

YOLACT:

Real-time Instance Segmentation

Daniel Bolya, Chong Zhou, Fanyi Xiao, Yong Jae Lee
[University of California, Davis]

ICCV 2019

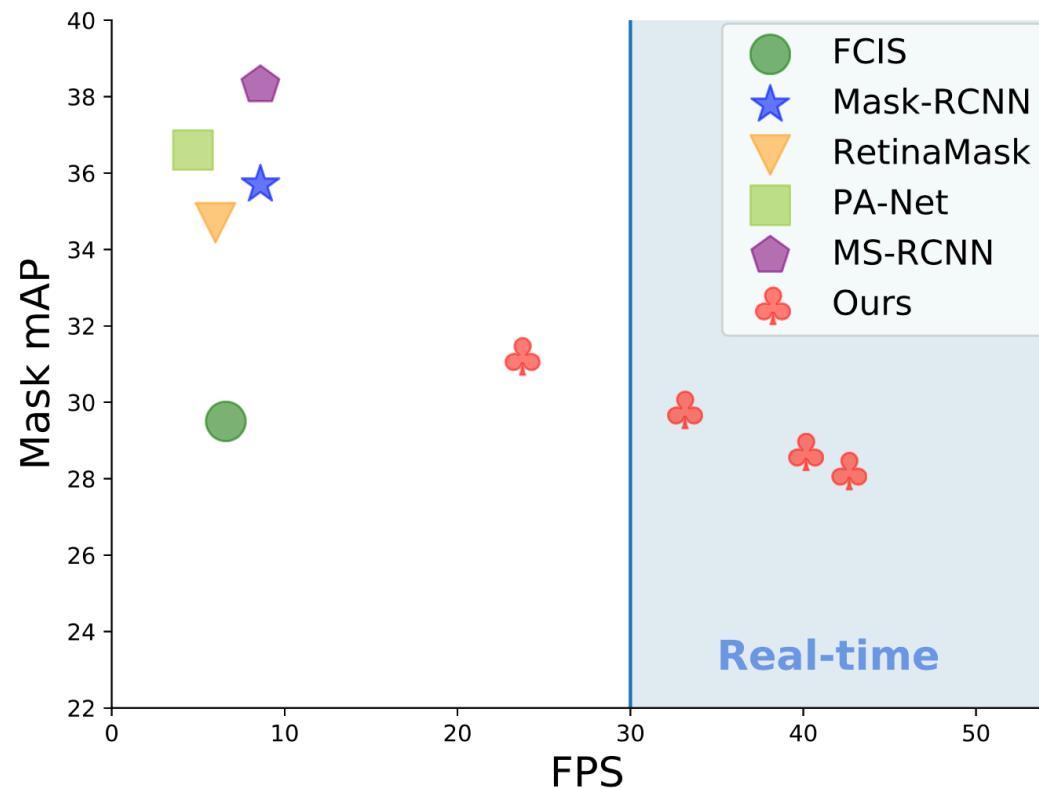
Sungman, Cho.



Introduction

Challenges

- Instance segmentation is hard – much harder than object detection.



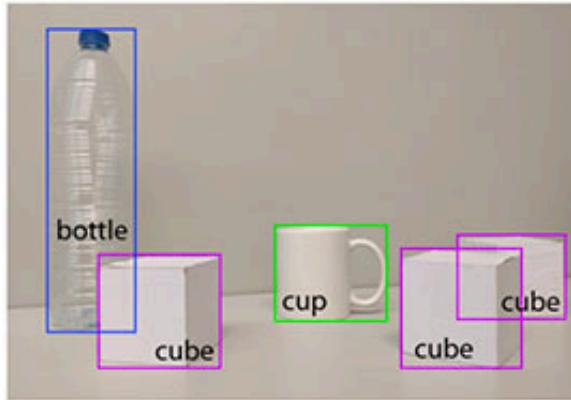
Contribution

- **Real-time Instance Segmentation.**
(29.8mAP on MS COCO at 33.5fps)
- Propose Fast NMS

