

# CV Assignment 4



Name 1: Abdelrahman Salem Mohamed

ID 1: 6309

Name 2: Reem Abdelhalim

ID 2: 6114

# Object Detection

In this assignment, you will work on Cocco dataset which is a large-scale object detection, segmentation, and captioning dataset. You are required to run 3 different object detectors on this dataset. We will learn to use and differentiate between the architectures.

Either PyTorch or TensorFlow are allowed to run your models.

## Dataset

We test models using coco 2017 validation set which consists of 5000 different images with their meta data from segmentation points and object bounding boxes.

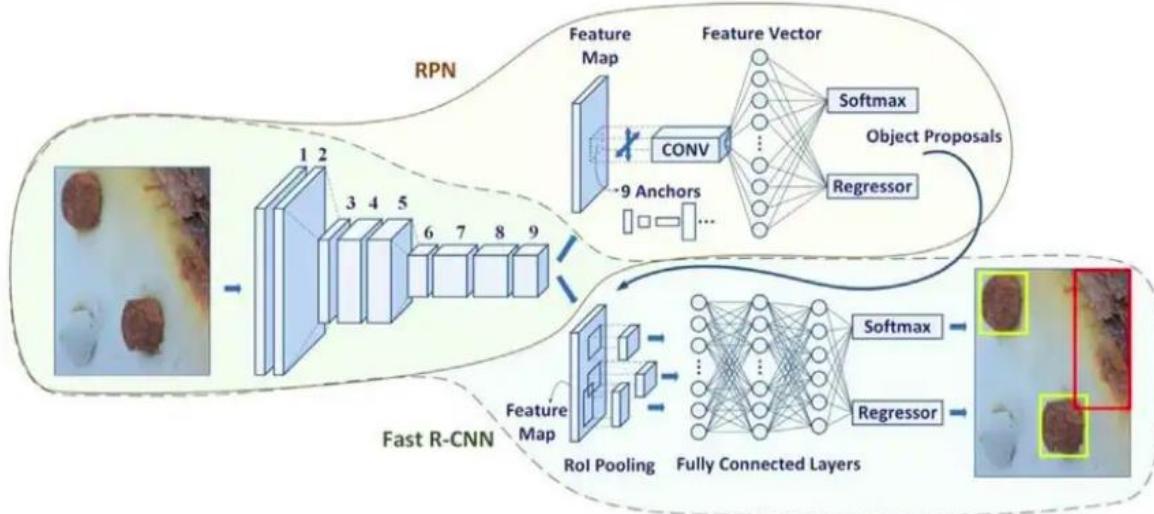
## Models

We used three different models:

- Faster R-CNN (Two stage Model)
- RetinaNet One stage Model)
- Single shot detector (One stage Model)

