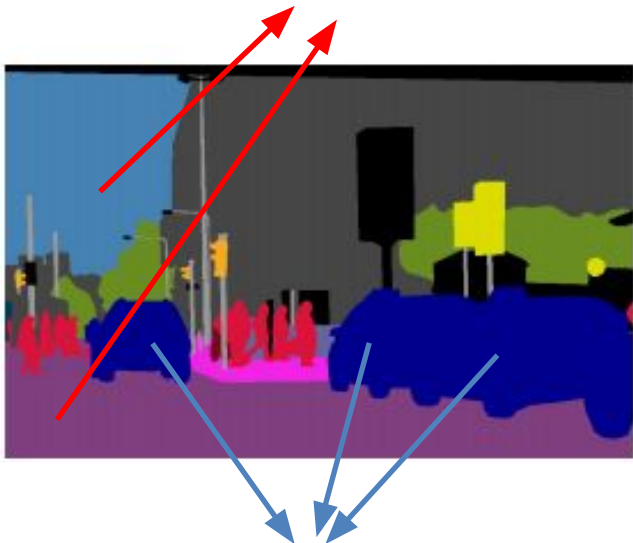


# Instance segmentation

# Semantic segmentation

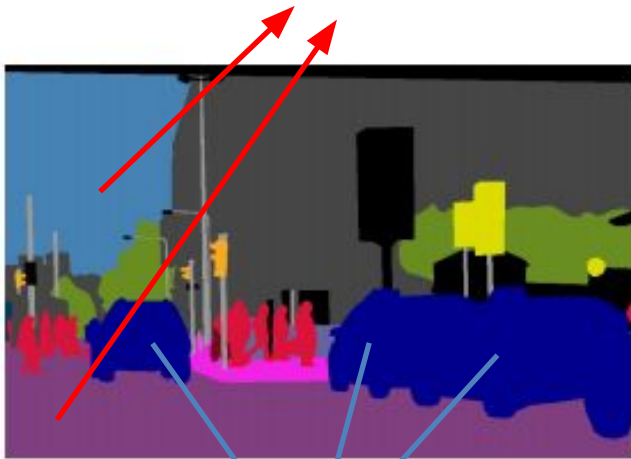
Label every pixel, including the background (sky, grass, road)



Do not differentiate between the pixels coming from instances of the same class

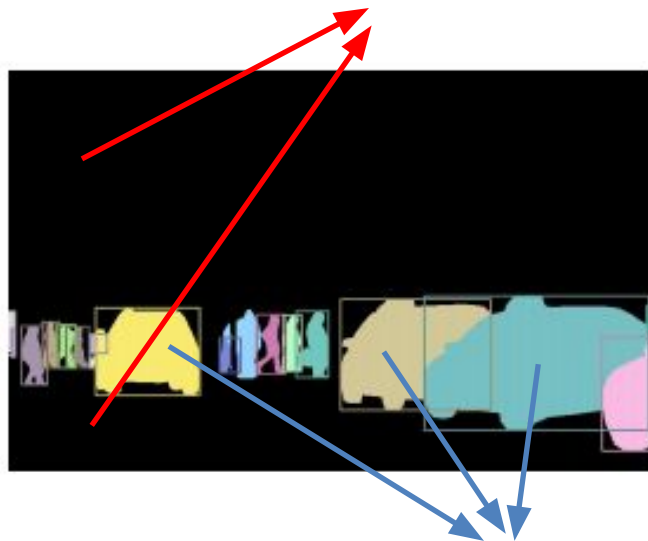
# Instance segmentation

Label every pixel, including the background (sky, grass, road)



Do not differentiate between the pixels coming from instances of the same class

Do not label pixels coming from uncountable objects (sky, grass, road)



Differentiate between the pixels coming from instances of the same class

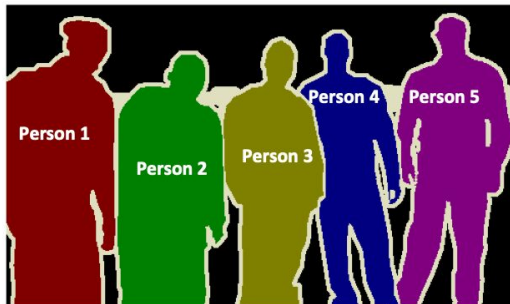
# Instance segmentation methods

Proposal-based

1. Proposals



2. Assign a class



FCN-based

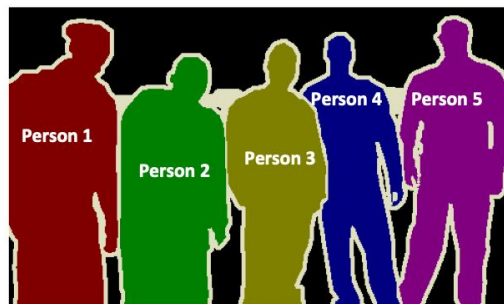
1. Semantic segmentation



vs.



2. Find instances



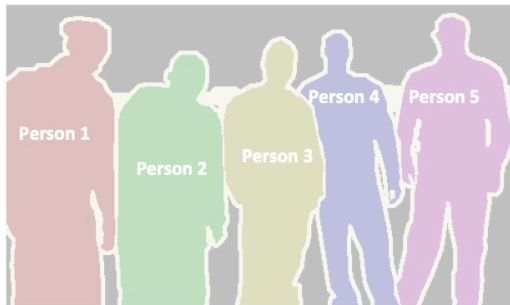
# Instance segmentation methods

Proposal-based

1. Proposals



2. Assign a class

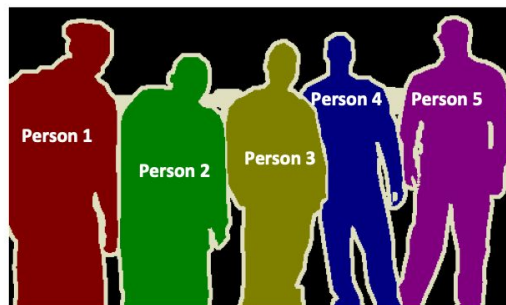


FCN-based

1. Semantic segmentation

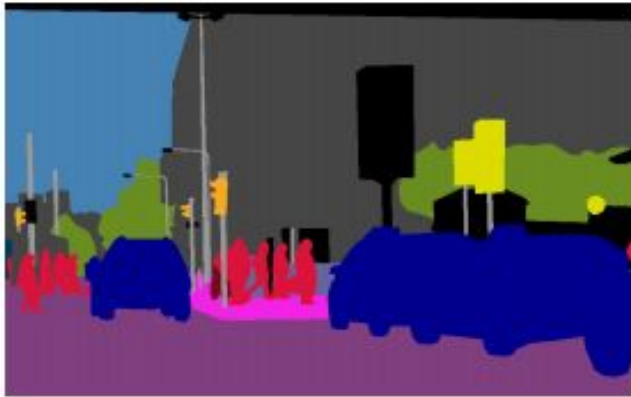


2. Find instances



vs.

# FCN-based methods

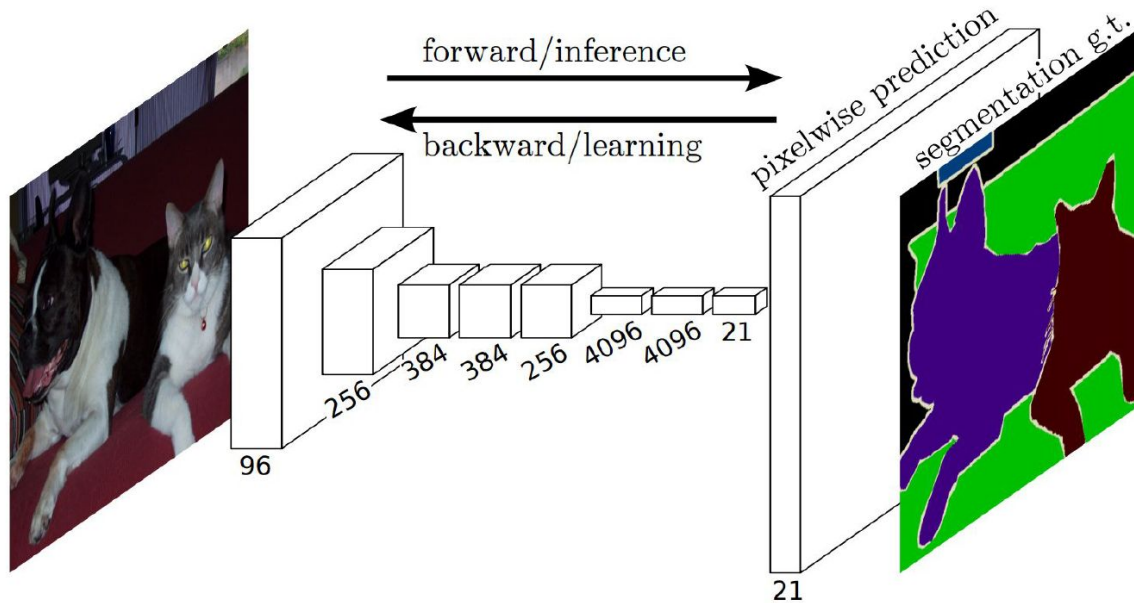


A semantic map

We already know how to obtain this!

# Why FCN-based?

- Fully Convolutional Networks for Semantic Segmentation



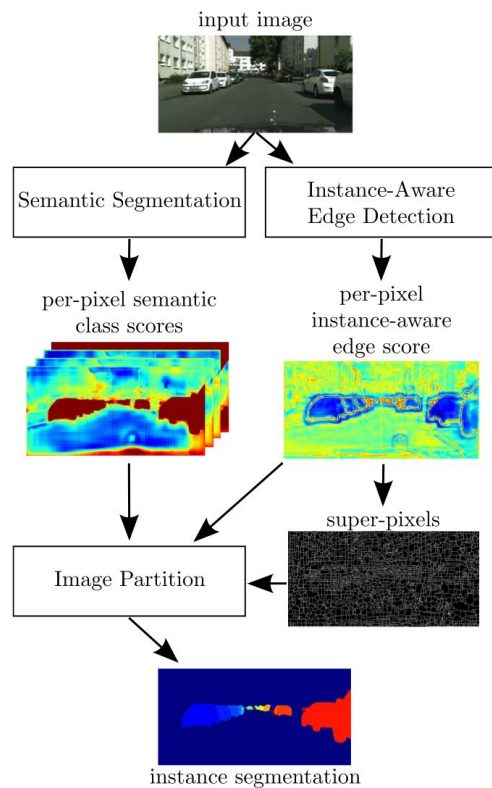
Long, Shelhamer, Darrell - Fully Convolutional Networks for Semantic Segmentation, CVPR 2015, PAMI 2016

# FCN-based methods

- X. Liang et al. "Proposal-free Network for Instance-level Object Segmentation". Arxiv 2015
- A. Kirillov et al. „InstanceCut: from Edges to Instances with MultiCut". CVPR 2017
- M. Bai and R. Urtasun "Deep Watershed Transform for Instance Segmentation ". CVPR 2017



# Instances through clustering



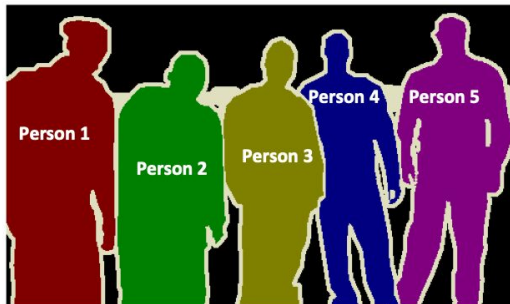
# Instance segmentation methods

Proposal-based

1. Proposals



2. Assign a class

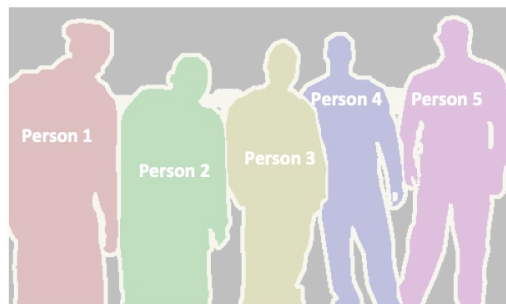


FCN-based

1. Semantic segmentation



2. Find instances



vs.

# Proposal-based methods

Bounding boxes.....

We already know how to obtain those!

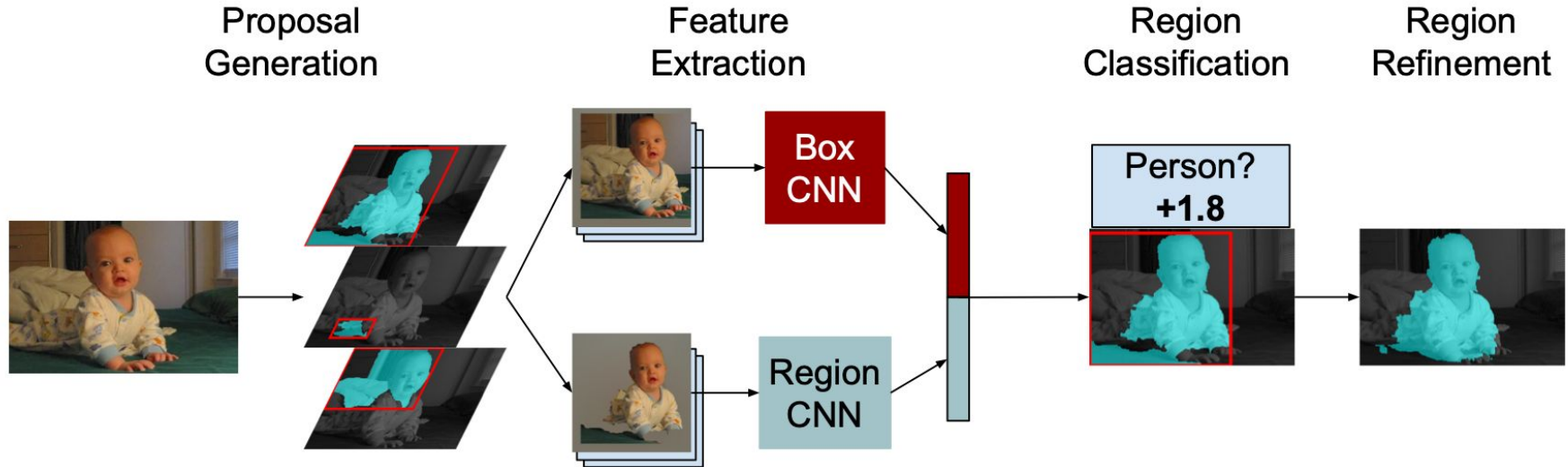


# Proposal-based methods

- B. Hariharan et al. "Simultaneous Detection and Segmentation". ECCV 2014
  - Follow-up work: B. Hariharan et al. "Hypercolumns for Object Segmentation and Fine-grained Localization ". CVPR 2015
- Dai et al. „Instance-aware Semantic Segmentation via Multi-task Network Cascades". CVPR 2016
  - Previous work: Dai et al. "Convolutional Feature Masking for Joint Object and Stuff Segmentation". CVPR 2015

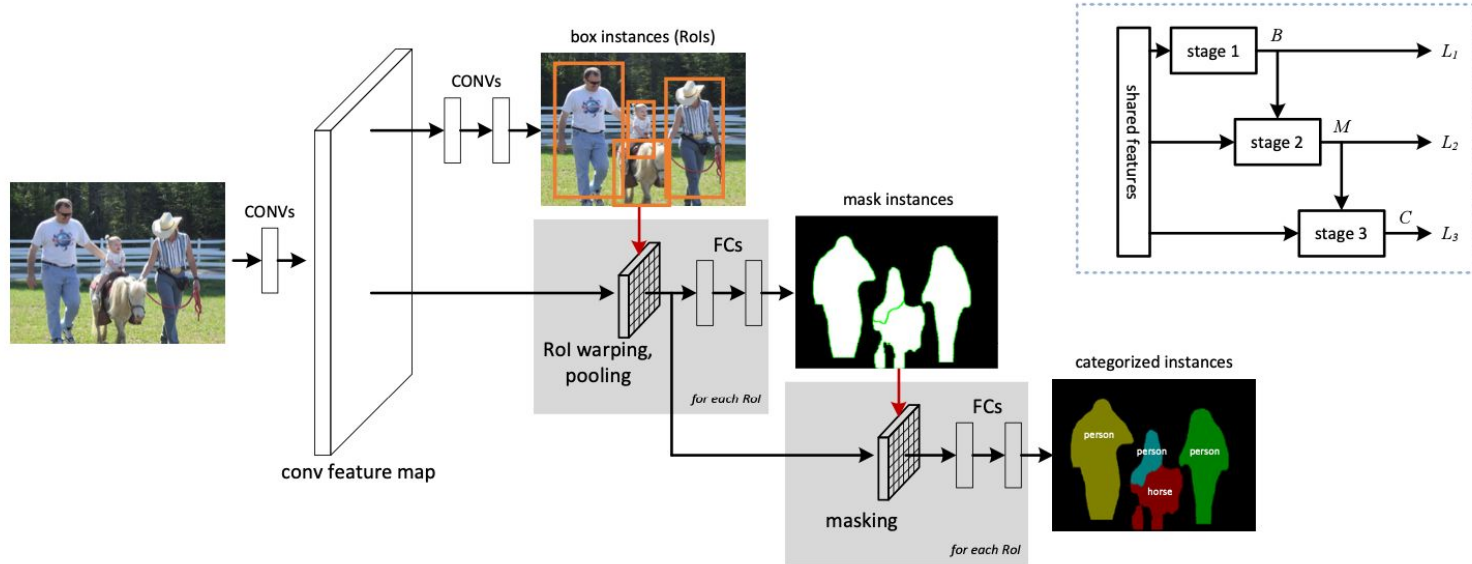
# SDS

- SDS: Simultaneous Detection and Segmentation



# MNC

- MNC: Multi-task network cascades



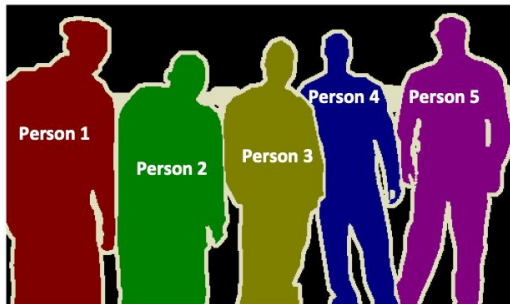
# IS: the best of both worlds

Proposal-based

1. Proposals



2. Assign a class

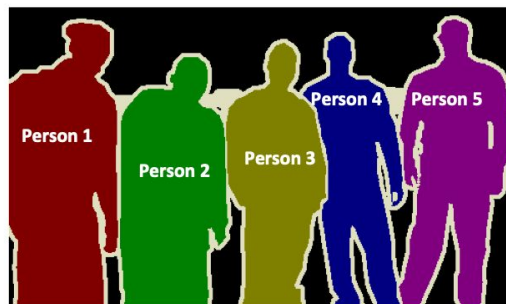


FCN-based

1. Semantic segmentation



2. Find instances

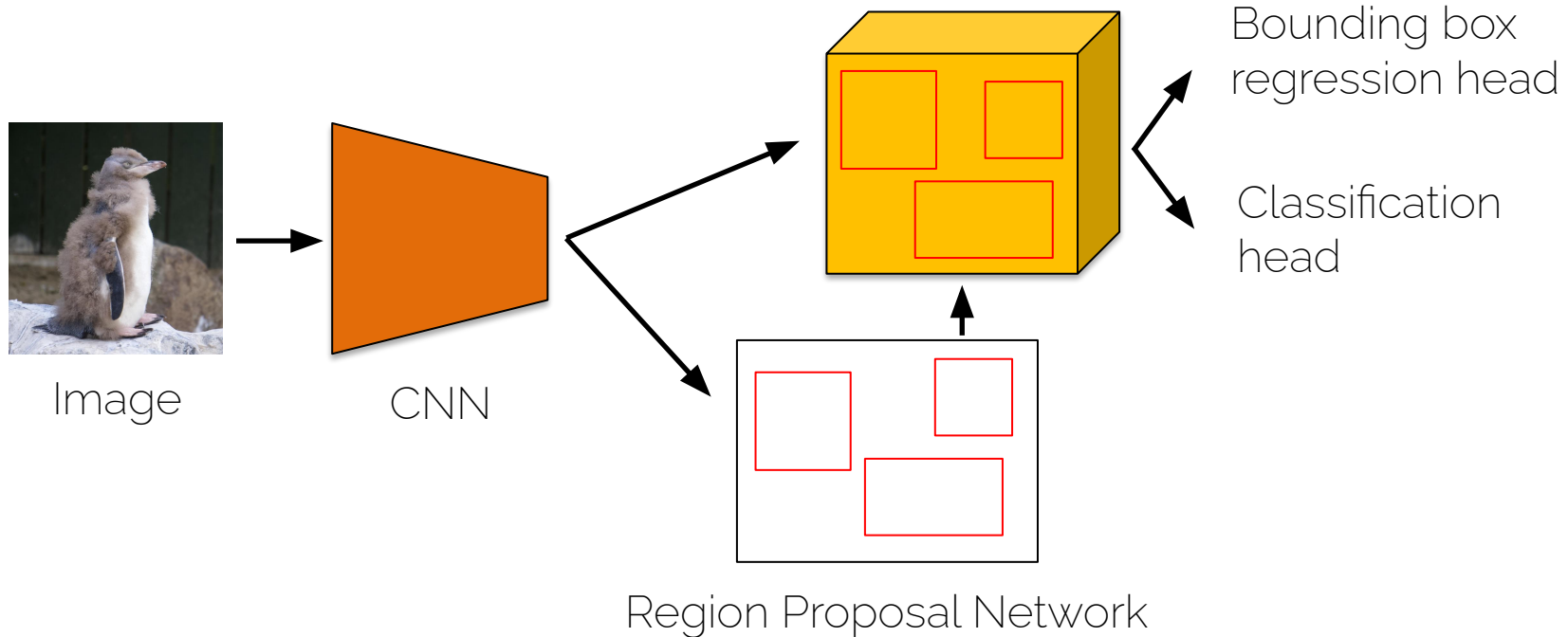


# Mask R-CNN



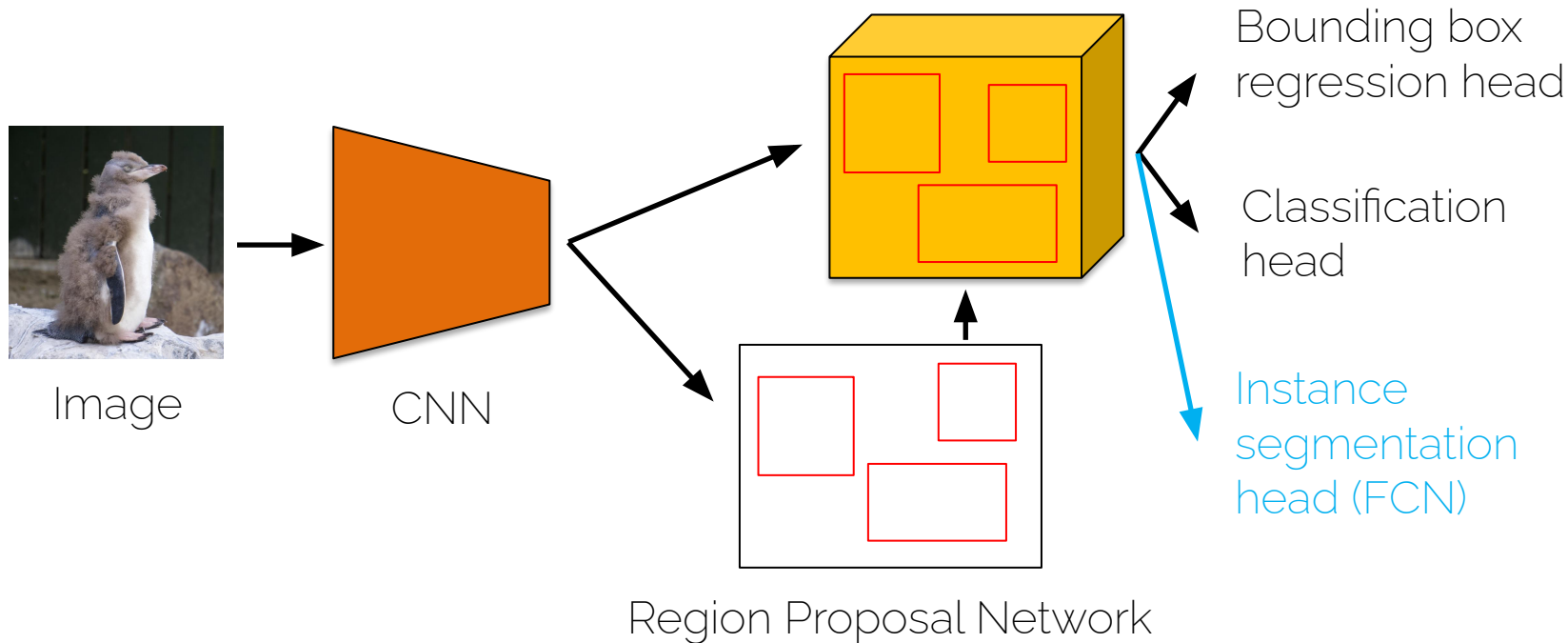
# What is Mask-RCNN?