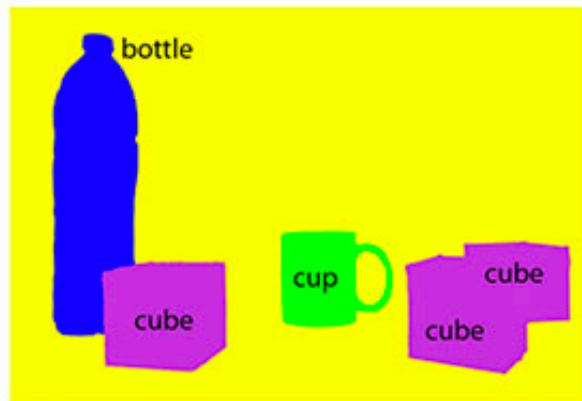
Related Works – Instance Segmentation



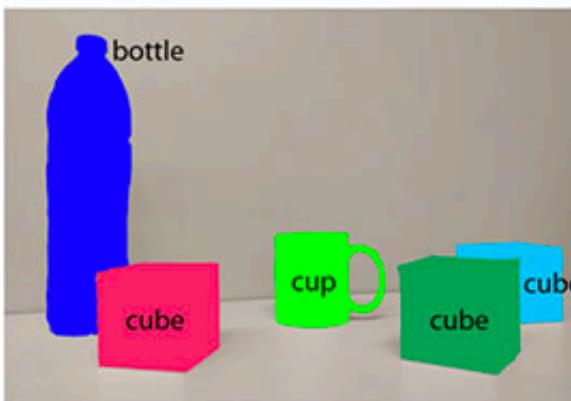
(a) Image classification



(b) Object localization



(c) Semantic segmentation

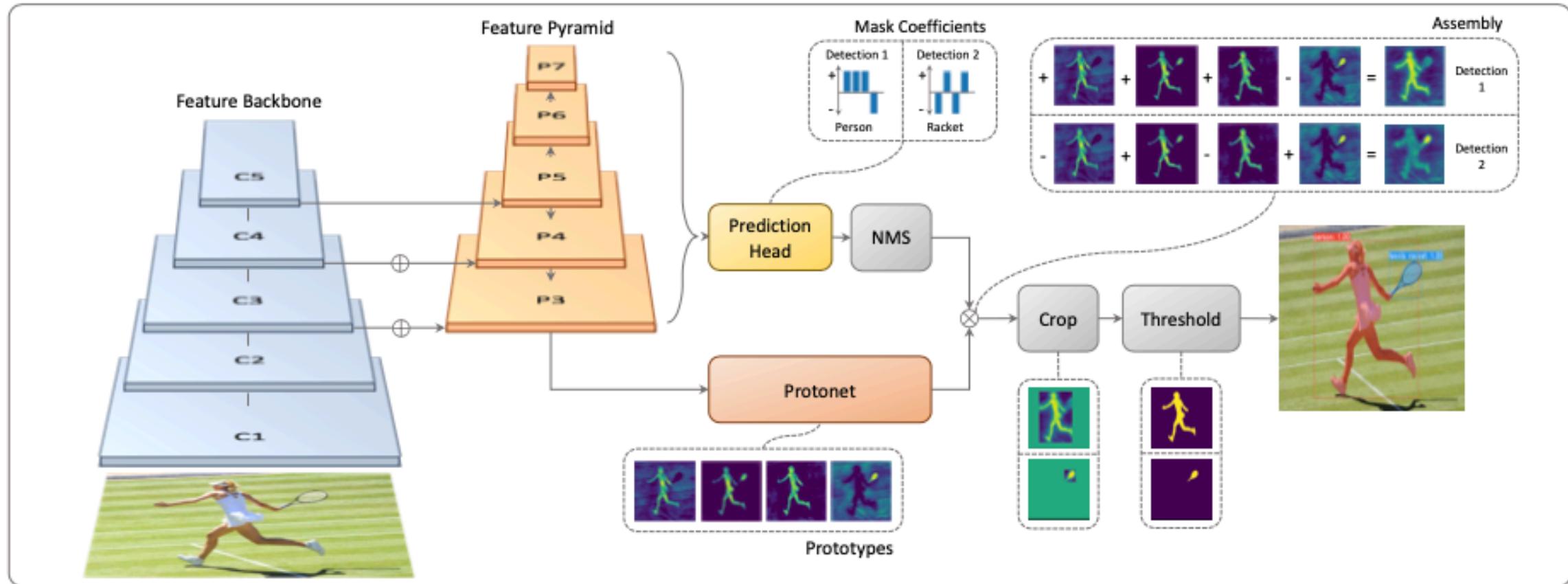


(d) Instance segmentation

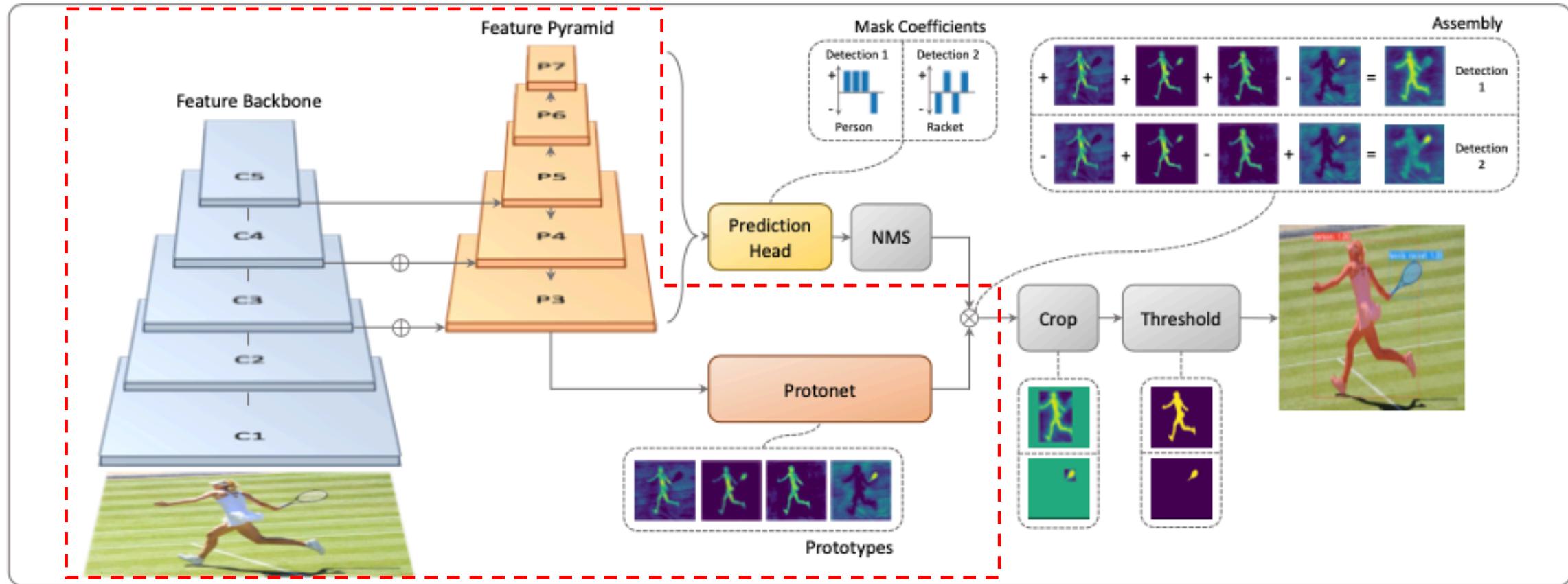
- Box2Pix : 30 fps on PASCAL SBD 2012
- Straight to Shapes : 10.9 fps on Cityscapes
- Mask R-CNN : 13.5 fps on 550x550 images

