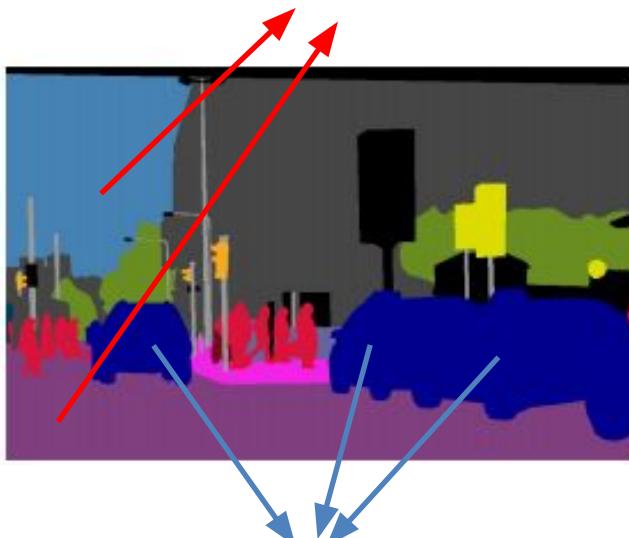


Instance segmentation

Semantic segmentation

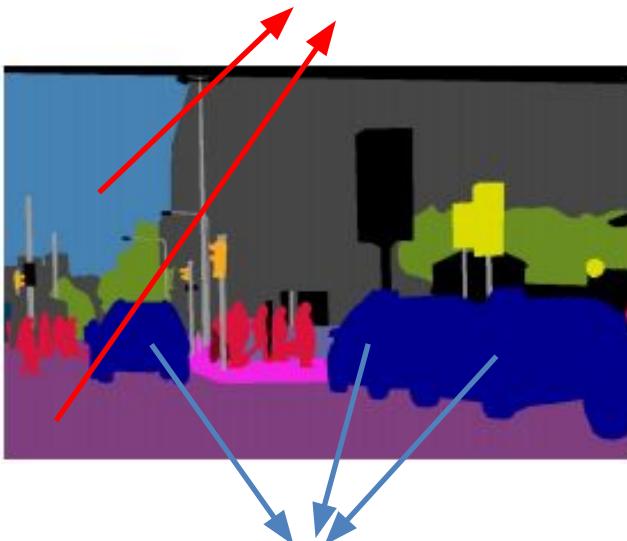
Label every pixel, including the background (sky, grass, road)



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

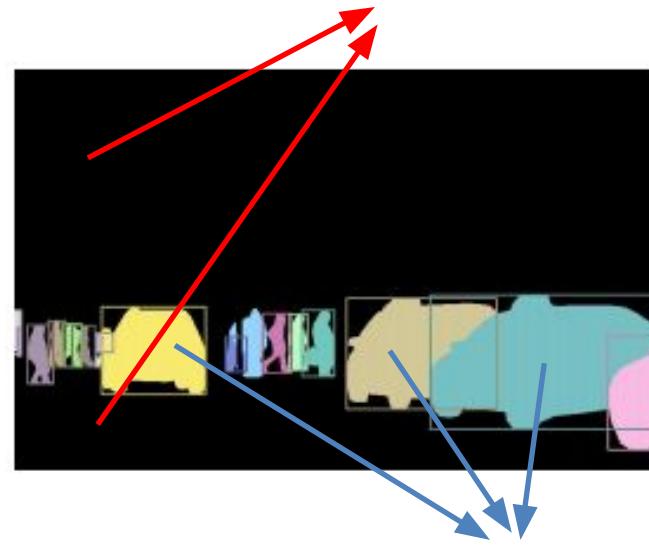
Instance segmentation

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

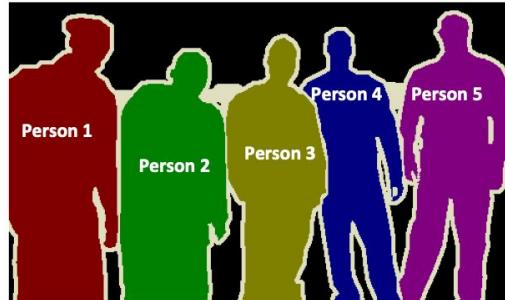
FCN-based



VS.

1. Semantic segmentation

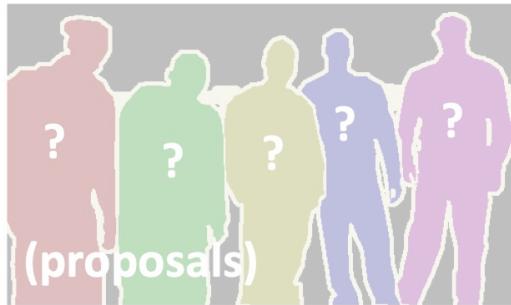
2. Find instances



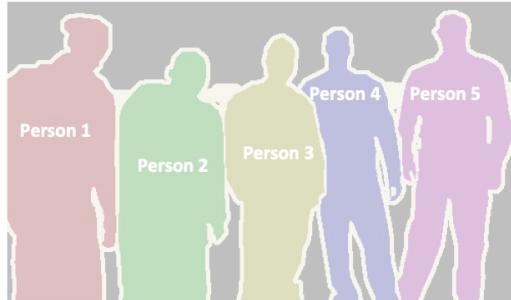
Instance segmentation methods

Proposal-based

1. Proposals

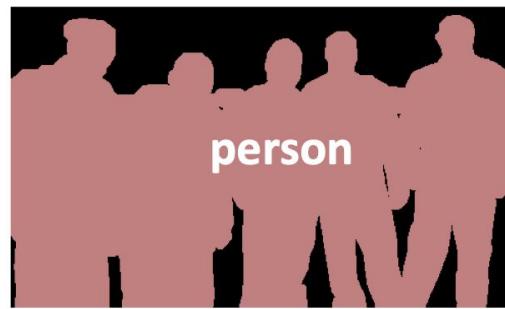


2. Assign a class

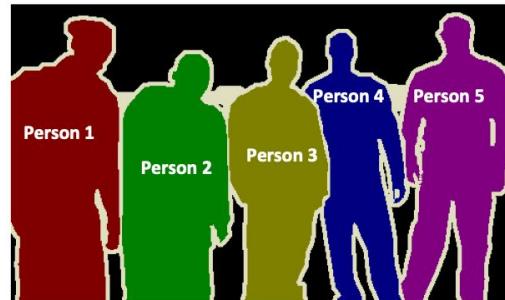


FCN-based

1. Semantic segmentation

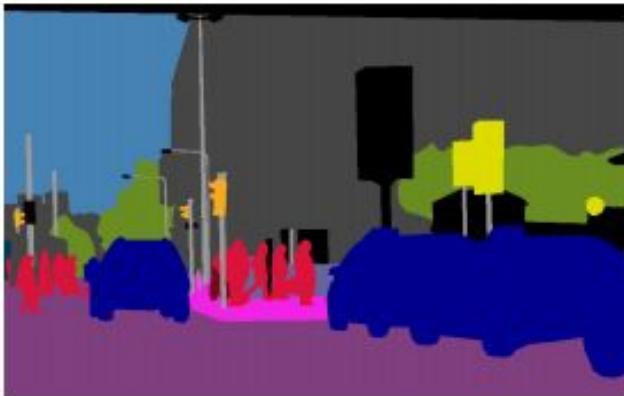


vs.



2. Find instances

FCN-based methods

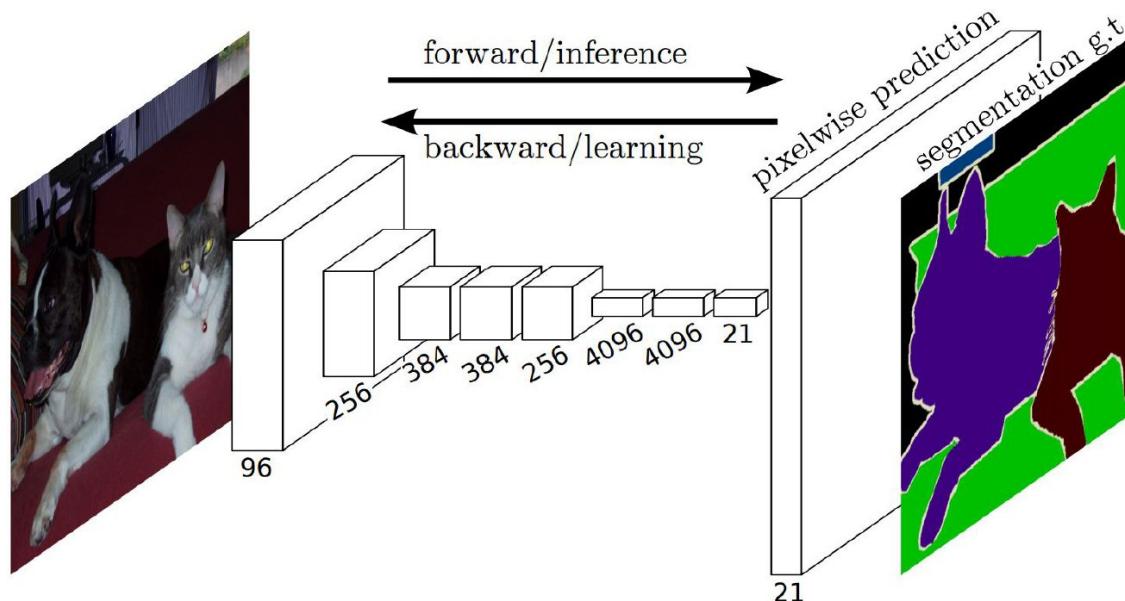


A semantic map...

We already know how to obtain this!

Why FCN-based?

- Fully Convolutional Networks for Semantic Segmentation

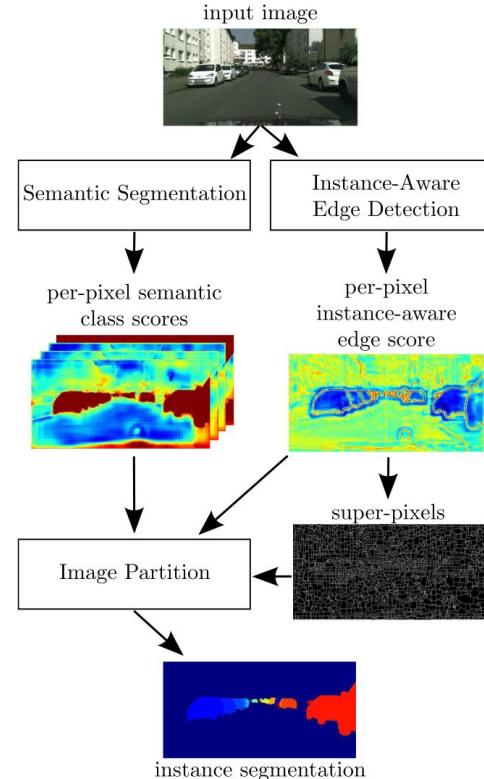


Long, Shelhamer, Darrell - Fully Convolutional Networks for Semantic Segmentation, CVPR 2015, PAMI 2016

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

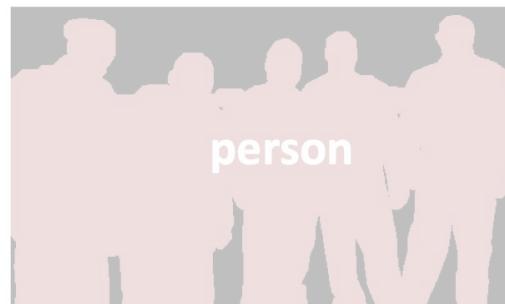
Proposal-based



1. Proposals

2. Assign a class

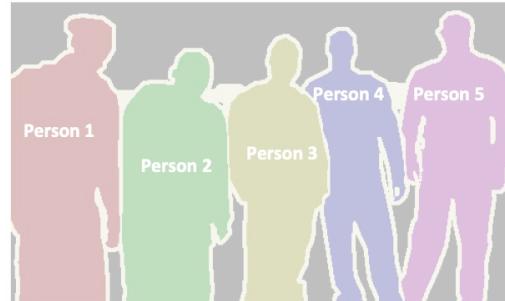
FCN-based



1. Semantic segmentation

2. Find instances

VS.



Proposal-based methods

Bounding boxes....

We already know how to obtain those!