- Starting from the Faster R-CNN architecture



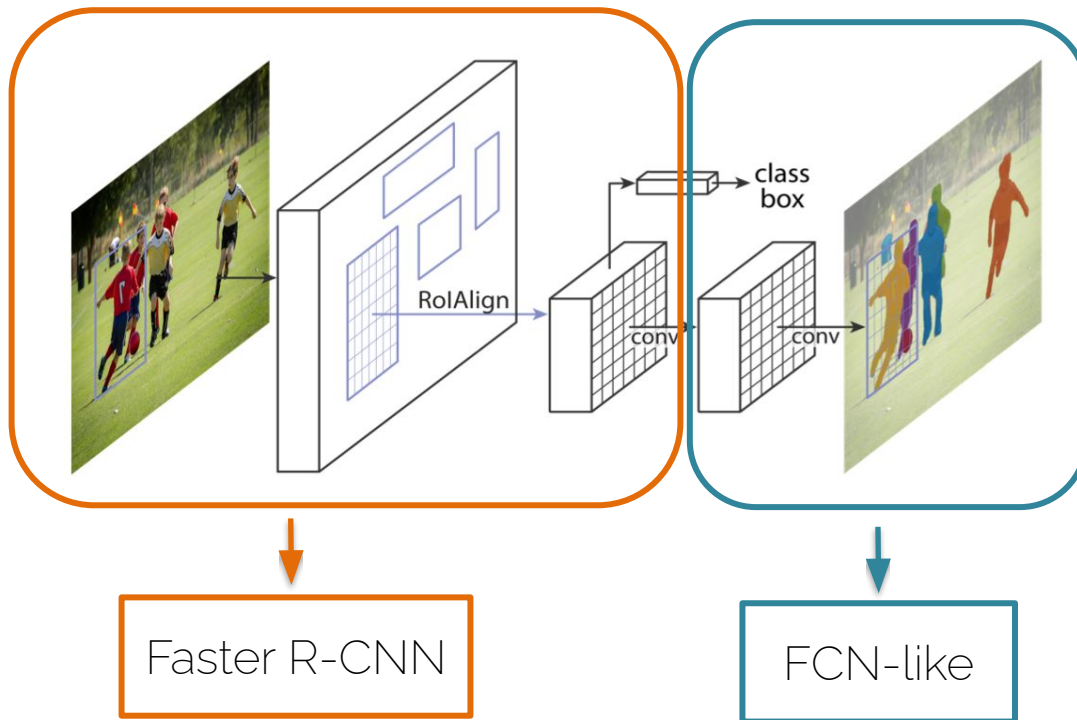
# What is Mask-RCNN?

- Faster R-CNN + FCN for segmentation



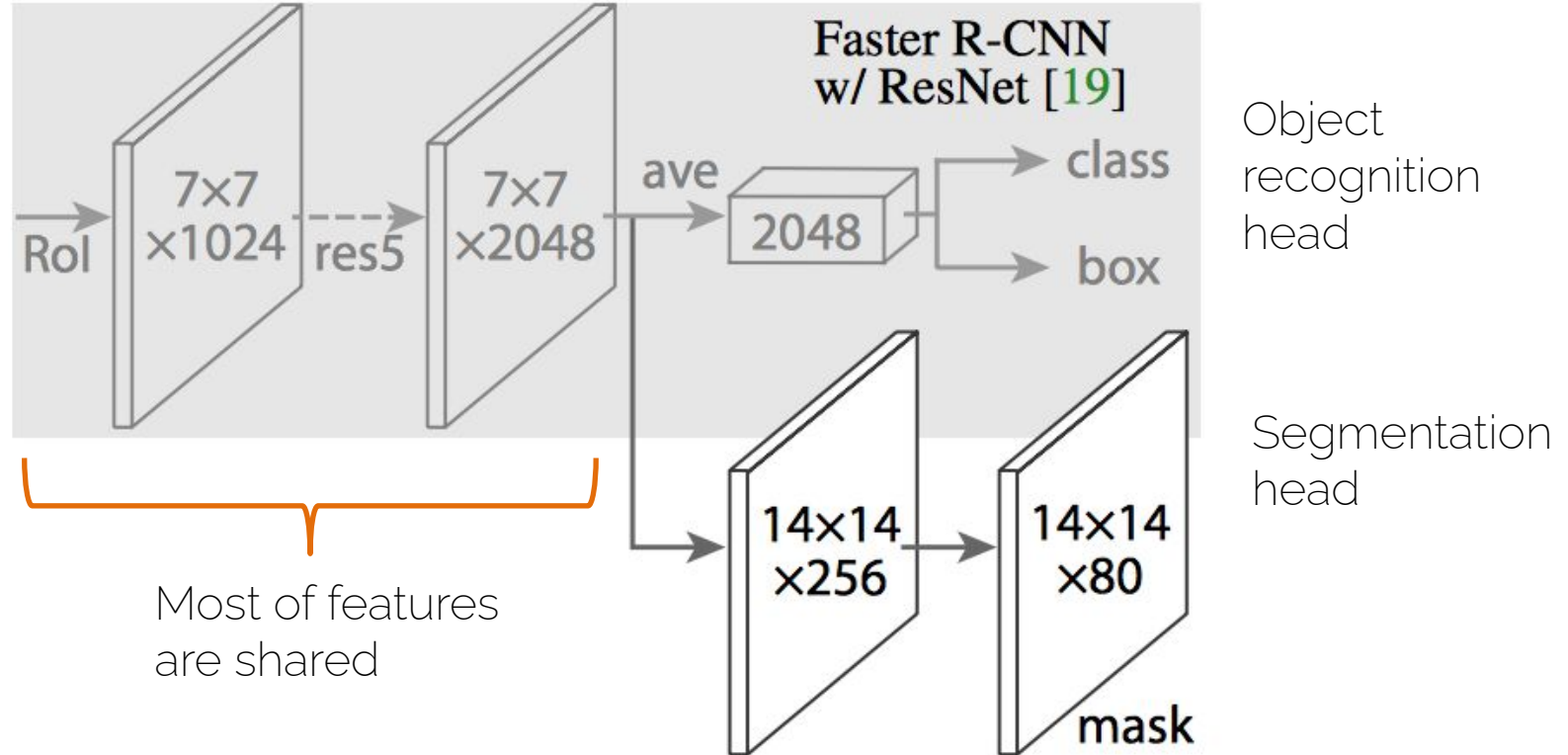
# What is Mask-RCNN?

- Faster R-CNN + FCN for segmentation

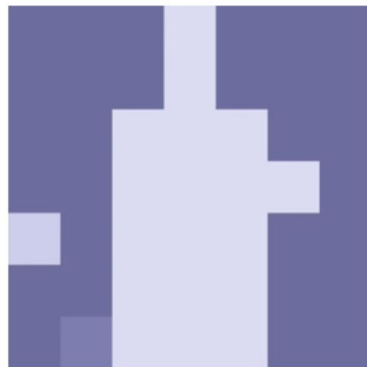


Mask loss =  
binary cross  
entropy per pixel  
for the  $k$  semantic  
classes

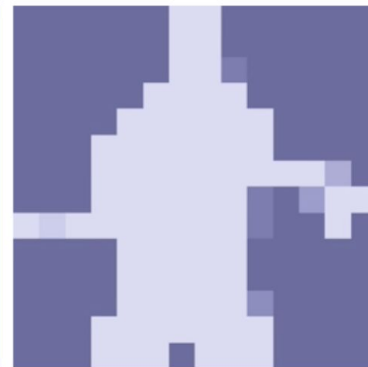
# Mask R-CNN



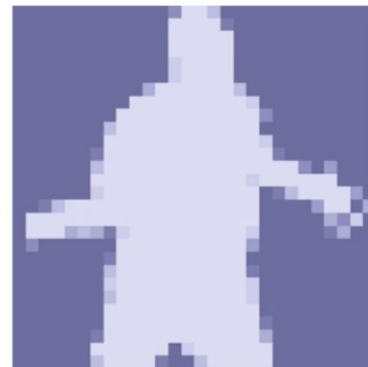
# Resolution vs. computation



7x7



14x14



28x28

output resolution is a tradeoff between computational cost and level of detail



56x56



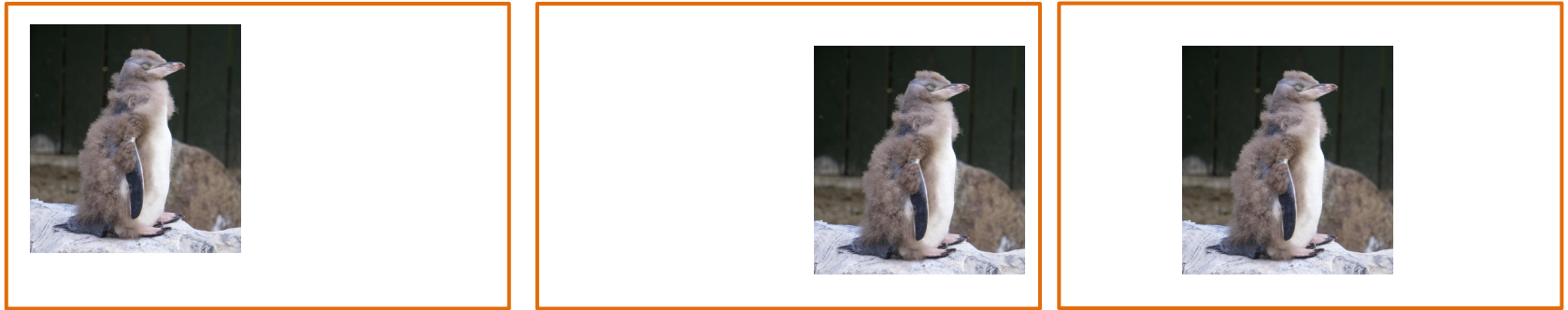
112x112



224x224

# Detection vs. segmentation

- Detection: for object classification, you require **invariant** representations



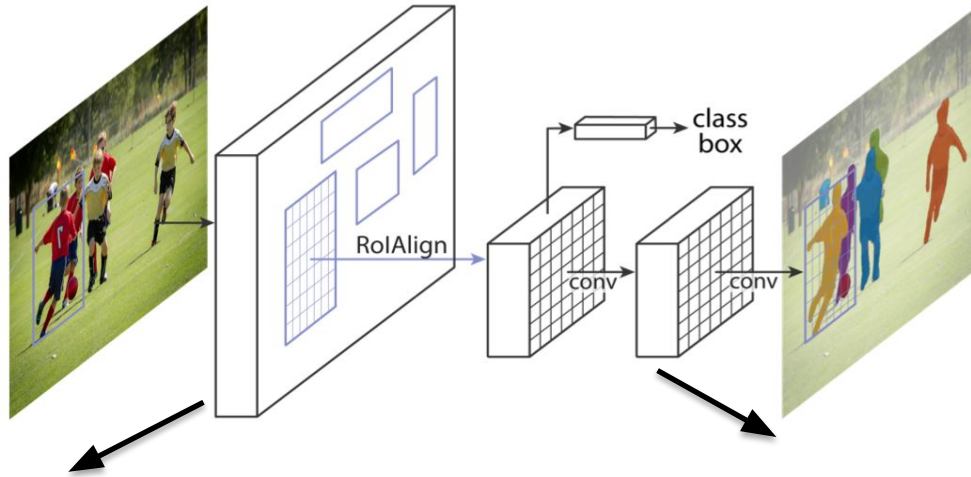
Translation invariance: wherever the penguin is in the image, I still want to have "penguin" as my classification output

# Detection vs. segmentation

- Detection: for object classification, you require **invariant** representations
- Segmentation: you require **equivariant** representations
  - Translated object  $\square$  Translated mask
  - Scaled object  $\square$  scaled mask
  - For semantic segmentation, small objects are less important (less pixels), but for instance segmentation, all objects (no matter the size) are equally important

# Mask-RCNN: operations

- What operations are equivariant?



Features extraction = convolutional layers equivariant

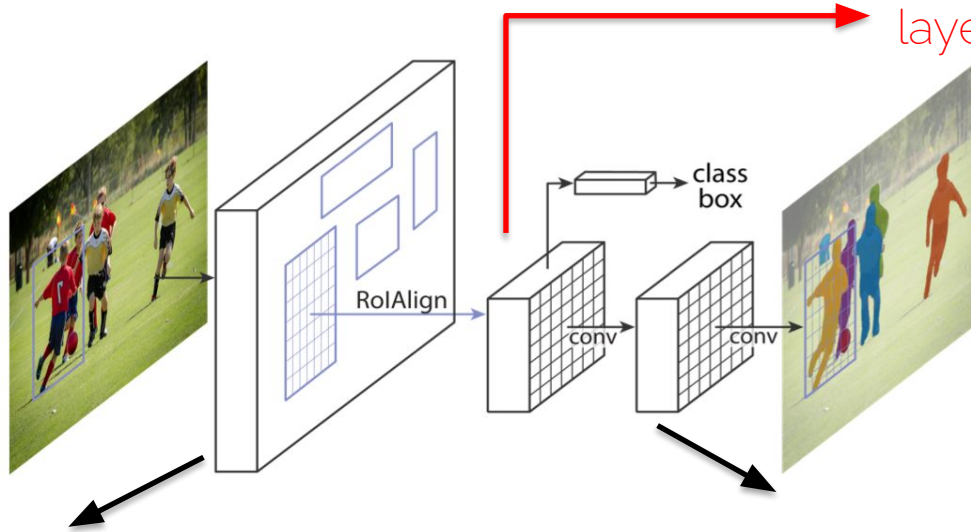
Segmentation head is a fully convolutional network equivariant



# Mask-RCNN: operations

- What operations are equivariant?

Fully connected layers  
and global pooling  
layers give invariance!

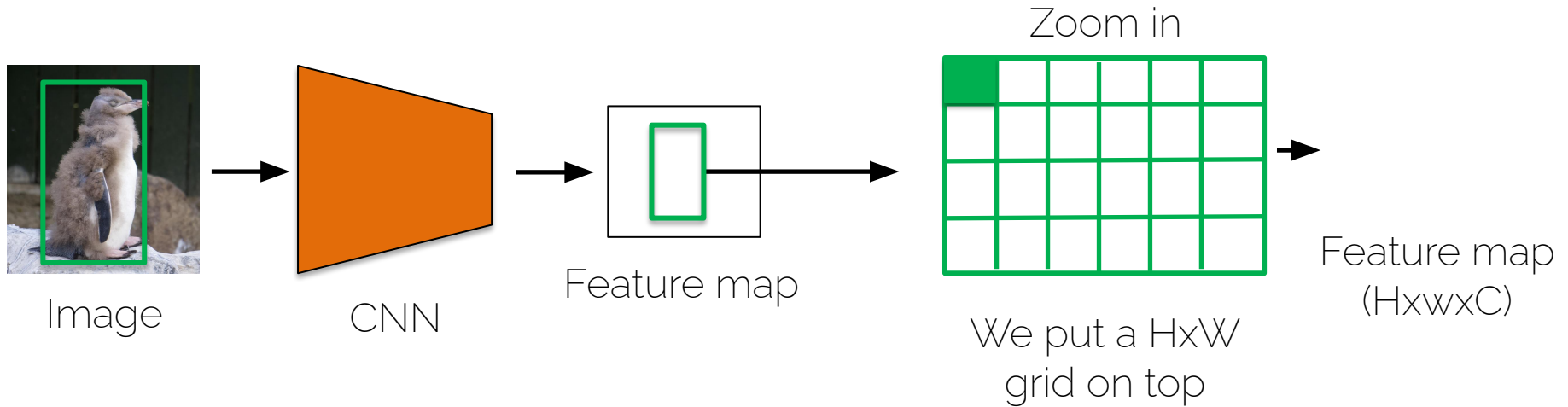


Features extraction = convolutional  
layers equivariant

Segmentation head is a fully  
convolutional network equivariant

# Recall: RoI pooling

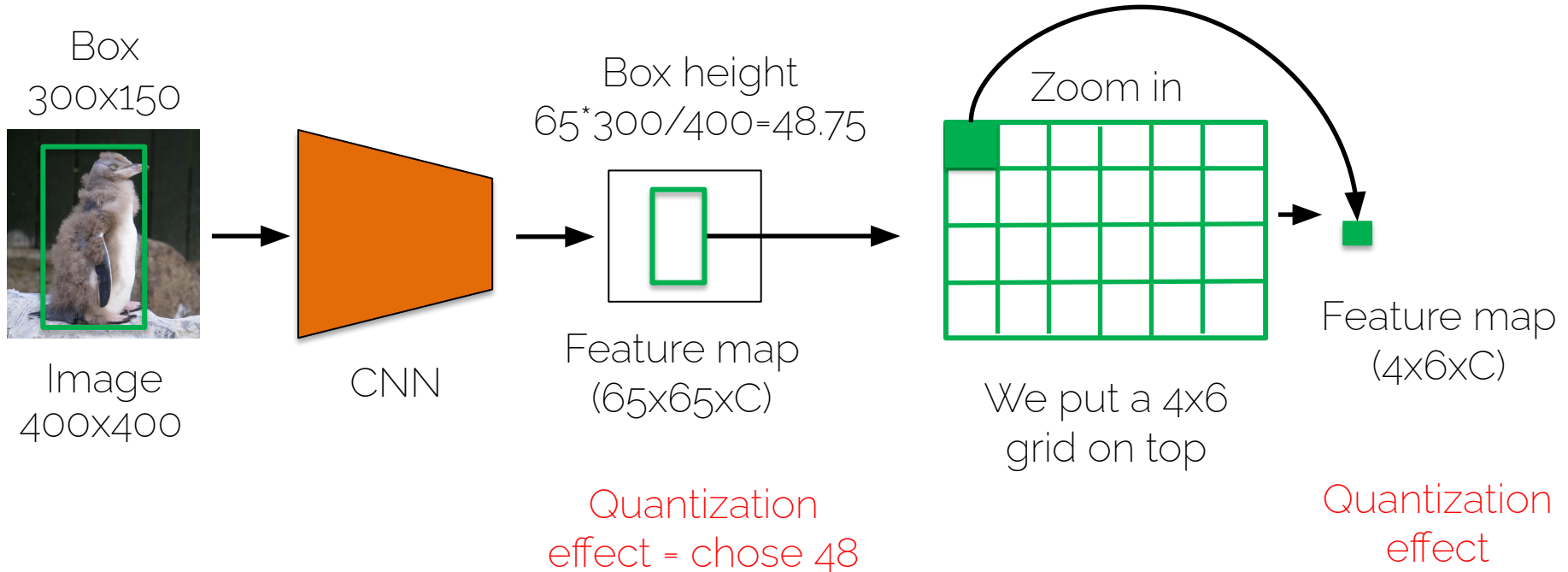
- Region of Interest Pooling: for every proposal