### 1- Faster R-CNN

Architecture diagram



### Discussion

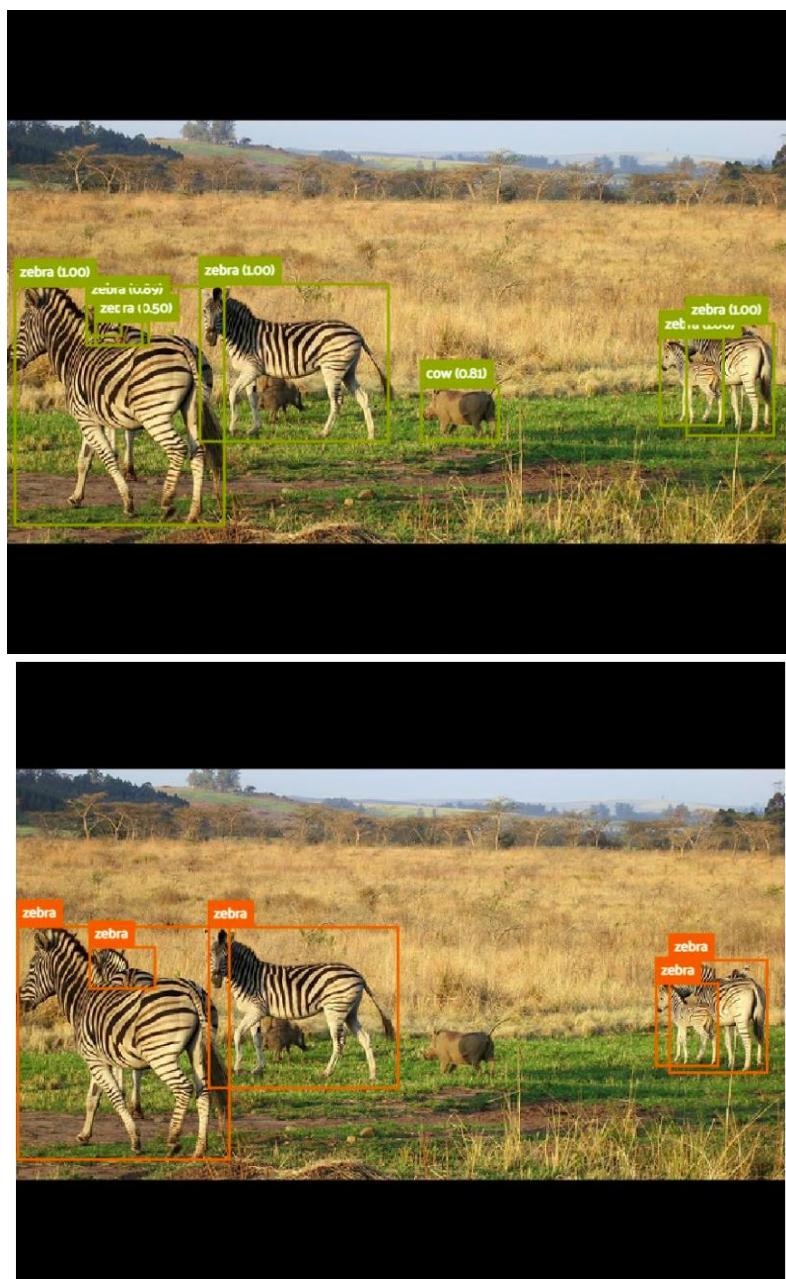
Two stage based model where, first the model learn how to extract Region of proposals, then for every Region of proposal fe perform the normal Fast R-CNN network from ROI

pooling then from fully connected layer compute the two losses one for object box (regression), and one for class type (classification).

Each network has its own loss functions and the whole system weights consider the both losses path, and during the test the network needs first to identify the ROI and apply Object detection technique, so Faster R-CNN consider slower than one stage models but have better mAP.

The backbone network is usually a dense convolutional network like ResNet or VGG16

