# Blog Post#2

Ezgi Gümüştekin & Emir Yorgun

Mask R-CNN is a deep neural network that aims to carry out Instance Segmentation. It can separate different objects in an image or video.
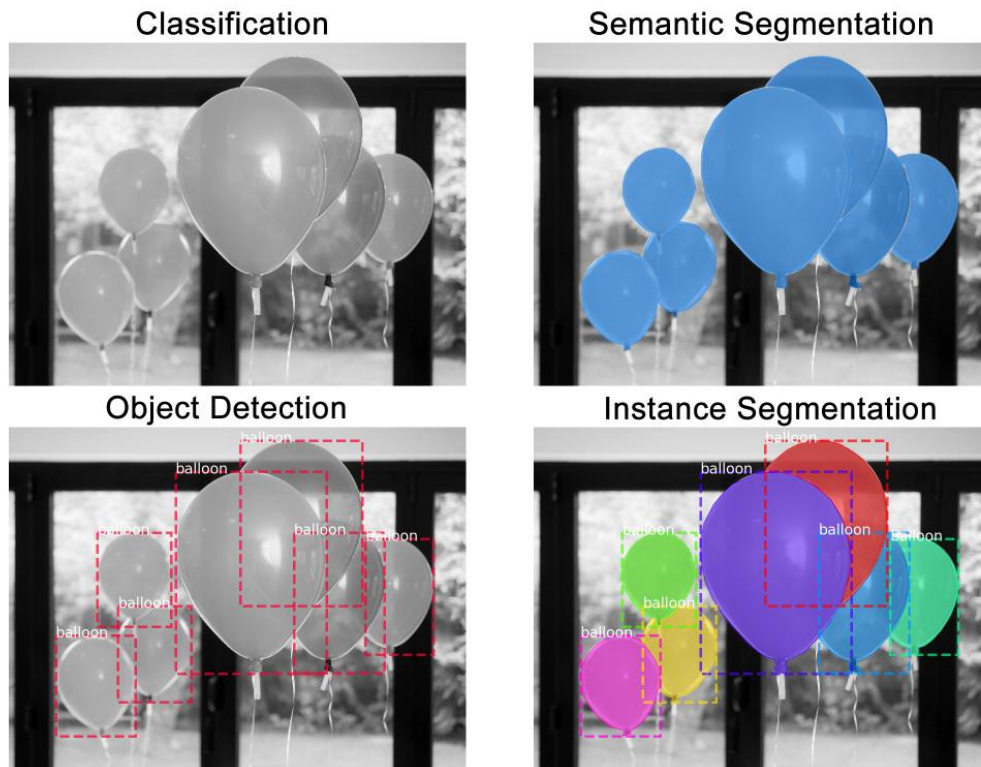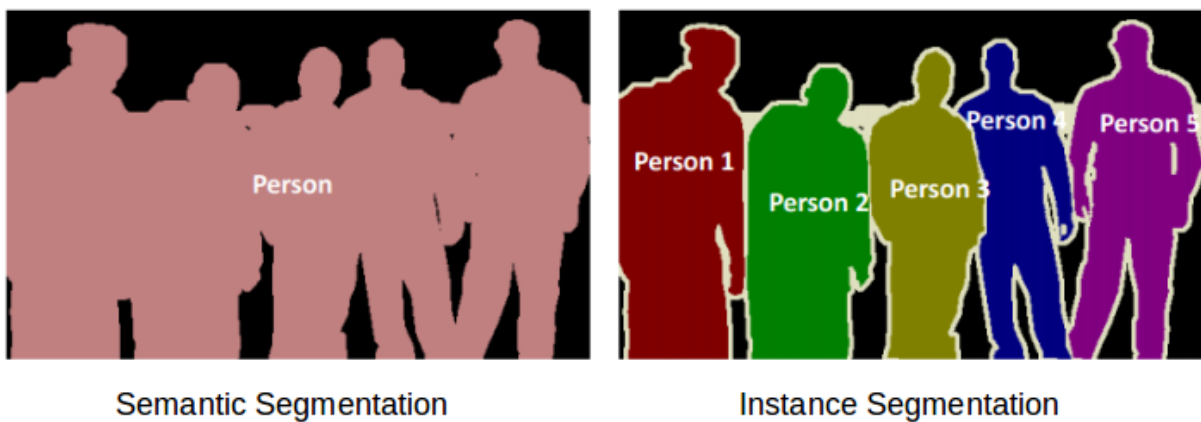


Fig. [1] - Instance Segmentation



Fig. [2] - Instance Segmentation

It comprises of; Backbone Network, Region Proposal Network, Mask Representation, RoI Align.