# Recall: RoI pooling

Not suitable to extract pixel-wise precise masks

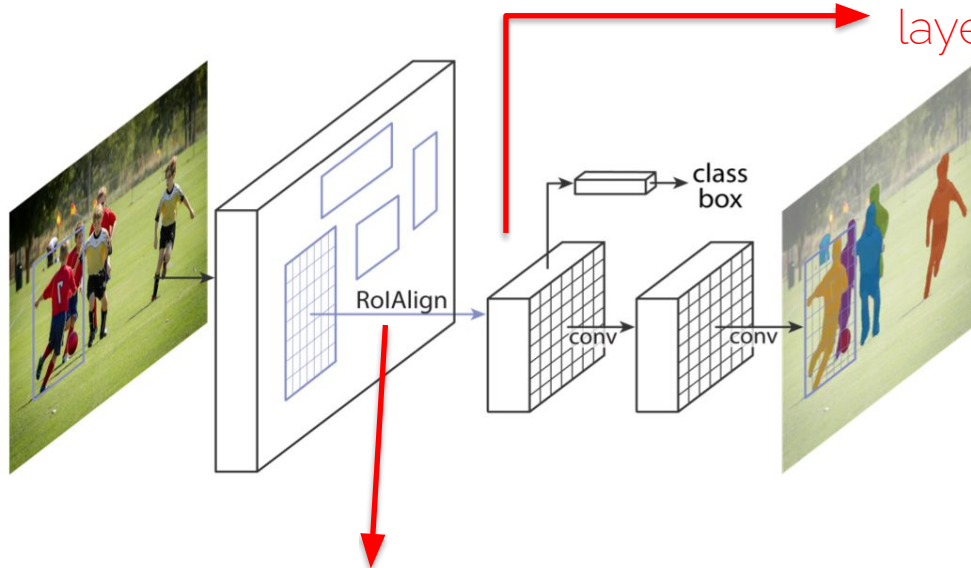
- Let us look at sizes



# Mask-RCNN: operations

- Make all operations equivariant

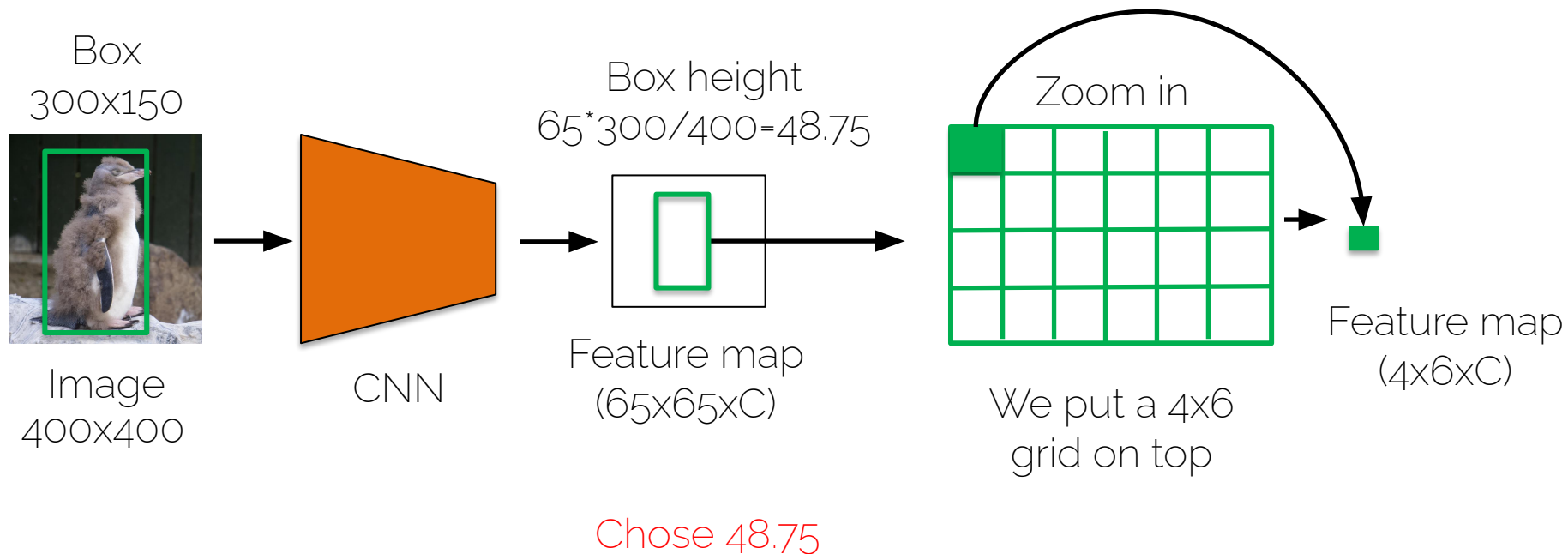
Fully connected layers  
and global pooling  
layers give invariance!



Exchange RoI pooling by an equivariant operation = RoI Align

# RoIAlign

- Erase quantization effects



# ROIAlign

To obtain the value  
use bilinear  
interpolation

Feature map

Each unit is  
sampled 4 times

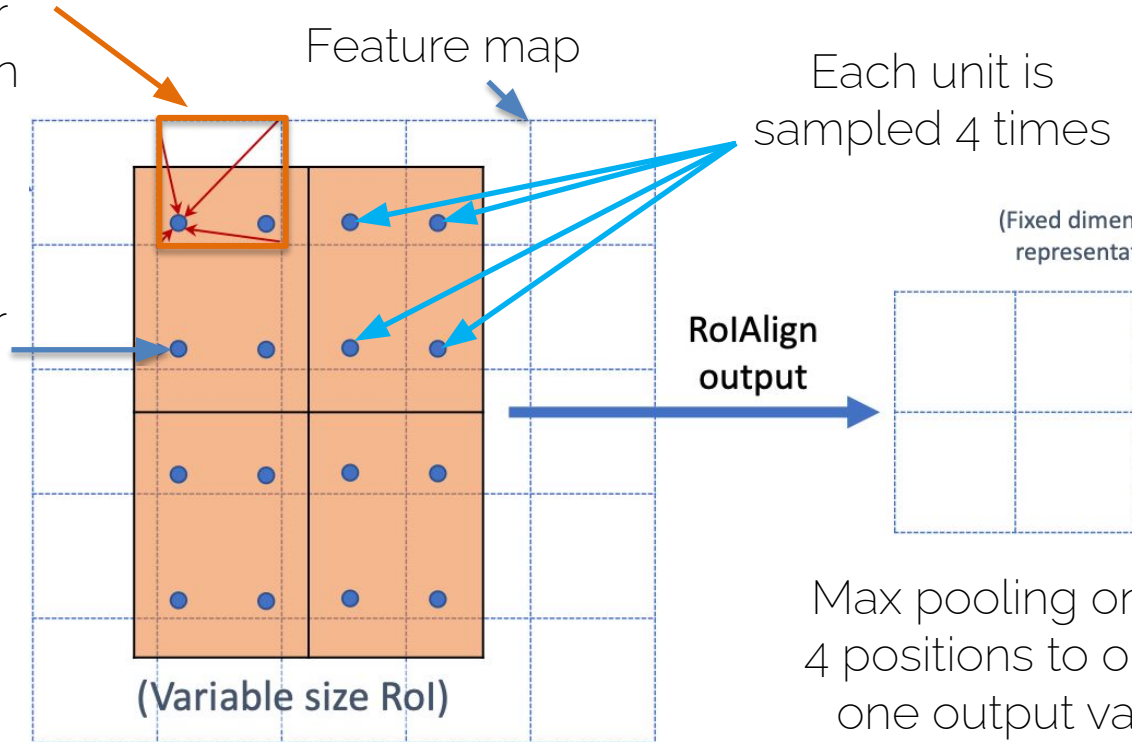
Grid points for  
bilinear  
interpolation

(Fixed dimensional  
representation)

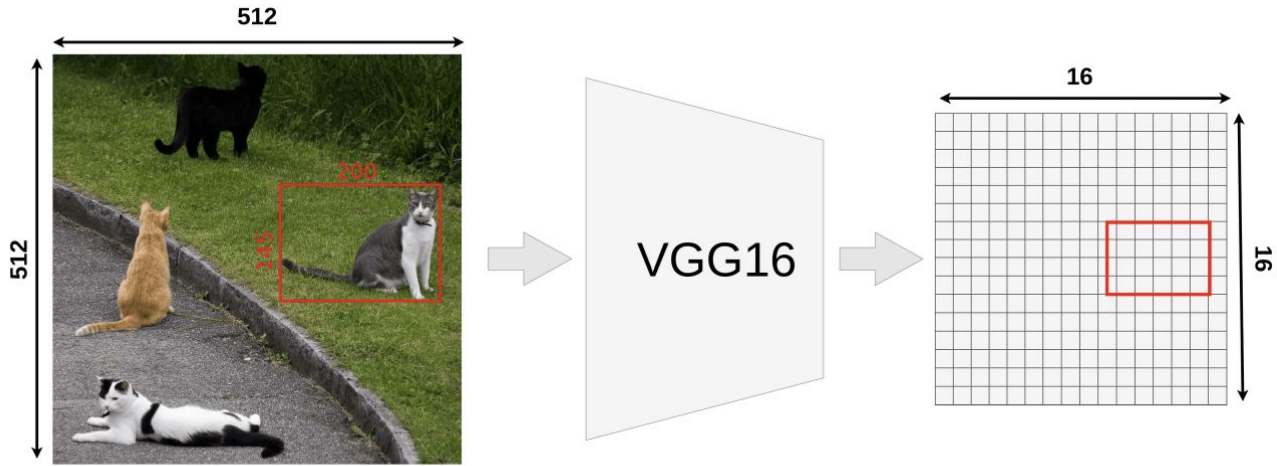
RoIAlign  
output

(Variable size RoI)

Max pooling on the  
4 positions to obtain  
one output value

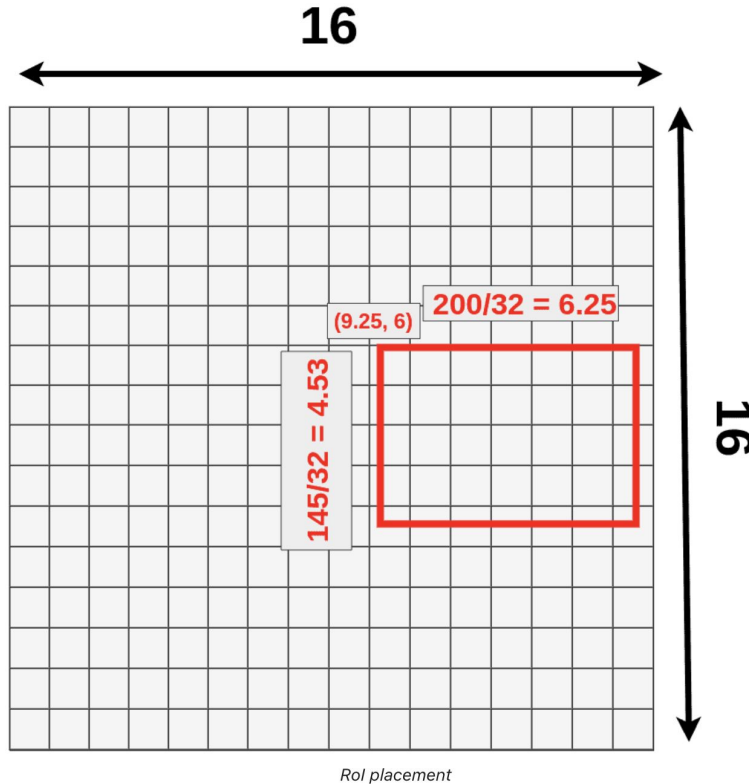


# RoI Pooling: Recall



*Model feature mapping process*

# RoI Pooling: Recall

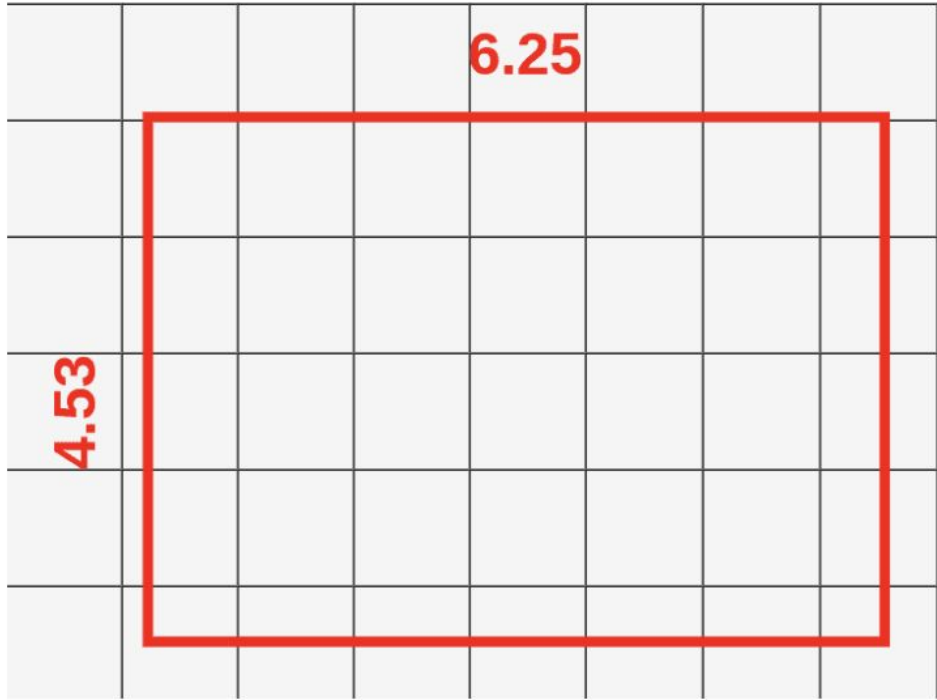


Next, we're using one of the proposed Rols (145x200 box) and try to map it onto the feature map. Because not all of our object dimensions can be divided by 32, we're placing RoI not align with our grid.

- (9.25,6) - top left corner
- 6.25 - width
- 4.53 - height

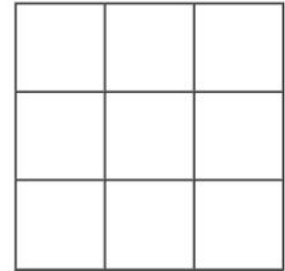


# RoI Pooling: Recall

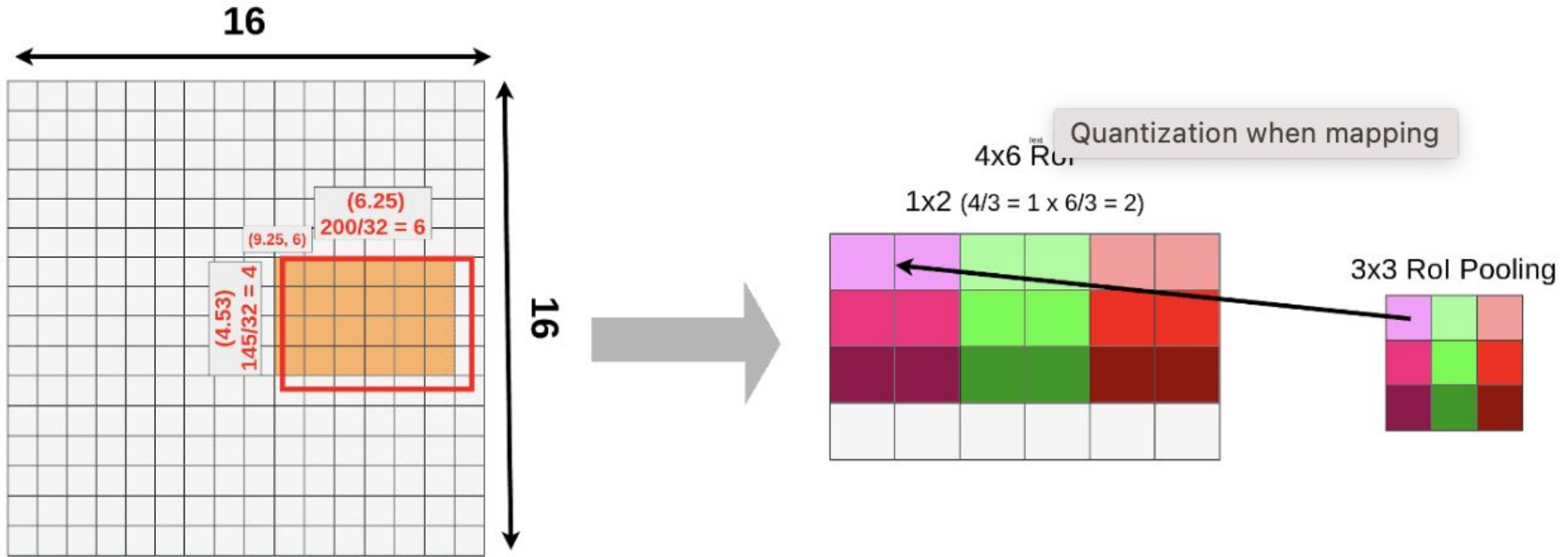


*Pooling layer*

3x3 RoI Pooling



# RoI Pooling: Recall



*Quantization when mapping and pooling*

# RoI Align



*RoI Box size*

Divide it into 9 boxes (because in our case the dimensions of our RoI Align are  $3 \times 3$ ).

# RoI Align



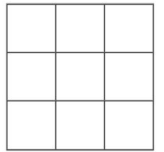
RoI Box size

Divide it into 9 boxes (because in our case the dimensions of our RoI Align are 3 x 3). That gives us a box with a height of 1.51 and a width of 2.08

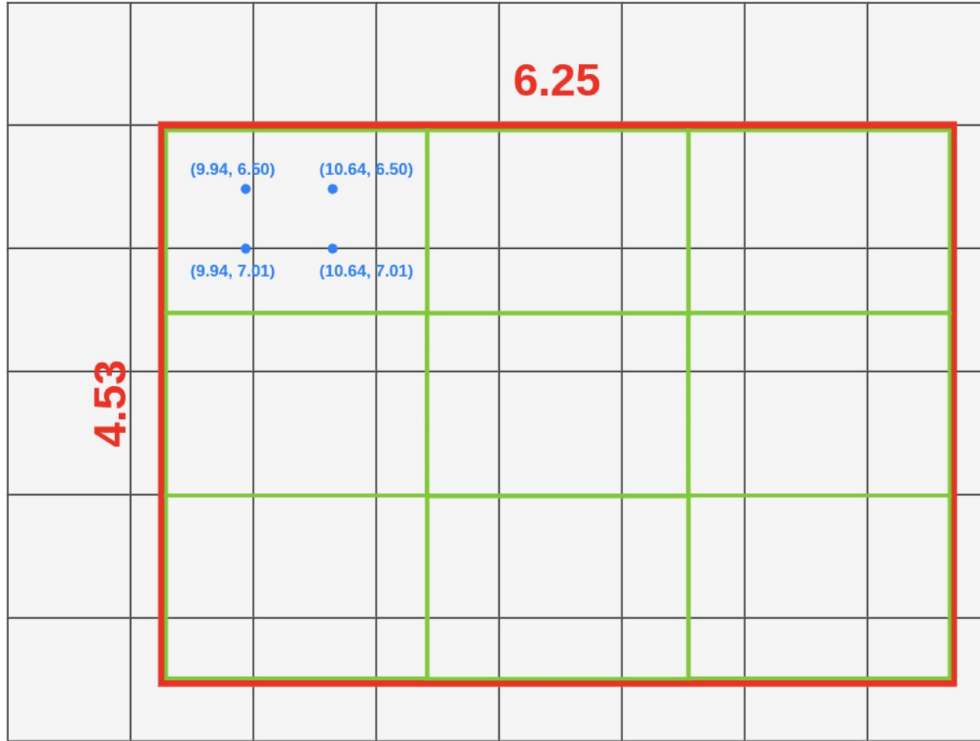


RoI divided into boxes

3x3 RoI Pooling



# RoI Align



Sampling points distribution

Get 4 sampling points. We get them by dividing the height and the width of the box by 3. For the first point:

- $X = X_{\text{box}} + (\text{width}/3) * 1 = 9.94$
- $Y = Y_{\text{box}} + (\text{height}/3) * 1 = 6.50$

where  $(X_{\text{box}}, Y_{\text{box}}) = (9.25, 6)$ .

For the second point:

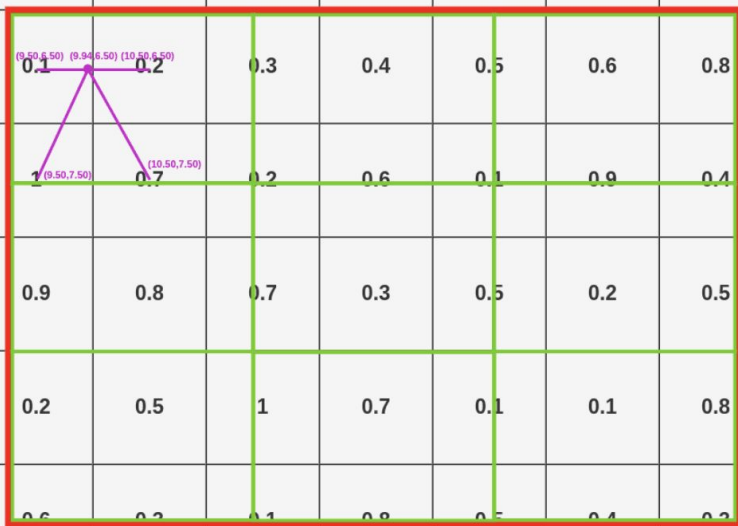
- $X = X_{\text{box}} + (\text{width}/3) * 1 = 9.94$
- $Y = Y_{\text{box}} + (\text{height}/3) * 2 = 7.01$

# RoI Align

$$P \approx \frac{7.5 - 6.5}{7.5 - 6.5} \left( \frac{10.5 - 9.94}{10.5 - 9.5} 0.1 + \frac{9.94 - 9.5}{10.5 - 9.5} 0.2 \right) + \frac{6.5 - 6.5}{7.5 - 6.5} \left( \frac{10.5 - 9.94}{10.5 - 9.5} 1 + \frac{9.94 - 9.5}{10.5 - 9.5} 0.7 \right)$$