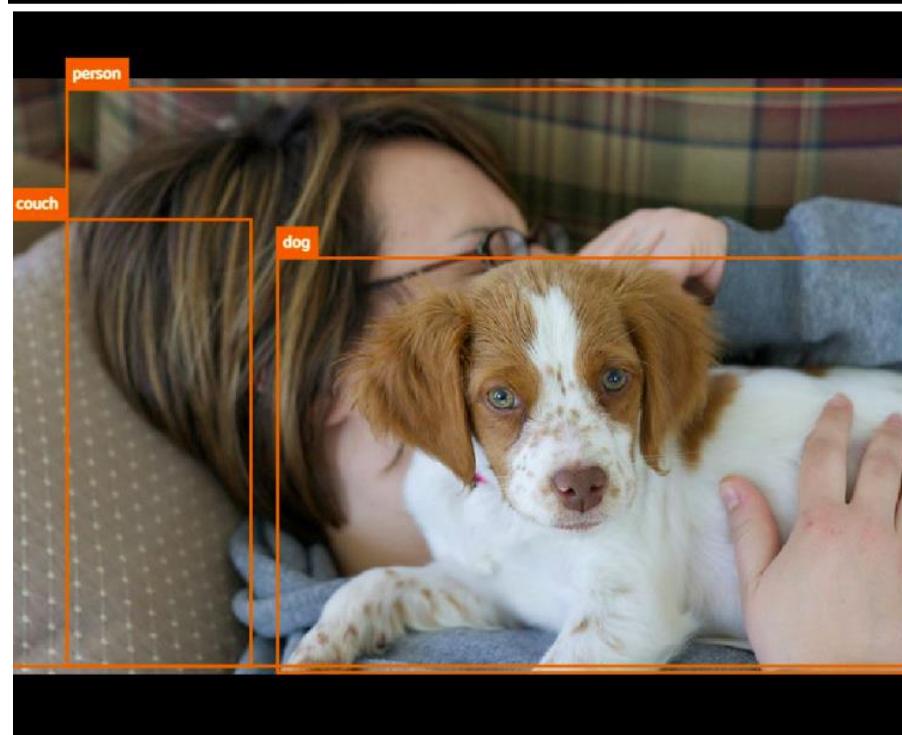
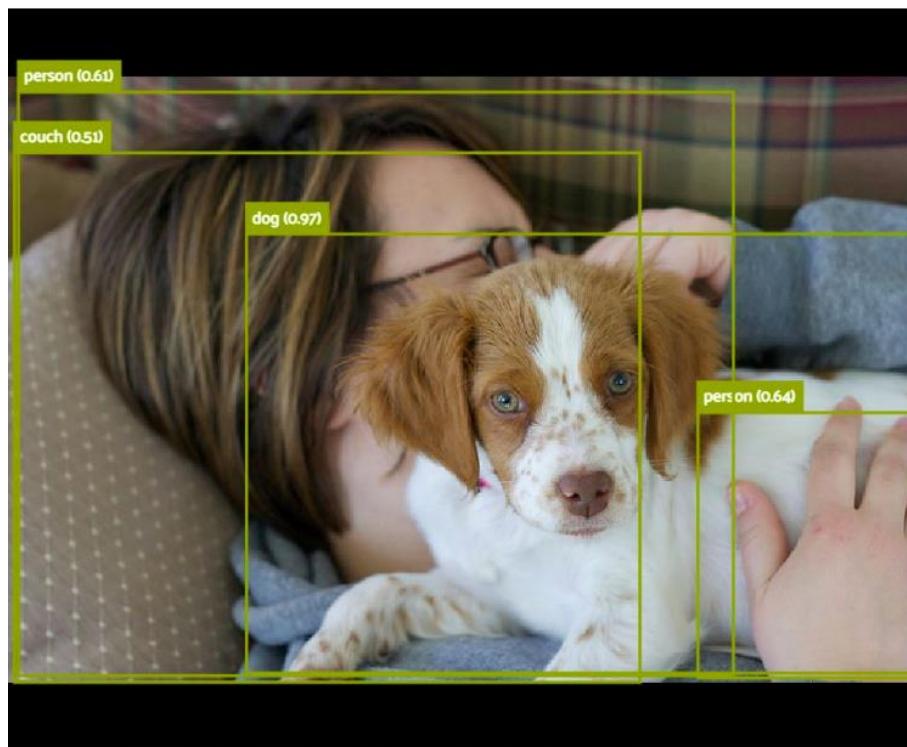
## Results on COCO dataset