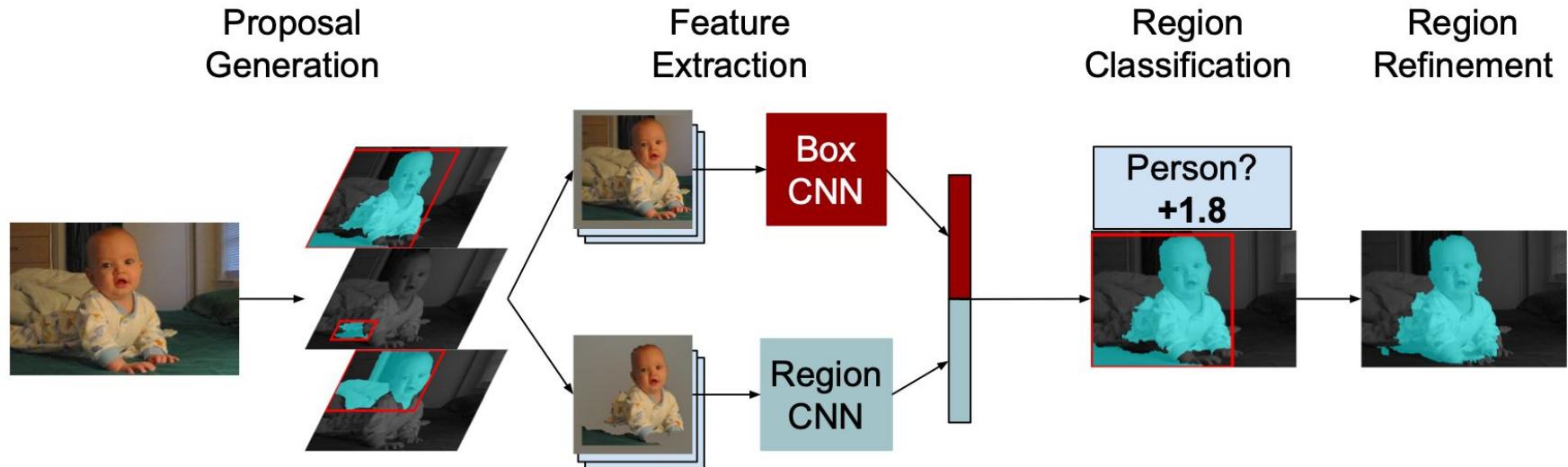


Proposal-based methods

- 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

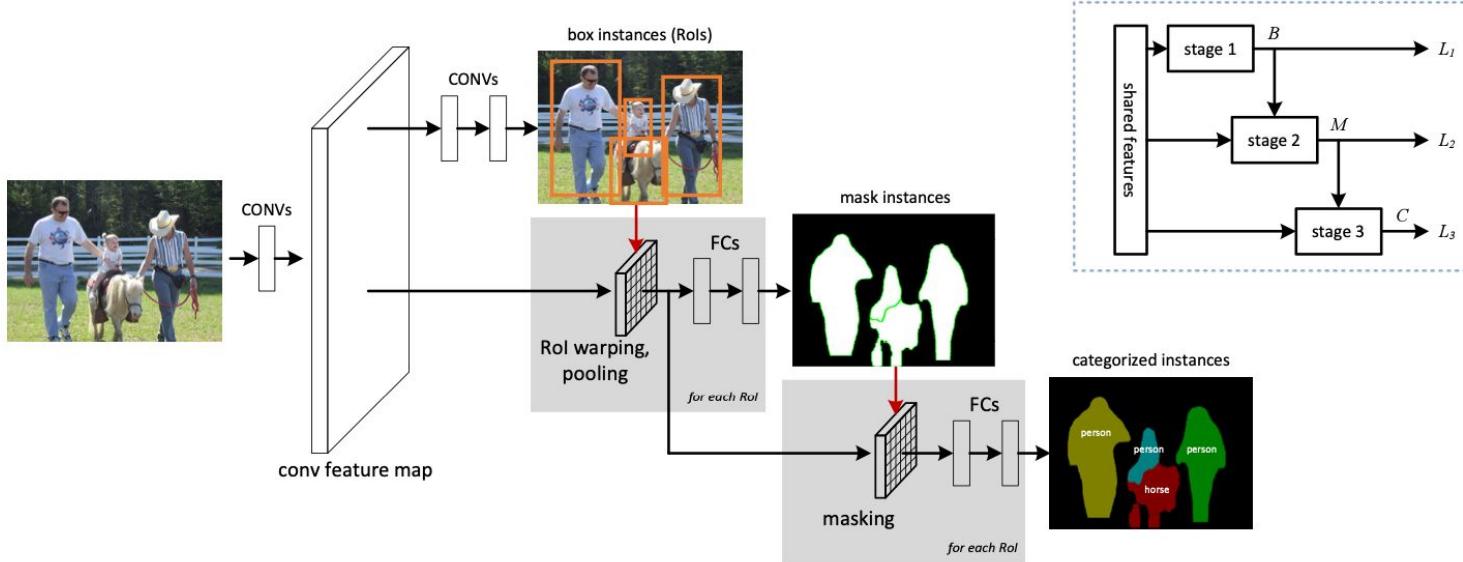
SDS

- SDS: Simultaneous Detection and Segmentation

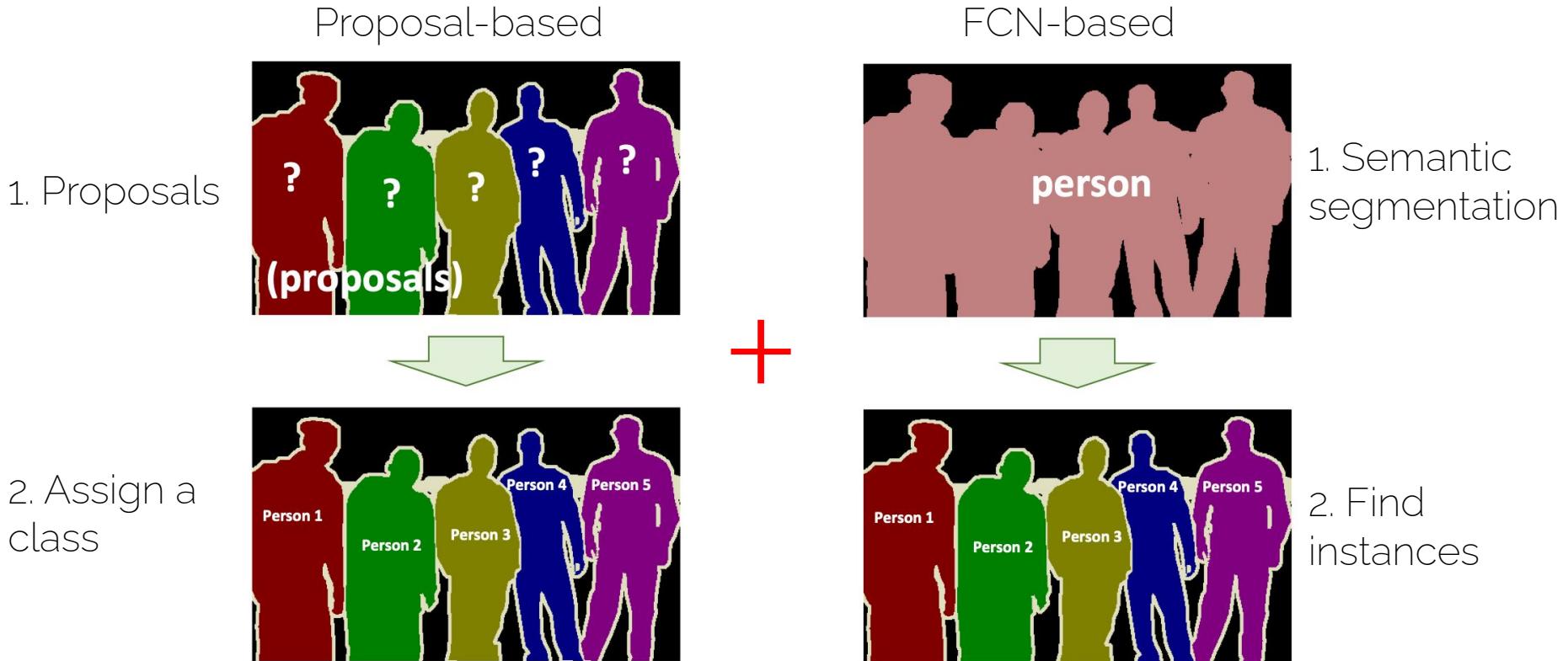


MNC

- MNC: Multi-task network cascades



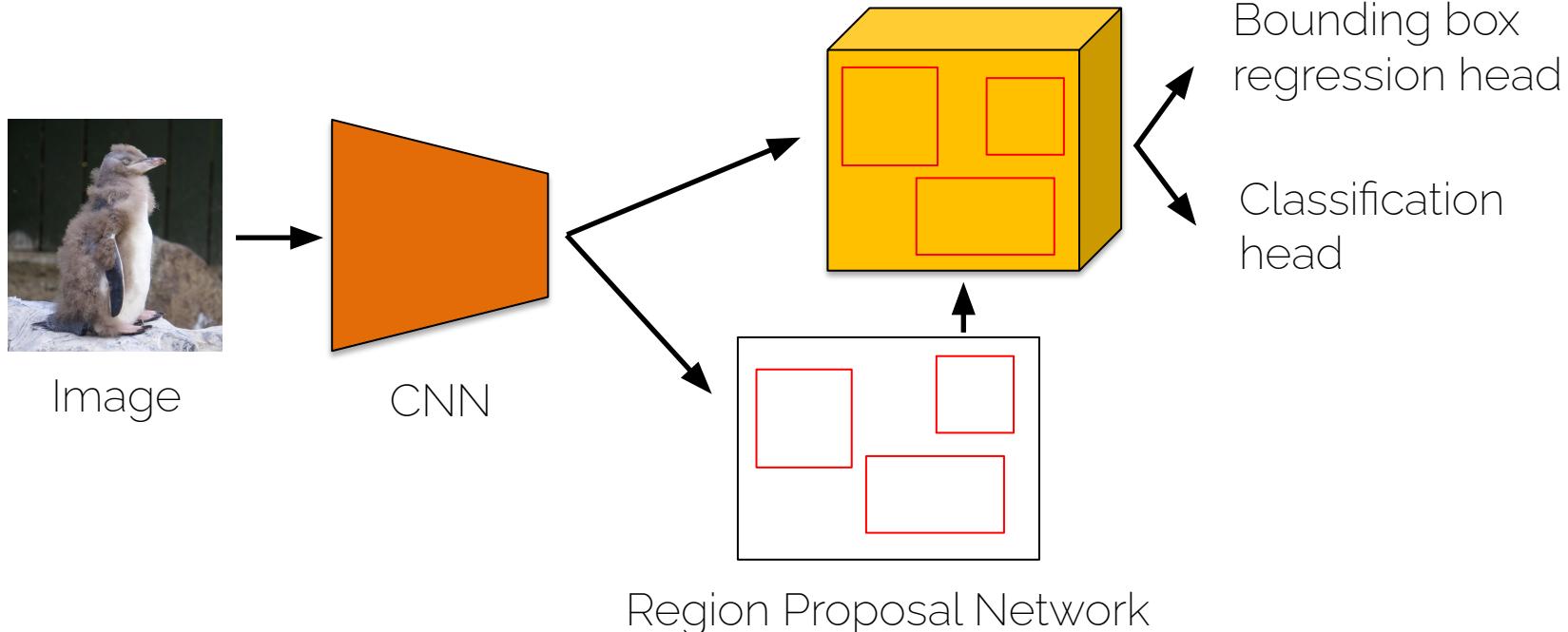
IS: the best of both worlds



Mask R-CNN

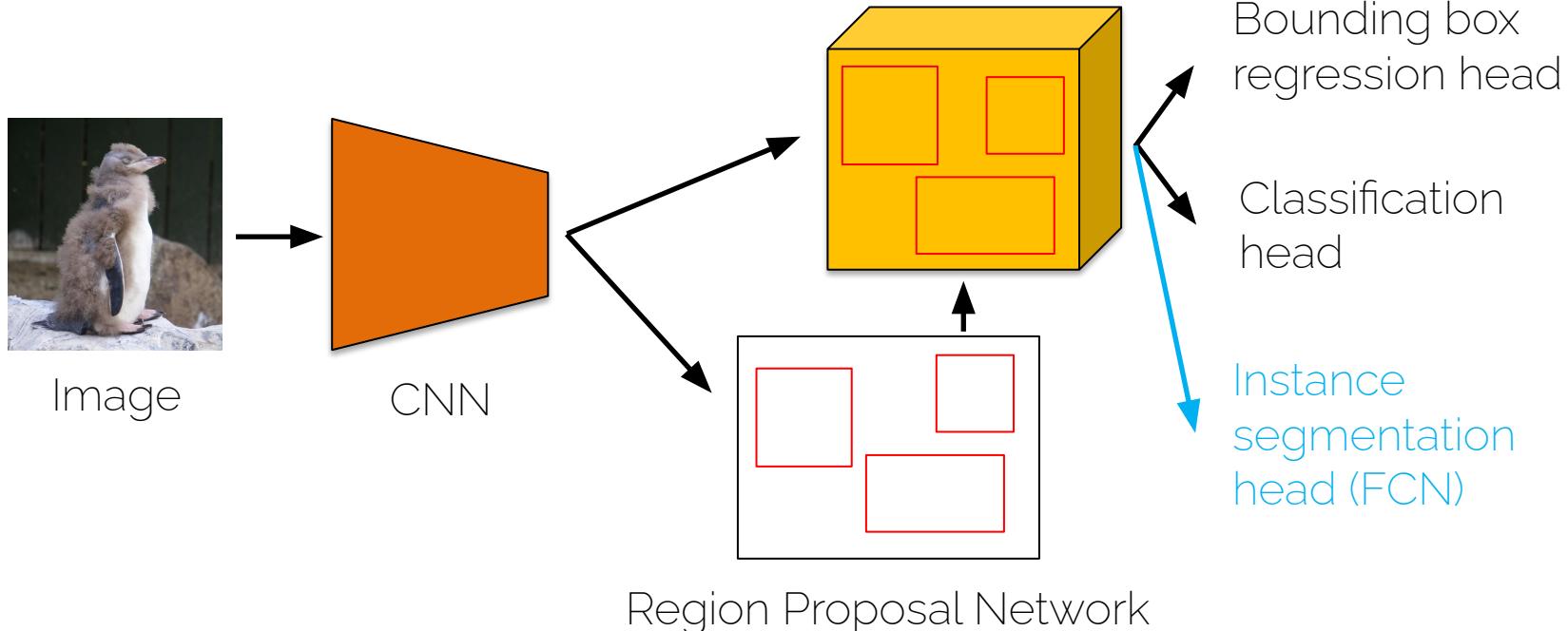
What is Mask-RCNN?

- Starting from the Faster R-CNN architecture



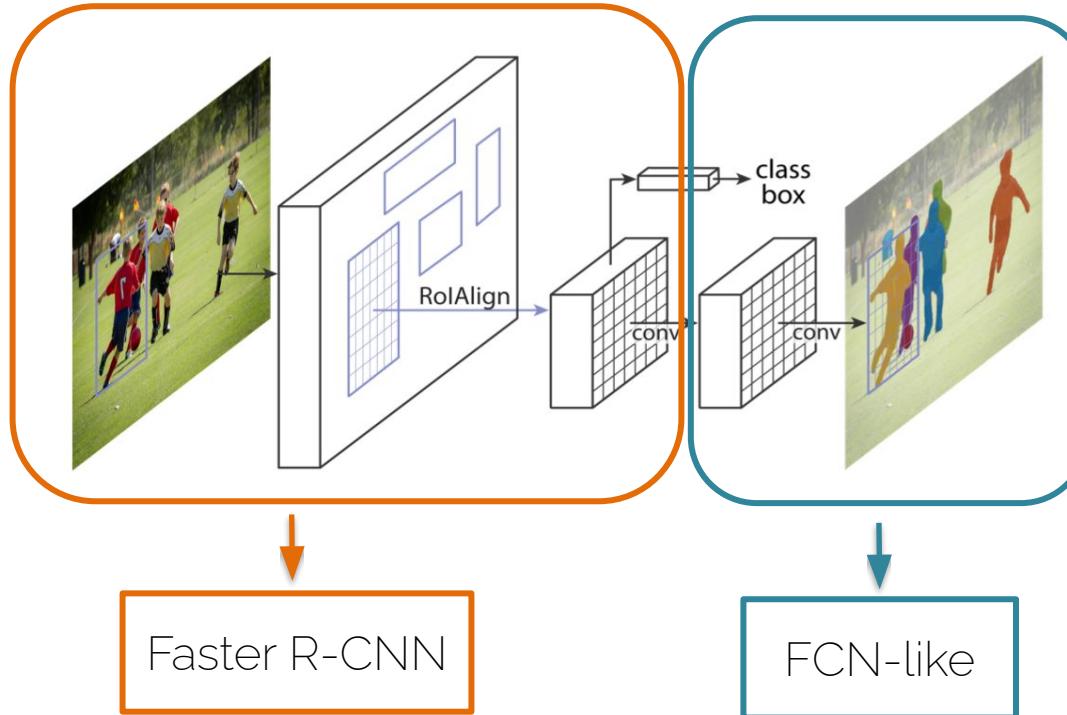
What is Mask-RCNN?

- Faster R-CNN + FCN for segmentation



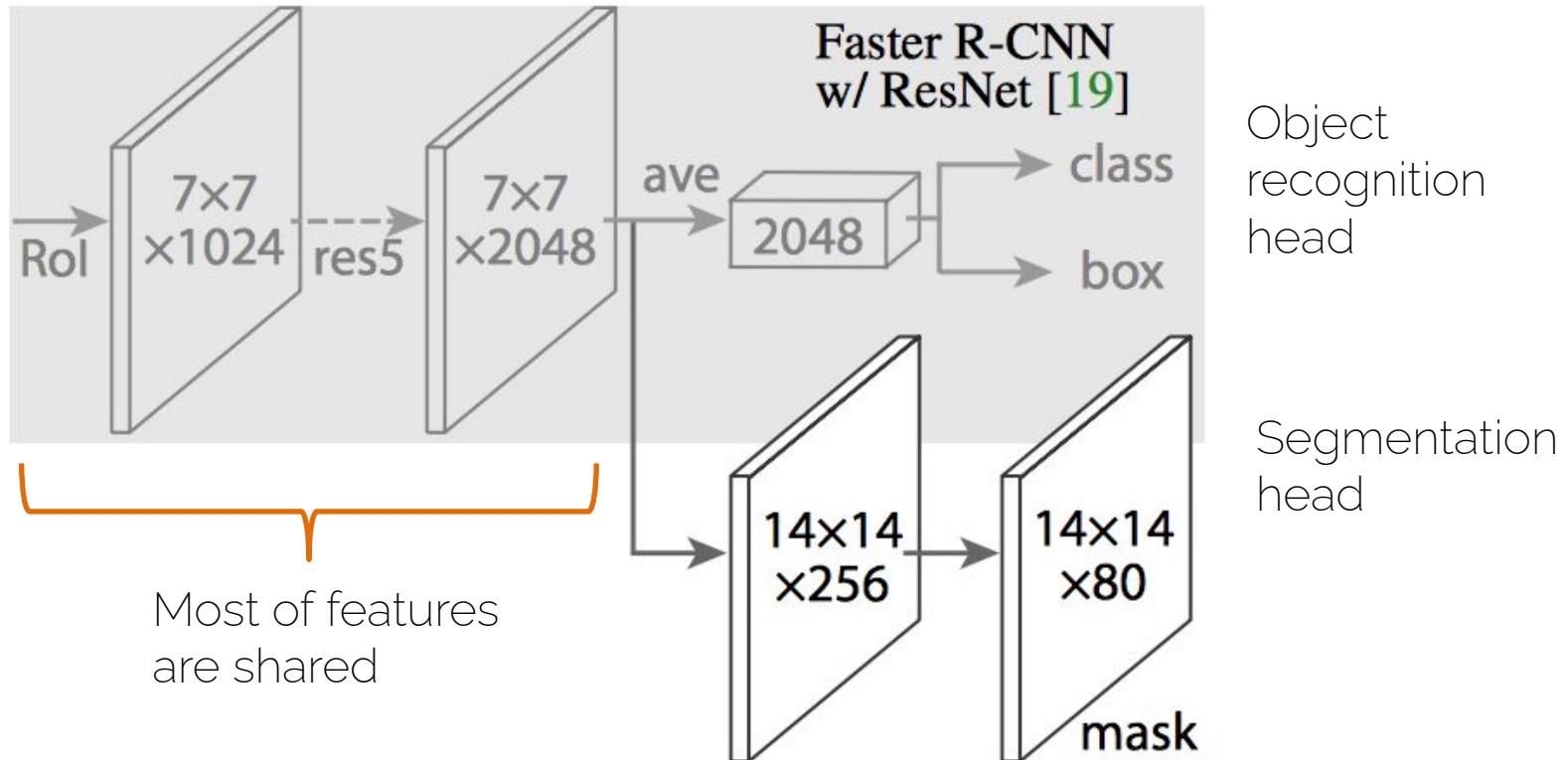
What is Mask-RCNN?

- Faster R-CNN + FCN for segmentation

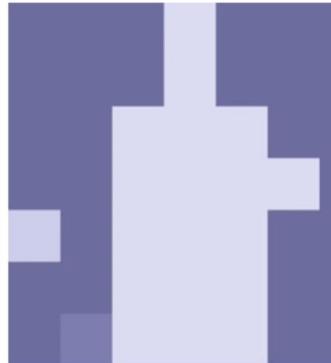


Mask loss =
binary cross
entropy per pixel
for the k semantic
classes

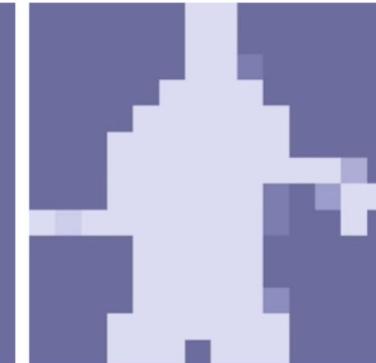
Mask R-CNN



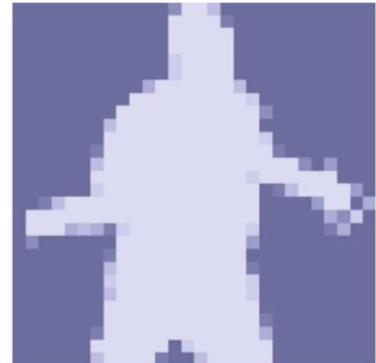
Resolution vs. computation



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



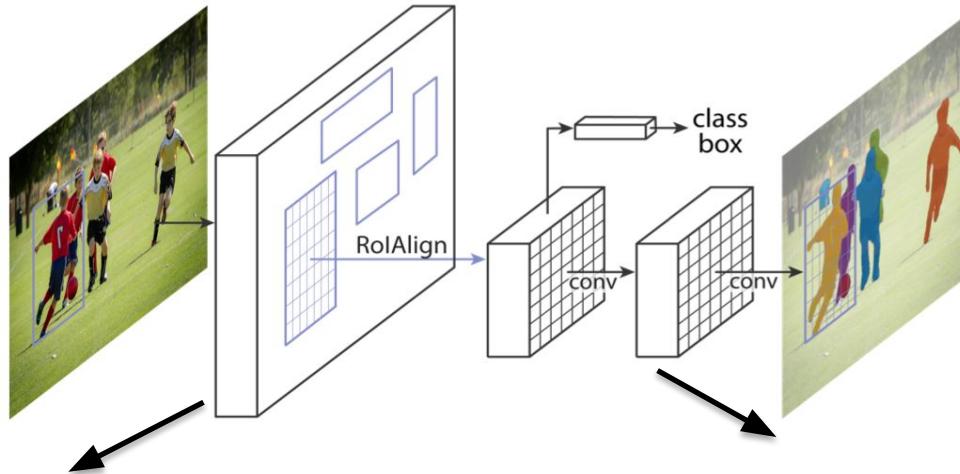
Translation invariance: wherever the penguin is in the image, I still want to have "penguin" as my classification output

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?

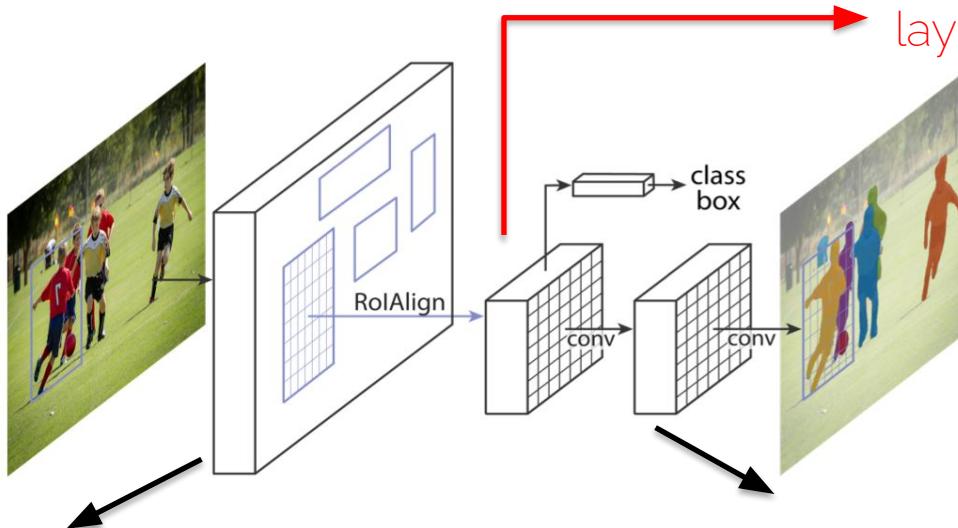


Features extraction = convolutional
layers \square equivariant

Segmentation head is a fully
convolutional network \square equivariant

Mask-RCNN: operations

- What operations are equivariant?