**6.25**

**4.53**



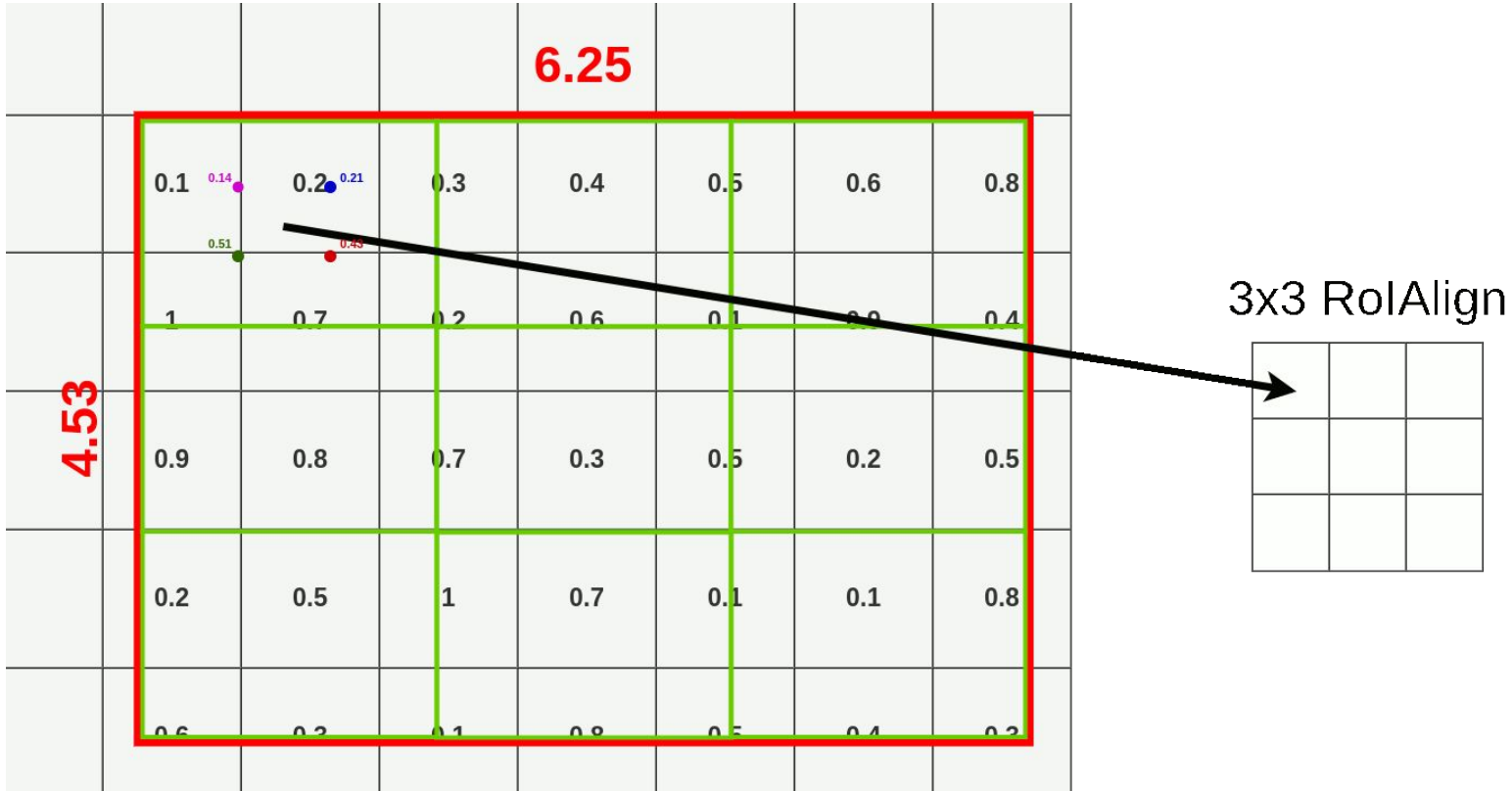
Bilinear Interpolation for the first point

Do bilinear interpolation following the equation:

$$P \approx \frac{y_2 - y}{y_2 - y_1} \left( \frac{x_2 - x}{x_2 - x_1} Q_{11} + \frac{x - x_1}{x_2 - x_1} Q_{21} \right) + \frac{y - y_1}{y_2 - y_1} \left( \frac{x_2 - x}{x_2 - x_1} Q_{12} + \frac{x - x_1}{x_2 - x_1} Q_{22} \right)$$

*Bilinear Interpolation equation*

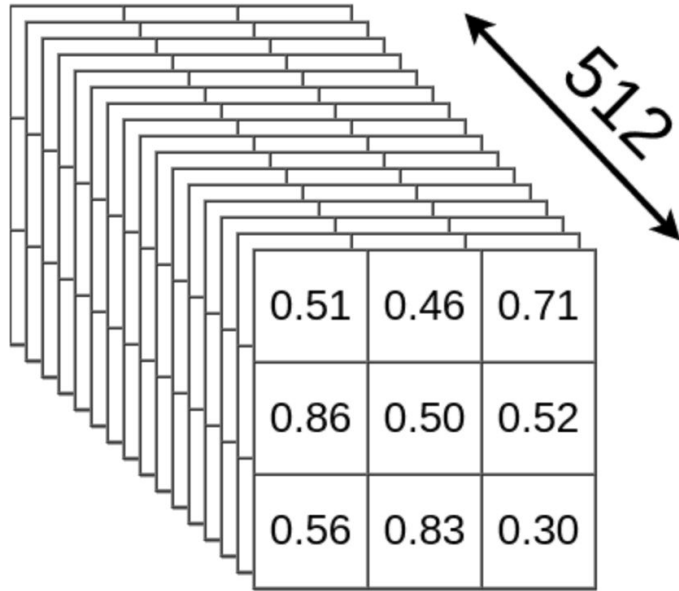
# RoI Align



# RoI Align

## 3x3 RoIAlign

Do this for every channel.

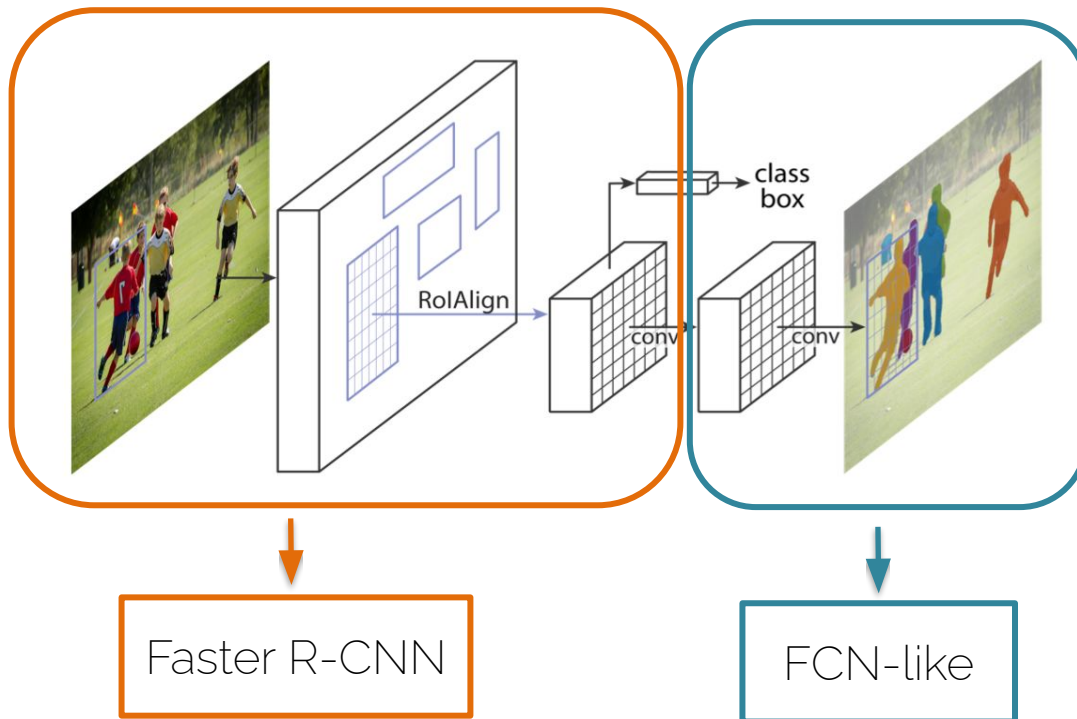


*RoIAlign full size*



# What is Mask-RCNN?

- Faster R-CNN + FCN for segmentation



Mask loss =  
binary cross  
entropy per pixel  
for the  $k$  semantic  
classes

# Mask R-CNN: qualitative results



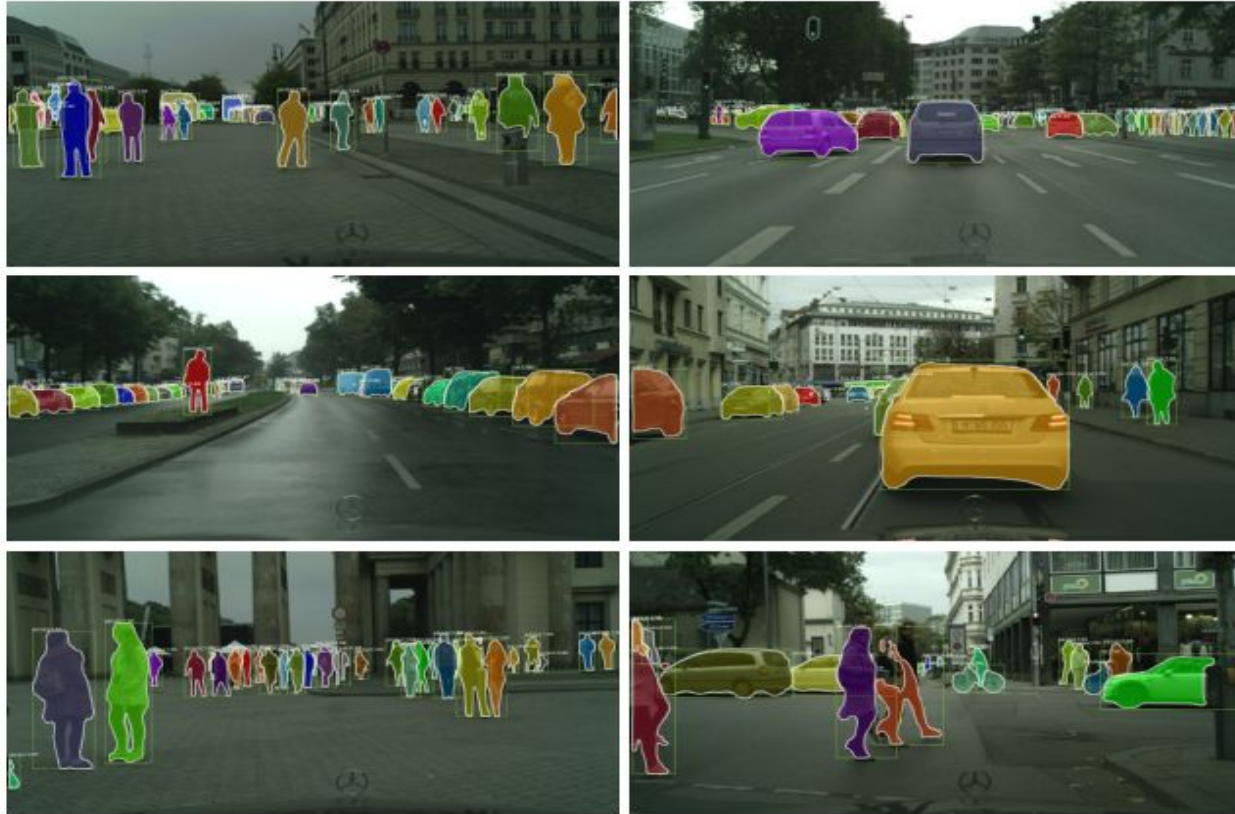
# Mask R-CNN: qualitative results



# Mask R-CNN: qualitative results



# Mask R-CNN: qualitative results



# Mask R-CNN: extended for joints



Model a keypoint's location as a one-hot mask, and adopt Mask R-CNN to predict  $K$  masks, one for each of  $K$  keypoint types (e.g., left shoulder, right elbow). This demonstrates the flexibility of Mask R-CNN.

# Improving Mask-RCNN

- One problem with Mask R-CNN is that the mask quality score is computed as the confidence score for the bounding box

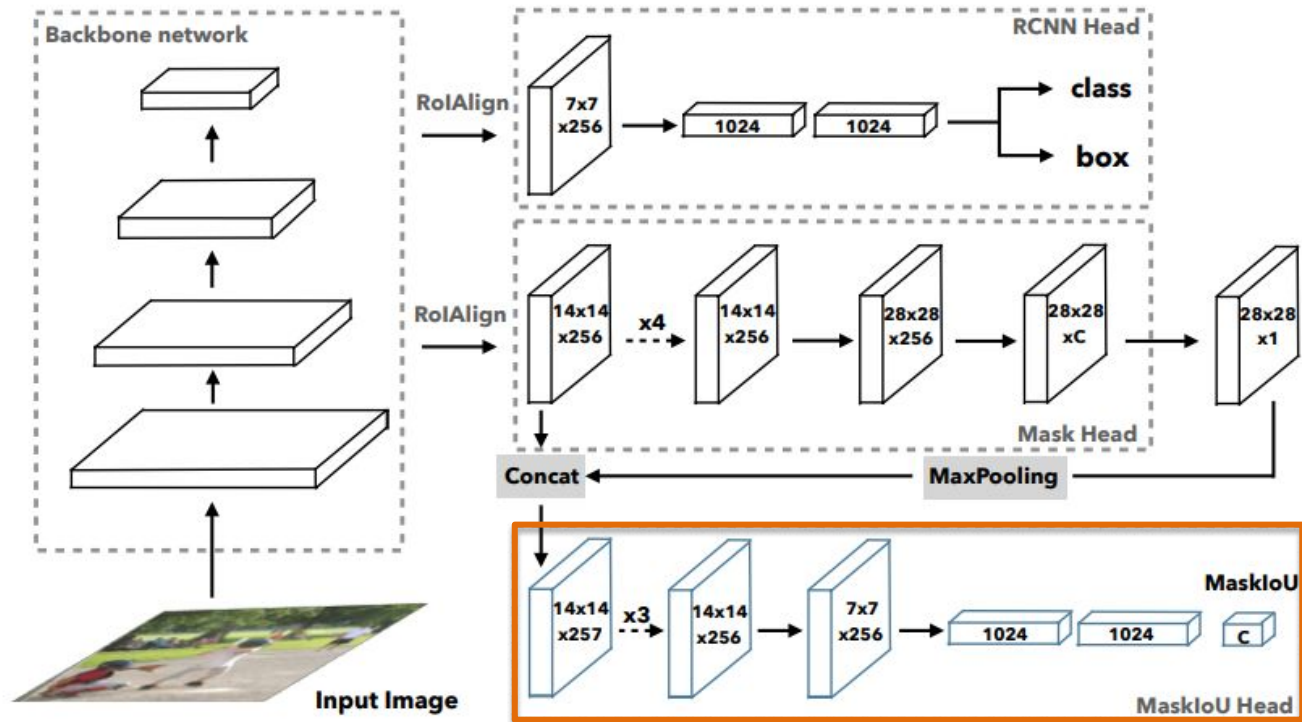
The only way the "instance" is evaluated is through the box loss



Recall the mask loss just evaluates if the pixels have the correct semantic class, not the correct instance!

Both instances have the same class = person

# Mask IoU head



Measure the intersection over union between the predicted mask and ground truth mask



# Mask confidence score



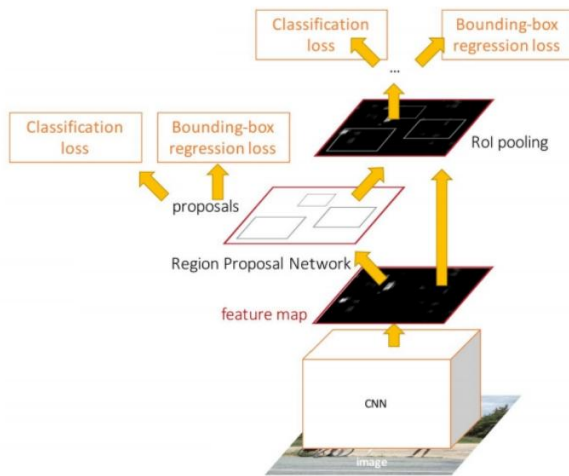
Typically, Mask scoring R-CNN gives lower confidence scores than Mask R-CNN, which corresponds to masks not being perfect ( $\text{IoU} < 1.0$ ).

This tiny modification achieves SOTA results.

Is one-stage vs  
two-stage also  
applicable to masks?

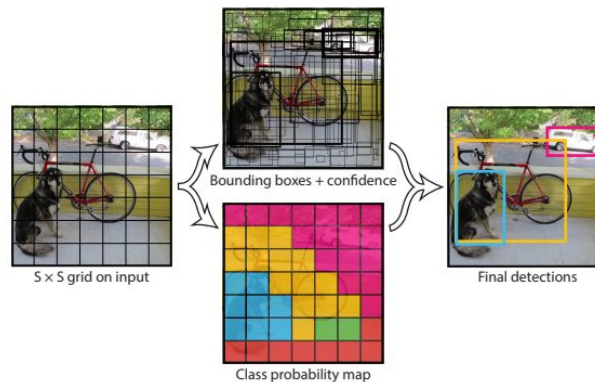
# One-stage vs two-stage detectors

## Faster R-CNN



Slower, but has higher performance

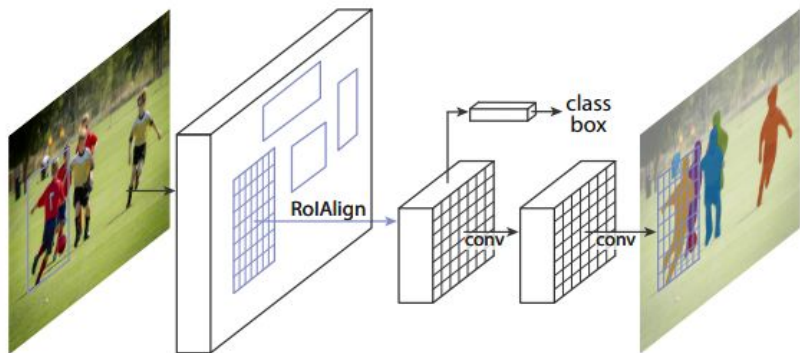
## YOLO



Faster, but has lower performance

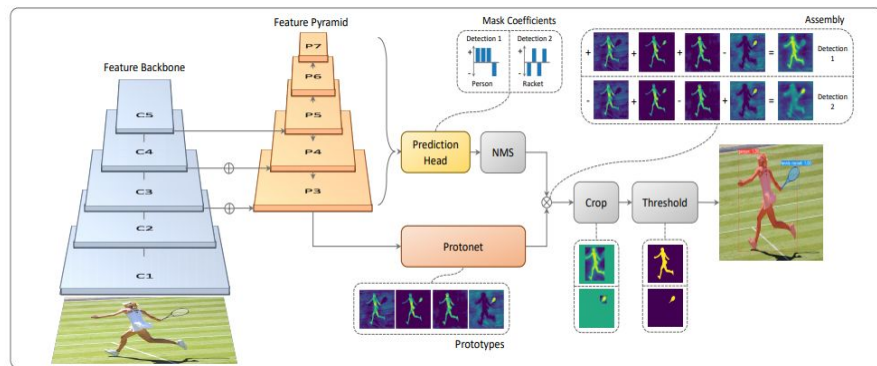
# One-stage vs two-stage instance segmenters

Mask R-CNN



Slower, but has higher performance

YOLOACT



Faster, but has lower performance

# YOLO with masks?

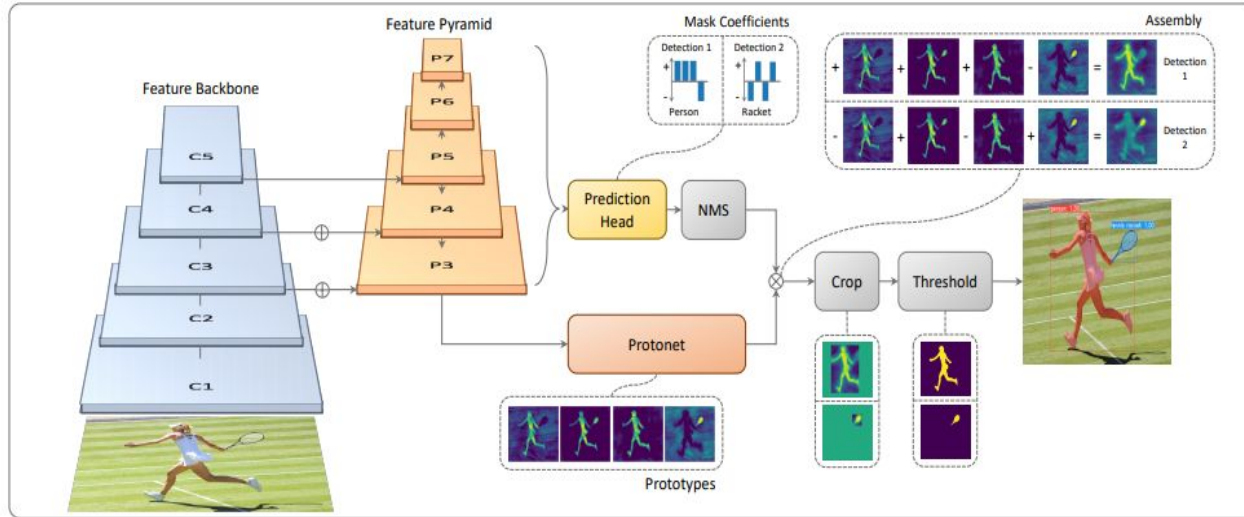
*"Boxes are stupid anyway though, I'm probably a true believer in masks except I can't get YOLO to learn them."*