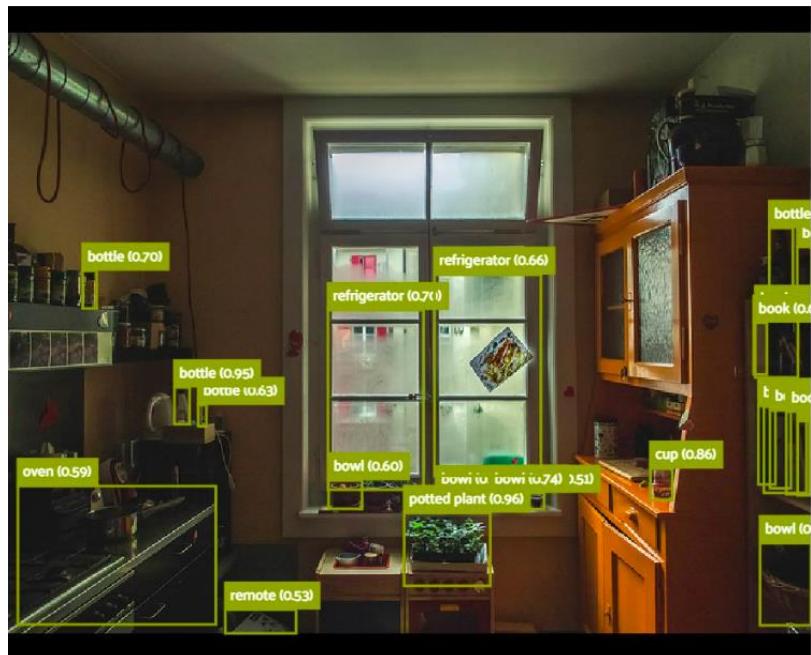
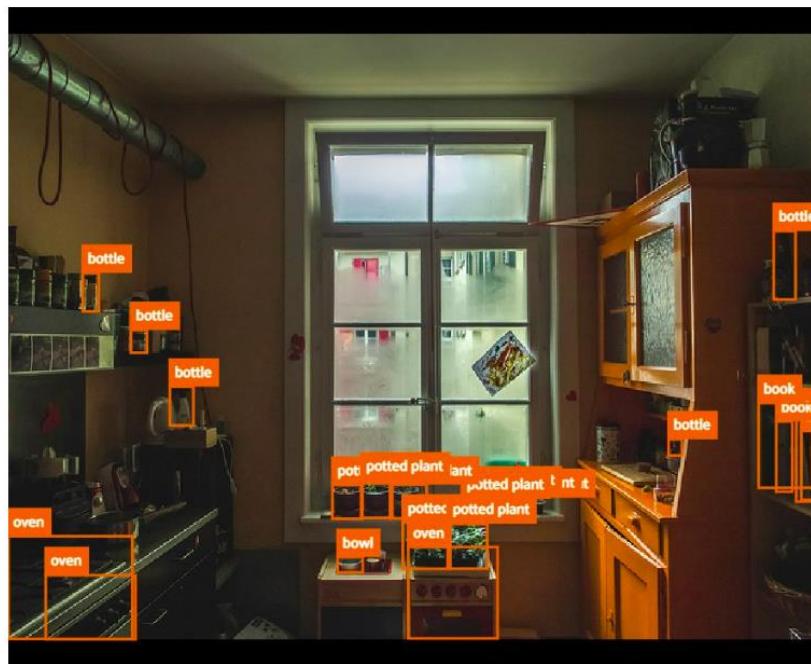
Note: The Model detected all objects and detected not annotated cow/pig