Methodology

Two branches

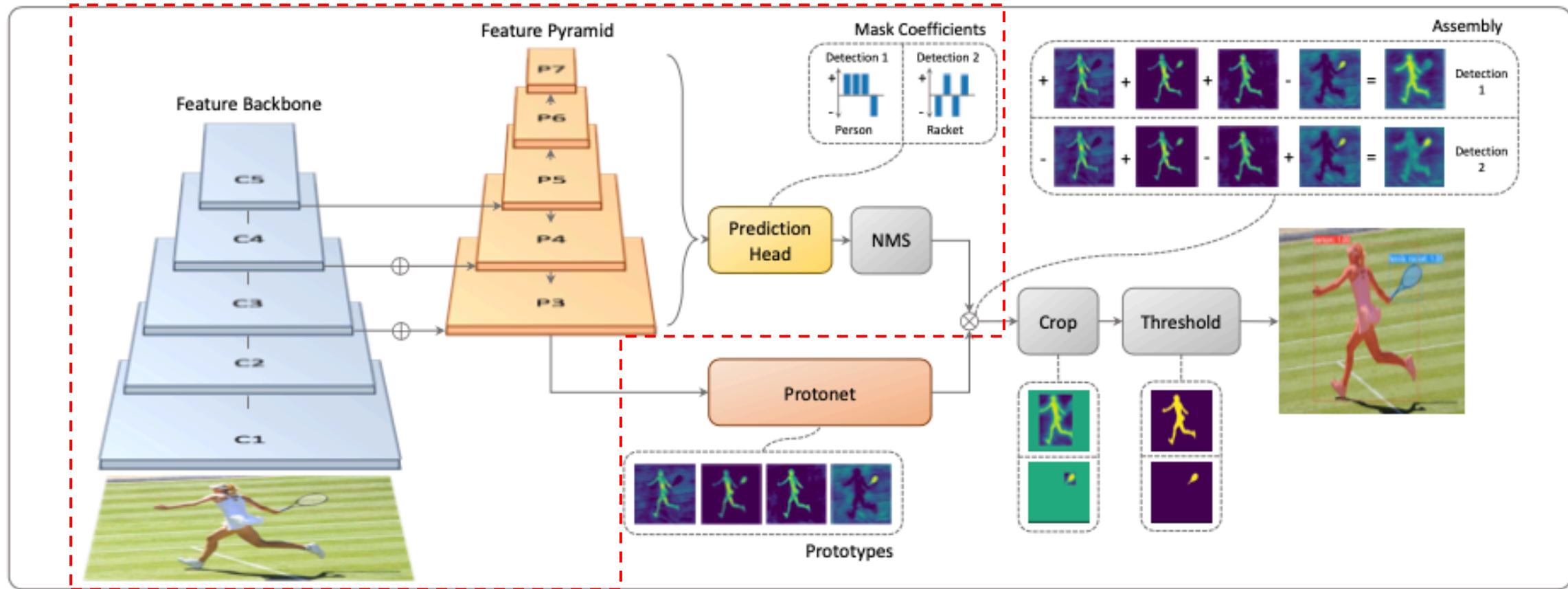


Two branches – FCN base.



Produce a set of image-sized “prototype masks” that do not depend on any one instance.

Two branches – Detection base.



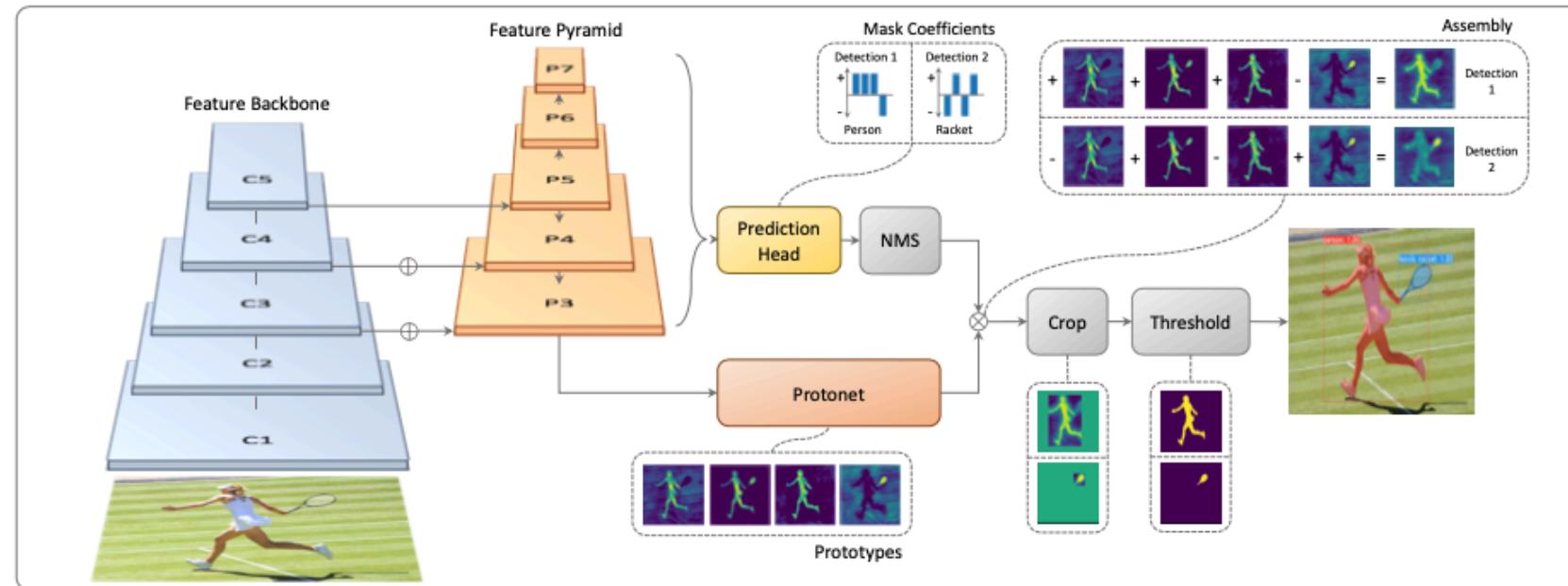
Predict a vector of “mask coefficients” for each anchor that encode an instance’s representation

Rationale

- **masks are spatially coherent.**
conv. layers naturally takes advantage of this coherence, a fc layer does not.
- Thus, making use of **fc layers**, which are good at **producing semantic vectors**,
conv. layers, which are good at producing **spatially coherent masks**, to produce the
" **mask coefficients** " and " **prototype masks** " , respectively.

