution

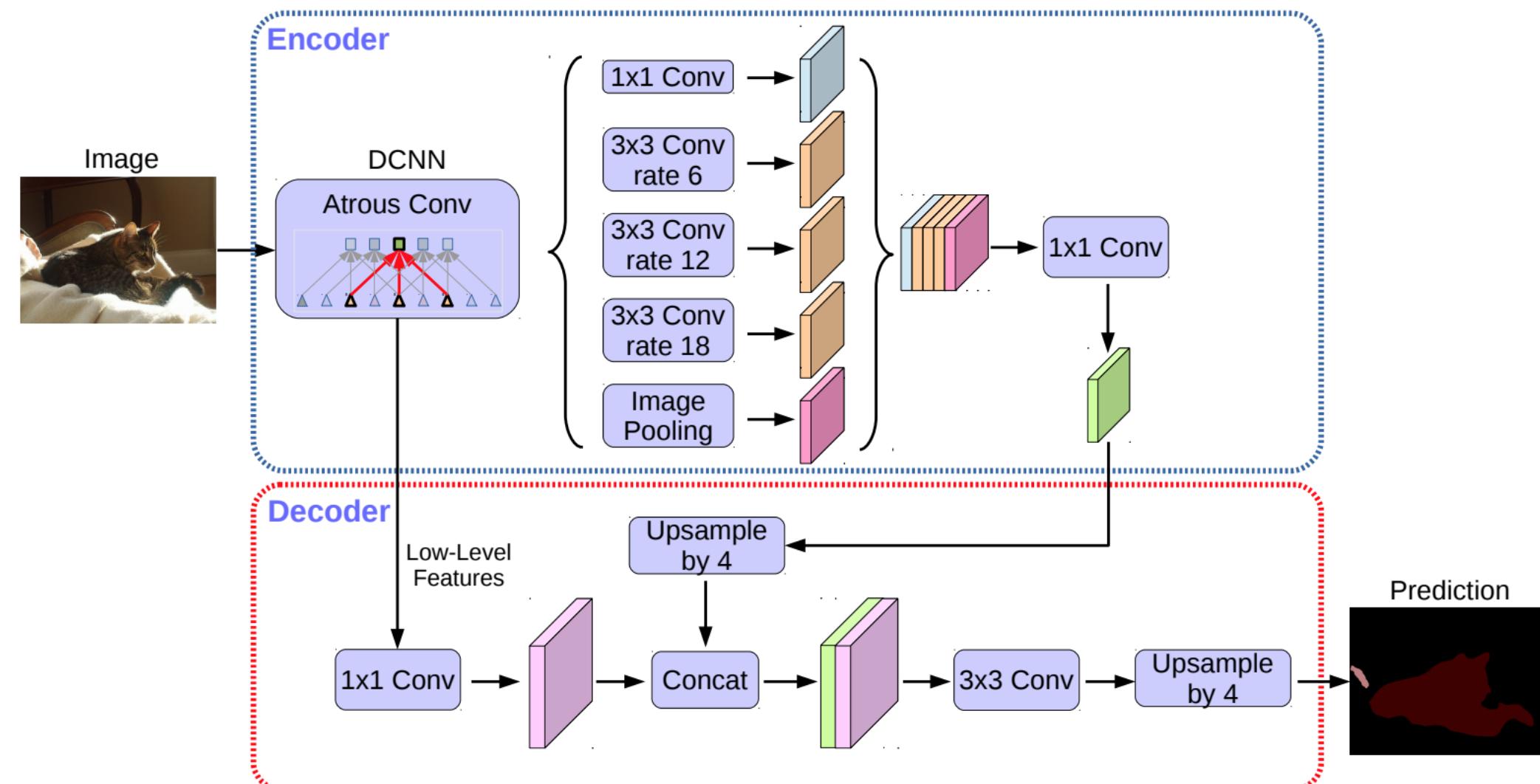


Liang-Chieh Chen et al. (2017) "DeepLab: Semantic Image Segmentation with Deep Convolutional Nets, Atrous Convolution, and Fully Connected CRFs". IEEE transactions on pattern analysis and machine intelligence, 2017, vol. 40, no 4, p. 834-848.

DeepLab



DeepLabV3+



DeepLabV3+



GRACIAS

Victor Flores Benites