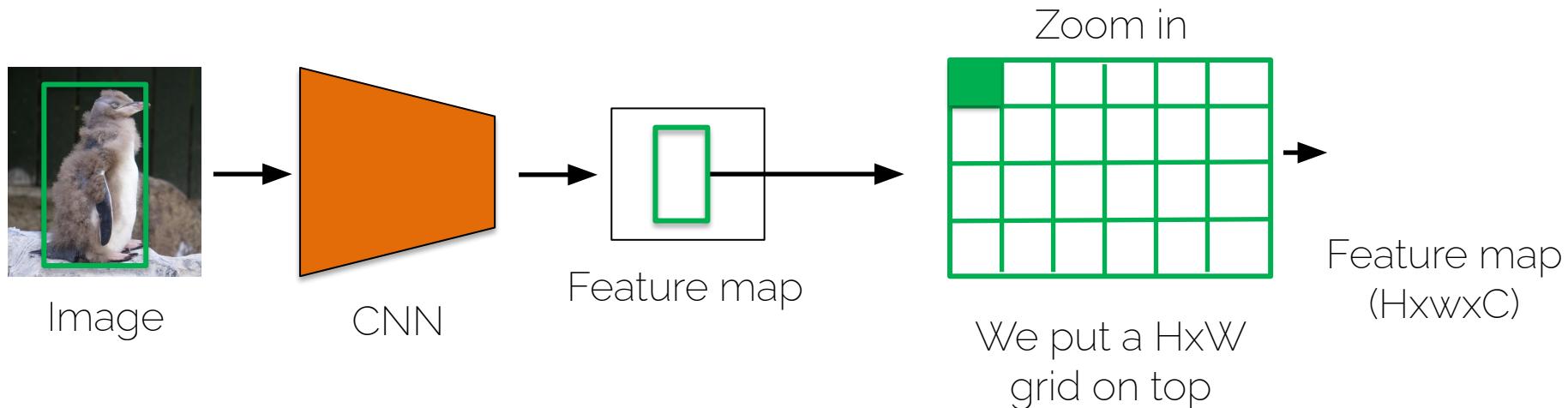
Features extraction = convolutional
layers □ equivariant

Fully connected layers
and global pooling
layers give invariance!

Segmentation head is a fully
convolutional network □ equivariant

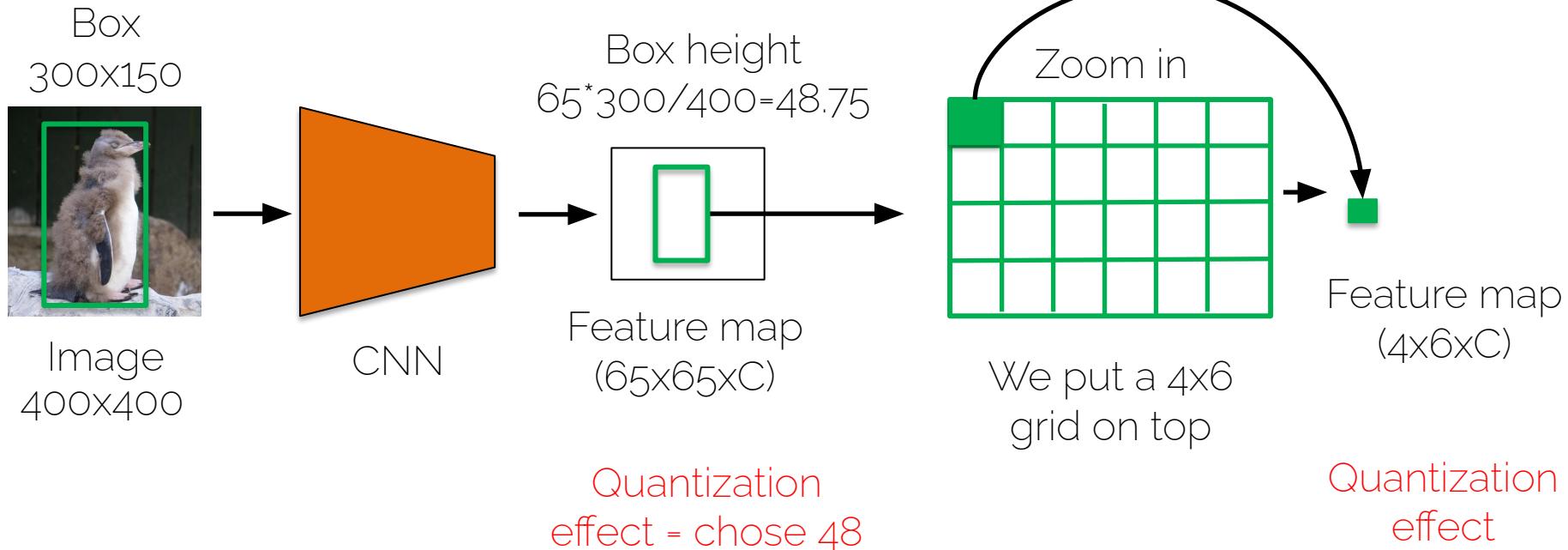
Recall: RoI pooling

- Region of Interest Pooling: for every proposal



Recall: RoI pooling

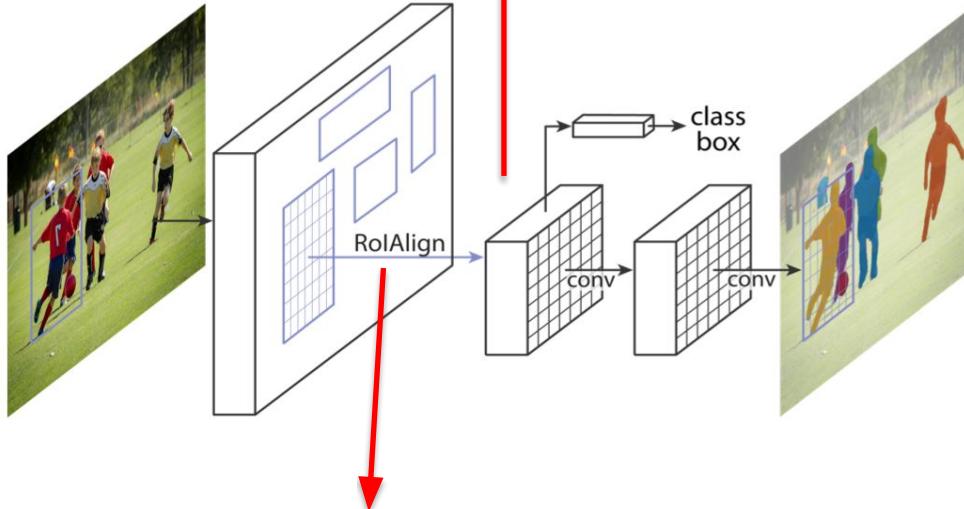
- Let us look at sizes



Mask-RCNN: operations

- Make all operations equivariant

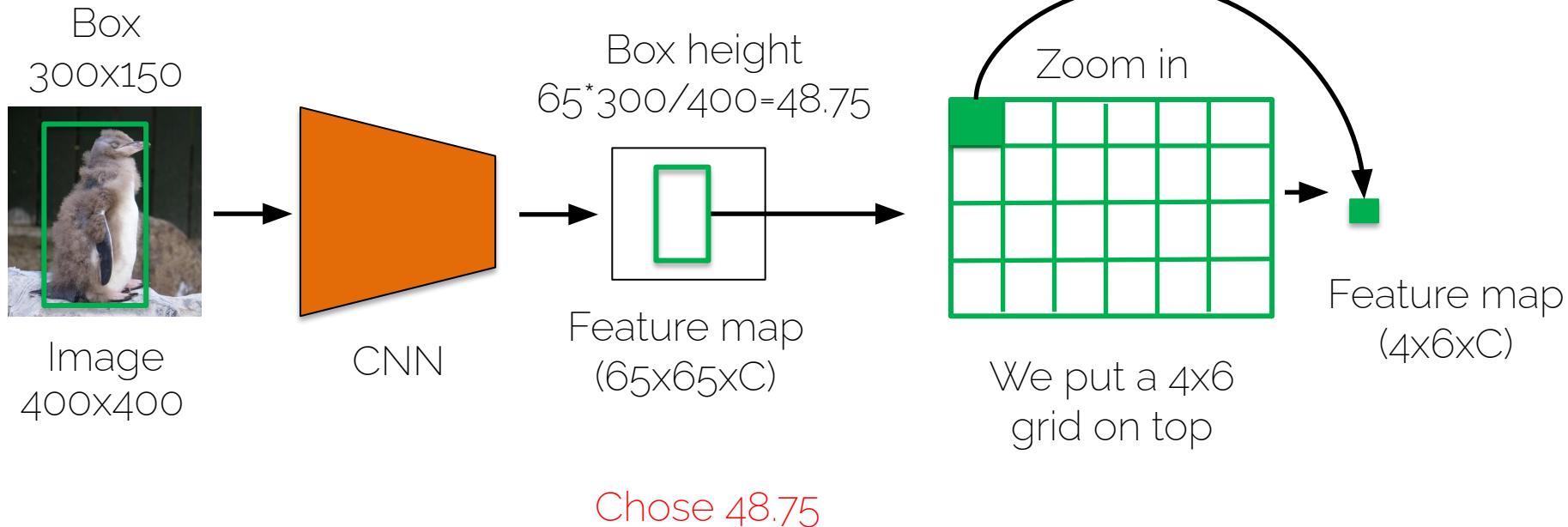
Fully connected layers
and global pooling
layers give invariance!



Exchange RoI pooling by an equivariant operation = RoI Align

RoIAlign

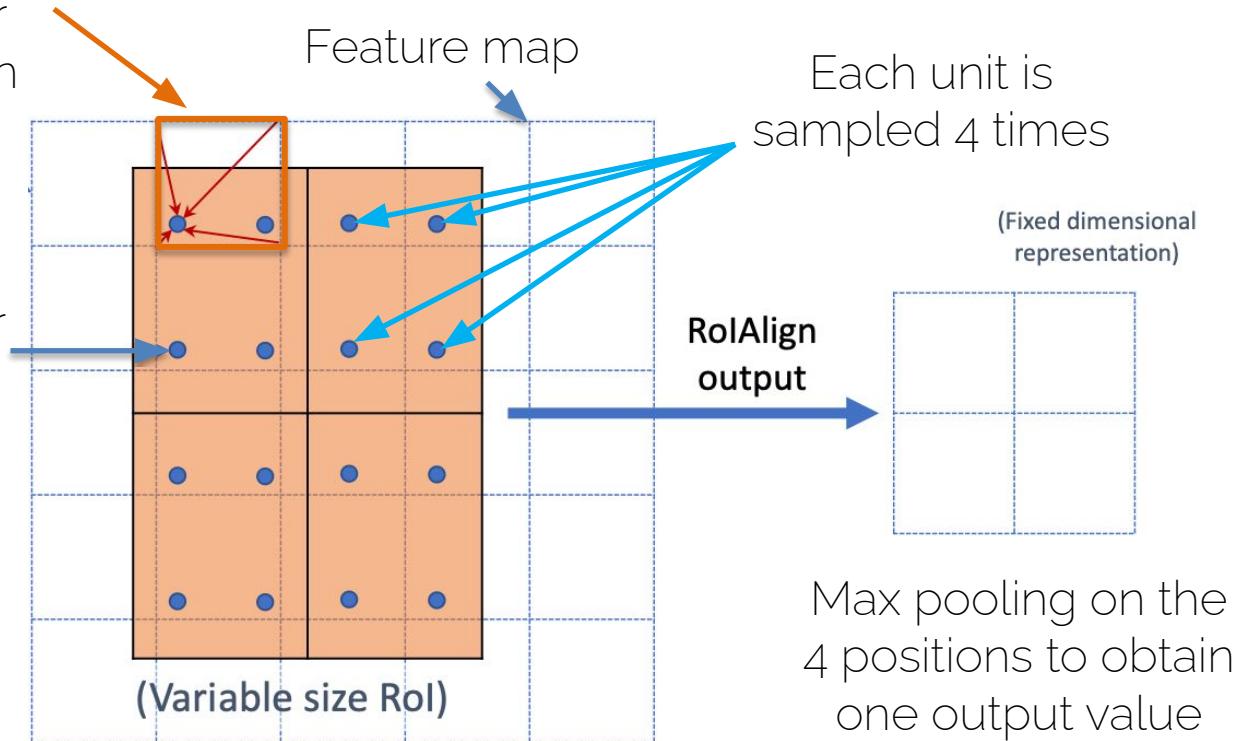
- Erase quantization effects



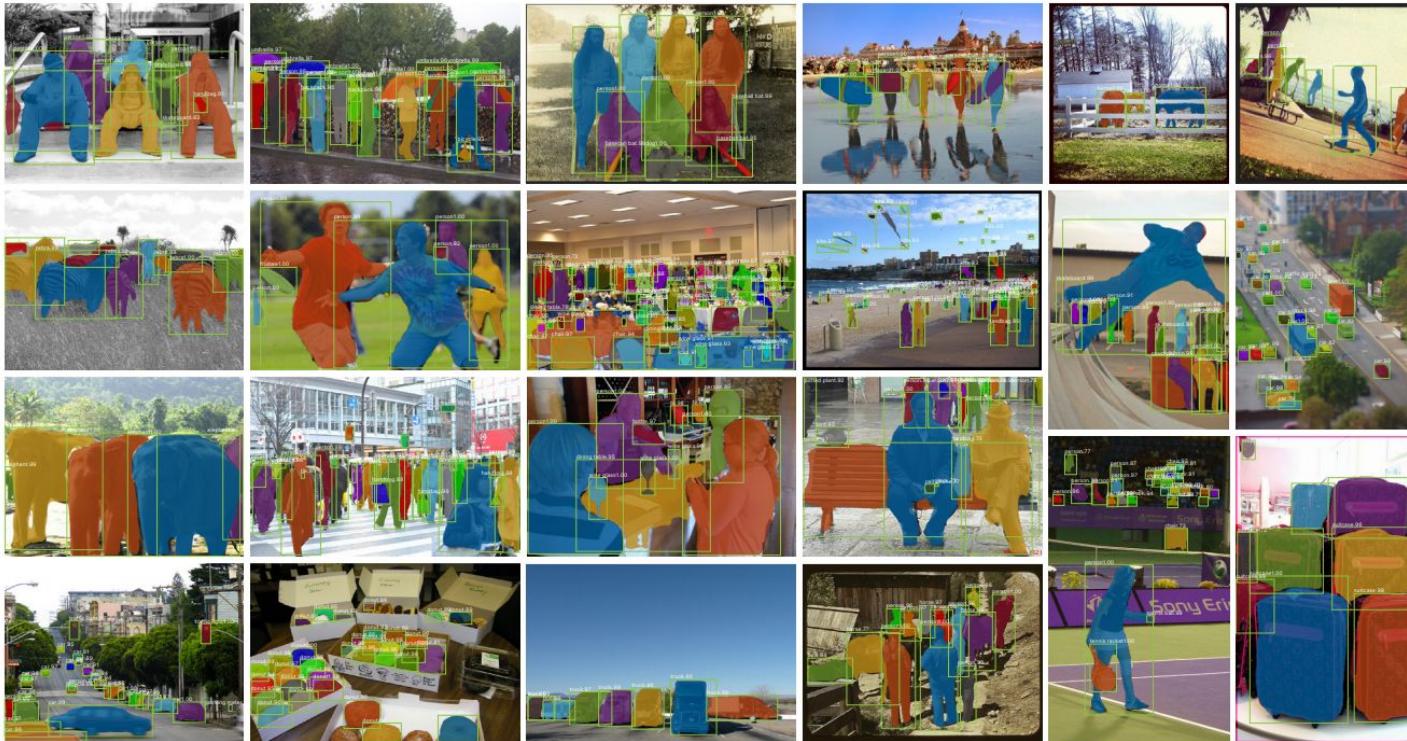
ROIAlign

To obtain the value
use bilinear
interpolation

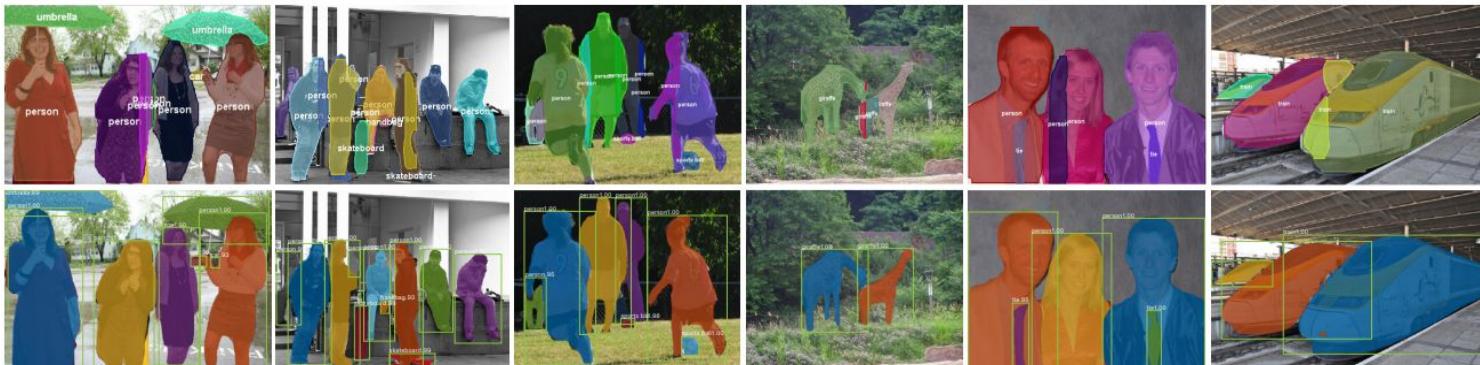
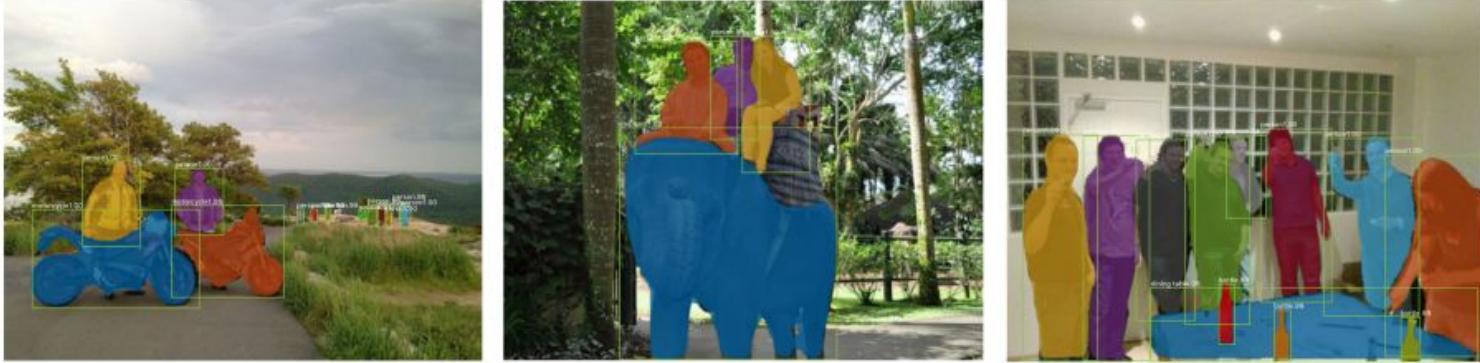
Grid points for
bilinear
interpolation