Rationale

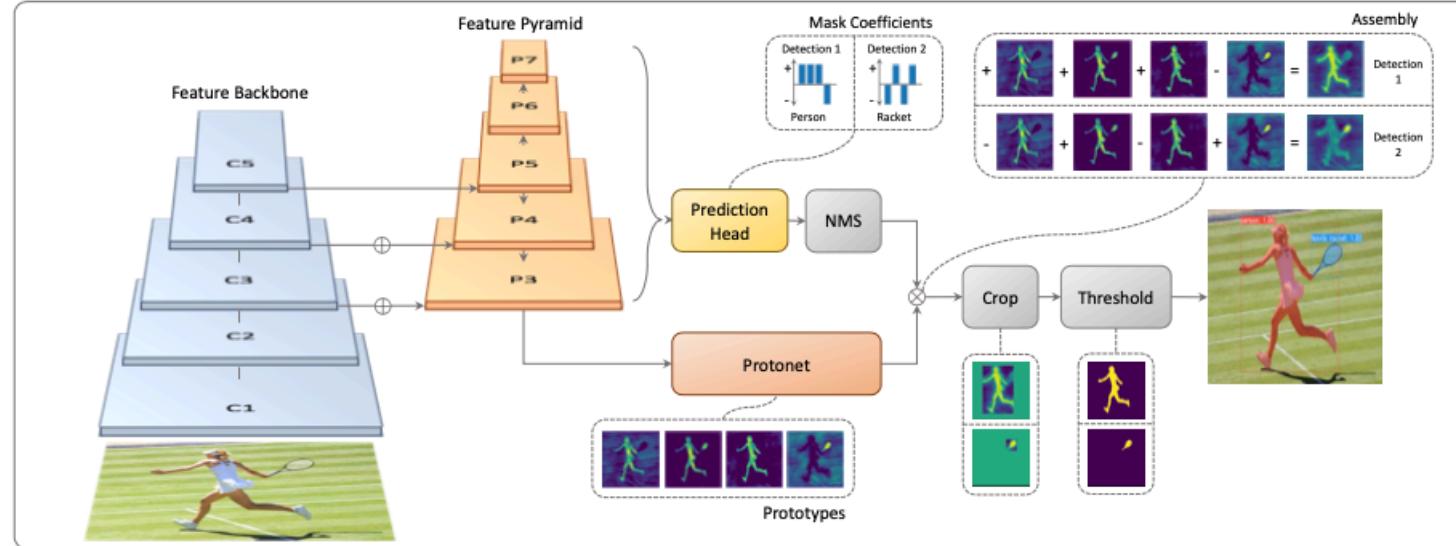
- Thus, making use of **fc layers**, which are good at **producing semantic vectors**, **conv. layers**, which are good at producing **spatially coherent masks**, to produce the "**mask coefficients**" and "**prototype masks**", respectively.



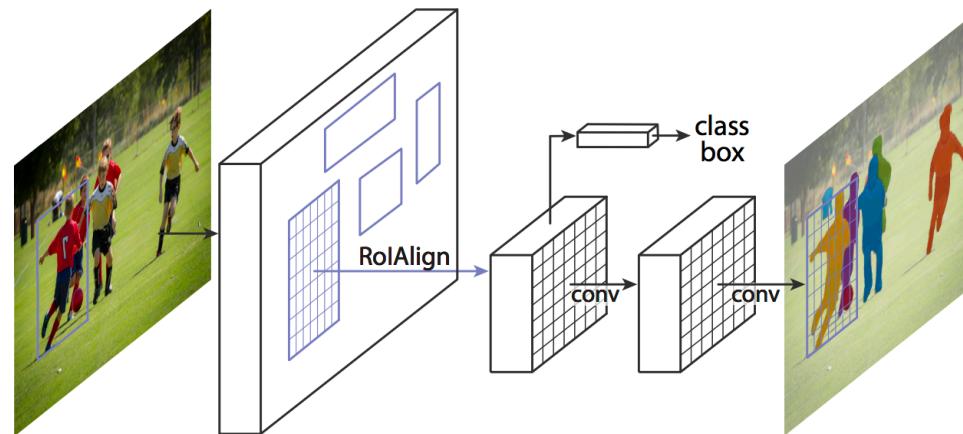
- Prototypes and mask coefficients can be computed independently.
- Assembly step can be implemented as a single matrix multiplication.