Note: The Model Detected the all object exactly as the GT with 100% confidence



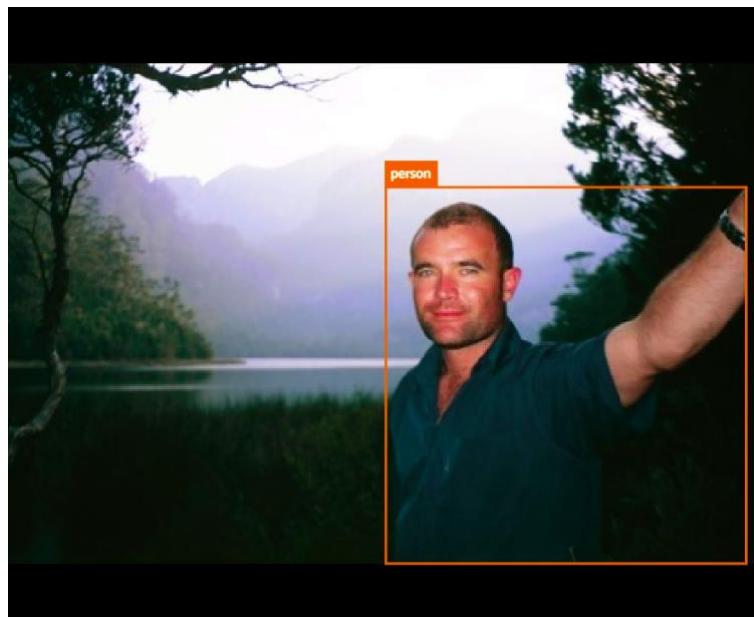
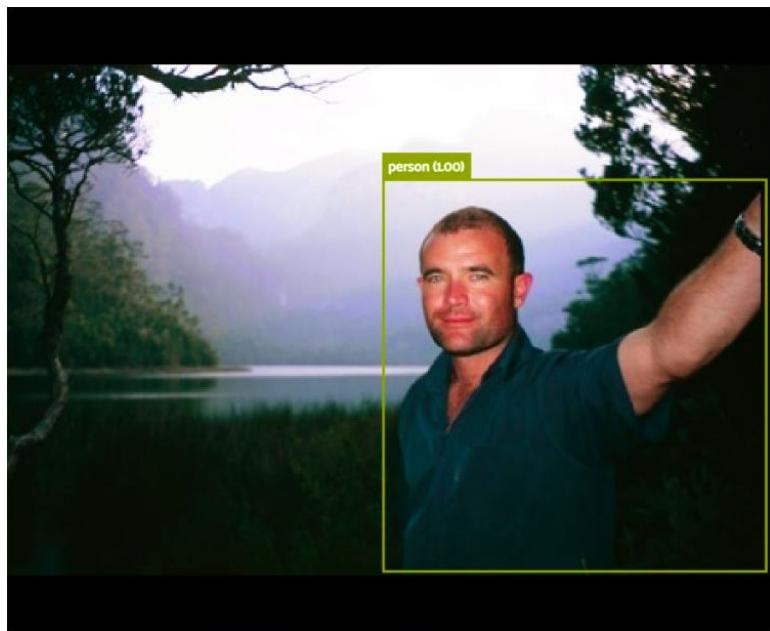
Note: The model detected the hand as human with is not exactly correct but detect the couch with low confidence and existence of person there