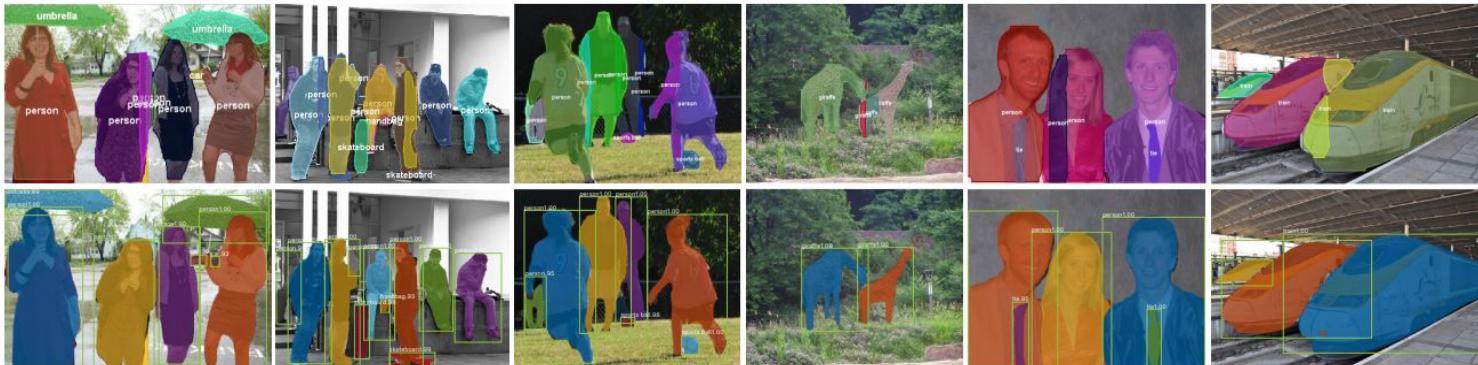
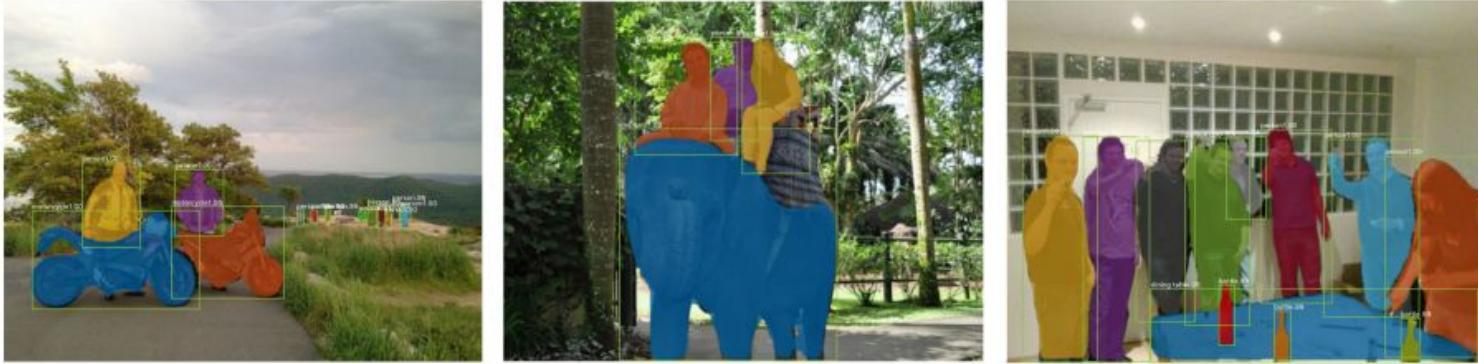
Mask R-CNN: qualitative results



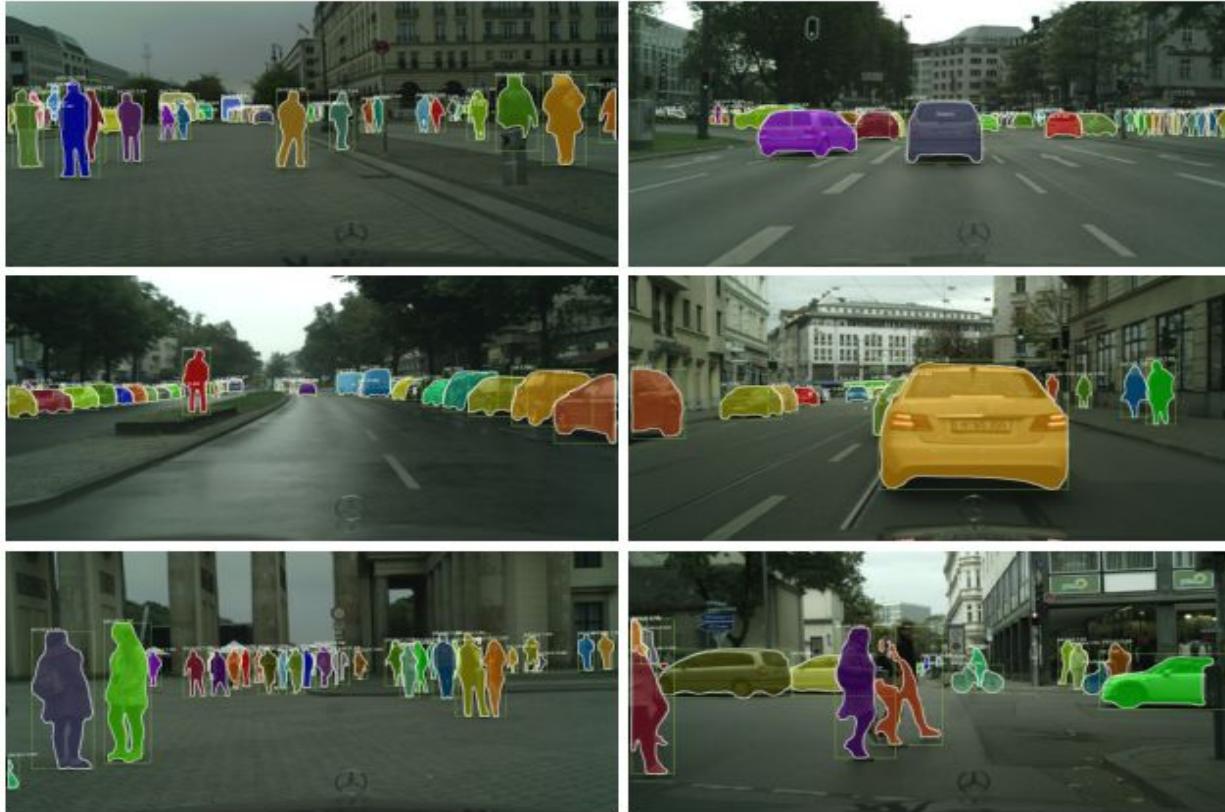
Mask R-CNN: qualitative results



Mask R-CNN: qualitative results



Mask R-CNN: qualitative results



Mask R-CNN: extended for joints

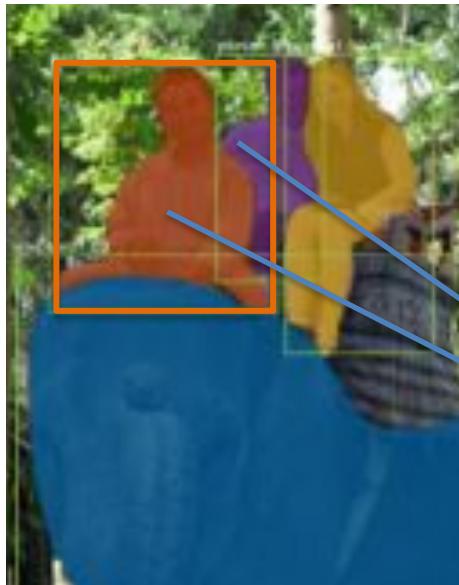


Model a keypoint's location as a one-hot mask, and adopt Mask R-CNN to predict K masks (which are in the end only 1 pixel), one for each of K keypoint types (e.g., left shoulder, right elbow). This demonstrates the flexibility of Mask R-CNN.

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

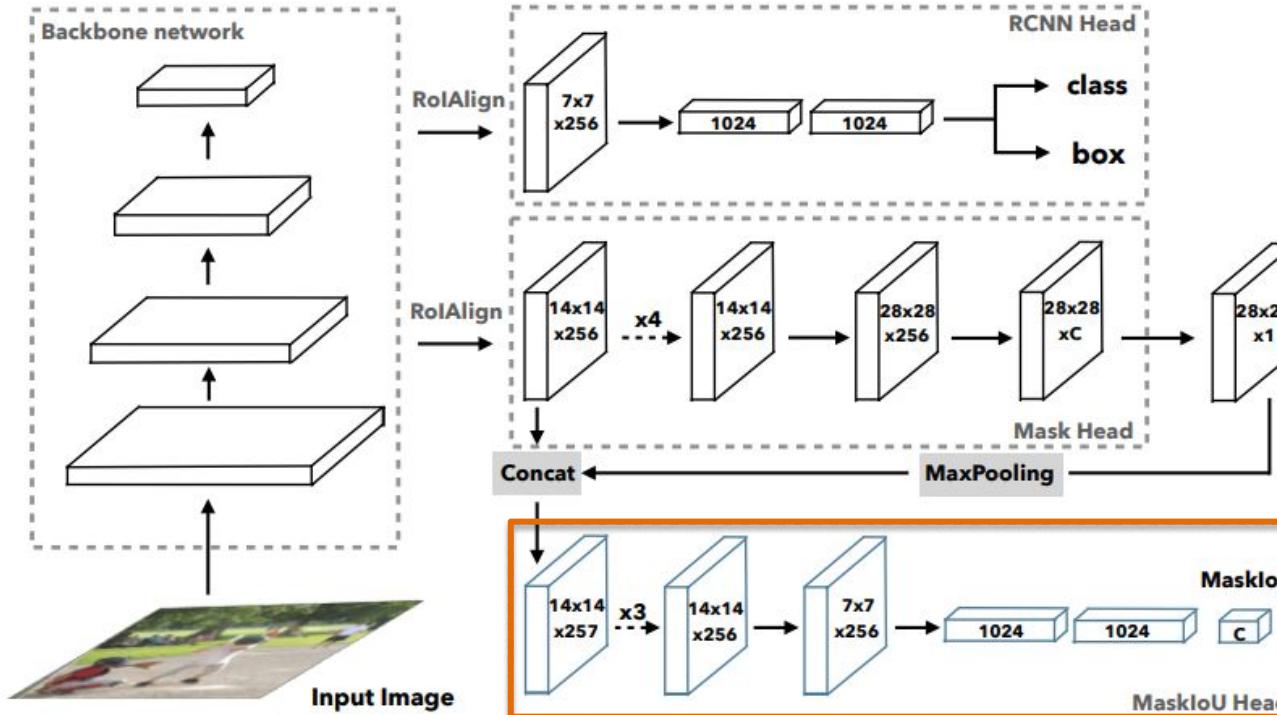
The only way the "instance" is evaluated is through the box loss



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



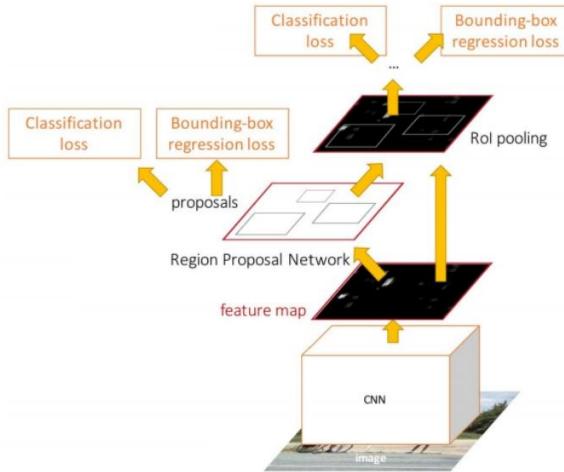
Typically, Mask scoring R-CNN gives lower confidence scores than Mask R-CNN, which corresponds to masks not being perfect ($\text{IoU} < 1.0$).

This tiny modification achieves SOTA results.

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

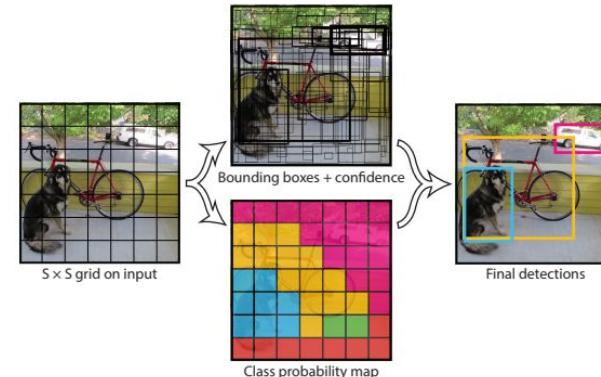
One-stage vs two-stage detectors

Faster R-CNN



Slower, but has higher performance

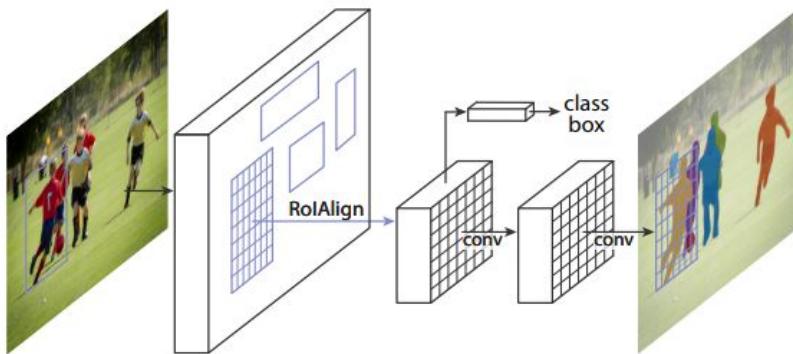
YOLO



Faster, but has lower performance

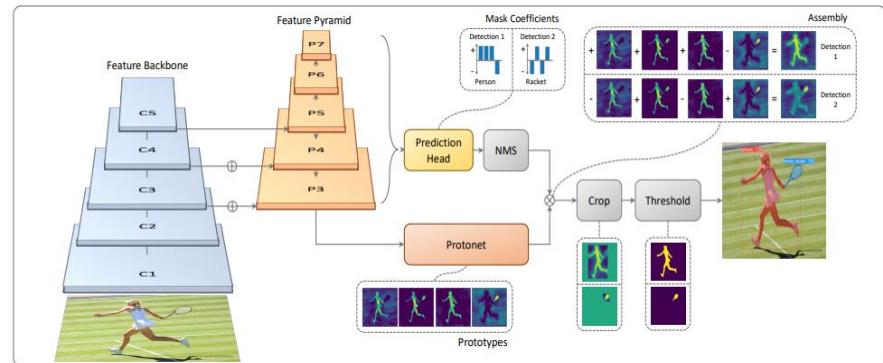
One-stage vs two-stage instance segmenters

Mask R-CNN



Slower, but has higher performance

YOLOACT



Faster, but has lower performance

YOLO with masks?

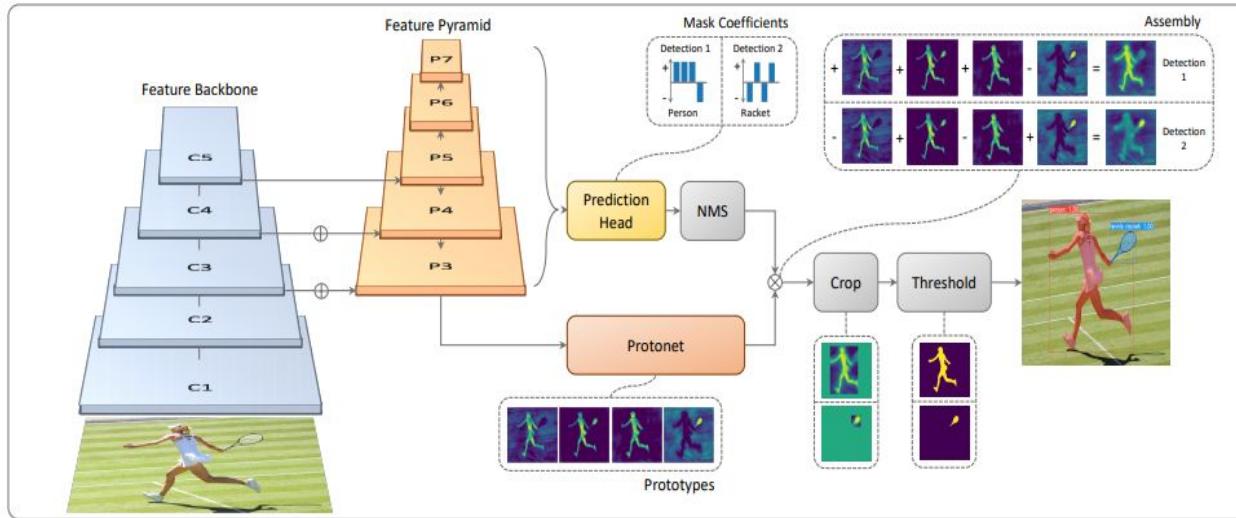
"Boxes are stupid anyway though, I'm probably a true believer in masks except I can't get YOLO to learn them."