} **Fast**

Rationale

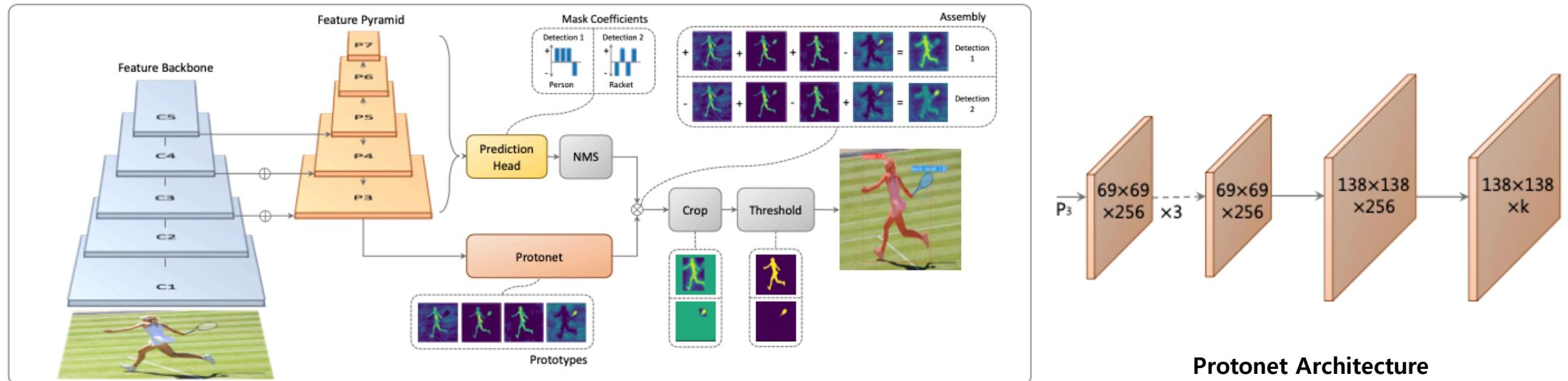


YOLOACT



Mask R-CNN

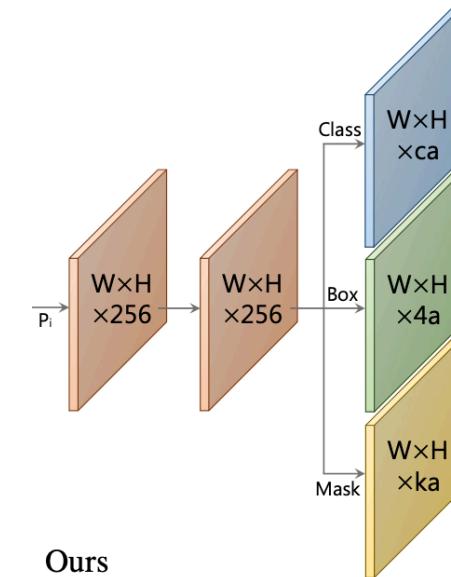
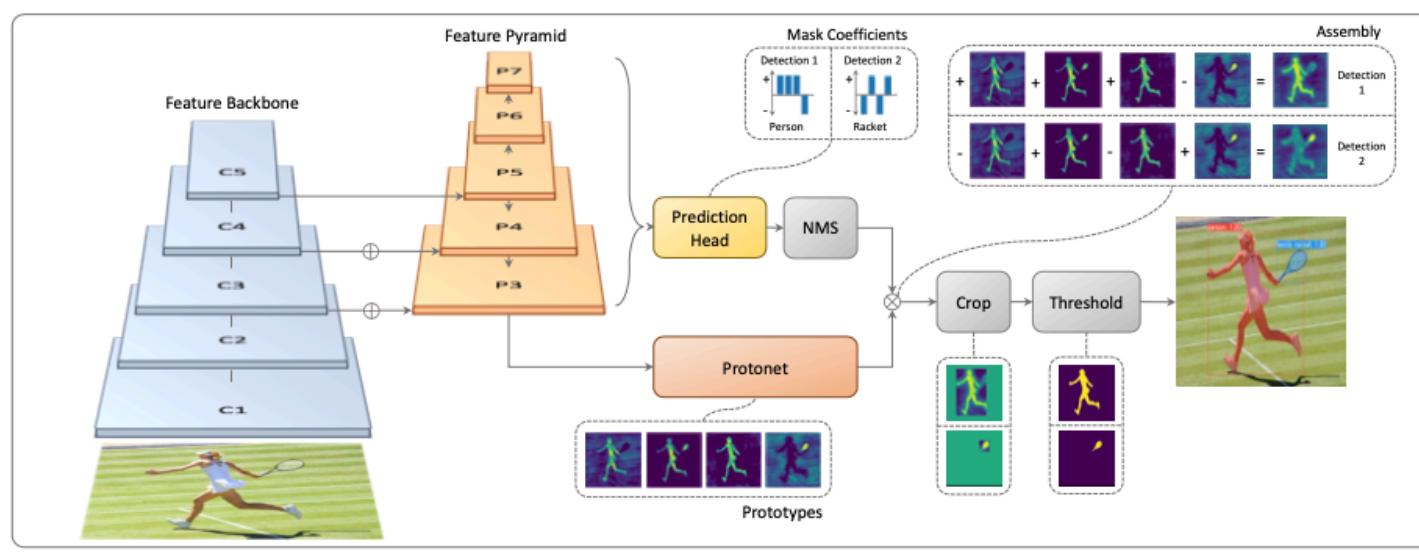
Prototype Generation



Protonet Architecture

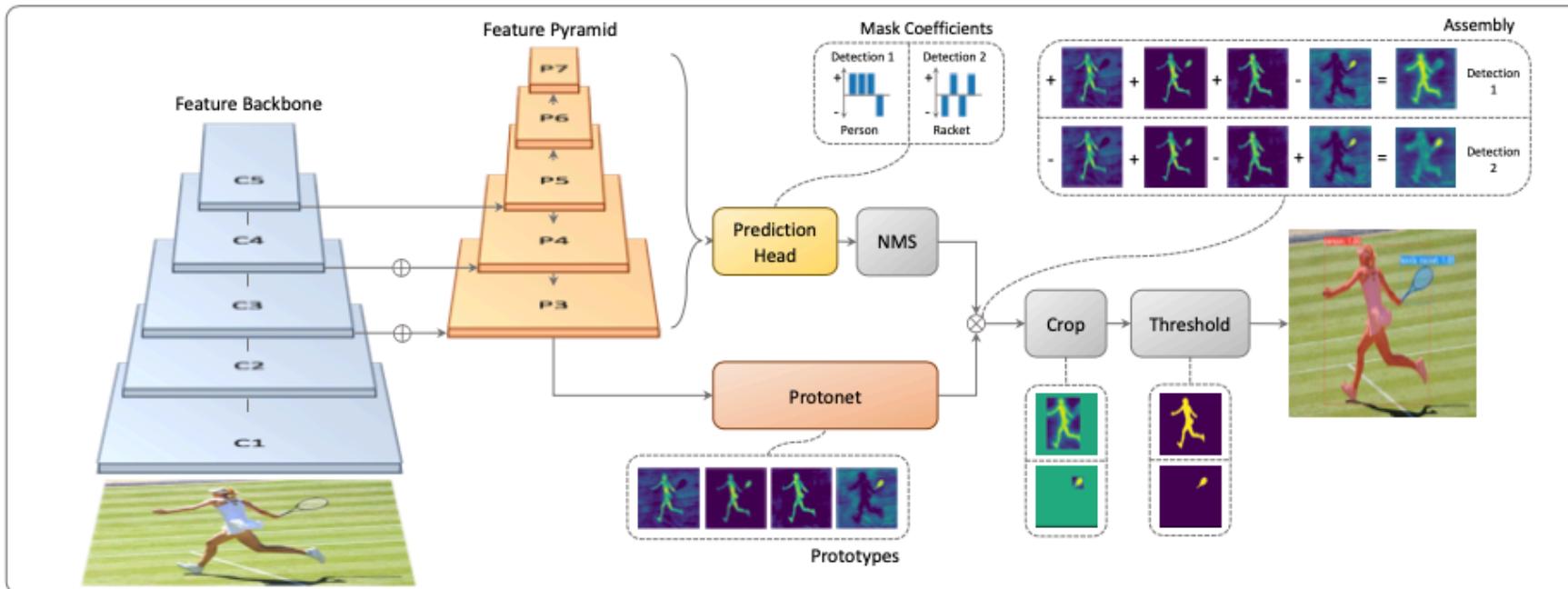
- Implement protonet as an **FCN** whose last layer has k channels.
- Use FPN because its largest feature layers are the deepest.
- Upsample P3 to one fourth the dimensions of the input image to increase performance
- **Protonet output to be unbounded**, overpowering activations for prototypes it is very confident about (e.g., obvious background) → we have the option ReLU / no nonlinearity.

Mask Coefficients



- For mask coefficient prediction, we simply add a third branch in parallel that predicts k mask coefficients, one corresponding to each prototype.
- For nonlinearity, we find it important to be able to subtract our prototypes from the final mask.