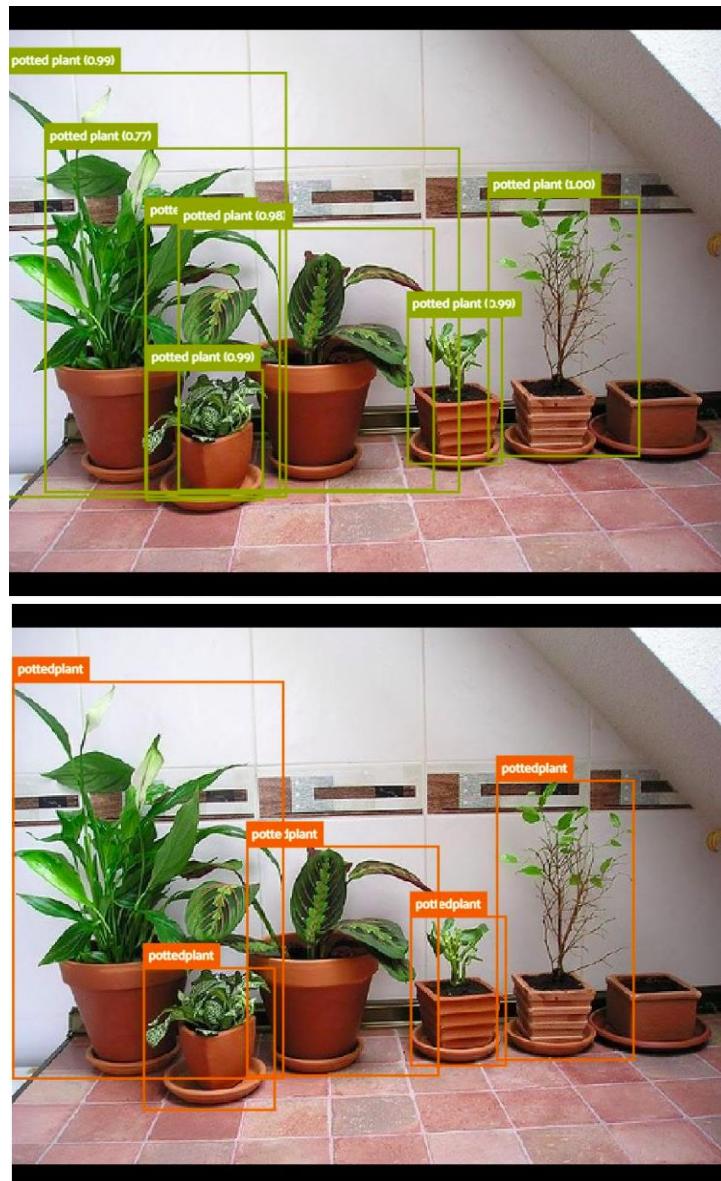


Note: As the image is so crowded, the model predicted most of the object but detected the window as refrigerator and some bottle as cup etc.