– Joseph Redmon, YOLOv3

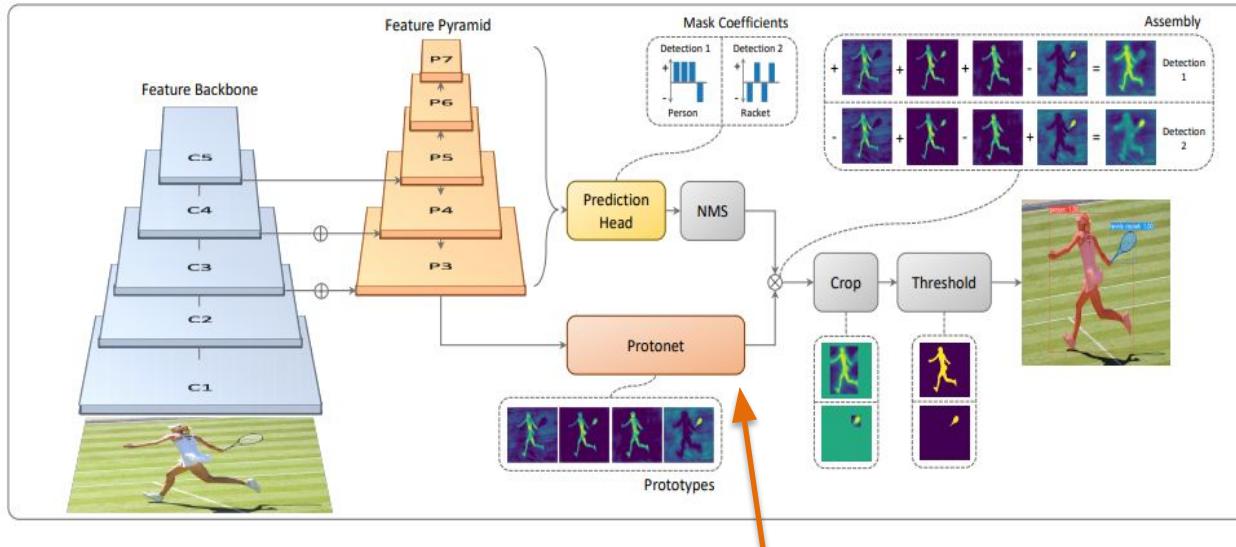
YolACT*

*You Only Look At Coefficients

YOLOCT: idea

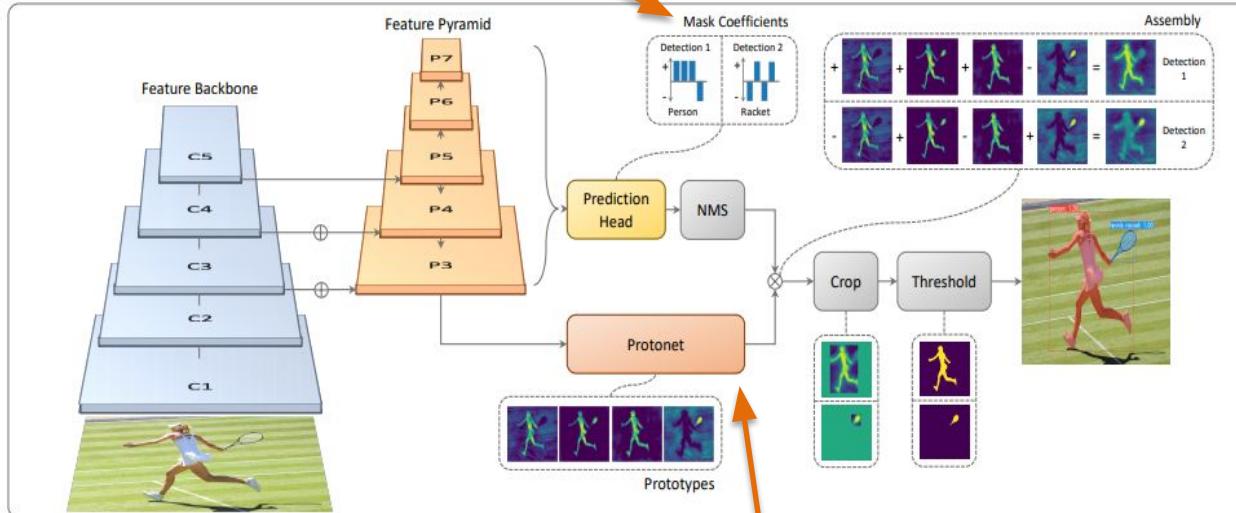


YOLOACT: idea



YOLOCT: idea

2) Generate mask coefficients

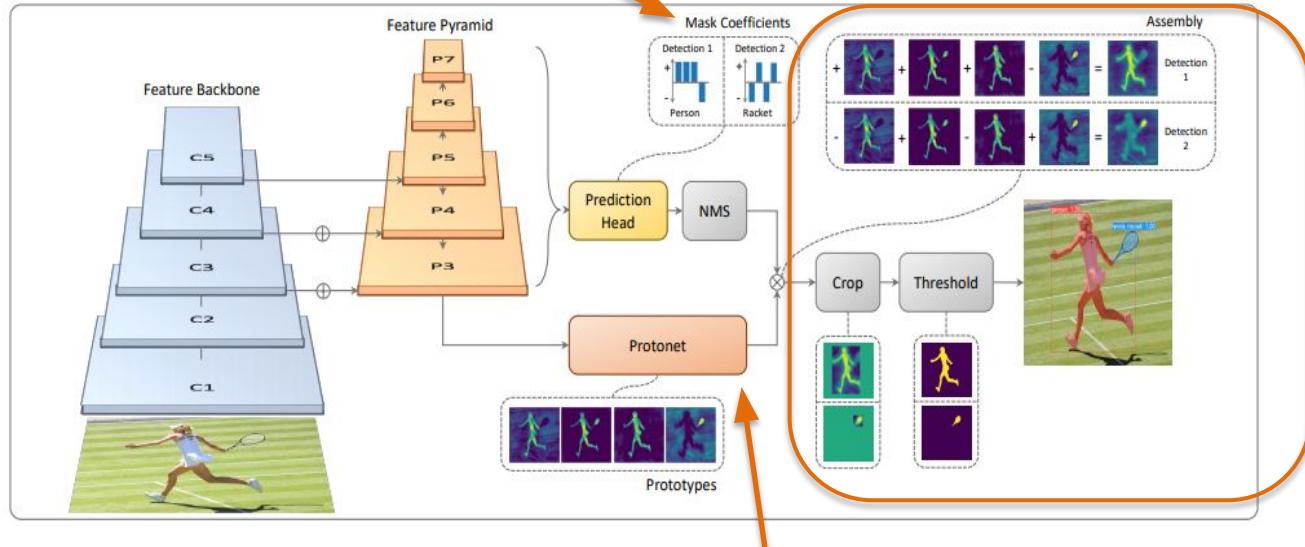


1) Generate mask prototypes

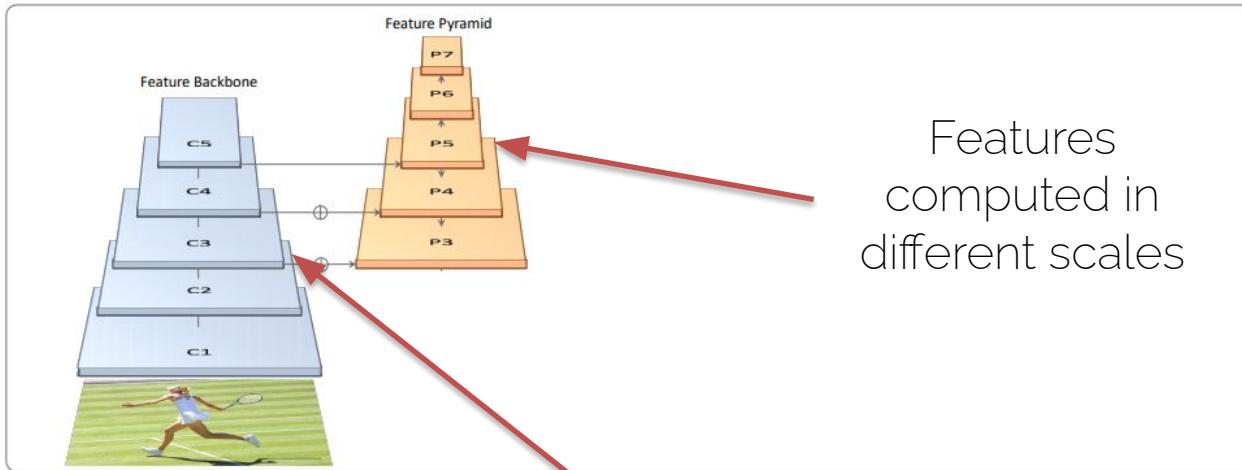
YOLOACT: idea

2) Generate mask coefficients

3) Combine (1) and (2)



YOLOCT: backbone

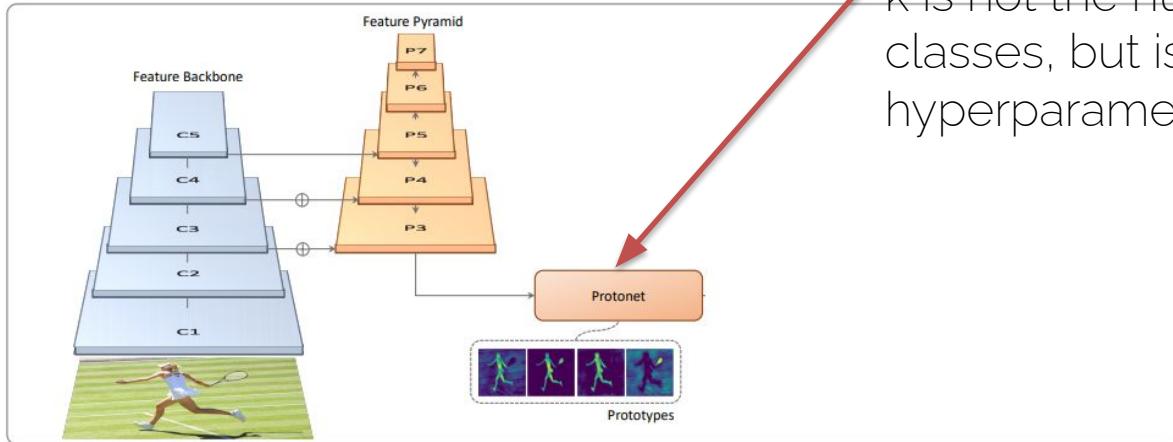


ResNet-101

YOLACT: protonet

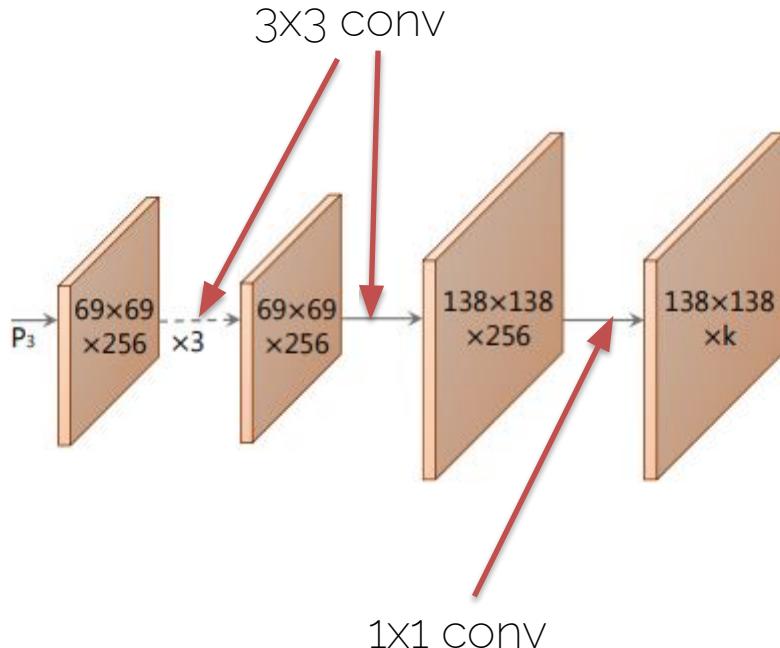
Generate k prototype masks.

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



YOLACT: protonet

- Fully convolutional network

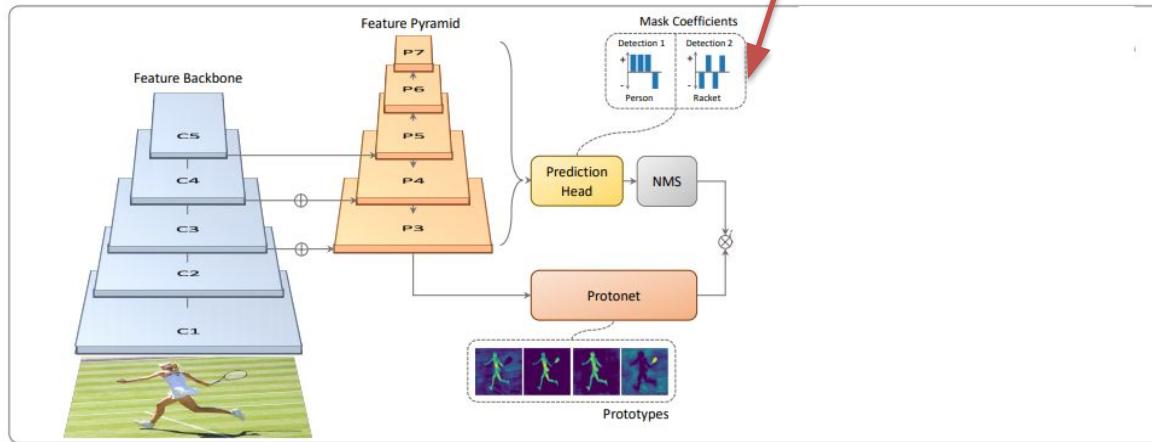


Similar to the mask branch in Mask R-CNN.

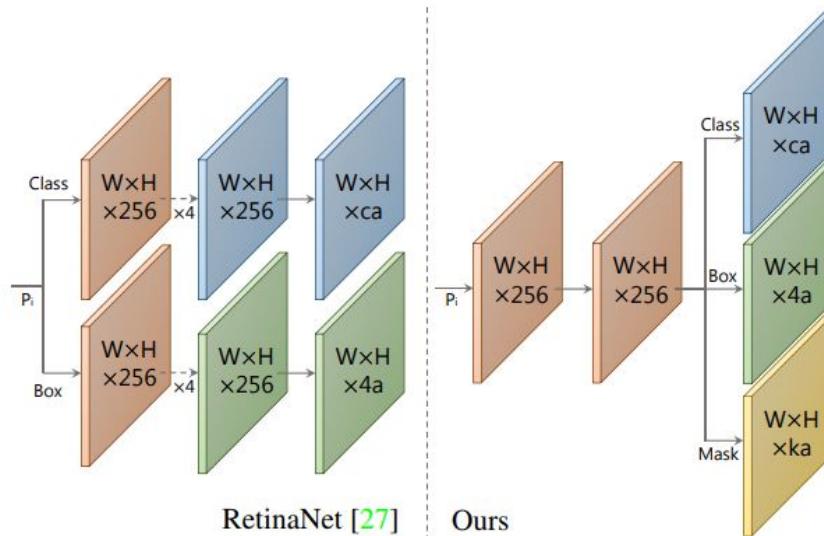
However, no loss function is applied on this stage.

YOLOCT: mask coefficients

Predict a coefficient for every predicted mask.



YOLOCT: mask coefficients



Predict one class per anchor box

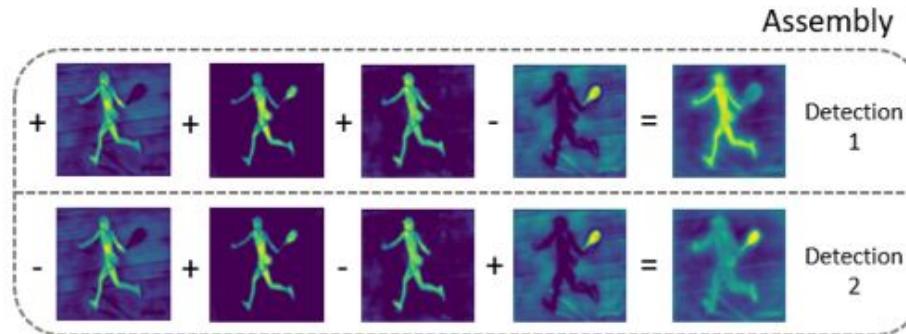
Predict the regression per anchor box

Predict k coefficients (one per prototype mask) per anchor

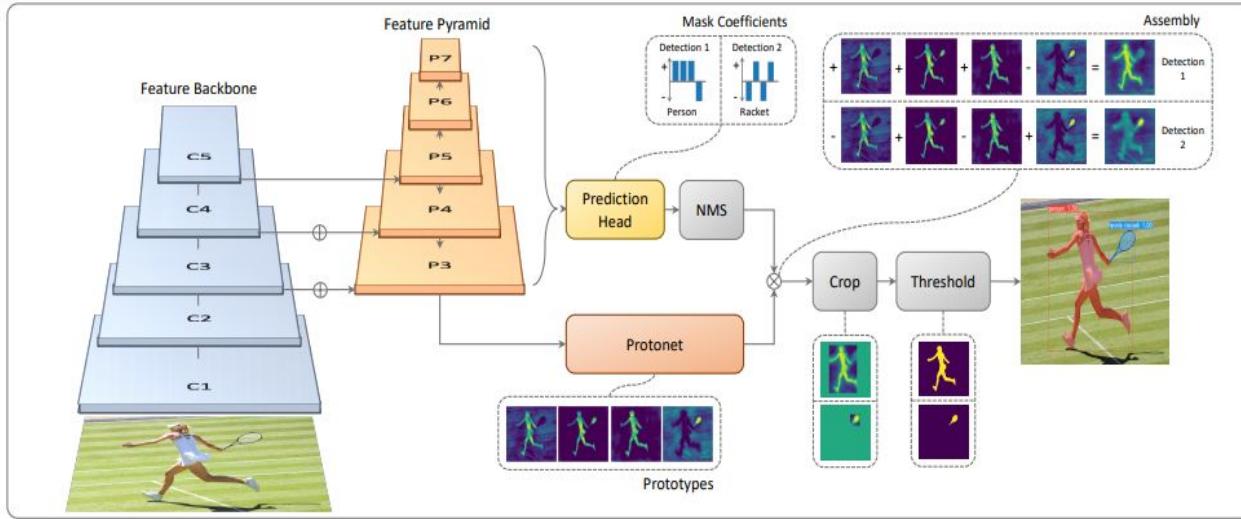
The network is similar but shallower than RetinaNet

YOLACT: 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 σ is a nonlinearity.



YOLACT: loss function

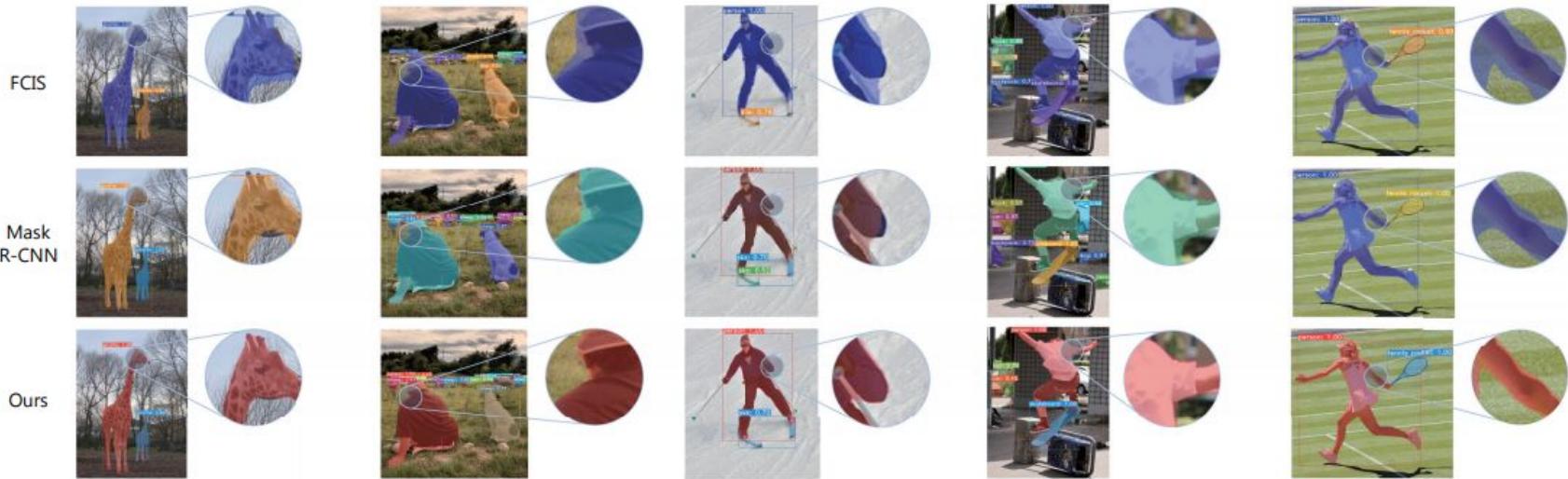


Cross-entropy between the assembled masks and the ground truth, in addition to the standard losses (regression for the bounding box, and classification for the class of the object/mask).

YOLACT: qualitative results

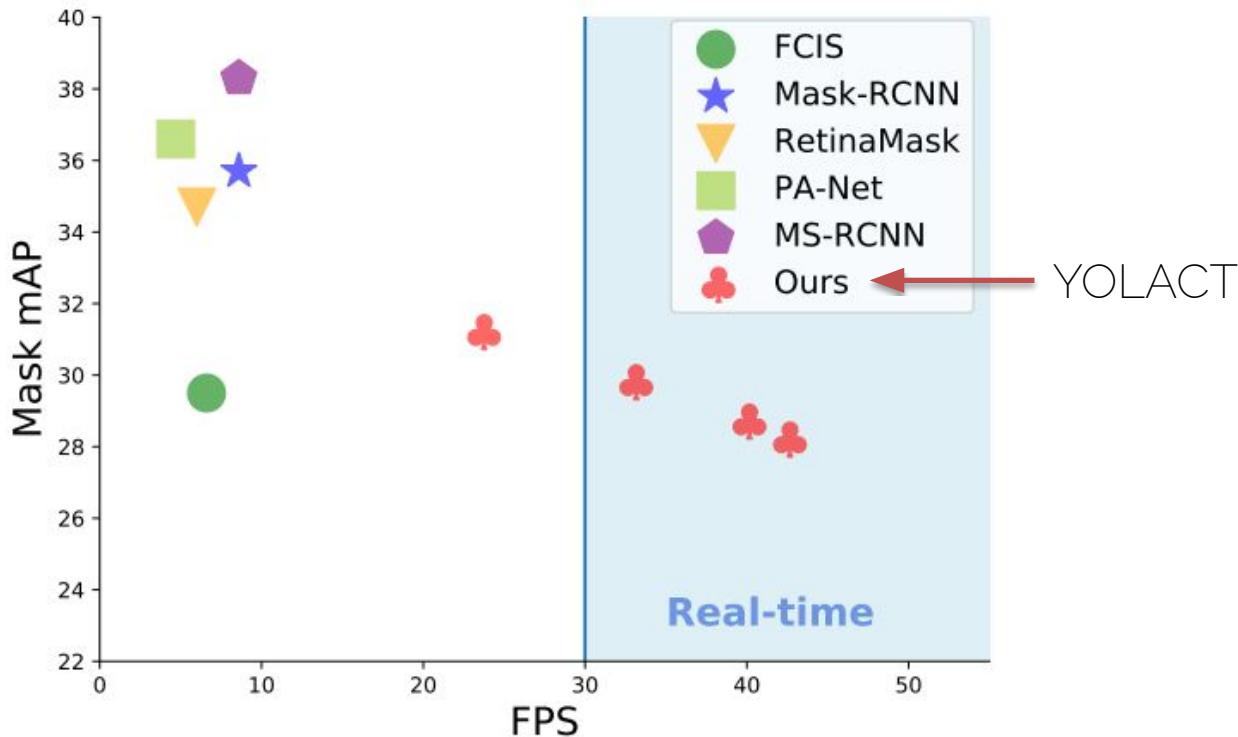


YOLACT: qualitative results



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

So, which segmenter to use?