Mask Assembly



$$M = \sigma(PC^T)$$

$P : h \times w \times k$ *Prototype mask*

$C : h \times w \times k$ *Mask coefficients*

Loss / Cropping Masks

$$L_{cls} : L_{box} : L_{mask} = 1 : 1.15 : 6.125$$

Loss / Cropping Masks

$$L_{cls} : L_{box} : L_{mask} = 1 : 1.15 : 6.125$$

SSD

Loss / Cropping Masks

$$L_{cls} : L_{box} : \boxed{L_{mask}} = 1 : 1.15 : 6.125$$

$$M_{gt}: L_{mask} = \text{BCE}(M, M_{gt}).$$

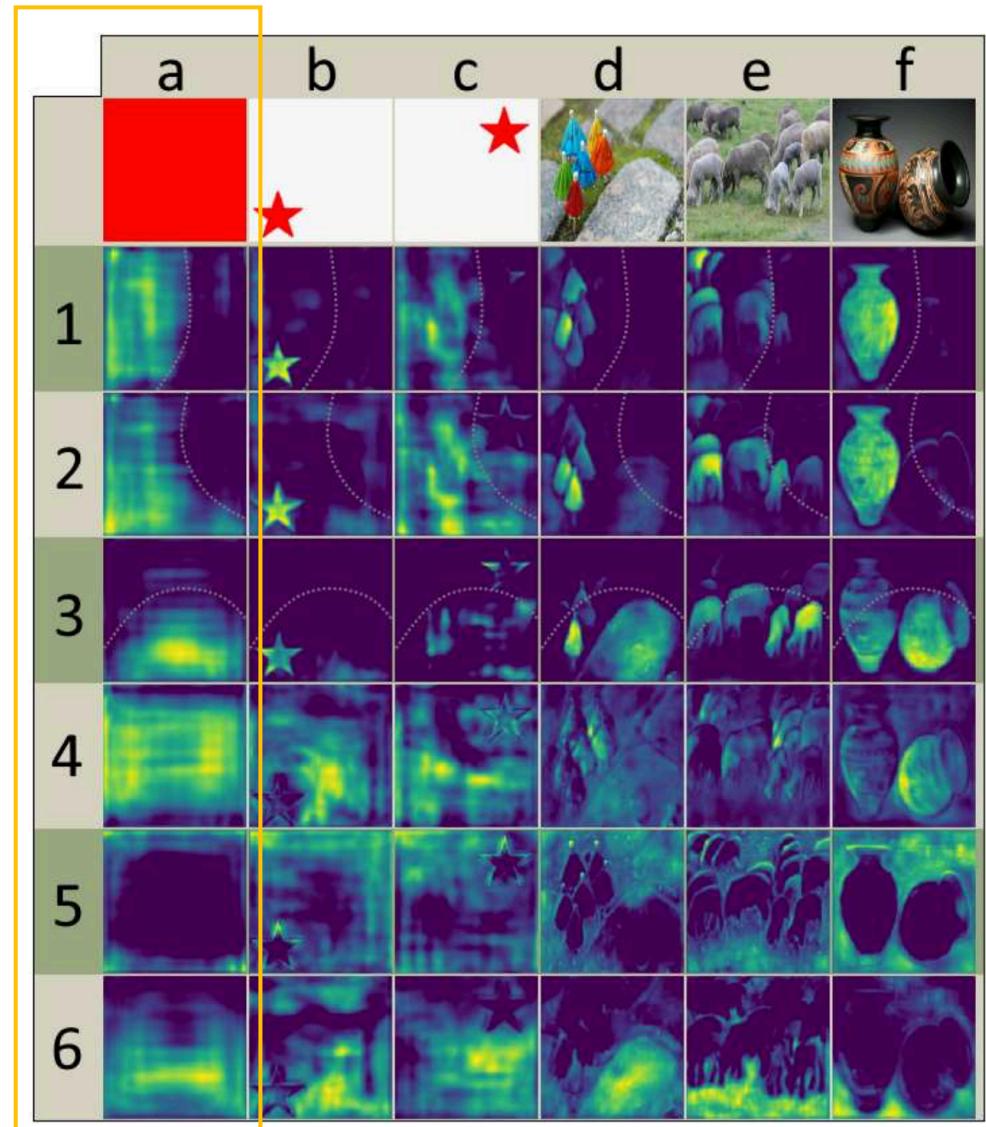
Loss / Cropping Masks

$$L_{cls} : L_{box} : L_{mask} = 1 : 1.15 : 6.125$$

Cropping Masks We crop the final masks with the predicted bounding box during evaluation. During training, we instead crop with the ground truth bounding box, and divide L_{mask} by the ground truth bounding box area to preserve small objects in the prototypes.

To be Translation Variant

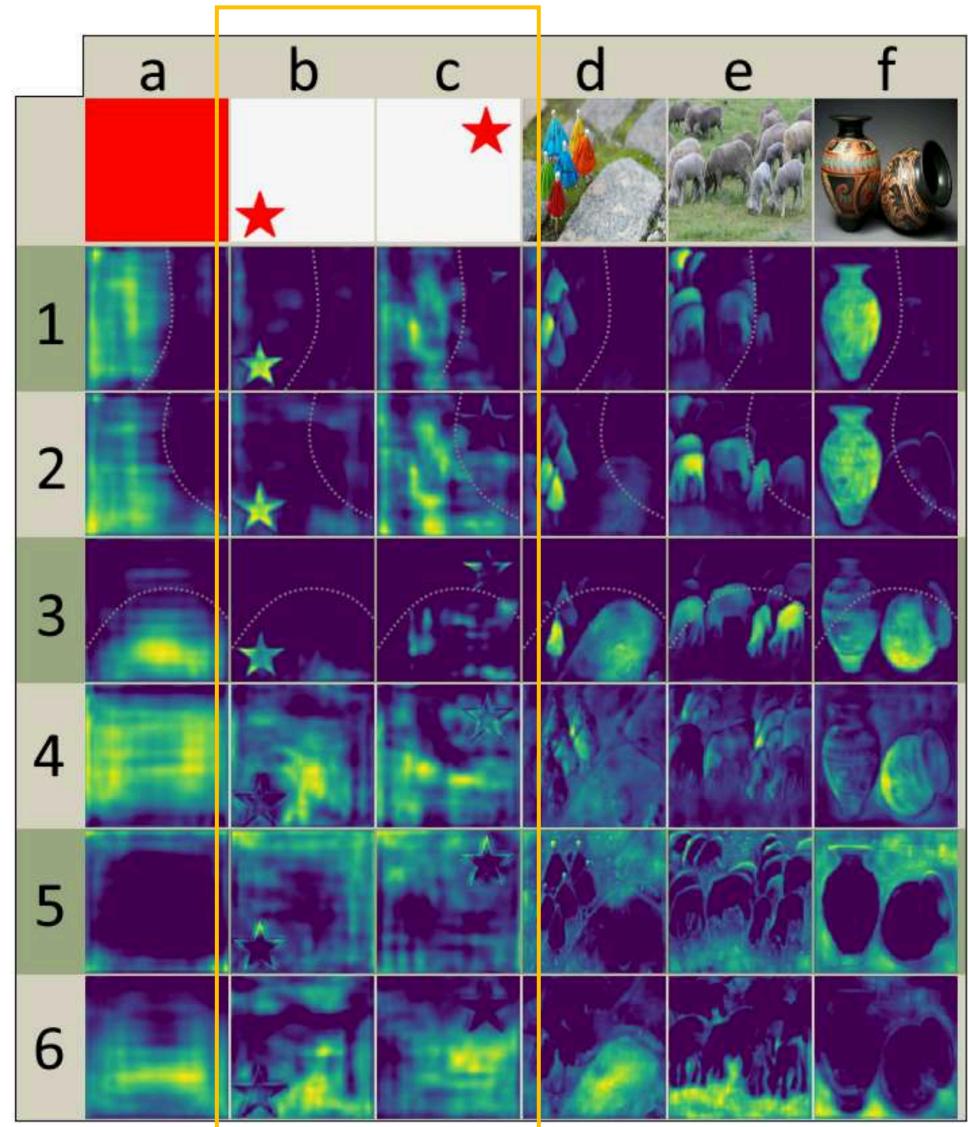
- FCN is translation invariant, the task needs translation variance. (FCIS, Mask R-CNN)
- YOLACT learns how to localize instances on its own via different activations in its prototypes.