This, Backbone Network, multi-layer feature pyramid network generates RoI of different scale which improves the accuracy of previous ResNet architecture.
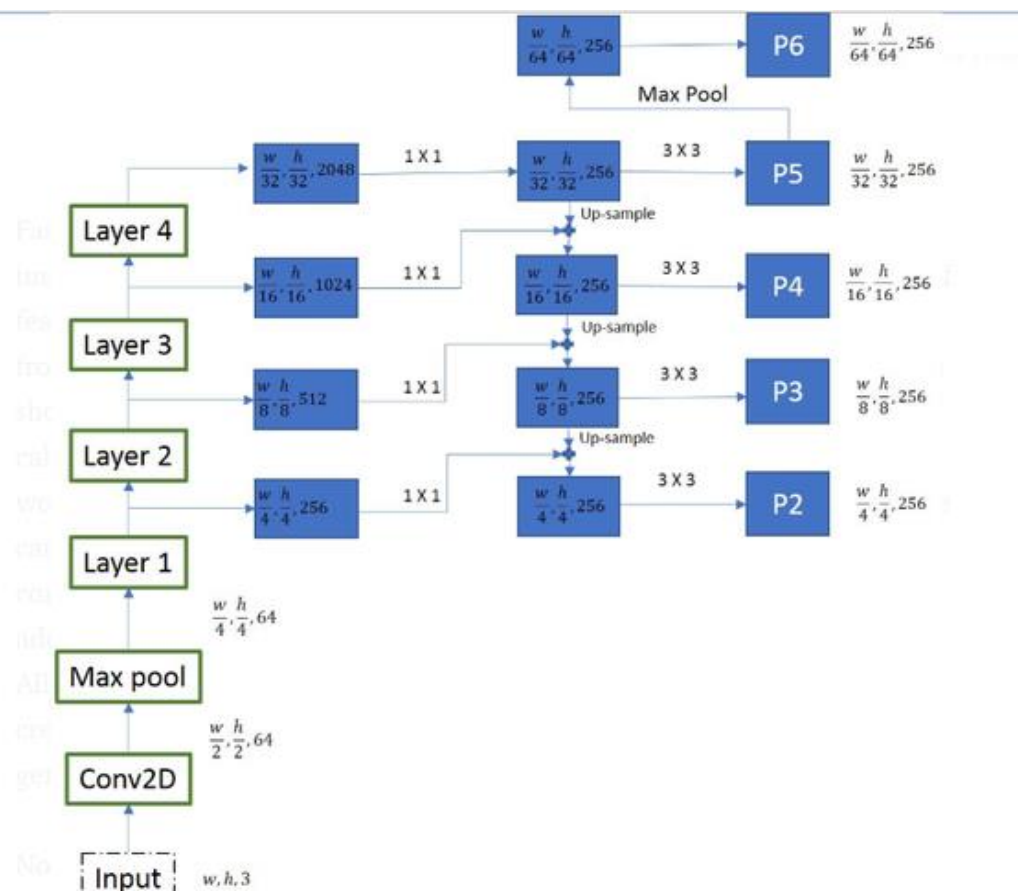


Fig. 3 - Mask R-CNN Backbone Architecture [4]