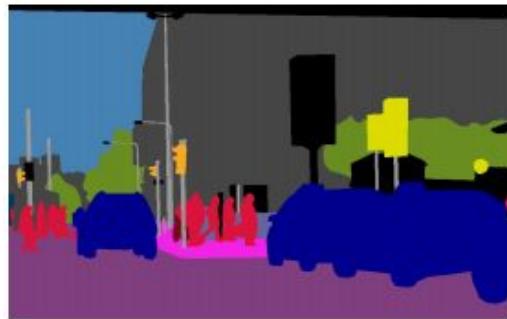
YOLOCT: improvements

- 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. "YOLOCT++: Better real-time instance segmentation". [arXiv:1912.06218](https://arxiv.org/abs/1912.06218) 2019

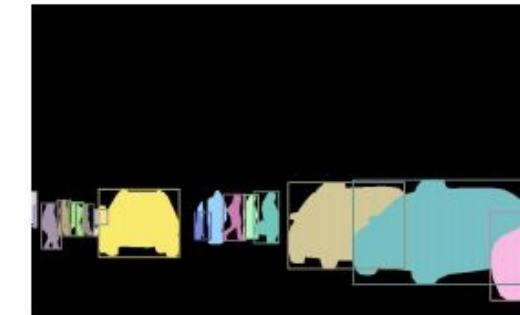
Panoptic segmentation

Panoptic segmentation

Semantic
segmentation



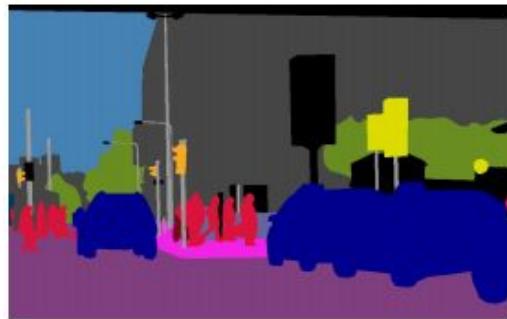
Instance
segmentation



+

Panoptic segmentation

Semantic
segmentation



FCN-like

Instance
segmentation

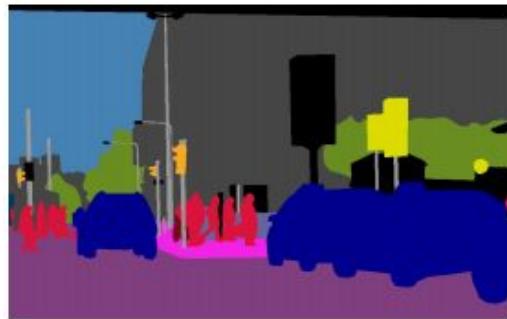


+

Mask R-CNN

Panoptic segmentation

Semantic
segmentation



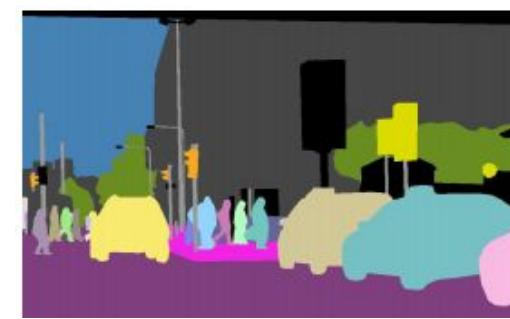
FCN-like

Instance
segmentation



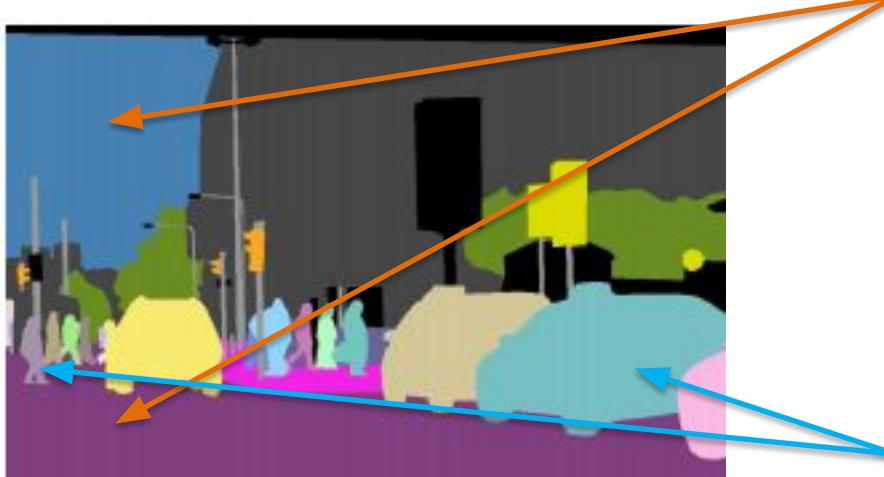
Mask R-CNN

Panoptic
segmentation



UPSNet

Panoptic segmentation



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

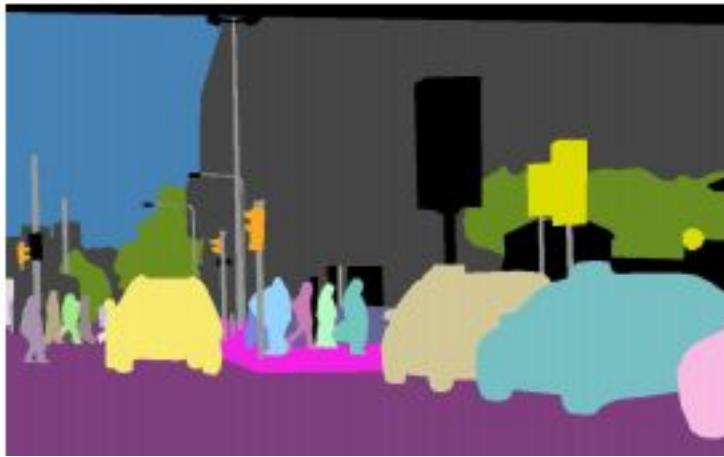
It differentiates between pixels coming from different instances of the same class (countable objects) called "things" (cars, pedestrians, etc).

Panoptic segmentation



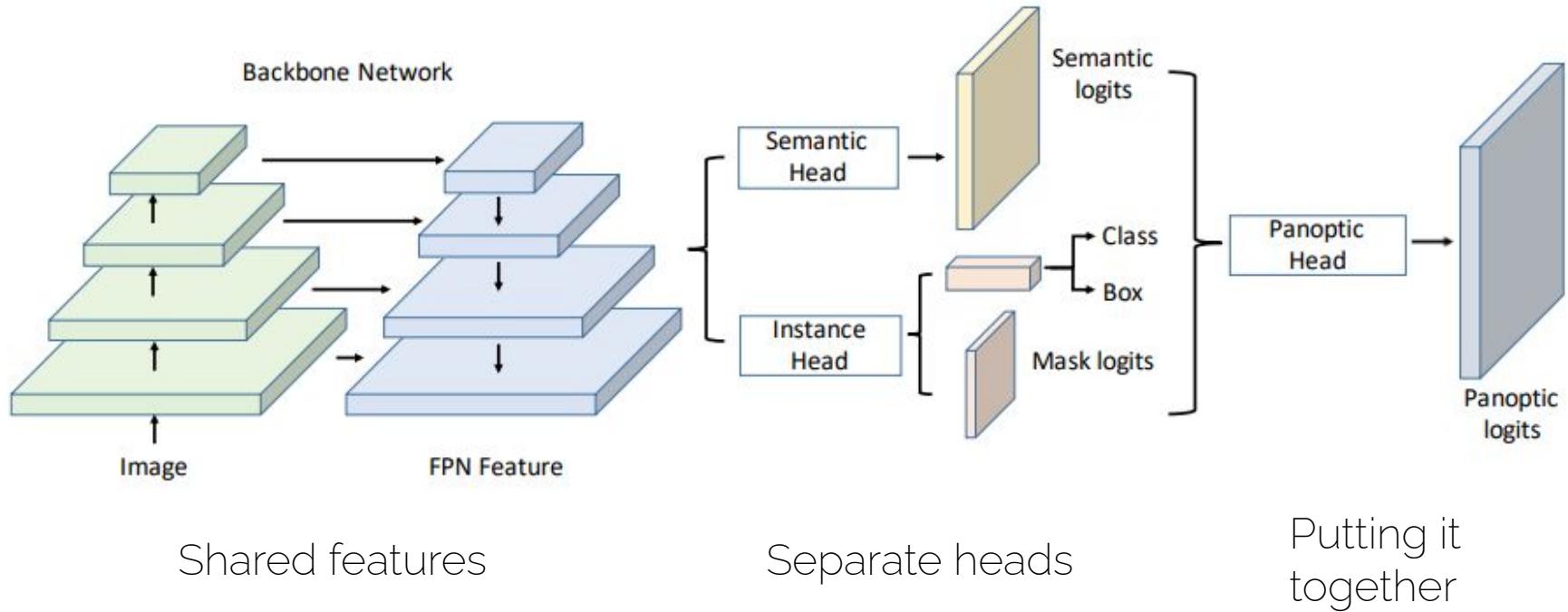
Problem: some pixels might get classified as stuff from FCN network, while at the same time being classified as instances of some class from Mask R-CNN (conflicting results)!

Panoptic segmentation

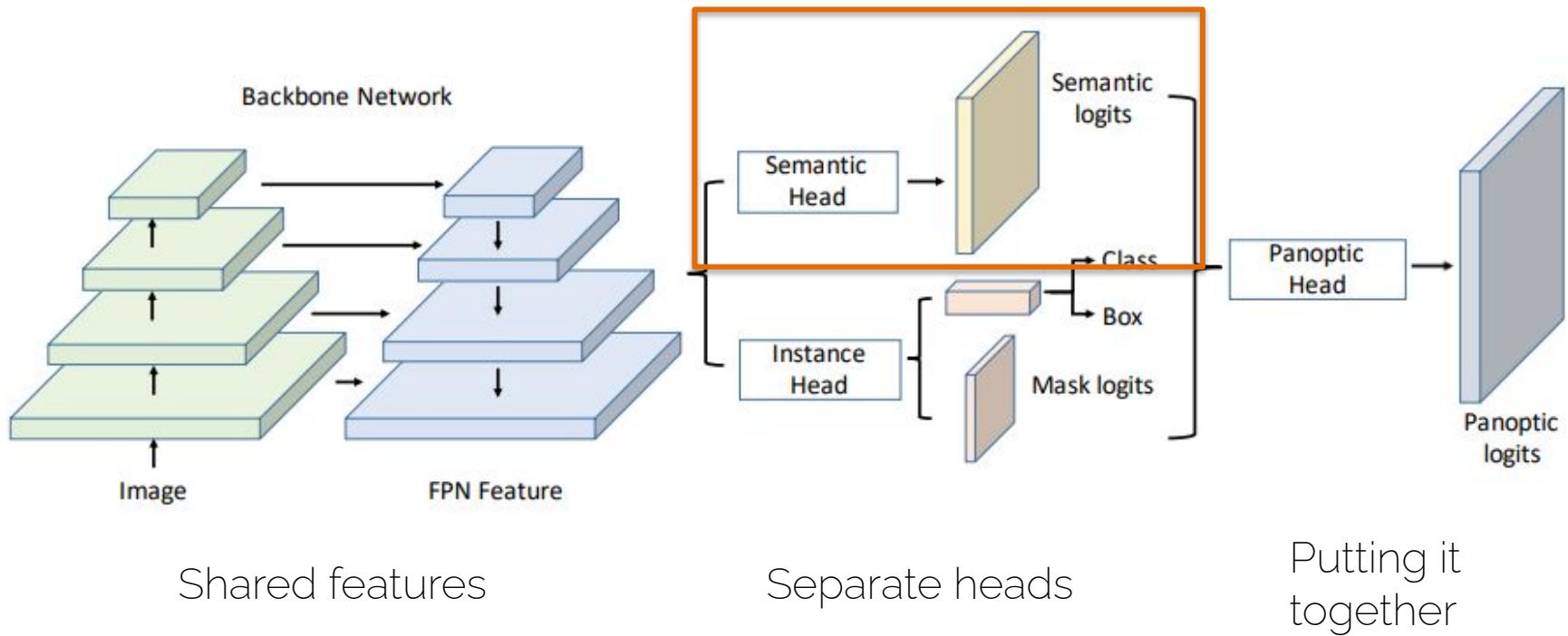


Solution: Parametric-free panoptic head which combines the information from the FCN and Mask R-CNN, giving final predictions.

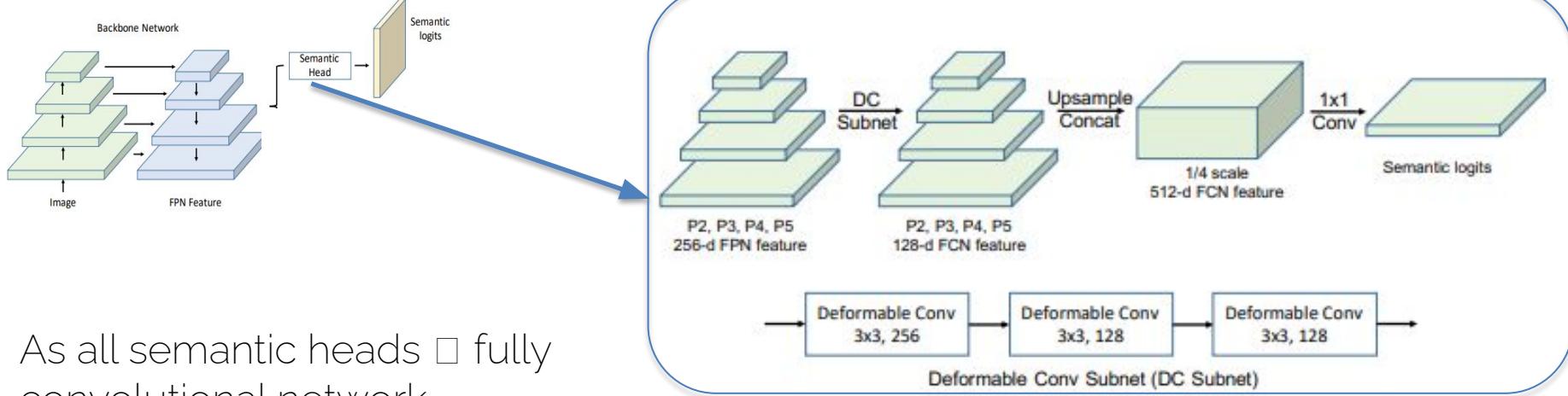
Network architecture



Network architecture



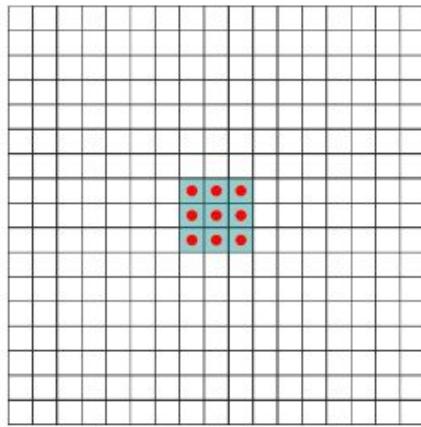
The semantic head



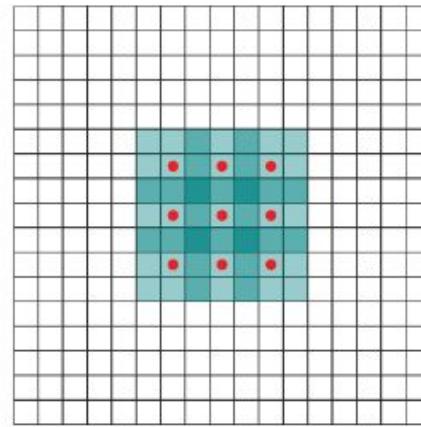
As all semantic heads □ fully convolutional network.

New: deformable convolutions!

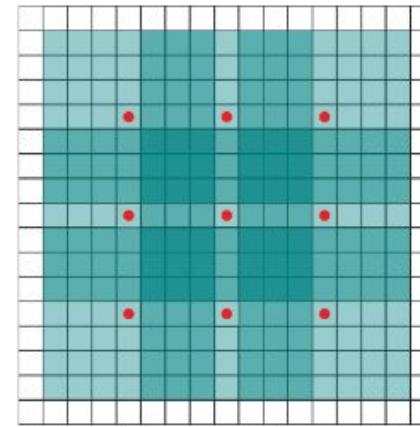
Recall: Dilated (atrous) convolutions 2D



(a)



(b)



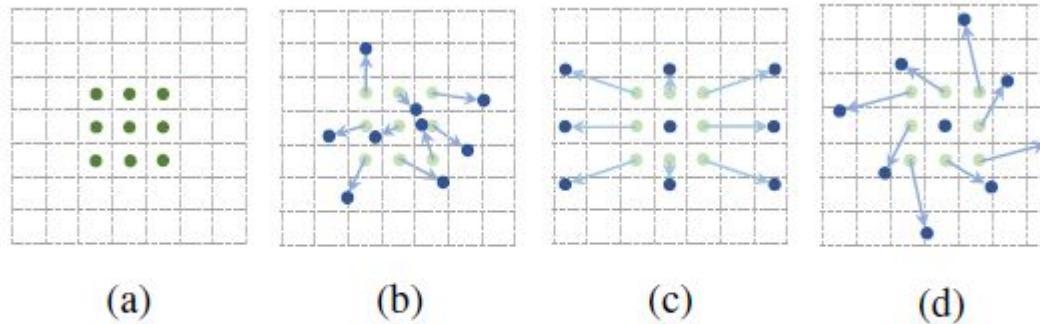
(c)

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

(b) the dilation parameter is 2, and each element produced by it has receptive field of 7×7 .

(c) the dilation parameter is 4, and each element produced by it has receptive field of 15×15 .