To be Translation Variant

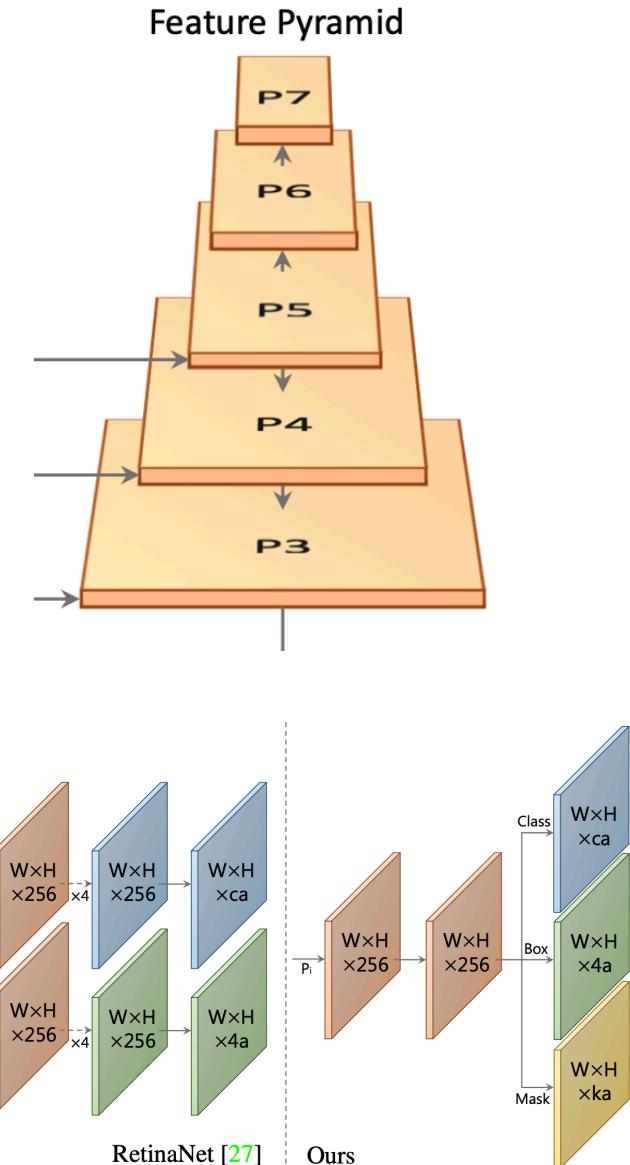
- FCN is translation invariant, the task needs translation variance. (FCIS, Mask R-CNN)
- YOLACT learns how to localize instances on its own via different activations in its prototypes.

conv output will be the same. On the other hand, the consistent rim of padding in modern FCNs like ResNet gives the network the ability to tell how far away from the image's edge a pixel is. Conceptually, one way it could accomplish this is to have multiple layers in sequence spread the padded 0's out from the edge toward the center (e.g., with a kernel like [1, 0]). This means ResNet, for instance, *is inherently translation variant*, and our method makes heavy use of that property (images b and c exhibit clear translation variance).



YOLACT Detector

- Backbone : ResNet-101
- Image size : 550 x 550
- Modify FPN by not producing P2 and producing P6, P7
- 3 anchors : [1, 0.5, 2]
- Smooth-L1 loss to bbox regression.
- Cross entropy with c positive labels and 1 background label.
- Using OHEM with a 3:1(neg:pos)
- Do not use focal loss.



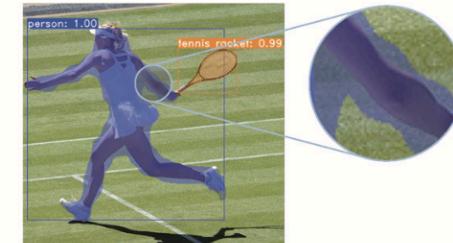
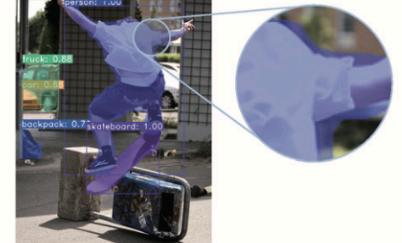
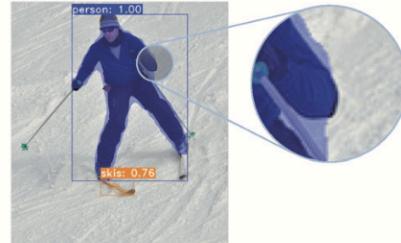
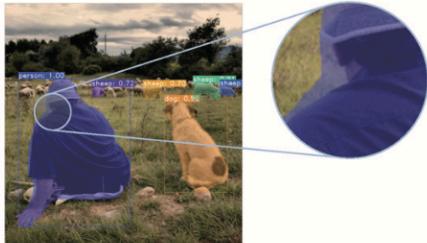
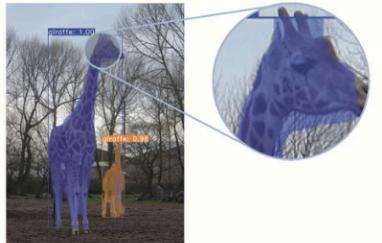
Fast NMS