– Joseph Redmon, YOLOv3

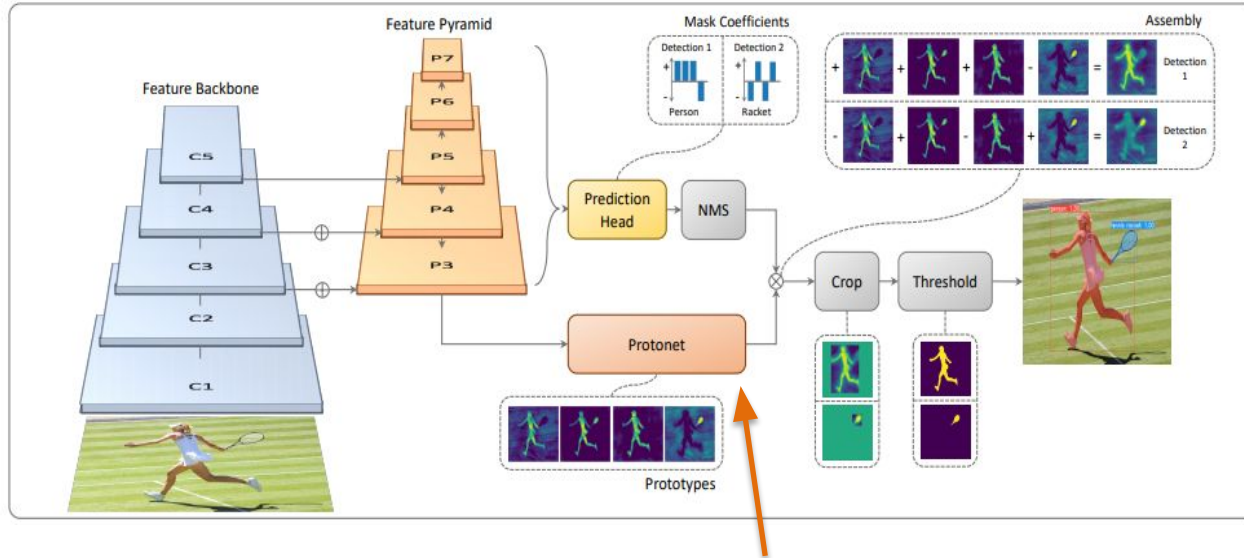
# YoLACT\*

\*You Only Look At CoefficientTs

# YOLOACT: idea



# YOLOACT: idea

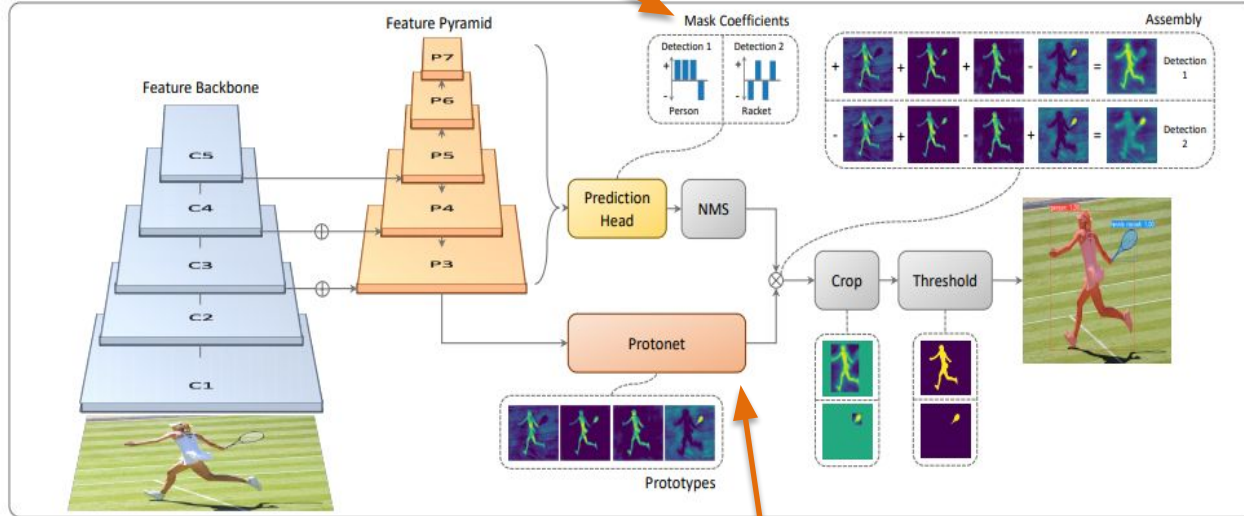


1) Generate mask prototypes



# YOLOACT: idea

2) Generate mask coefficients

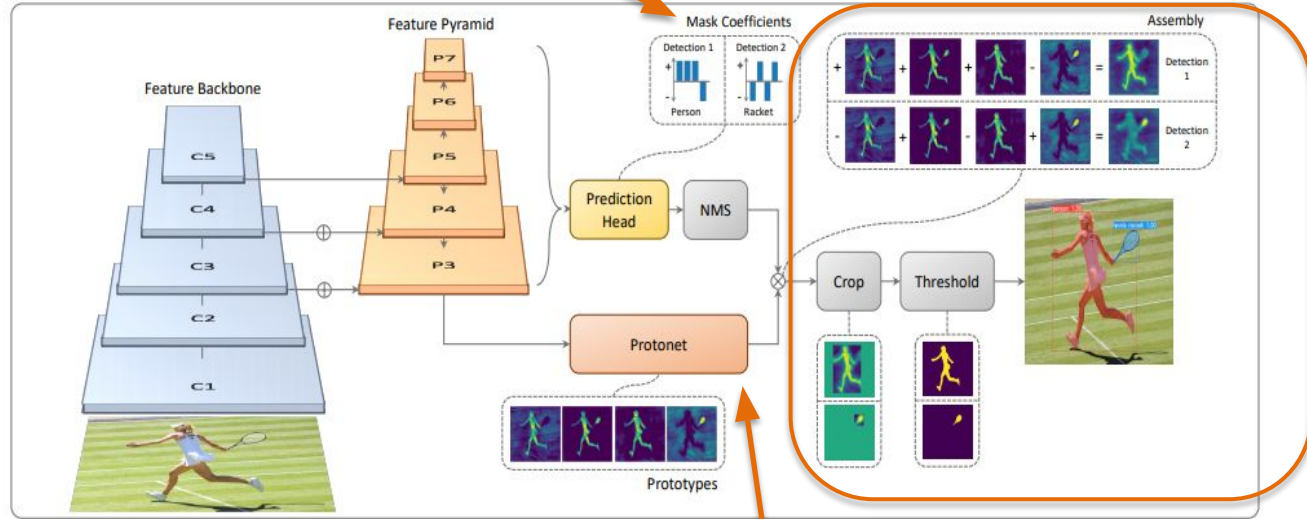


1) Generate mask prototypes

# YOLOACT: idea

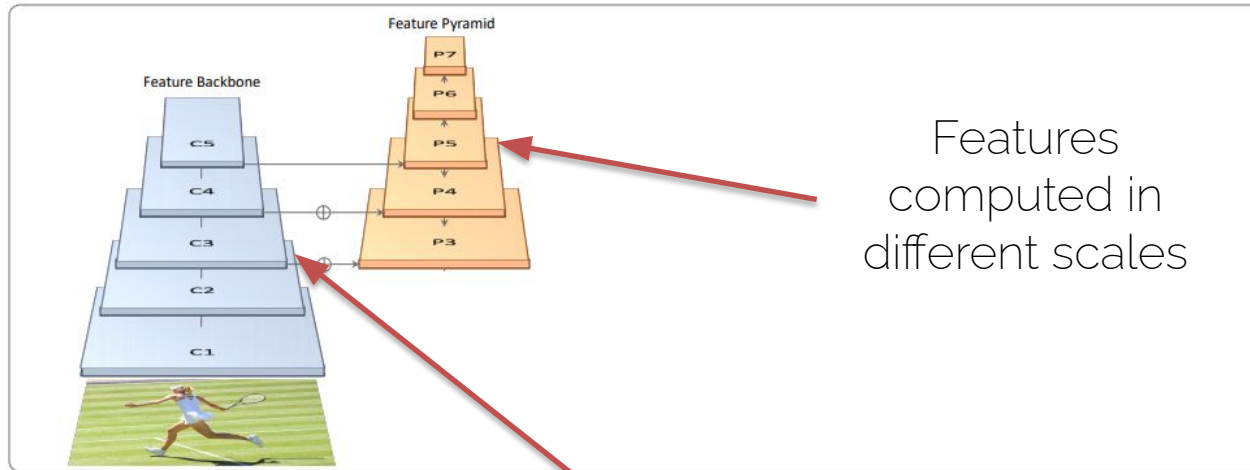
2) Generate mask coefficients

3) Combine (1) and (2)



1) Generate mask prototypes

# YOLOACT: backbone



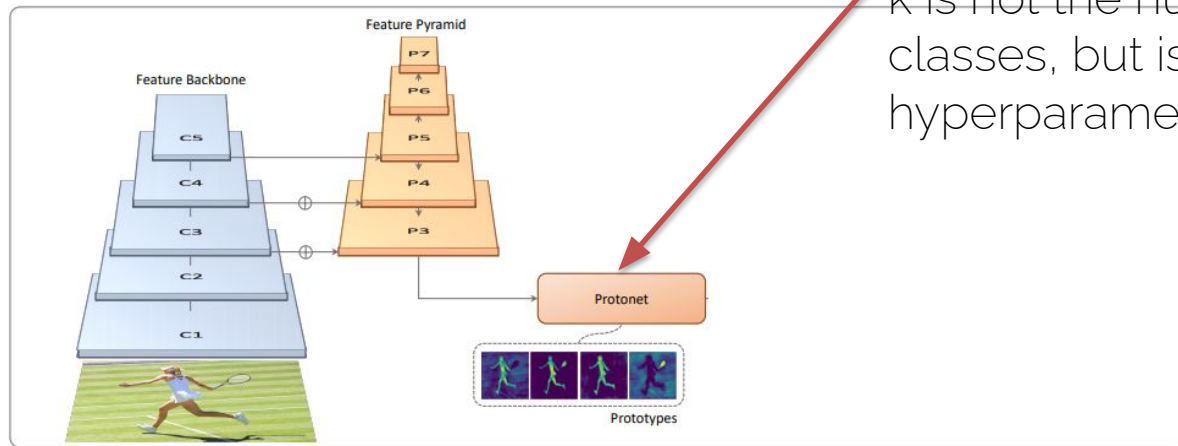
Features  
computed in  
different scales

ResNet-101

# YOLOACT: protonet

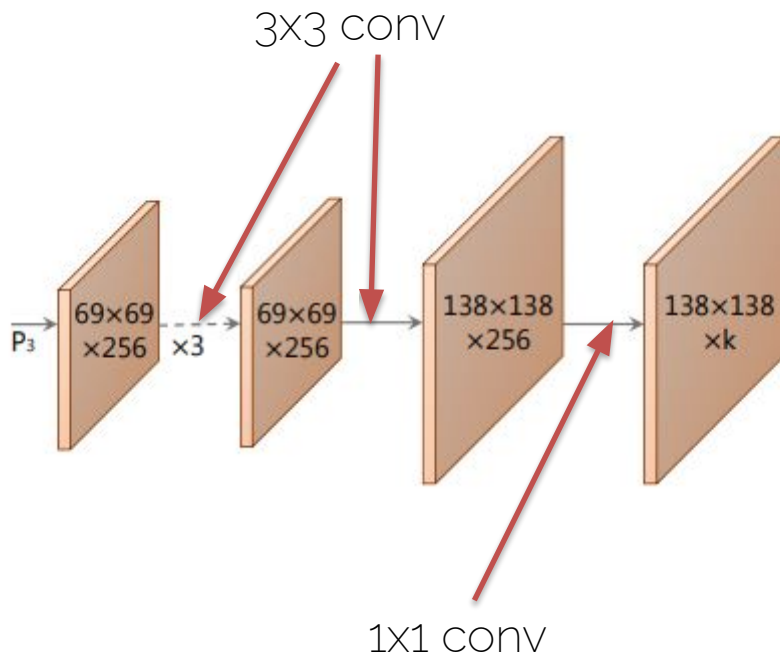
Generate k prototype masks.

k is not the number of classes, but is a hyperparameter.



# YOLOACT: protonet

- Fully convolutional network

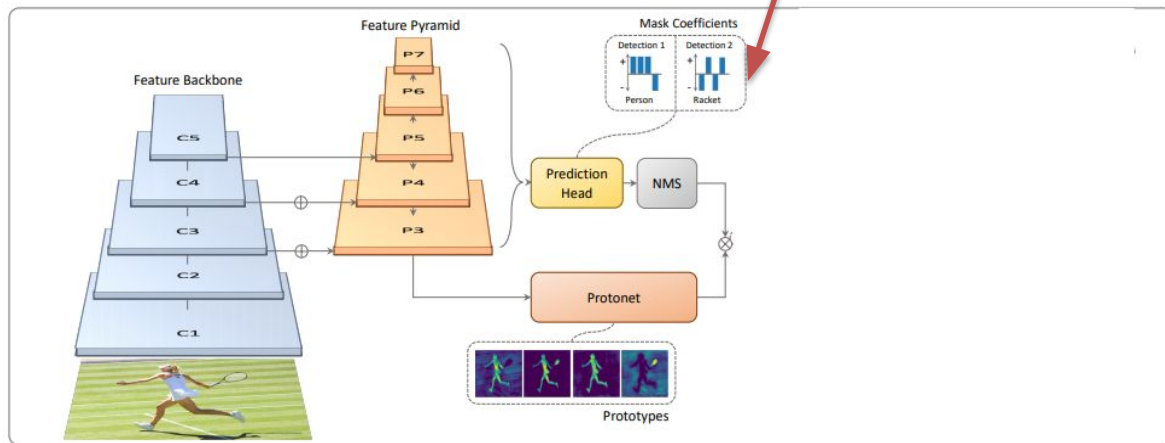


Similar to the mask branch in Mask R-CNN.

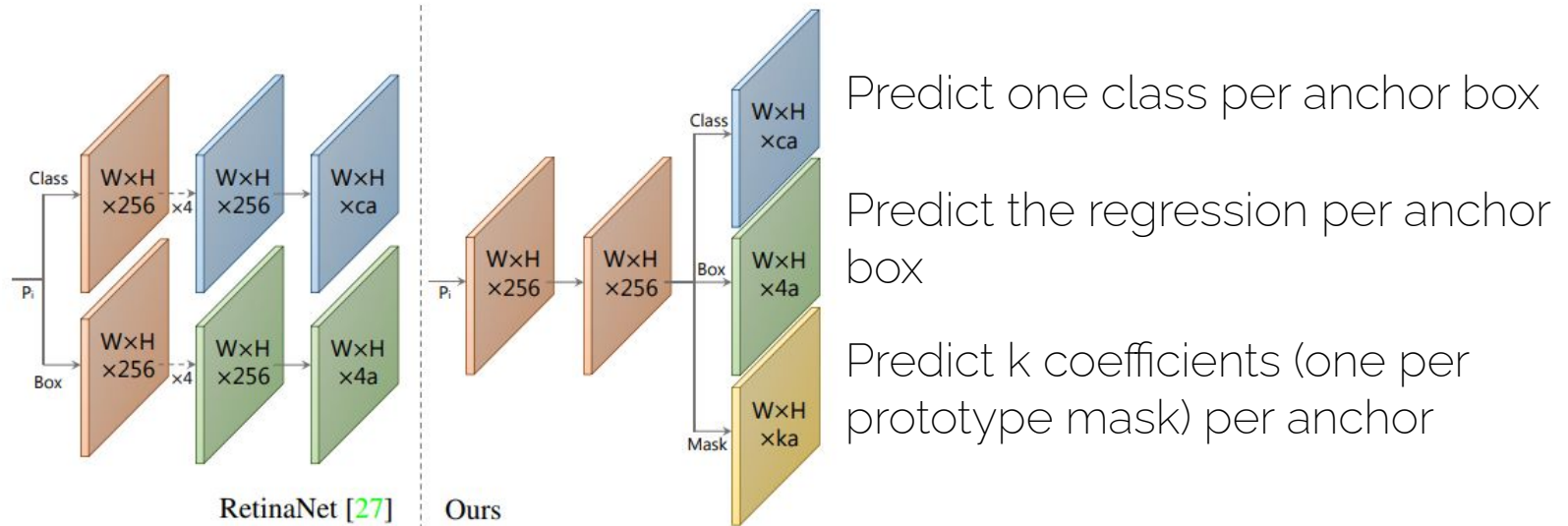
However, no loss function is applied on this stage.

# YOLOACT: mask coefficients

Predict a coefficient for every predicted mask.



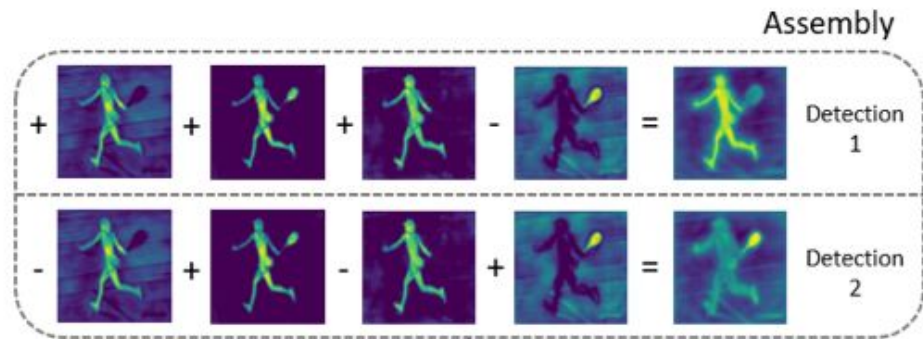
# YOLOACT: mask coefficients



The network is similar but shallower than RetinaNet

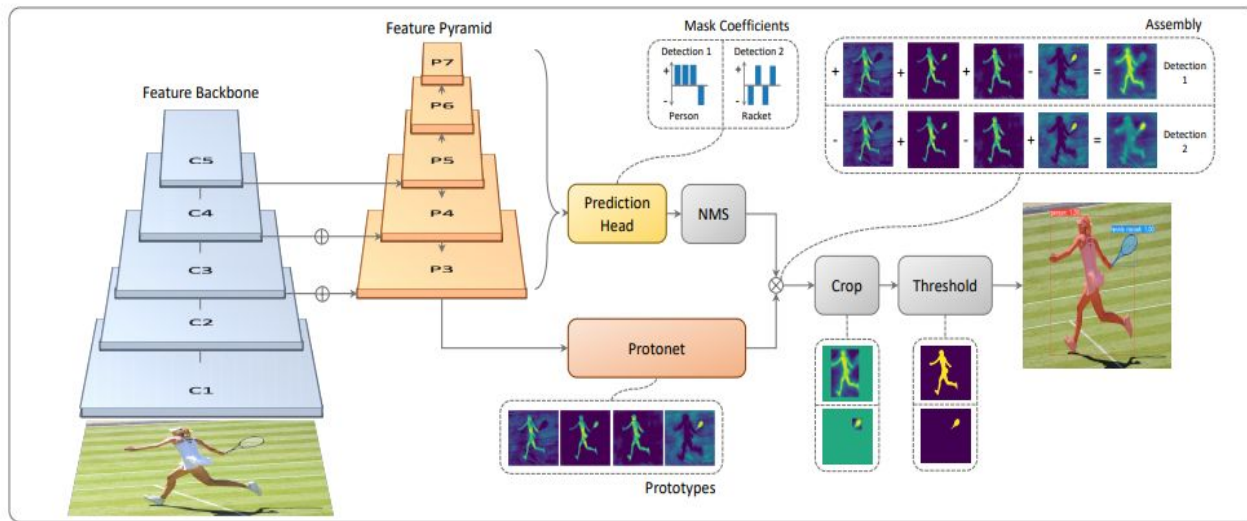
# YOLOACT: mask assembly

1. Do a linear combination between the mask coefficients and the mask prototypes.
2. Predict the mask as  $M = \sigma(PC^T)$  where  $P$  is a  $(H \times W \times K)$  matrix of prototype masks,  $C$  is a  $(N \times K)$  matrix of mask coefficients surviving NMS, and  $\sigma$  is a nonlinearity.





# YOLOACT: loss function

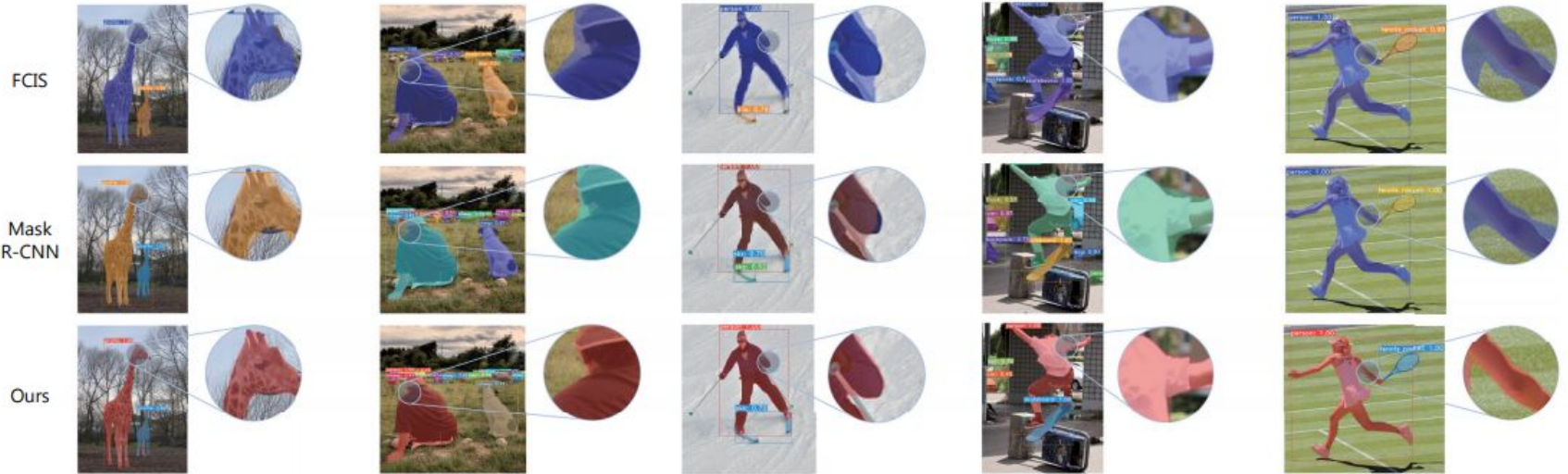


Cross-entropy between the assembled masks and the ground truth, in addition to the standard losses (regression for the bounding box, and classification for the class of the object/mask).