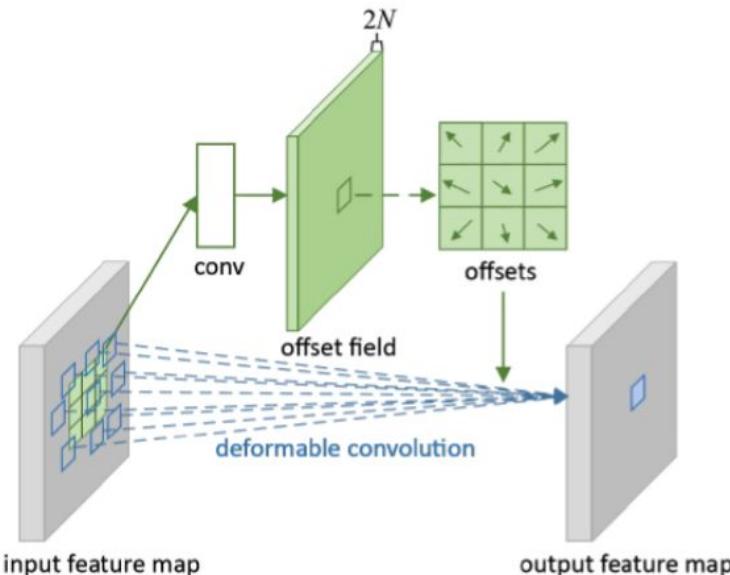
Deformable convolutions



(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

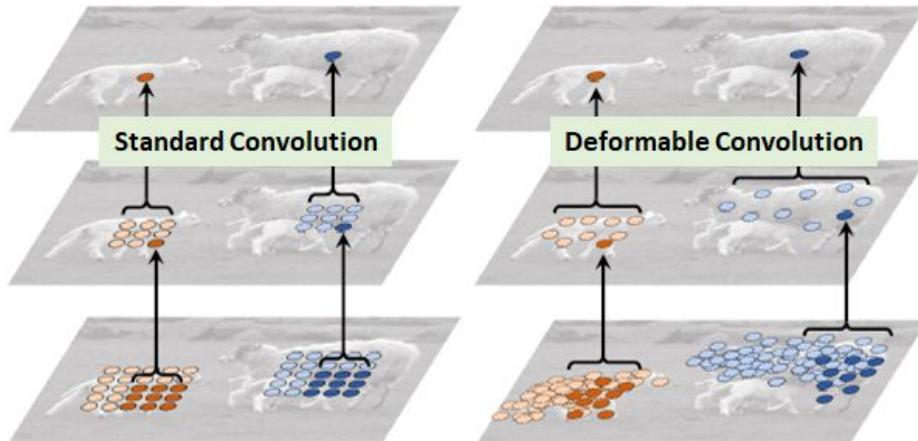
$$y(p_0) = \sum_{p_n \in \mathcal{R}} w(p_n) \cdot x(p_0 + p_n)$$

Deformable convolution

$$y(p_0) = \sum_{p_n \in \mathcal{R}} w(p_n) \cdot x(p_0 + p_n + \Delta p_n)$$

where Δp_n is generated by a sibling branch of regular convolution

Deformable convolutions



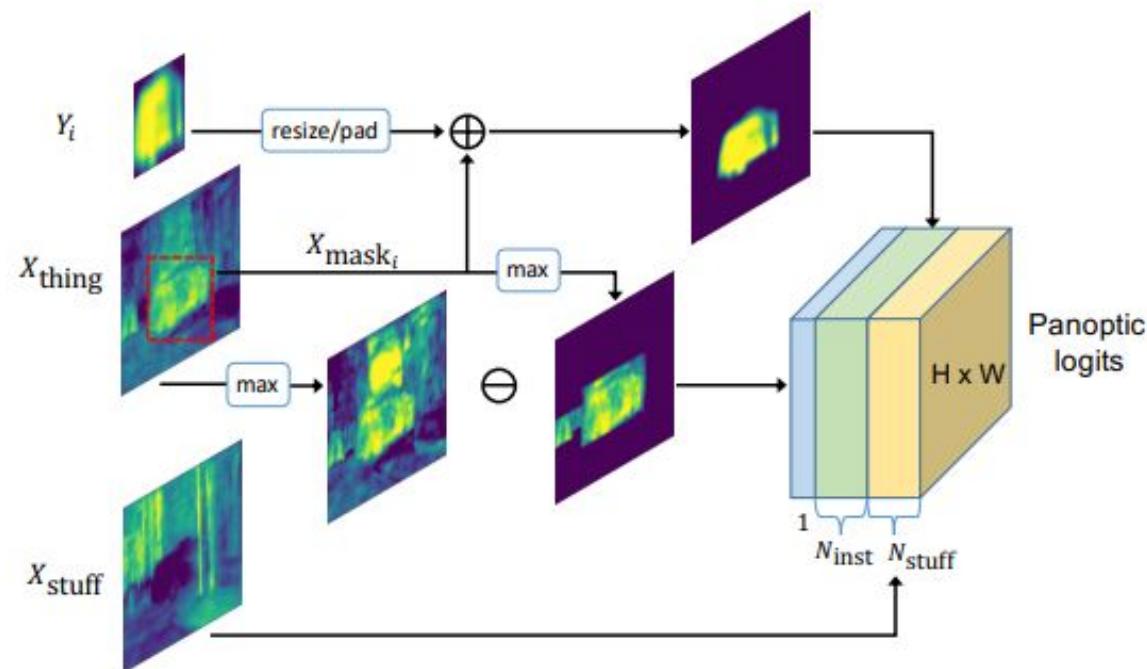
The deformable convolution will pick the values at different locations for convolutions conditioned on the input image of the feature maps.

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)



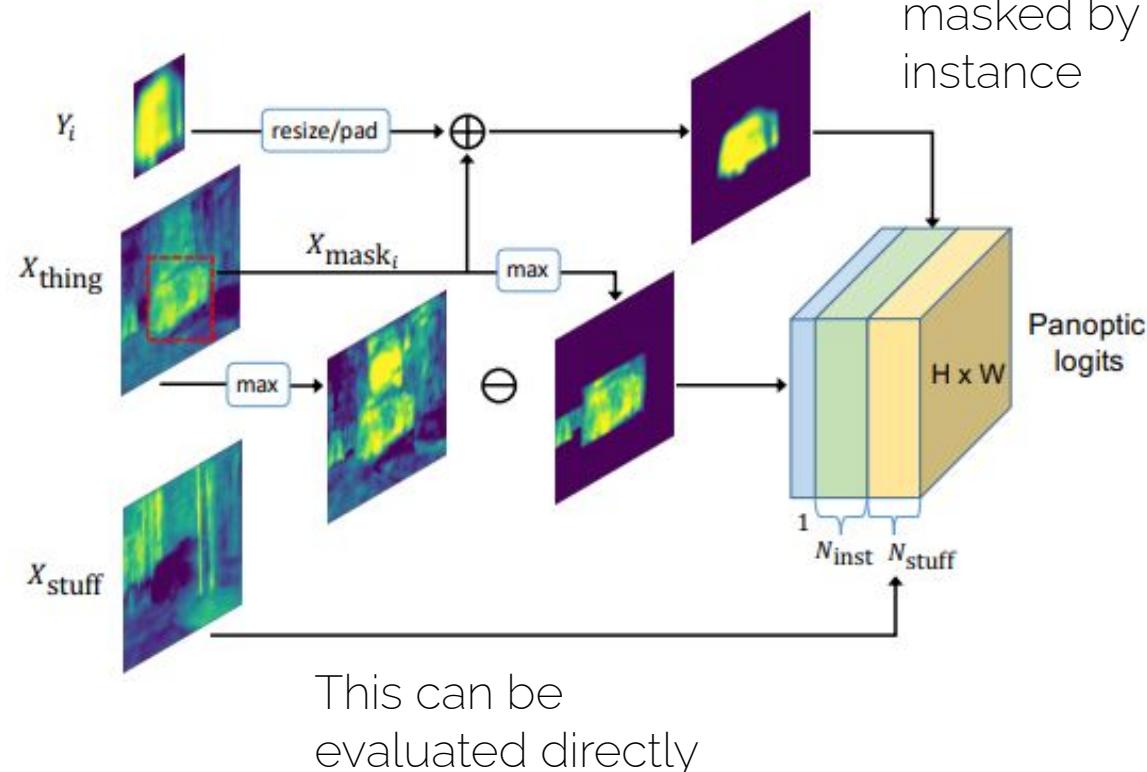
The Panoptic head

Mask logits from the instance head

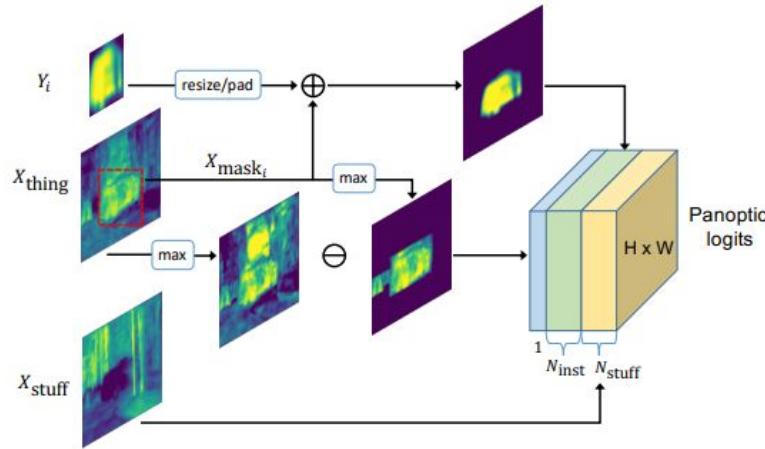
Object logits coming from the semantic head (e.g., car)

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

Objects need to be masked by the instance



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



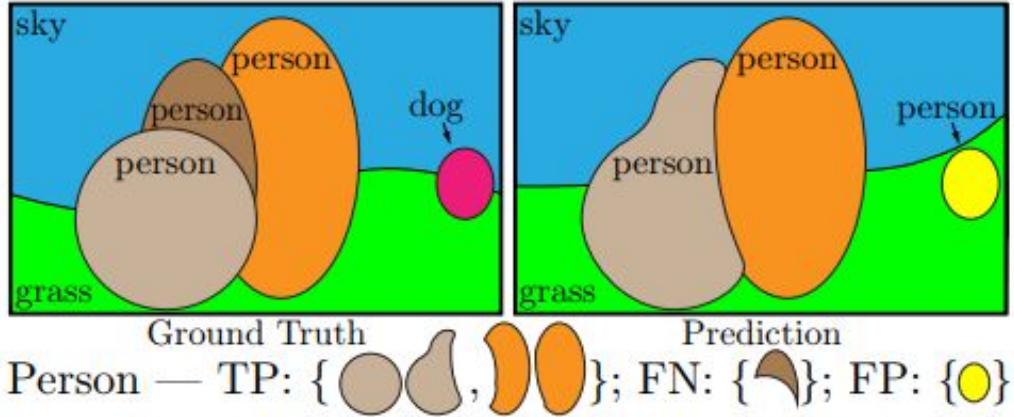
Xiong et al., "UPSNNet: A Unified Panoptic Segmentation Network". CVPR 2019

Metrics

Panoptic quality

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

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

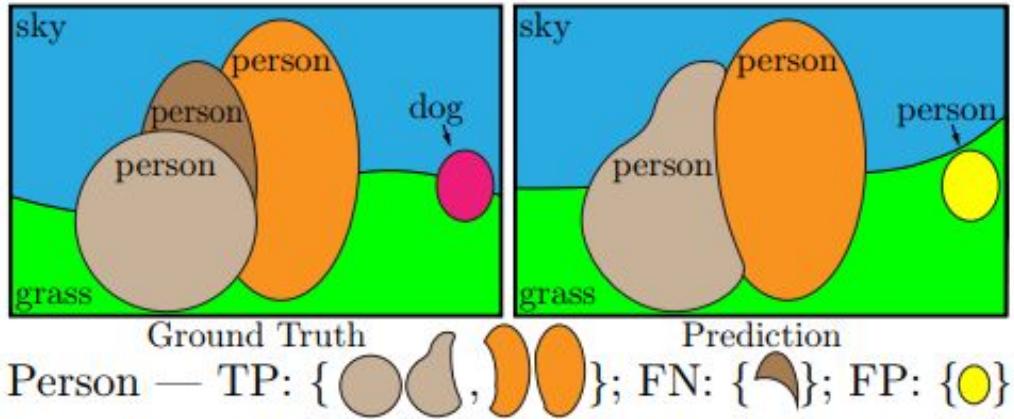


- SQ: Segmentation Quality = how close the predicted segments are to the ground truth segment
(does not take into account bad predictions!)

Panoptic quality

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

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



- 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).

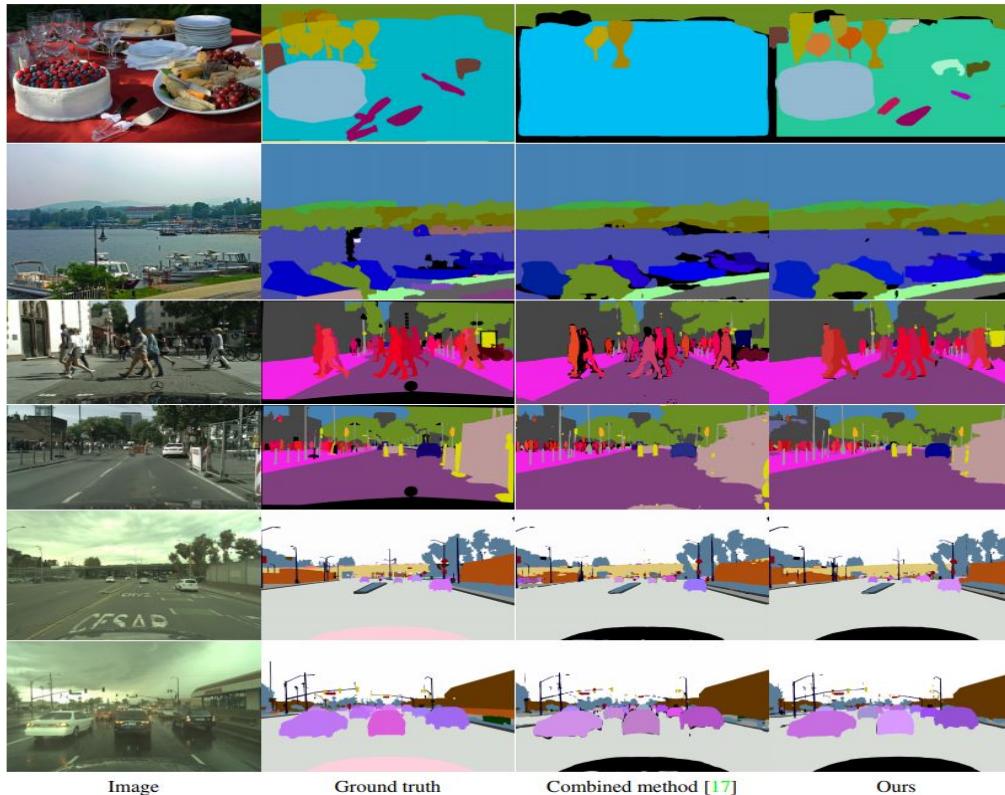
Panoptic quality

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

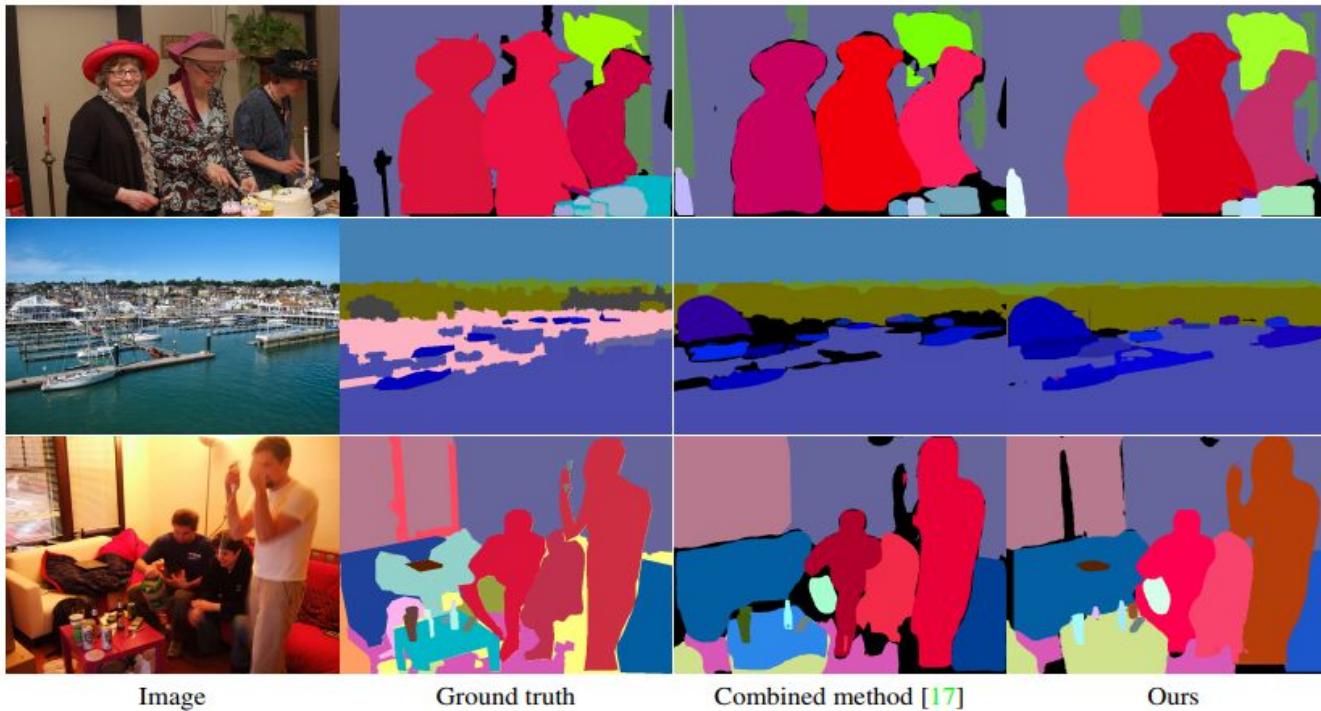


- Segment is matched if $\text{IoU} > 0.5$. No pixel can belong to two predicted segments.