- A version of NMS where every instance can be decided to be kept or discarded in parallel.
(Fast NMS entirely in standard GPU-accelerated matrix operations)
 1. Compute $c \times n \times n$ pairwise IoU matrix X .
 2. Batched sorting on the GPU is readily available and computing IoU can be easily vectorized.
 3. Remove detections (threshold t)
 - I. Setting the lower triangle and diagonal of X to 0: $X_{kij} = 0, \forall k, j, i \geq j$
 - II. Taking column-wise max $K_{kj} = \max_i(X_{kij}) \quad \forall k, j$
 - III. thresholding with t ($K < t$)

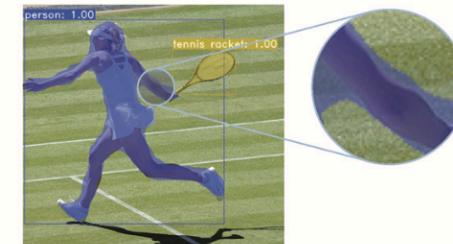
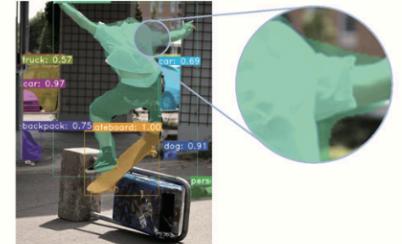
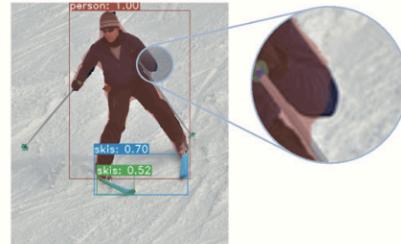
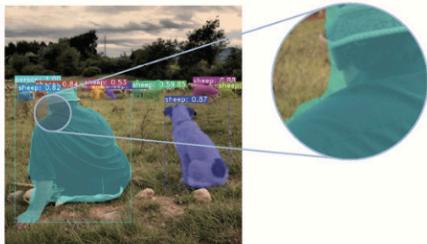
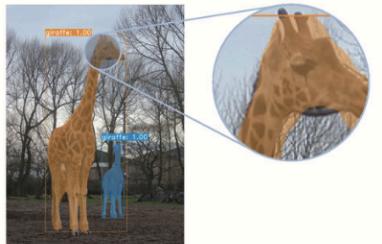
Fast NMS is 15.0 ms faster than traditional NMS with a performance loss of only 0.3 mAP.
(in Mask R-CNN)

Results

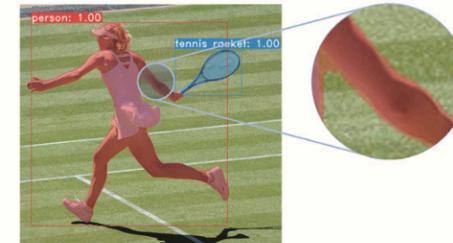
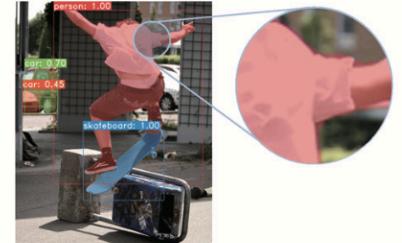
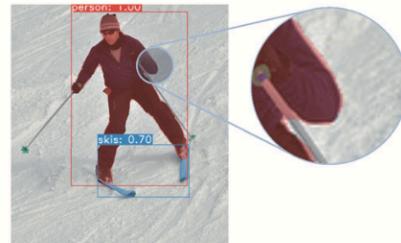
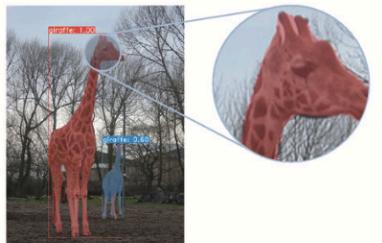
FCIS



Mask R-CNN



Ours



Results

Method	Backbone	FPS	Time	AP	AP ₅₀	AP ₇₅	AP _S	AP _M	AP _L
PA-Net [29]	R-50-FPN	4.7	212.8	36.6	58.0	39.3	16.3	38.1	53.1
RetinaMask [14]	R-101-FPN	6.0	166.7	34.7	55.4	36.9	14.3	36.7	50.5
FCIS [24]	R-101-C5	6.6	151.5	29.5	51.5	30.2	8.0	31.0	49.7
Mask R-CNN [18]	R-101-FPN	8.6	116.3	35.7	58.0	37.8	15.5	38.1	52.4
MS R-CNN [20]	R-101-FPN	8.6	116.3	38.3	58.8	41.5	17.8	40.4	54.4
YOLACT-550	R-101-FPN	33.5	29.8	29.8	48.5	31.2	9.9	31.3	47.7
YOLACT-400	R-101-FPN	45.3	22.1	24.9	42.0	25.4	5.0	25.3	45.0
YOLACT-550	R-50-FPN	45.0	22.2	28.2	46.6	29.2	9.2	29.3	44.8
YOLACT-550	D-53-FPN	40.7	24.6	28.7	46.8	30.0	9.5	29.6	45.5
YOLACT-700	R-101-FPN	23.4	42.7	31.2	50.6	32.8	12.1	33.3	47.1

Results

Method	NMS	AP	FPS	Time
YOLACT	Standard	30.0	24.0	41.6
	Fast	29.9	33.5	29.8
Mask R-CNN	Standard	36.1	8.6	116.0
	Fast	35.8	9.9	101.0

k	AP	FPS	Time
8	26.8	33.0	30.4
16	27.1	32.8	30.5
*32	27.7	32.4	30.9
64	27.8	31.7	31.5
128	27.6	31.5	31.8
256	27.7	29.8	33.6