# YOLOACT: qualitative results

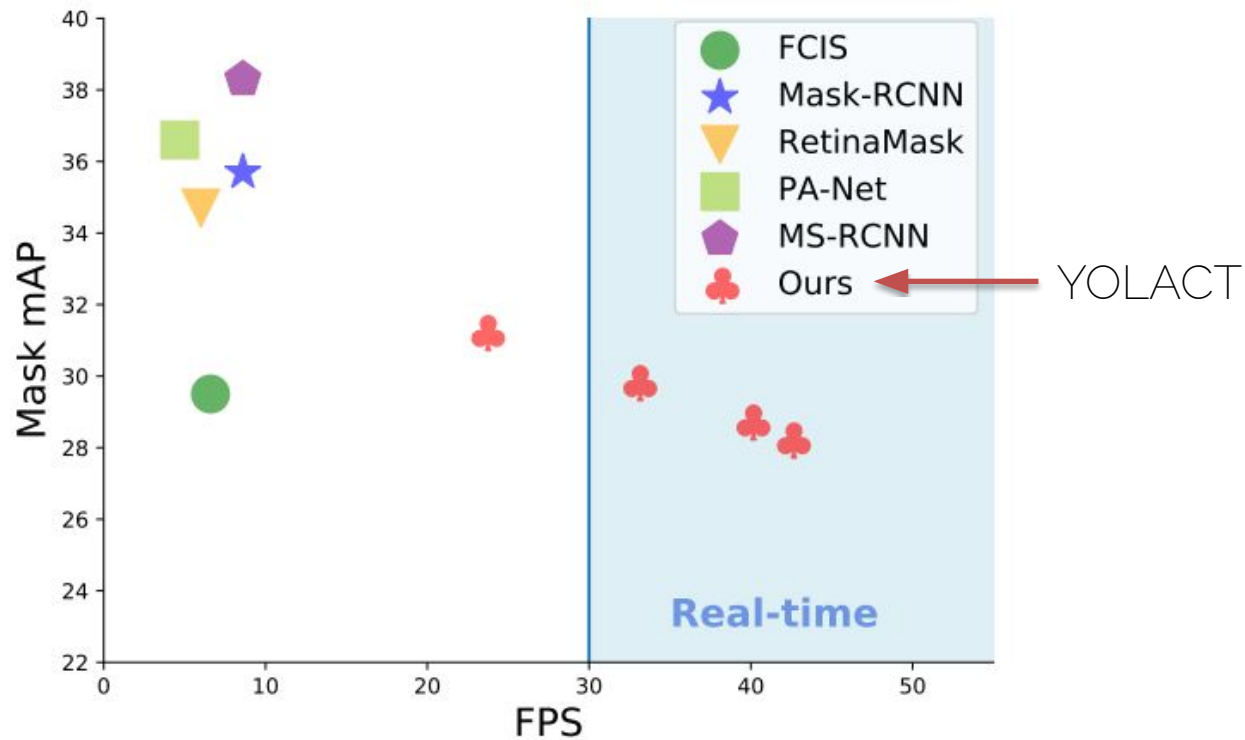


# YOLOACT: qualitative results



For large objects, the quality of the masks is even better than those of two-stage detectors

# So, which segmenter to use?



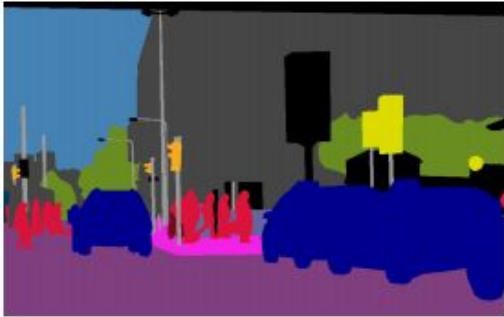
# YOLOACT++: improvements

- A specially designed version of NMS, in order to make the procedure faster.
- An auxiliary semantic segmentation loss function performed on the final features of the FPN. The module is not used during the inference stage.
- D. Boyla et al. "YOLOACT++: Better real-time instance segmentation". [arXiv:1912.06218](https://arxiv.org/abs/1912.06218) 2019

# Panoptic segmentation

# Panoptic segmentation

Semantic segmentation

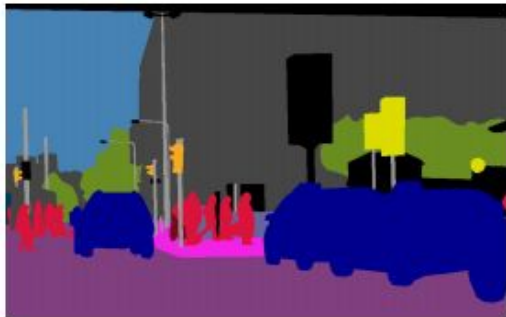


Instance segmentation



# Panoptic segmentation

Semantic segmentation



FCN-like

Instance segmentation

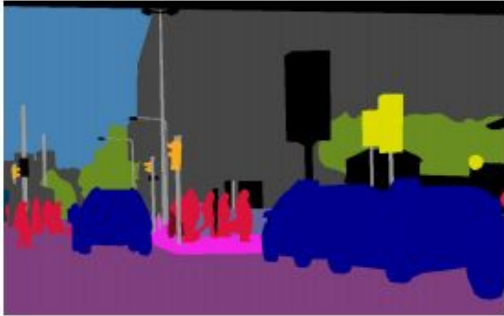


Mask R-CNN



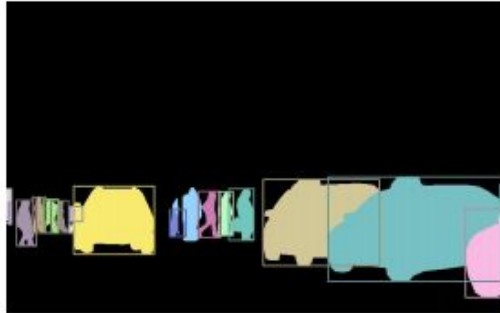
# Panoptic segmentation

Semantic segmentation



FCN-like

Instance segmentation



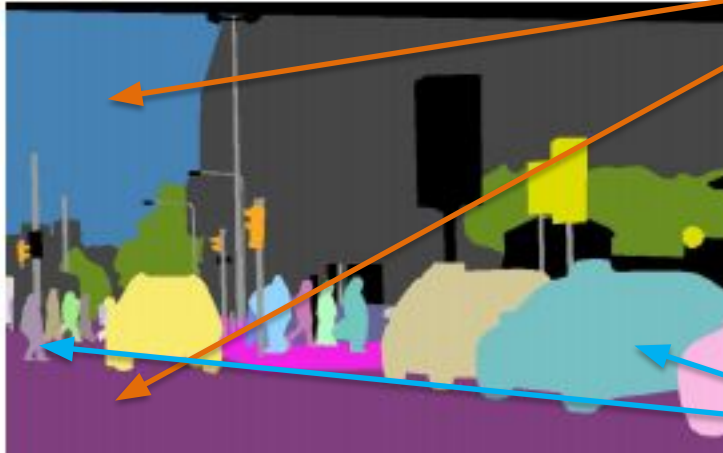
Mask R-CNN

Panoptic segmentation



UPSNet

# Panoptic segmentation



It gives labels to uncountable objects called "stuff" (sky, road, etc), similar to FCN-like networks.

It differentiates between pixels coming from different instances of the same class (countable objects) called "things" (cars, pedestrians, etc).

# Panoptic segmentation



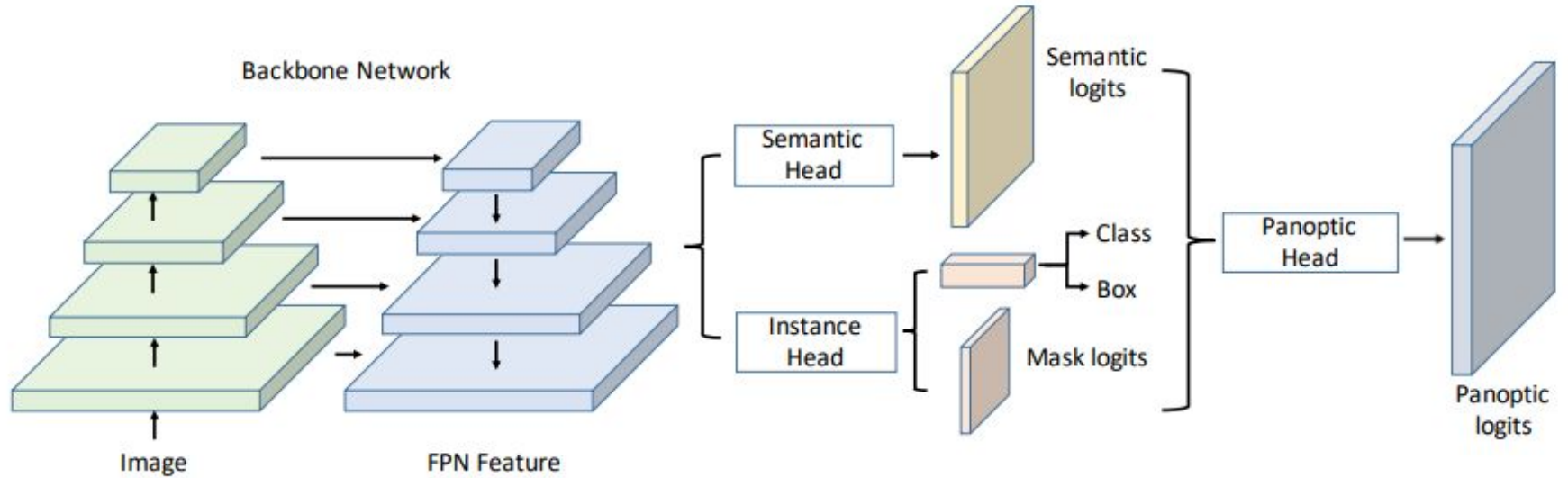
Problem: some pixels might get classified as stuff from FCN network, while at the same time being classified as instances of some class from Mask R-CNN (conflicting results)!

# Panoptic segmentation



Solution: Parametric-free panoptic head which combines the information from the FCN and Mask R-CNN, giving final predictions.

# Network architecture

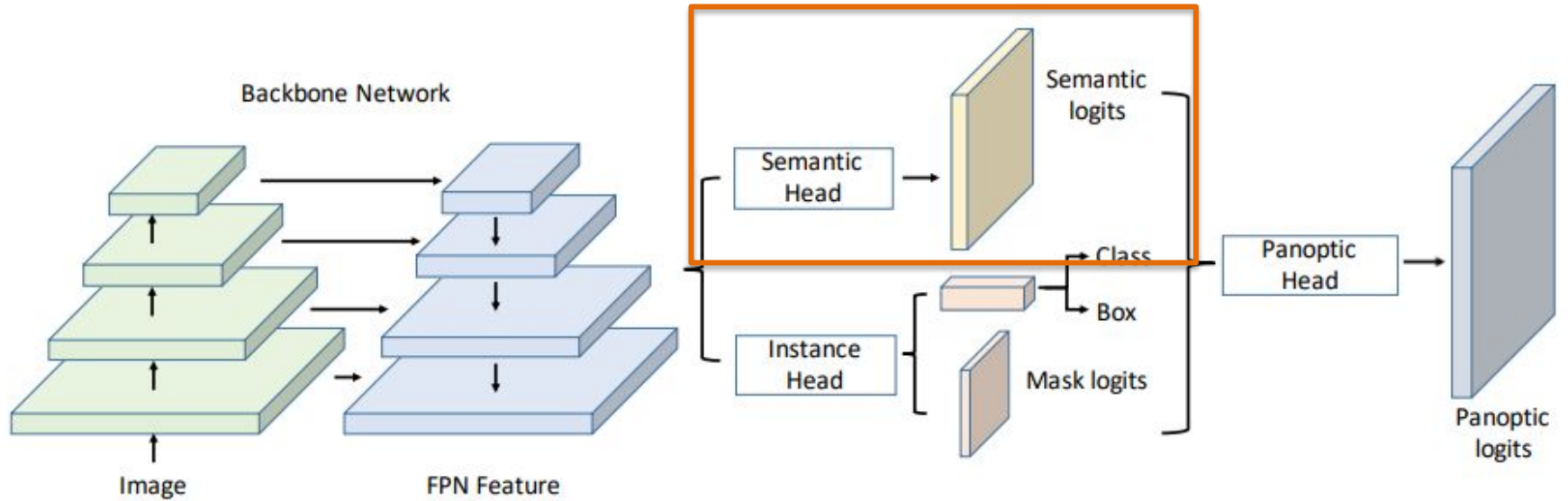


Shared features

Separate heads

Putting it together

# Network architecture

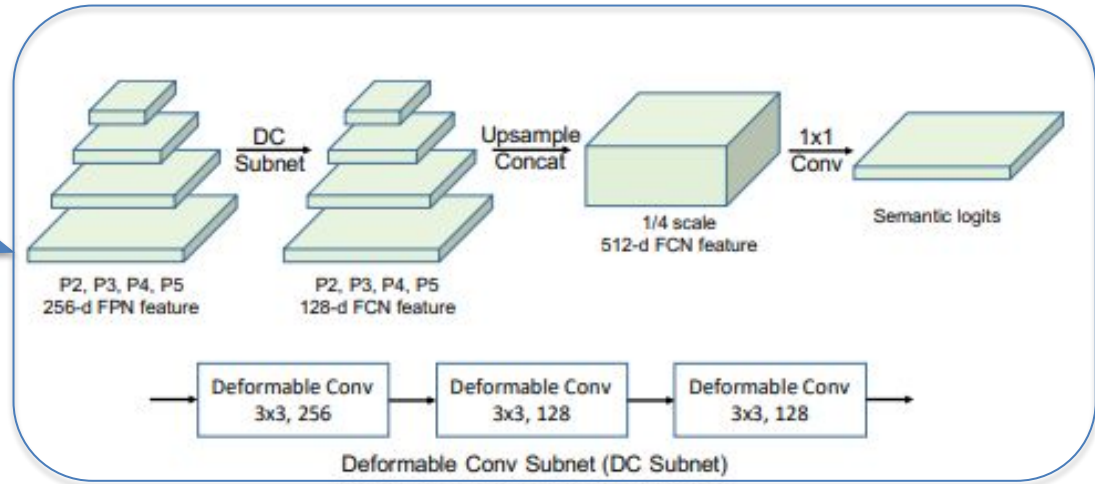
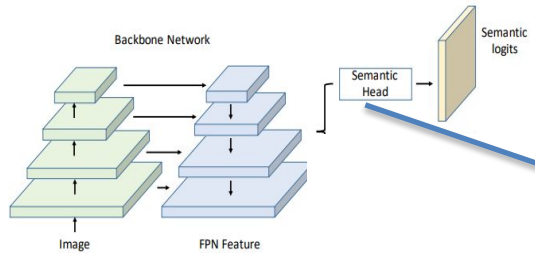


Shared features

Separate heads

Putting it together

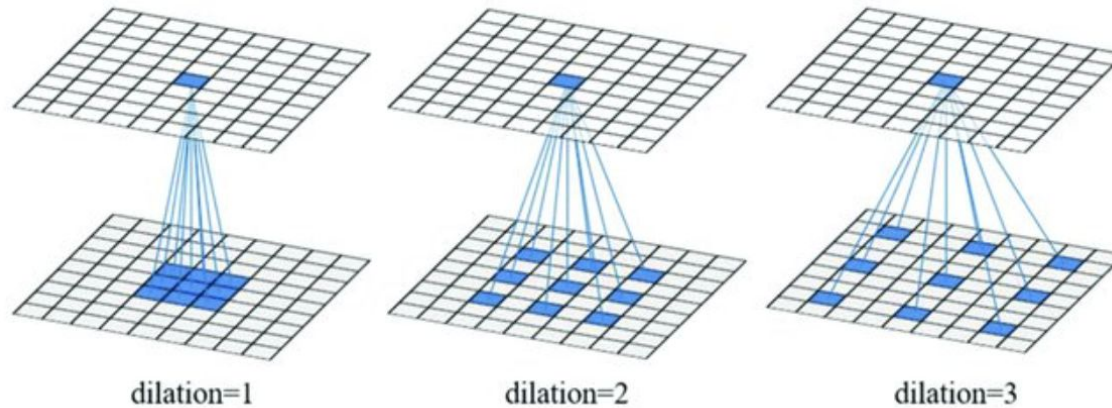
# The semantic head



As all semantic heads, fully convolutional network.

New: deformable convolutions!

# Dilated (atrous) convolutions 2D



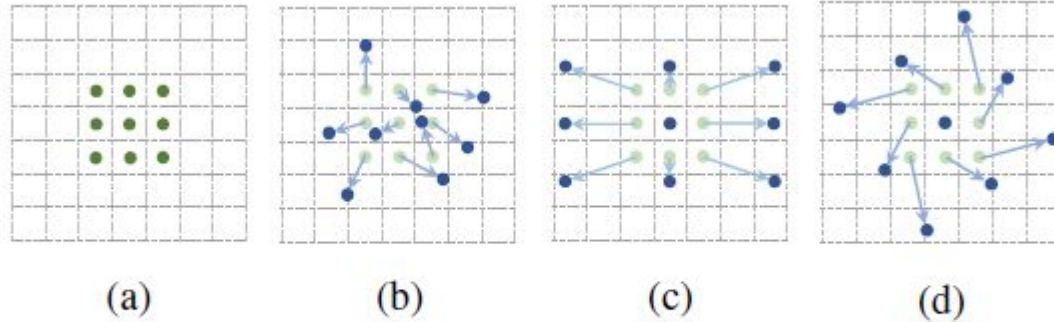
(a) the dilation parameter is 1, and each element produced by this filter has receptive field of  $3 \times 3$ .

(b) the dilation parameter is 2, and each element produced by it has receptive field of  $5 \times 5$ .

(c) the dilation parameter is 3, and each element produced by it has receptive field of  $7 \times 7$ .



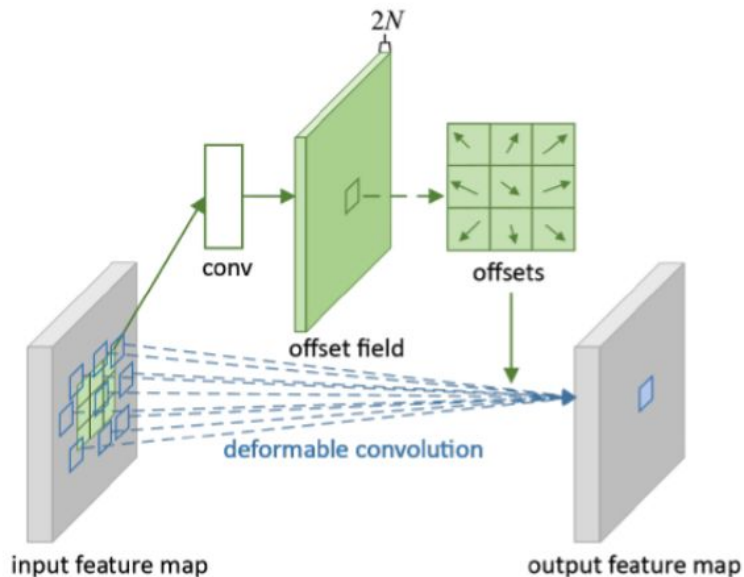
# Deformable convolutions



**(a) Conventional Convolution, (b) Deformable Convolution, (c) Special Case of Deformable Convolution with Scaling, (d) Special Case of Deformable Convolution with Rotation**

Deformable convolutions: generalization of dilated convolutions when you learn the offset

# Deformable convolutions



Regular convolution

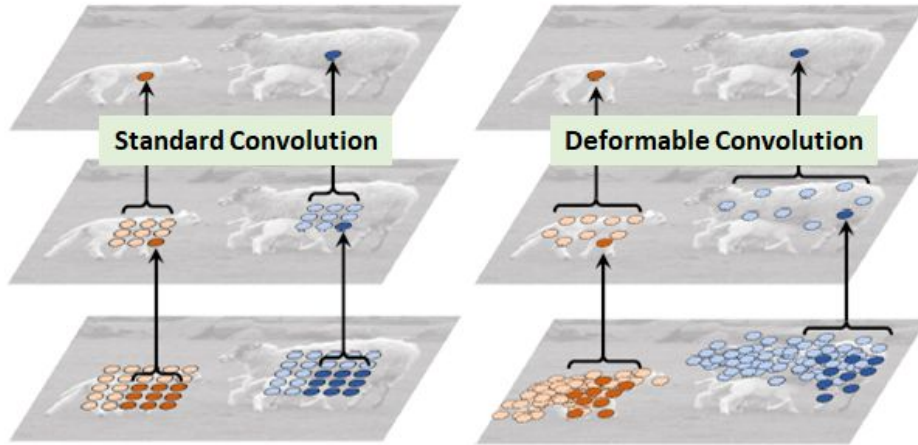
$$y(\mathbf{p}_0) = \sum_{\mathbf{p}_n \in \mathcal{R}} w(\mathbf{p}_n) \cdot x(\mathbf{p}_0 + \mathbf{p}_n)$$

Deformable convolution

$$y(\mathbf{p}_0) = \sum_{\mathbf{p}_n \in \mathcal{R}} w(\mathbf{p}_n) \cdot x(\mathbf{p}_0 + \mathbf{p}_n + \Delta \mathbf{p}_n)$$

where  $\Delta \mathbf{p}_n$  is generated by a sibling branch of regular convolution

# Deformable convolutions



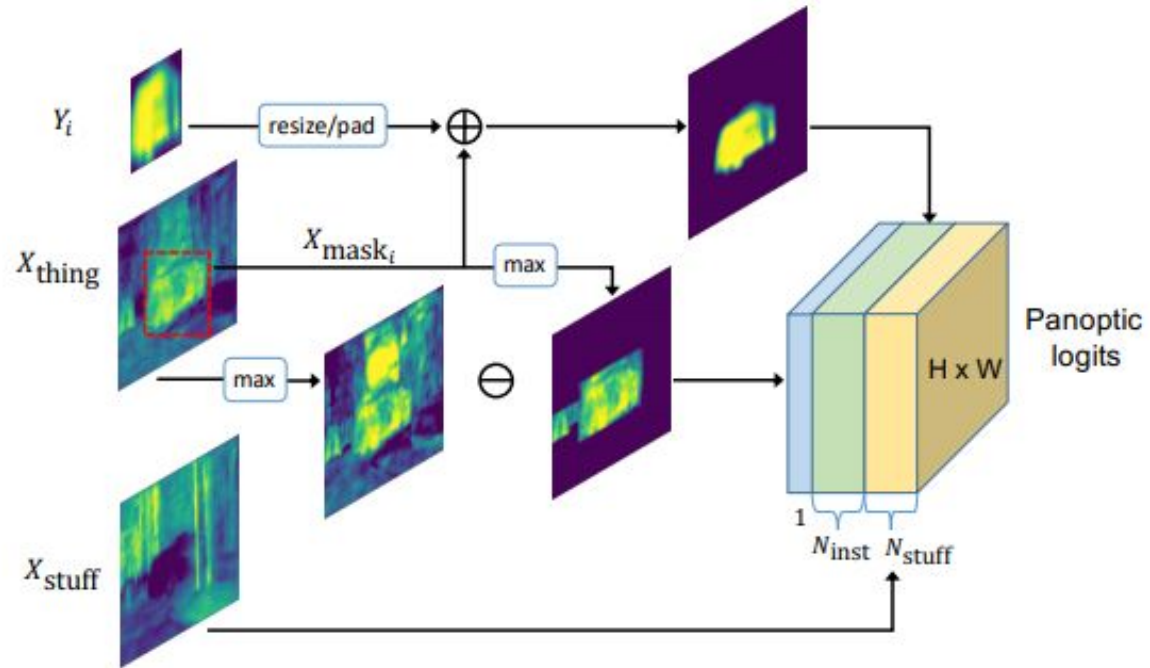
The deformable convolution will pick the values at different locations for convolutions conditioned on the input image of the feature maps.