Panoptic segmentation: qualitative



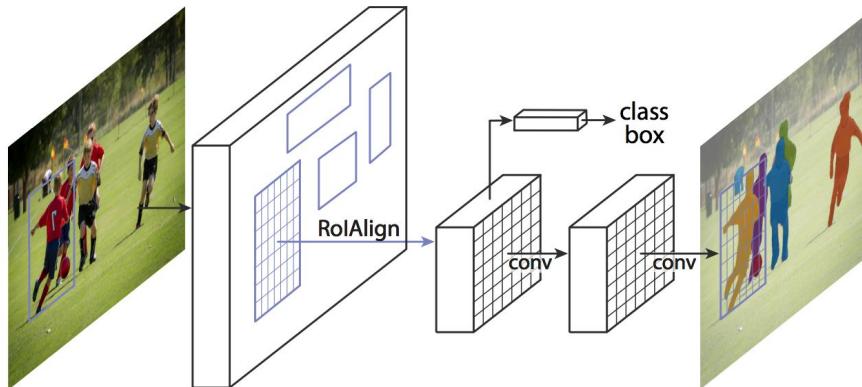
Panoptic segmentation: qualitative



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

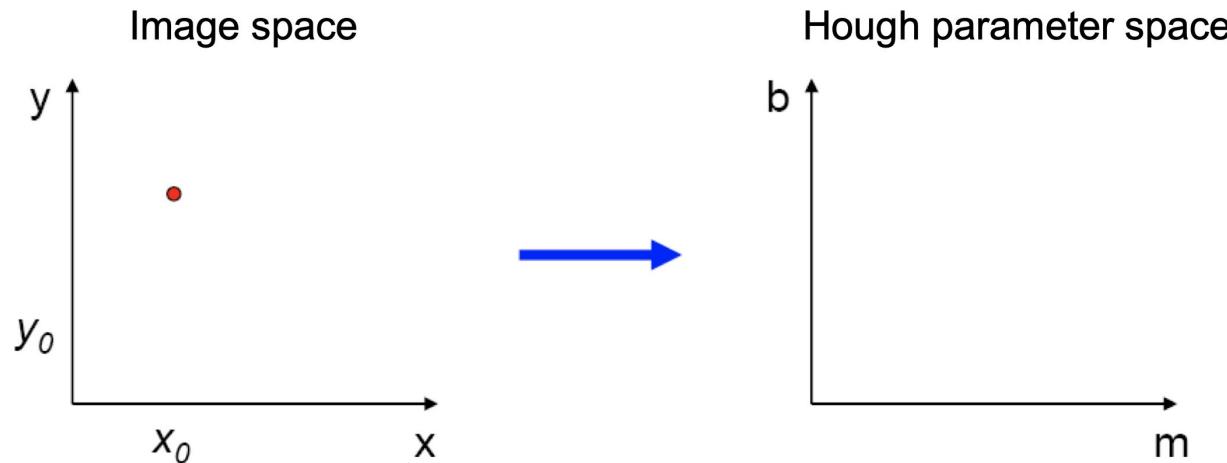


Hough Voting

- 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

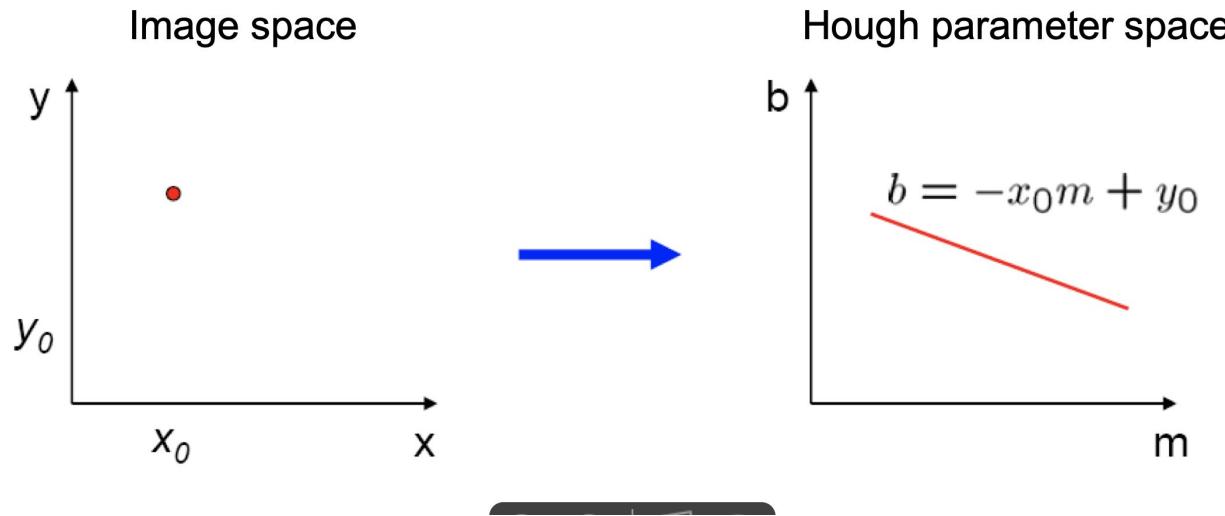
Example: Line Detection

- Each edge point in image space casts a vote



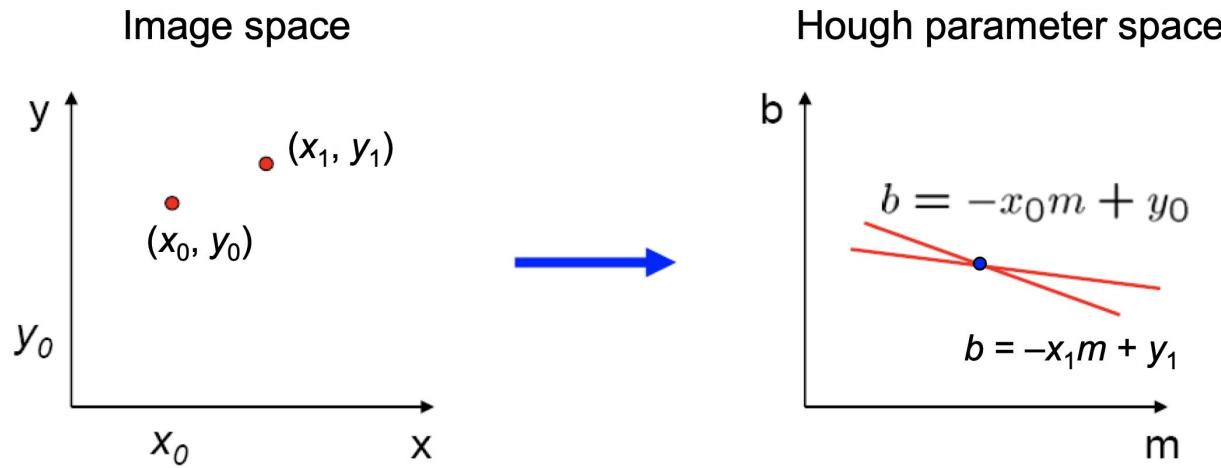
Example: Line Detection

- 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



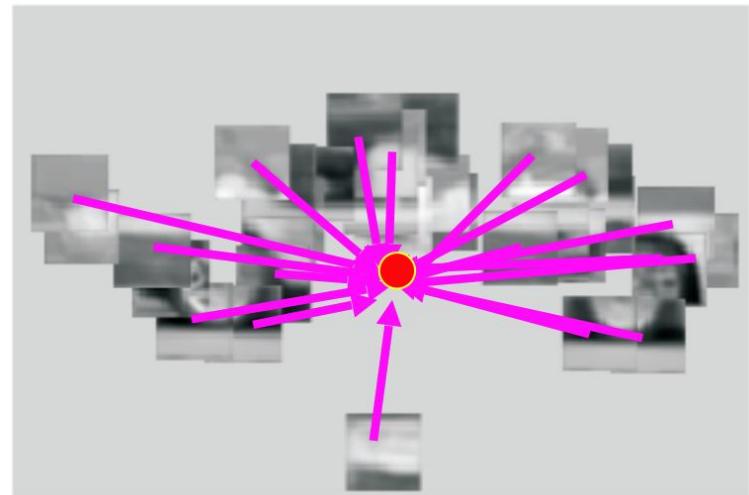
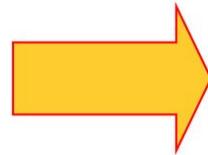
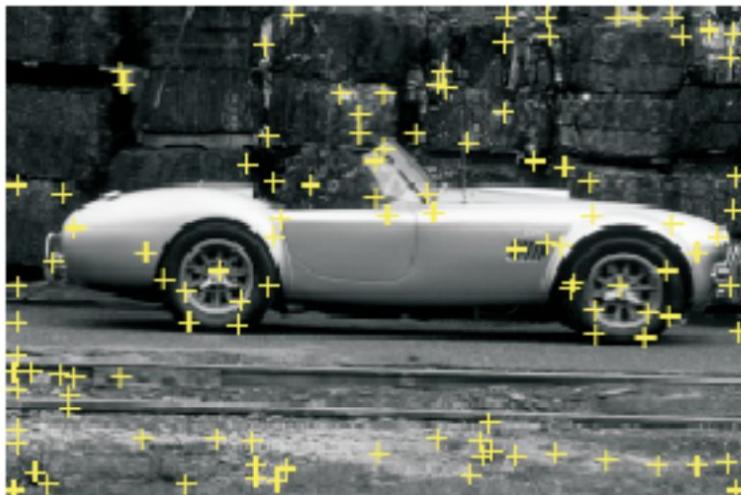
Object Detection as Voting

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



Object Detection

- Training

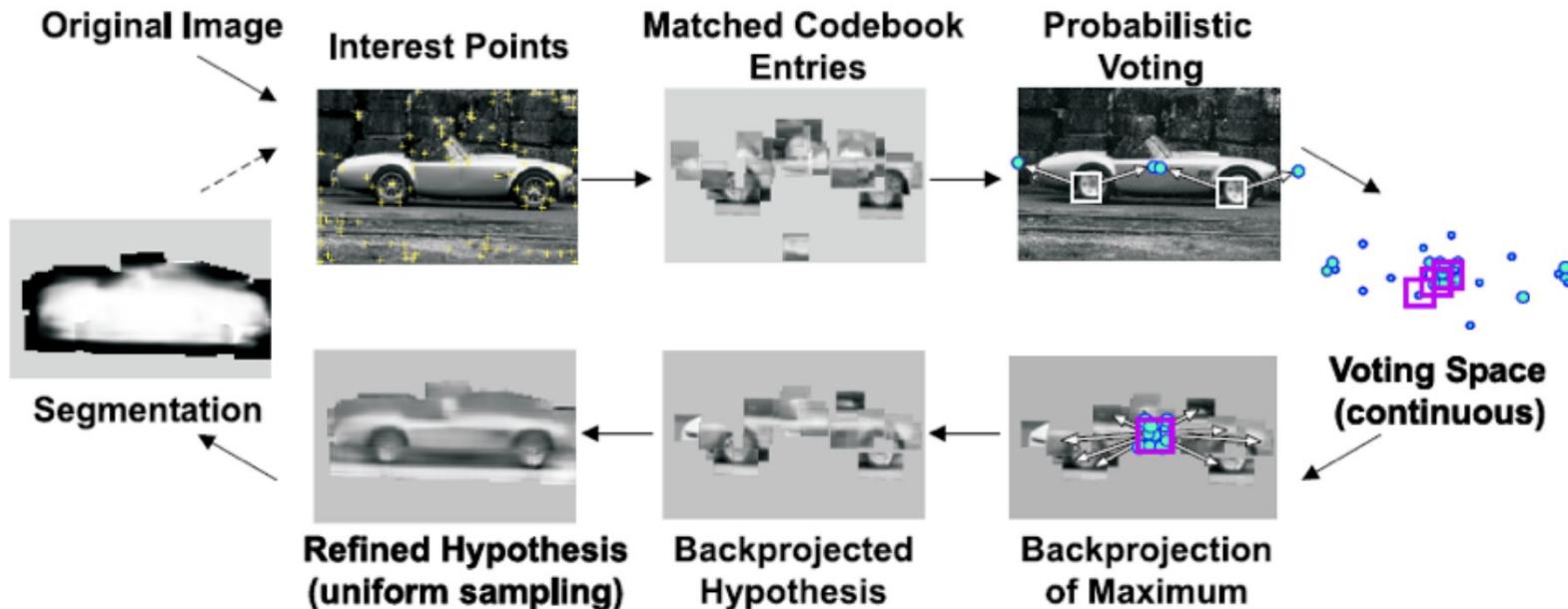


Interest point detection
(SIFT, SURF)

Center point voting

Object Detection

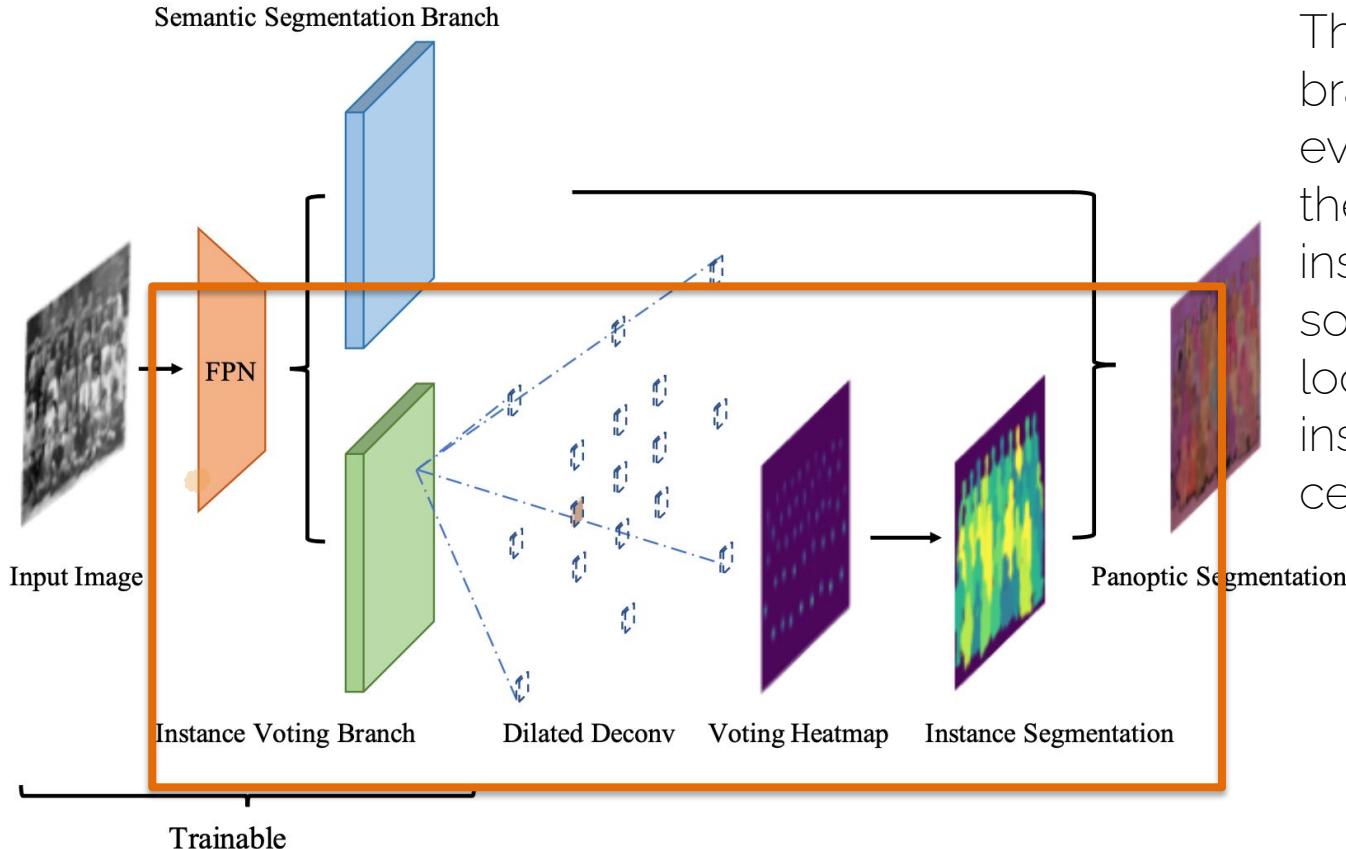
- Inference (test time)



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.



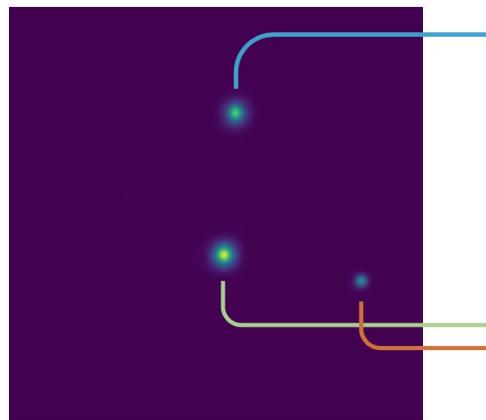
Discretization & Classification

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



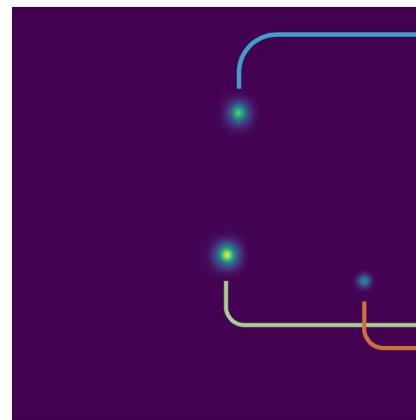
Voting as Transposed Convolution

In a Nutshell

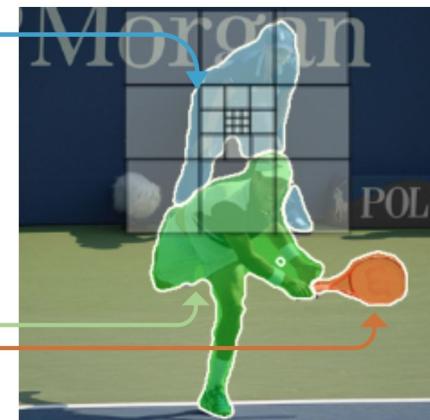
5. Back-projection of 'peaks' back to the image to get an instance masks
6. Category information provided by the parallel semantic segmentation head



Discretization & Classification



Voting as Transposed Convolution



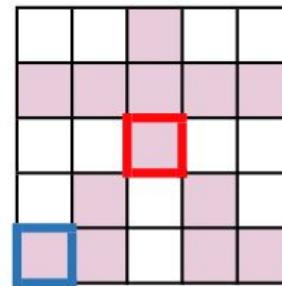
Backprojection as Filtering

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

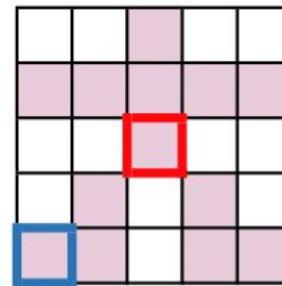
Ground Truth Assignment

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.

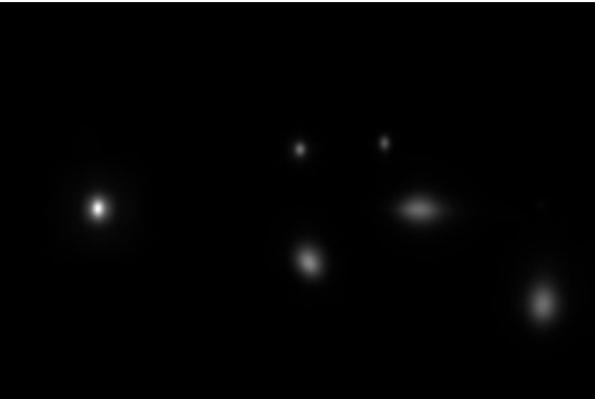
Voting - Implementation

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

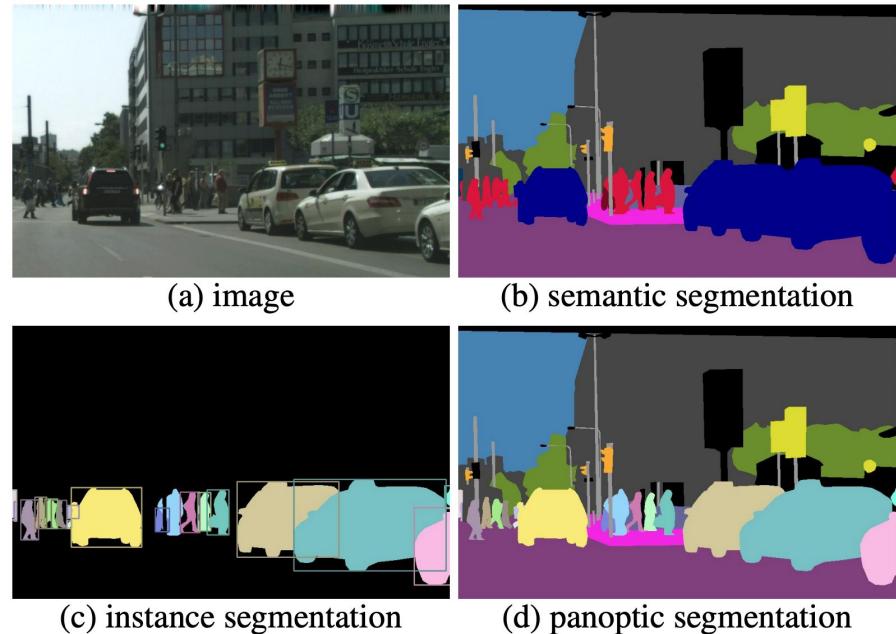
My center
is at pixel
8!

Qualitative Results



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.



Instance segmentation