Region Proposal Network, all the convolution feature map that is generated by the previous layer is passed through a 3*3 convolution layer. The output of this is then passed into two parallel branches that determine the objectness score and regress the bounding box coordinates.
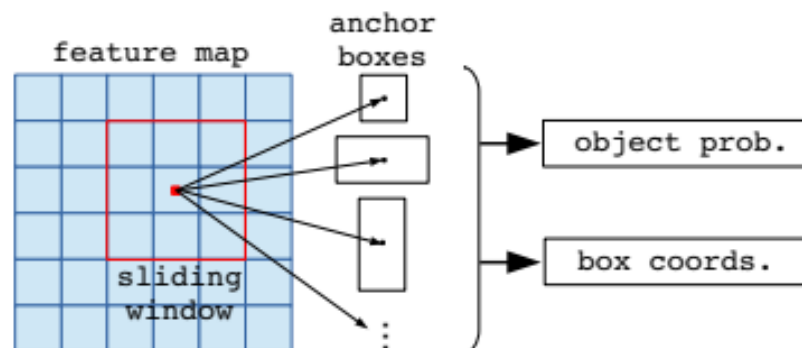


Fig 4. - Region Proposal Network - Anchor Generation Mask R-CNN [4]

Mask Representation, is a small Fully Convolutional Network (FCN) applied to each Region of Interest (RoI). The FCN outputs a fixed-size mask for each RoI and each class. The mask is used to predict the pixel-level segmentation of the object instances in the test images.

RoI Align, is a technique that improves the accuracy of object detection and instance segmentation fmodels in RoI Pooling.
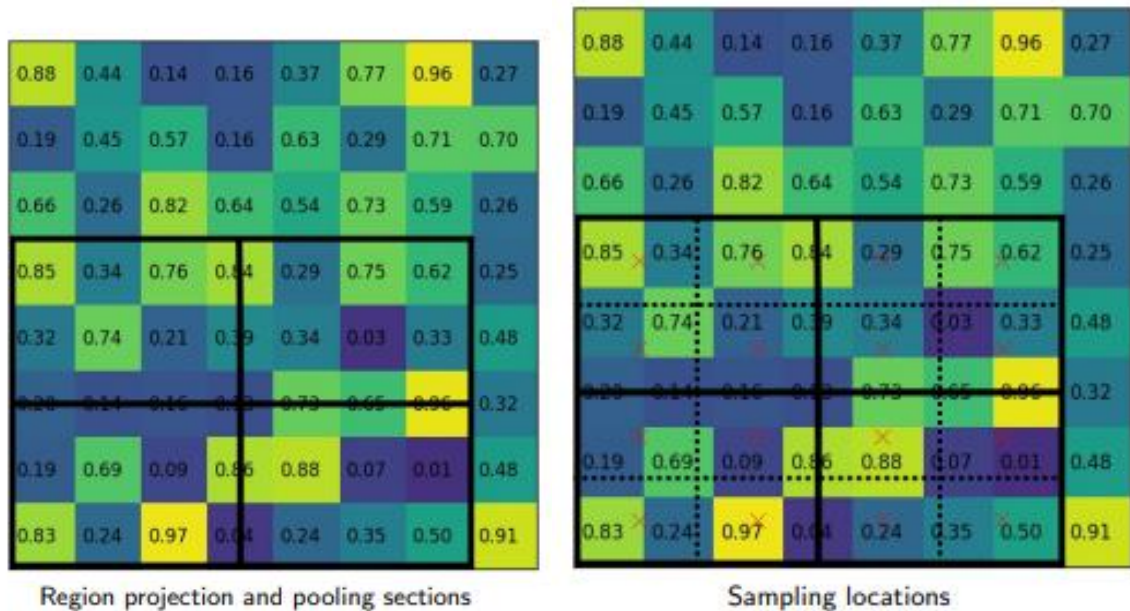


Fig. 5 - ROI Align [4]

**References**

[1] https://merveelifsarac.medium.com/cnn-r-cnn-fast-r-cnn-mask-r-cnn

[2] https://www.analyticsvidhya.com/blog/2019/07/computer-vision-implementing-mask-r-cnn-image-segmentation

[3] https://www.geeksforgeeks.org/r-cnn-vs-fast-r-cnn-vs-faster-r-cnn-ml

[4] https://www.geeksforgeeks.org/mask-r-cnn-ml

[5] https://arxiv.org/abs/1703.06870