# The Panoptic head

Mask logits from the instance head

Object logits coming from the semantic head (e.g., car)

Stuff logits coming from the semantic head (e.g., sky)

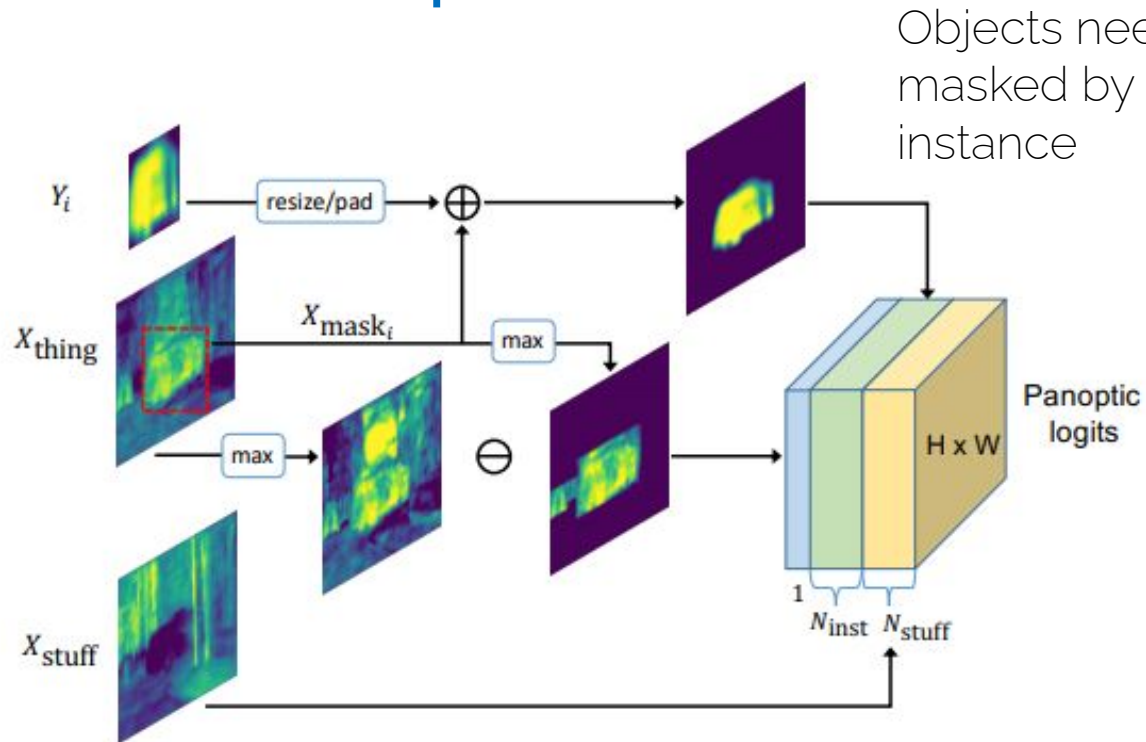


# The Panoptic head

Mask logits from the instance head

Object logits coming from the semantic head (e.g., car)

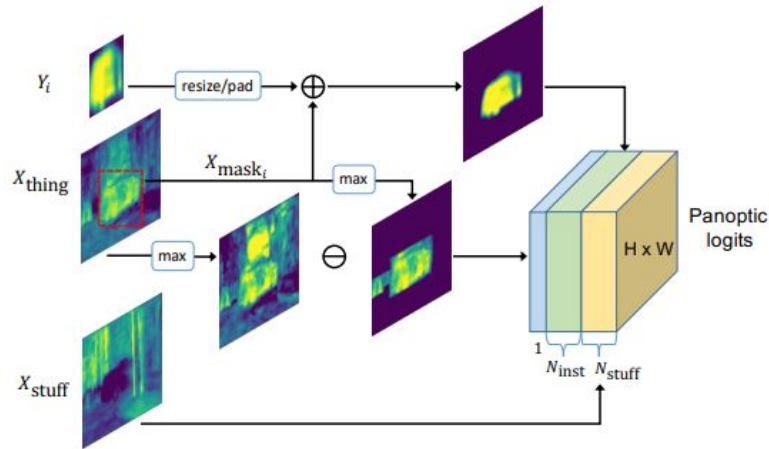
Stuff logits coming from the semantic head (e.g., sky)



Objects need to be masked by the instance

This can be evaluated directly

# The Panoptic head



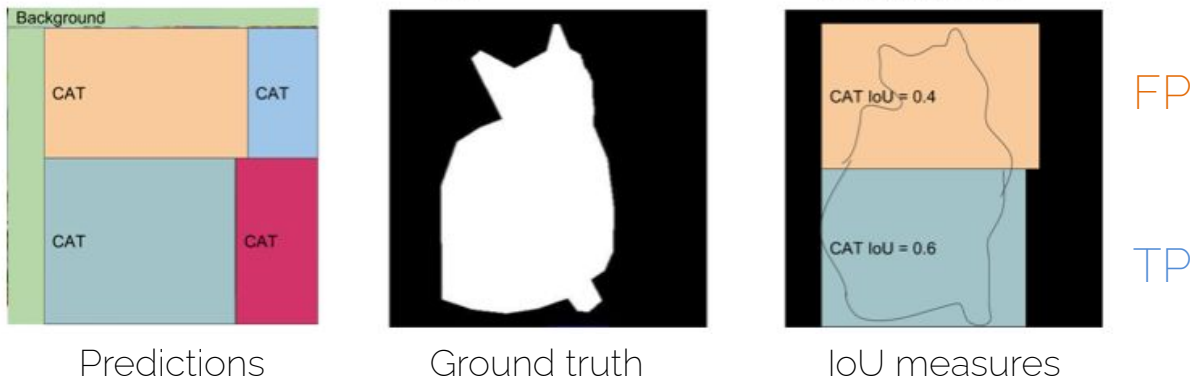
Perform softmax over the panoptic logits. If the maximum value falls into the first stuff channels, then it belongs to one of the stuff classes. Otherwise the index of the maximum value tells us the instance ID the pixel belongs to.

Read the details on how to use the unknown class



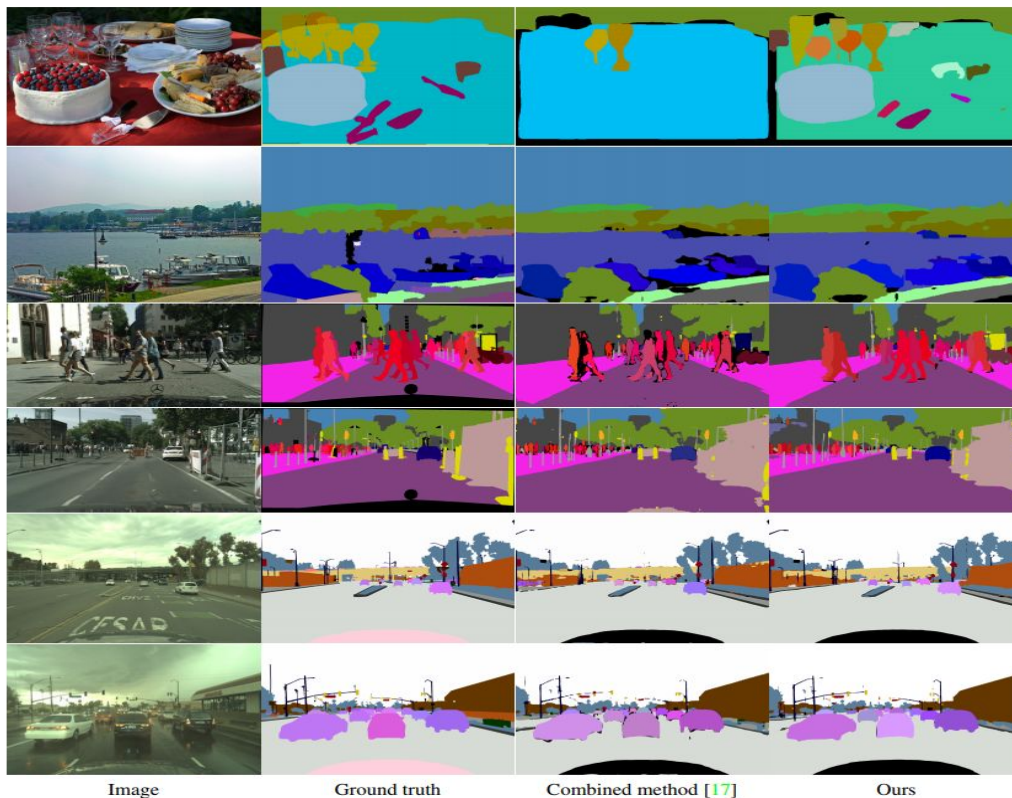
# Panoptic quality

- As in detection, we have to “match ground truth and predictions. In this case we have segment matching.



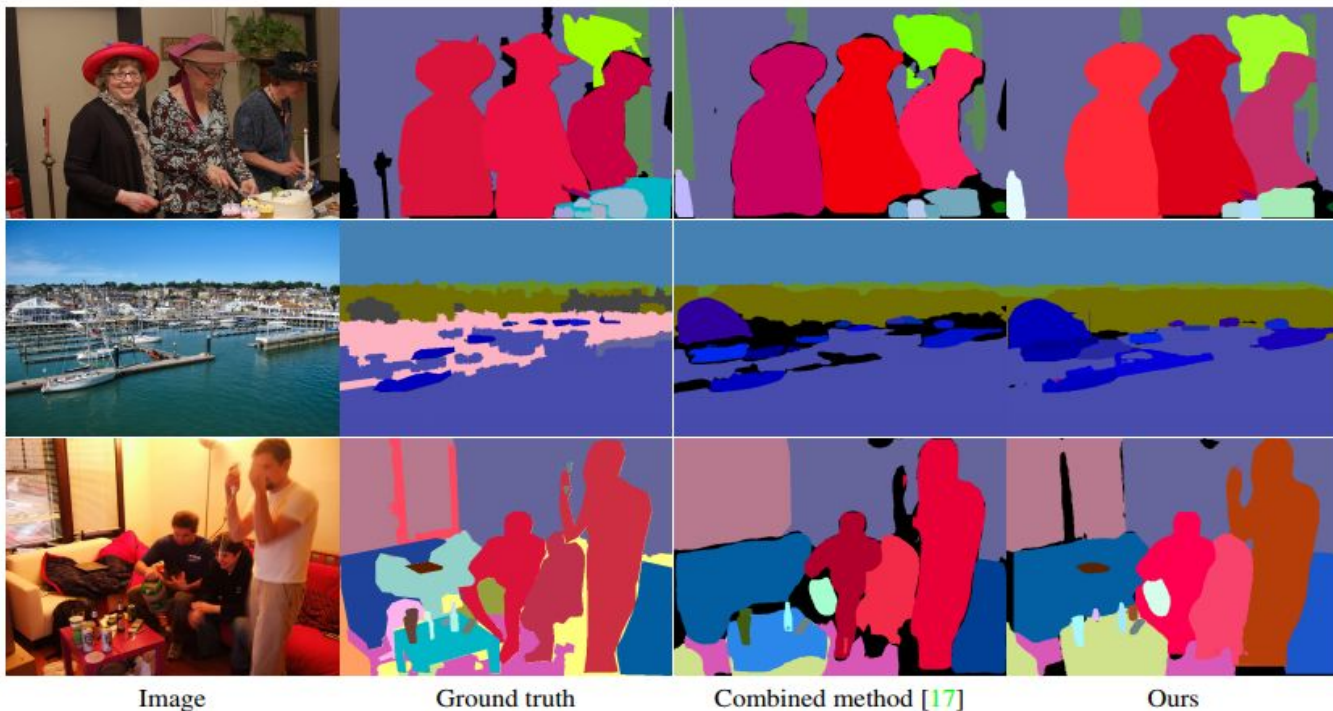
- Segment is matched if  $\text{IoU} > 0.5$ . No pixel can belong to two predicted segments.

# Panoptic segmentation: qualitative

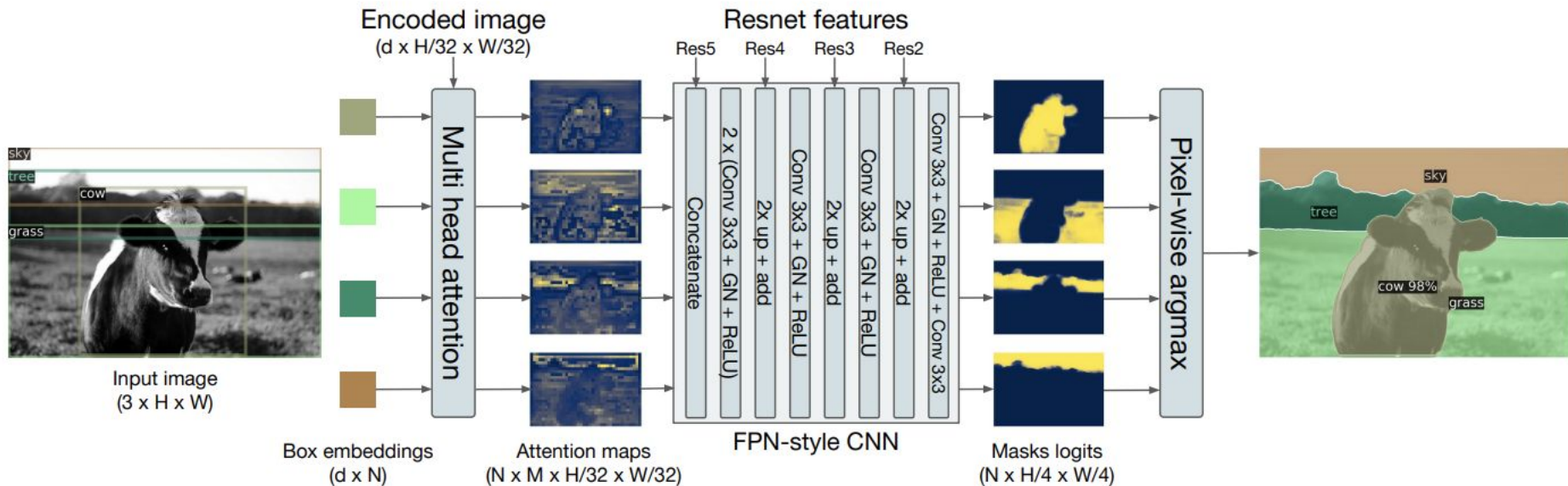




# Panoptic segmentation: qualitative



# Of course, it can be done with Transformers



- A binary mask is generated in parallel for each detected object, then the masks are merged using pixel-wise argmax.

# DETR Panoptic Segmentation – results

Model	Backbone	PQ	SQ	RQ	PQ <sup>th</sup>	SQ <sup>th</sup>	RQ <sup>th</sup>	PQ <sup>st</sup>	SQ <sup>st</sup>	RQ <sup>st</sup>	AP
PanopticFPN++	R50	42.4	79.3	51.6	49.2	82.4	58.8	32.3	74.8	40.6	37.7
UPSnet	R50	42.5	78.0	52.5	48.6	79.4	59.6	33.4	75.9	41.7	34.3
UPSnet-M	R50	43.0	79.1	52.8	48.9	79.7	59.7	34.1	78.2	42.3	34.3
PanopticFPN++	R101	44.1	79.5	53.3	<b>51.0</b>	<b>83.2</b>	60.6	33.6	74.0	42.1	<b>39.7</b>
DETR	R50	43.4	79.3	53.8	48.2	79.8	59.5	36.3	78.5	45.3	31.1
DETR-DC5	R50	44.6	79.8	55.0	49.4	80.5	60.6	<b>37.3</b>	<b>78.7</b>	<b>46.5</b>	31.9
DETR-R101	R101	<b>45.1</b>	<b>79.9</b>	<b>55.5</b>	50.5	80.9	<b>61.7</b>	37.0	78.5	46.0	33.0

# DETR Panoptic Segmentation – results

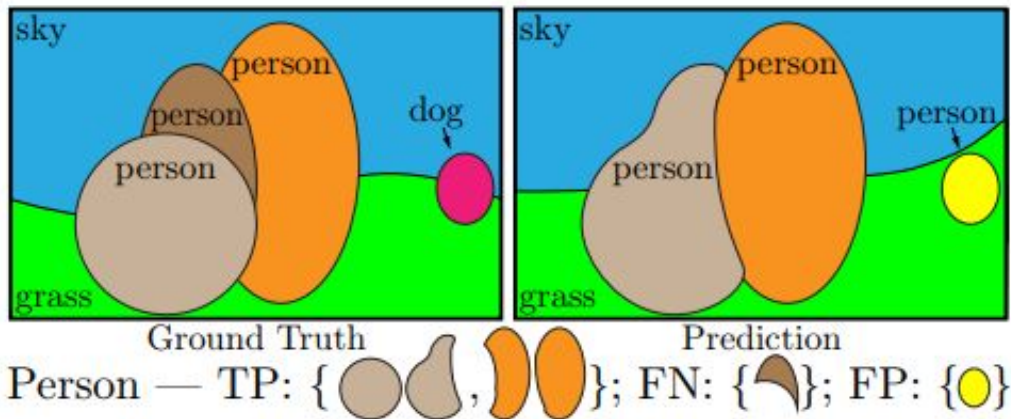


# Metrics

# Panoptic quality

TP = True positive, FN = False negative, FP = false positive

$$PQ = \underbrace{\frac{\sum_{(p,g) \in TP} \text{IoU}(p,g)}{|TP|}}_{SQ} \underbrace{\frac{|TP|}{|TP| + \frac{1}{2}|FP| + \frac{1}{2}|FN|}}_{RQ}$$

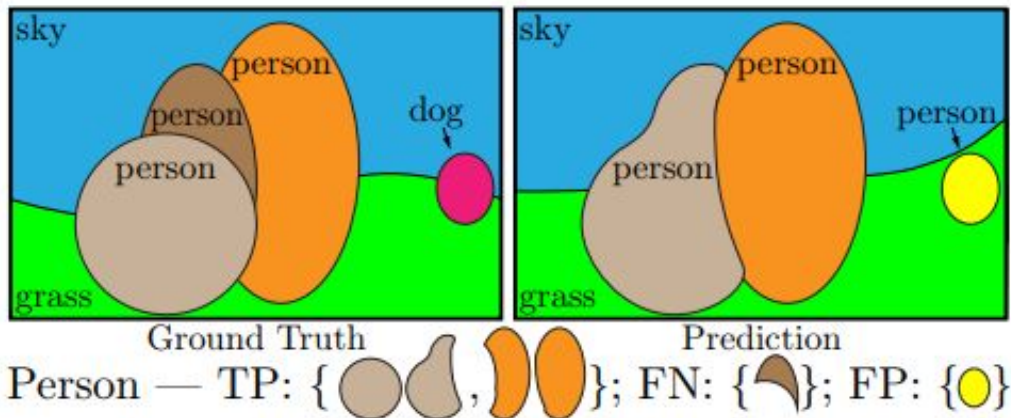


- SQ: Segmentation Quality = how close the predicted segments are to the ground truth segment (does not take into account bad predictions!)

# Panoptic quality

TP = True positive, FN = False negative, FP = false positive

$$PQ = \underbrace{\frac{\sum_{(p,g) \in TP} \text{IoU}(p,g)}{|TP|}}_{SQ} \underbrace{\frac{|TP|}{|TP| + \frac{1}{2}|FP| + \frac{1}{2}|FN|}}_{RQ}$$



- RQ: Recognition Quality = just like for detection, we want to know if we are missing any instances (FN) or we are predicting more instances (FP).

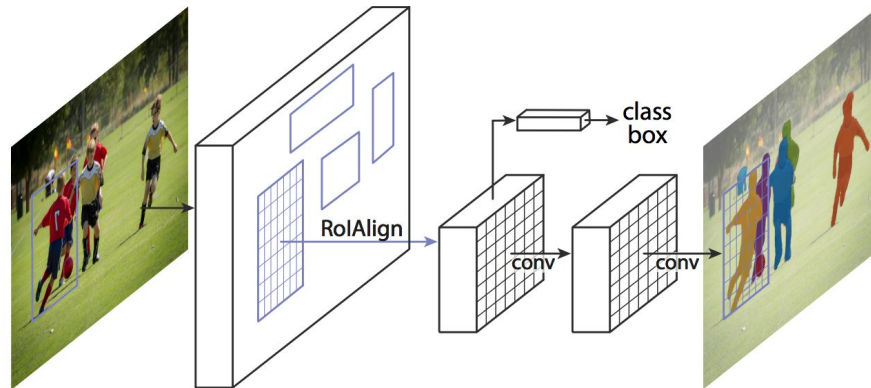
# Instance segmentation



# Object Instance Segmentation as Voting

# Sliding Window Approach

- DPM, RCNN families
- Densely enumerate box proposals + classify
- Tremendously successful paradigm, very well engineered
- SOTA methods are still based on this paradigm



# Generalized Hough Transform

Before DPM, RCNN dominance: detection-as-voting

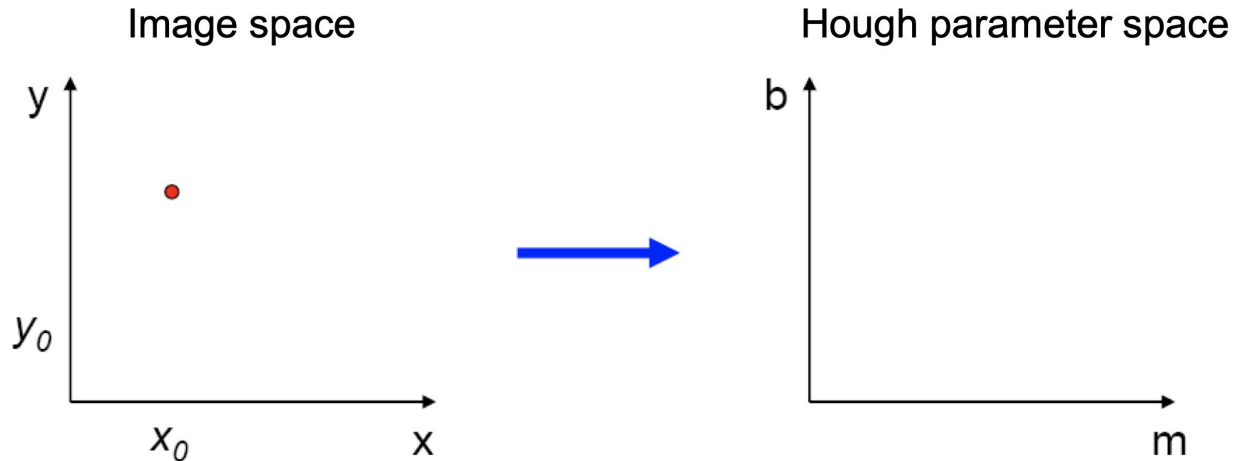


# Hough Voting

- Detect analytical shapes (e.g., lines) as peaks in the dual parametric space
- Each pixel casts a vote in this dual space
- Detect peaks and 'back-project' them to the image space

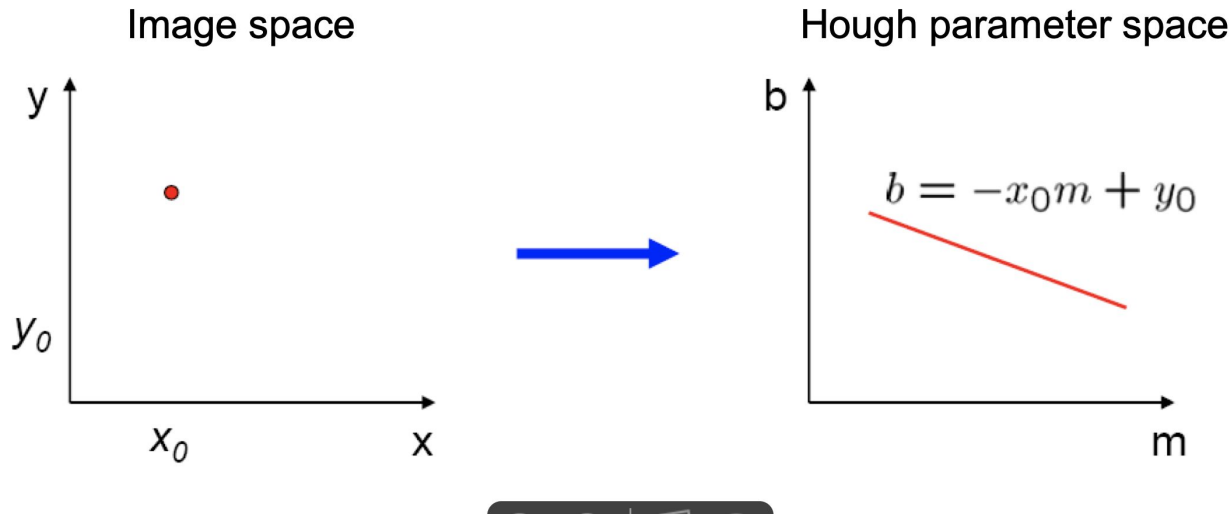
# Example: Line Detection

- Each edge point in image space casts a vote



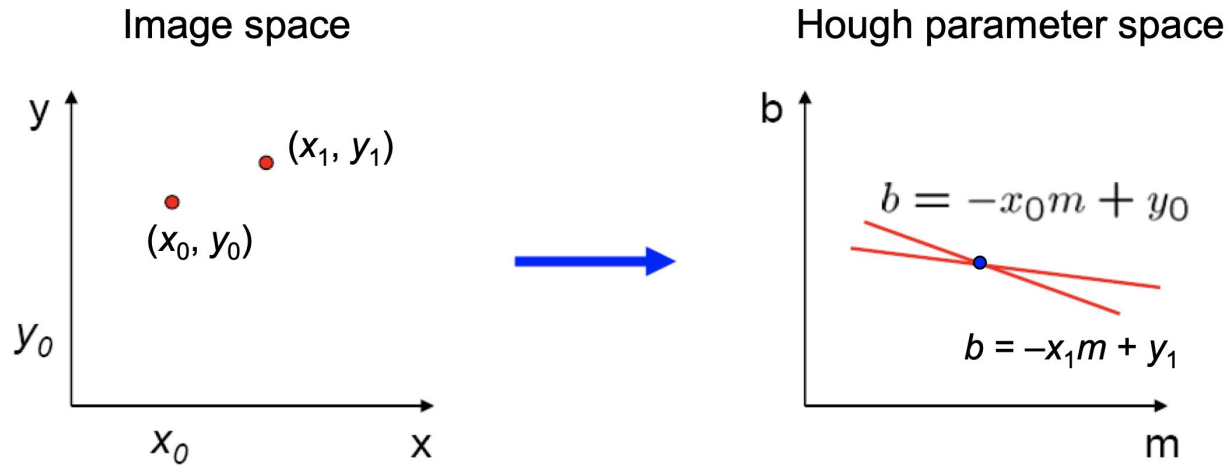
# Example: Line Detection

- Each edge point in image space casts a vote
- The vote is in the form of a line that crosses the point



# Example: Line Detection

- Accumulate votes from different points in (discretized) parameter space
- Read-out maxima (peaks) from the accumulator



# Object Detection as Voting

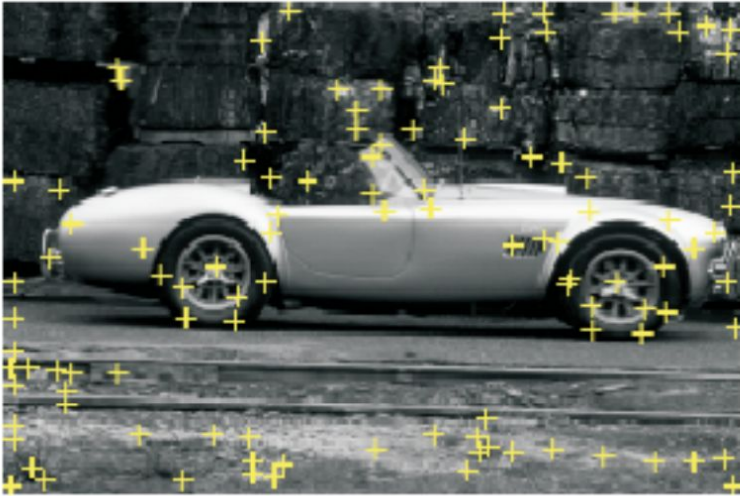
- Idea: Objects are detected as consistent configurations of the observed parts (visual words)



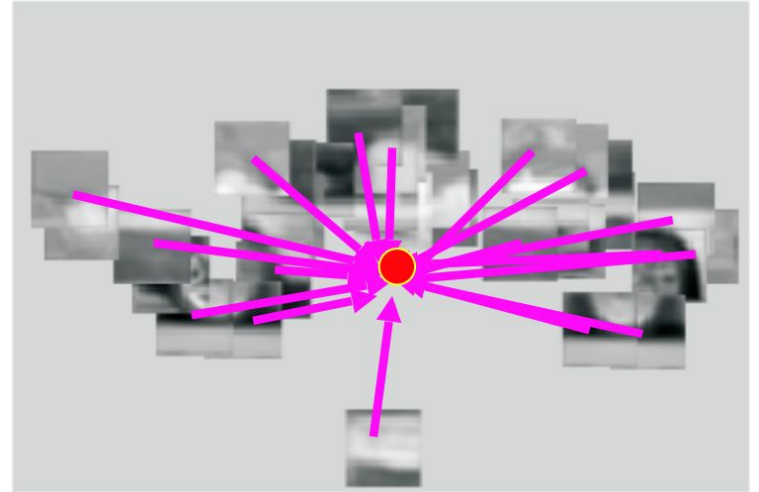


# Object Detection

- Training



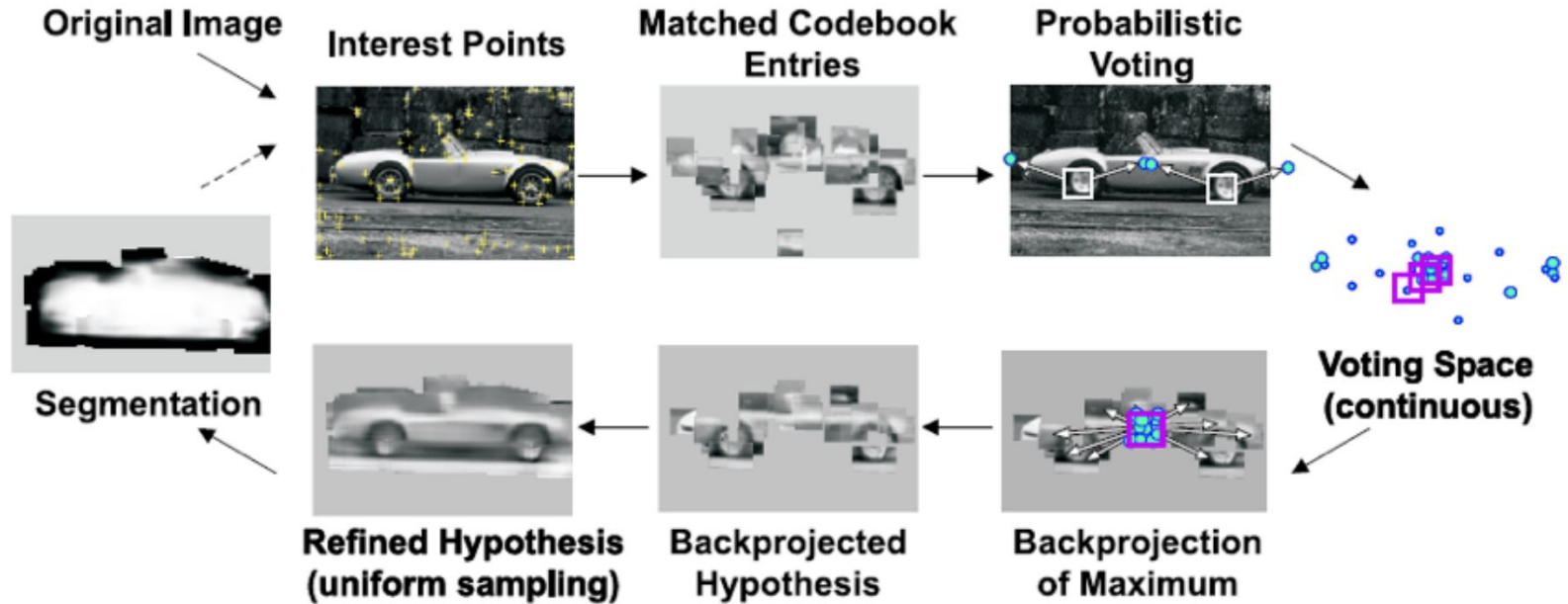
Interest point detection  
(SIFT, SURF)



Center point voting

# Object Detection

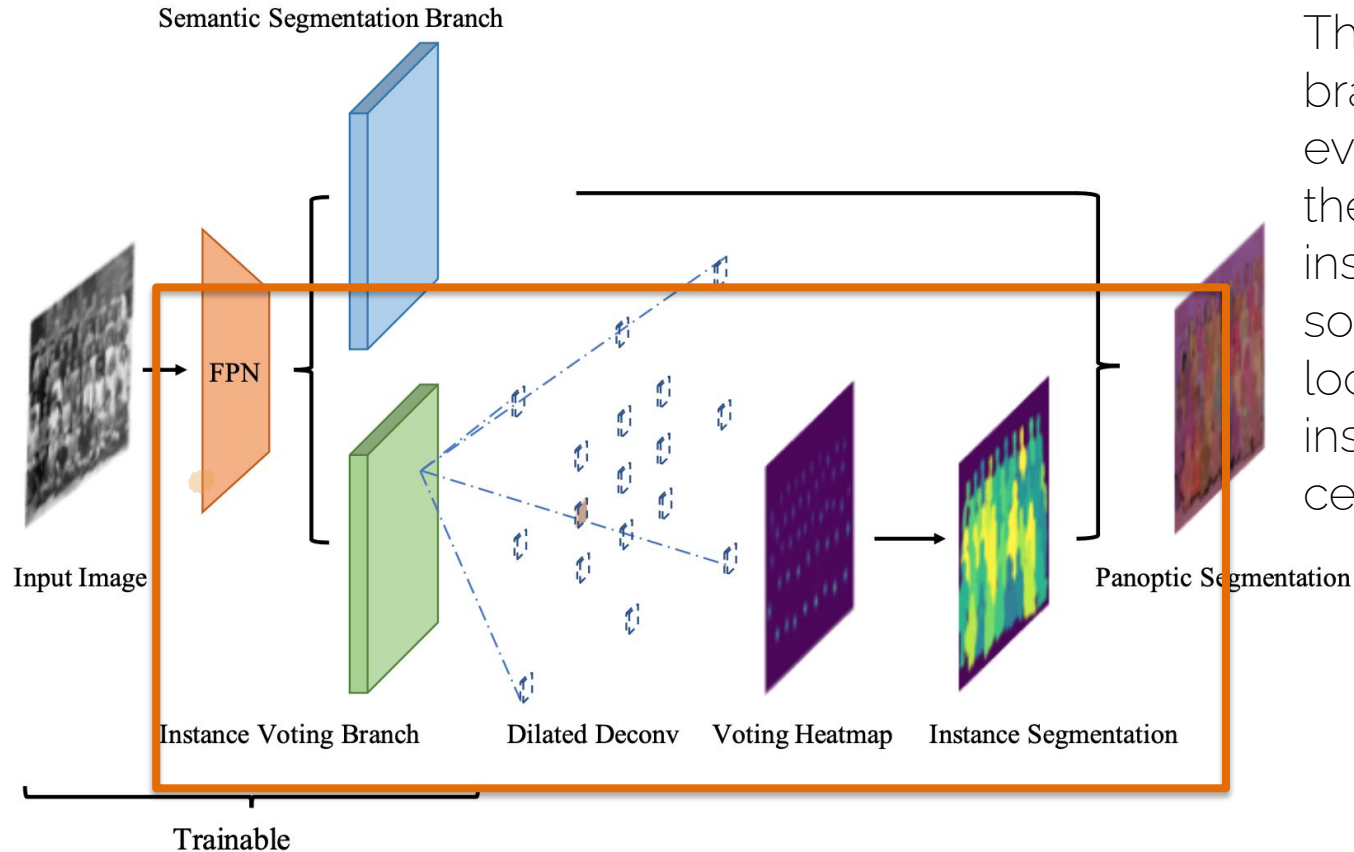
- Inference (test time)



# Back to the future

- Back to 2020...
- We can use pixel consensus voting for panoptic segmentation (CVPR 20)

# Overview



The instance voting branch predicts for every pixel whether the pixel is part of an instance mask, and if so, the relative location of the instance mask centroid.

# In a Nutshell

1. Discretize regions around each pixel.
2. Every pixel votes for a centroid (or no centroid for “stuff”) over a set of grid cells.



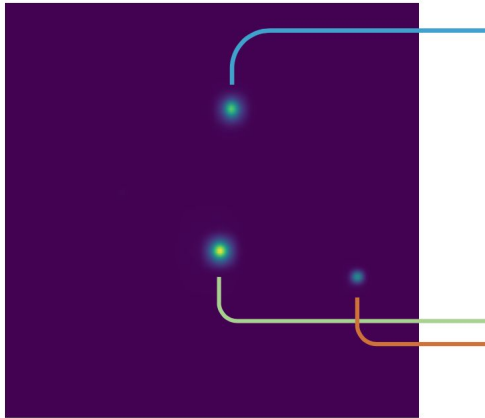
Discretization & Classification

# In a Nutshell

3. Vote aggregation probabilities at each pixel are cast to accumulator space via (dilated) transposed convolutions
4. Detect objects as 'peaks' in the accumulator space



Discretization & Classification



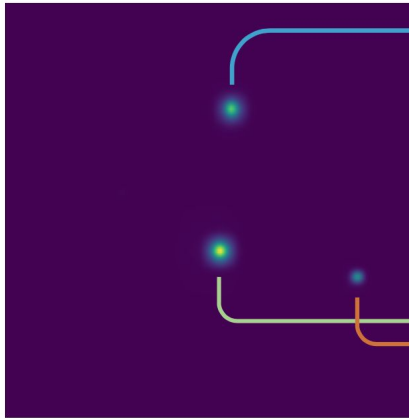
Voting as Transposed Convolution

# In a Nutshell

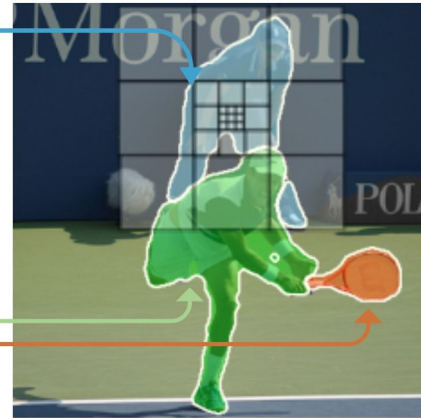
5. Back-projection of 'peaks' back to the image to get an instance masks
6. Category information provided by the parallel semantic segmentation head



Discretization & Classification



Voting as Transposed Convolution



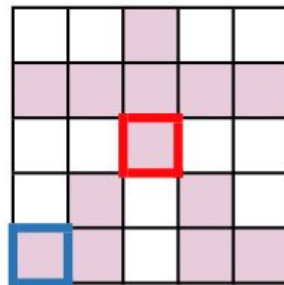
Backprojection as Filtering

# Voting Lookup Table

- Discretize region around the pixel:  $M \times M$  cells converted into  $K=17$  indices.

10	10	10	9	9	9	16	16	16
10	10	10	9	9	9	16	16	16
10	10	10	9	9	9	16	16	16
11	11	11	2	1	8	15	15	15
11	11	11	3	0	7	15	15	15
11	11	11	4	5	6	15	15	15
12	12	12	13	13	13	14	14	14
12	12	12	13	13	13	14	14	14
12	12	12	13	13	13	14	14	14

Voting filter



Instance Mask

10	10	10	9	9	9	16	16	16
10	10	10	9	9	9	16	16	16
10	10	10	9	9	9	16	16	16
11	11	11	2	1	8	15	15	15
11	11	11	3	0	7	15	15	15
11	11	11	4	5	6	15	15	15
12	12	12	13	13	13	14	14	14
12	12	12	13	13	13	14	14	14
12	12	12	13	13	13	14	14	14

Ground Truth Assignment

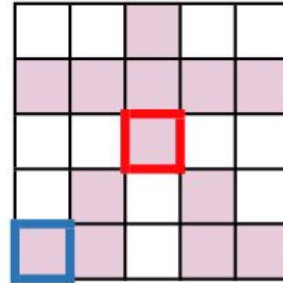


# Voting Lookup Table

- The vote should be cast to the center, which is the red pixel, which corresponds to position 16.

10	10	10	9	9	9	16	16	16
10	10	10	9	9	9	16	16	16
10	10	10	9	9	9	16	16	16
11	11	11	2	1	8	15	15	15
11	11	11	3	0	7	15	15	15
11	11	11	4	5	6	15	15	15
12	12	12	13	13	13	14	14	14
12	12	12	13	13	13	14	14	14
12	12	12	13	13	13	14	14	14

Voting filter



Instance Mask

10	10	10	9	9	9	16	16	16
10	10	10	9	9	9	16	16	16
10	10	10	9	9	9	16	16	16
11	11	11	2	1	8	15	15	15
11	11	11	3	0	7	15	15	15
11	11	11	4	5	6	15	15	15
12	12	12	13	13	13	14	14	14
12	12	12	13	13	13	14	14	14
12	12	12	13	13	13	14	14	14

Ground Truth Assignment

# Voting

- At inference, instance voting branch provides tensor of size  $[H,W,K+1]$
- Softly accumulate votes in the voting accumulator. **How?**

10	10	10	9	9	9	16	16	16
10	10	10	9	9	9	16	16	16
10	10	10	9	9	9	16	16	16
11	11	11	2	1	8	15	15	15
11	11	11	3	0	7	15	15	15
11	11	11	4	5	6	15	15	15
12	12	12	13	13	13	14	14	14
12	12	12	13	13	13	14	14	14
12	12	12	13	13	13	14	14	14

Example: for the blue pixel, we get a vote for index 16 with 0.9 probability (softmax output)

- Transfer 0.9 to cell 16 -- (dilated)  
transposed convolution
- Evenly distribute among pixels, each gets 0.1 -- average pooling

# Transposed Convolutions

- Take a single value in the input
- Multiply with a kernel and *distribute* in the output map
  - Kernel *defines* the amount of the input value that is being distributed to each of the output cells
- For the purpose of vote aggregation, however, we fix the kernel parameters to 1-hot across each channel that marks the target location.

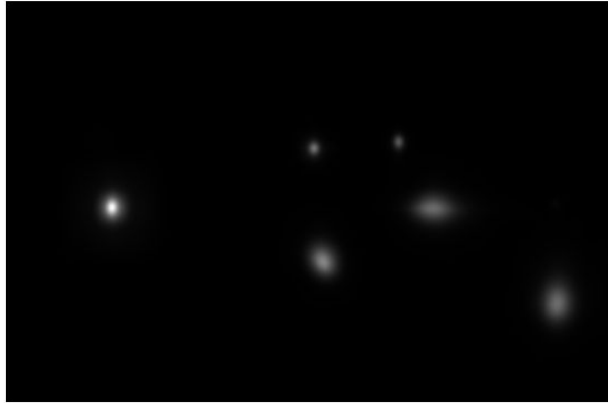
# Voting - Implementation

- Output tensor:  $[H, W, K+1]$
- Example: 9 inner, 8 outer bins,  $K=17$
- Split the output tensor to two tensors:  $[H, W, 9], [H, W, 8]$ 
  - Apply two transposed convolutions, with kernel of size  $[3, 3, 9]$ ,  $\text{stride}=1$  and  $[3, 3, 8]$ ,  $\text{stride}=3$
  - Pre-fixed kernel parameters; 1-hot across each channel that marks the target location
  - Dilation => spread votes to the **outer ring**
- Smooth votes evenly via average pooling

10	10	10	9	9	9	16	16	16
10	10	10	9	9	9	16	16	16
10	10	10	9	9	9	16	16	16
11	11	11	2	1	8	15	15	15
11	11	11	3	0	7	15	15	15
11	11	11	4	5	6	15	15	15
12	12	12	13	13	13	14	14	14
12	12	12	13	13	13	14	14	14
12	12	12	13	13	13	14	14	14

# Object Detection

- Peaks in the heatmap -- consensus detections
- Thresholding + connected components



# Object Localization

- Vote back-projection
  - For every peak, determine pixels that favor this region above all others



# Object Localization

- Idea: determine which pixels could have voted for a specific object center
  - Query filter
- Examine votes
  - Vote argmax
- Find “consensus”
  - Equality test

12	11	10	9	24
13	2	1	8	23
14	3	0	7	22
15	4	5	6	21
16	17	18	19	20

Voting filter

Bottom-left pixel should have voted for '8' if I'm the instance center!

Spatial Inversion



9	18	17	16	
11	6	5	4	15
22	7	0	3	14
23	8	1	2	13
24	9	10	11	12

Query filter

My center is at pixel 8!



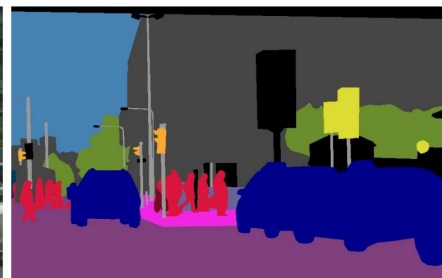


# Fine-grained Scene Interpretation

- Individual objects, surfaces (things and stuff)
- Mobile robots
  - Reason about the drivability of surfaces.
  - The type of objects and obstacles.
  - The intent of other agents in the vicinity.



(a) image



(b) semantic segmentation



(c) instance segmentation



(d) panoptic segmentation

# Instance segmentation