## Results on Pascal VOC 2012 dataset

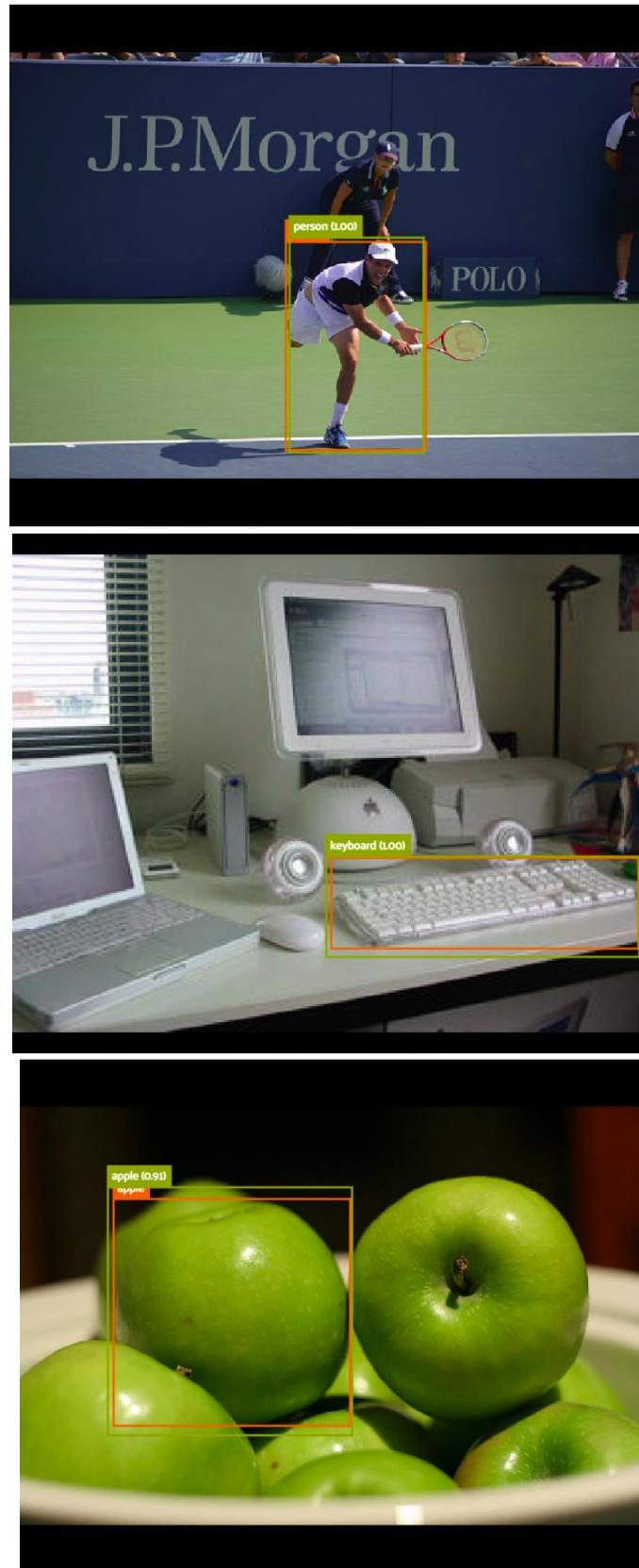


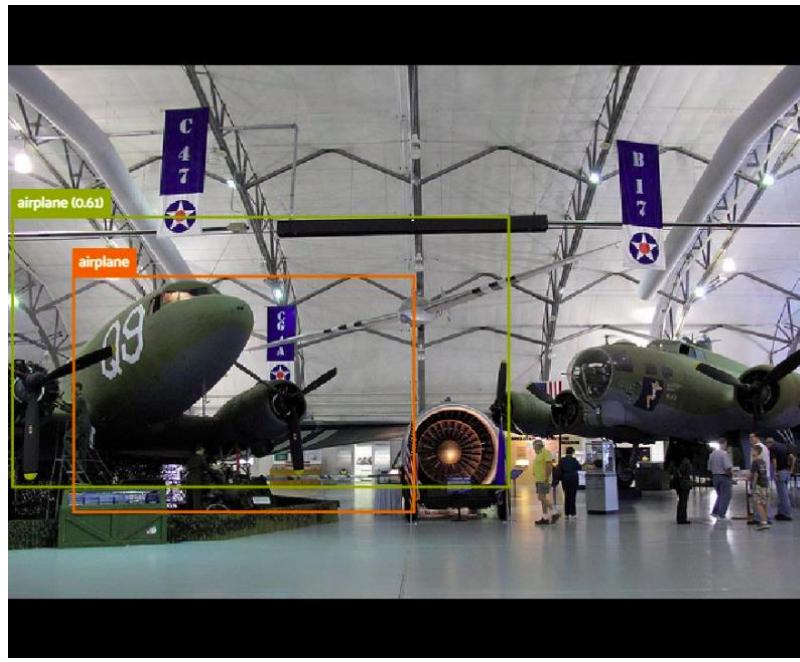




Note: Some bounding boxes are overlapped in the last image but the previous ones are correct with good match

## COCO IOU



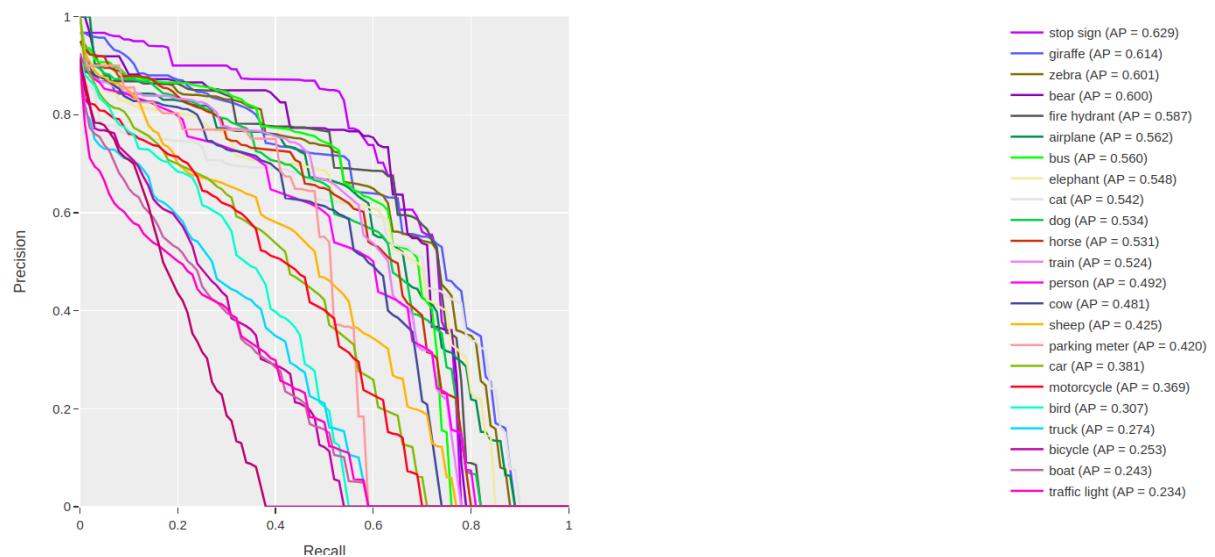


Note: for specific boxes and their GT we compare the IOU, as shown most of them are exactly matched except the last one the box is badly matched .

## mAP and IOU

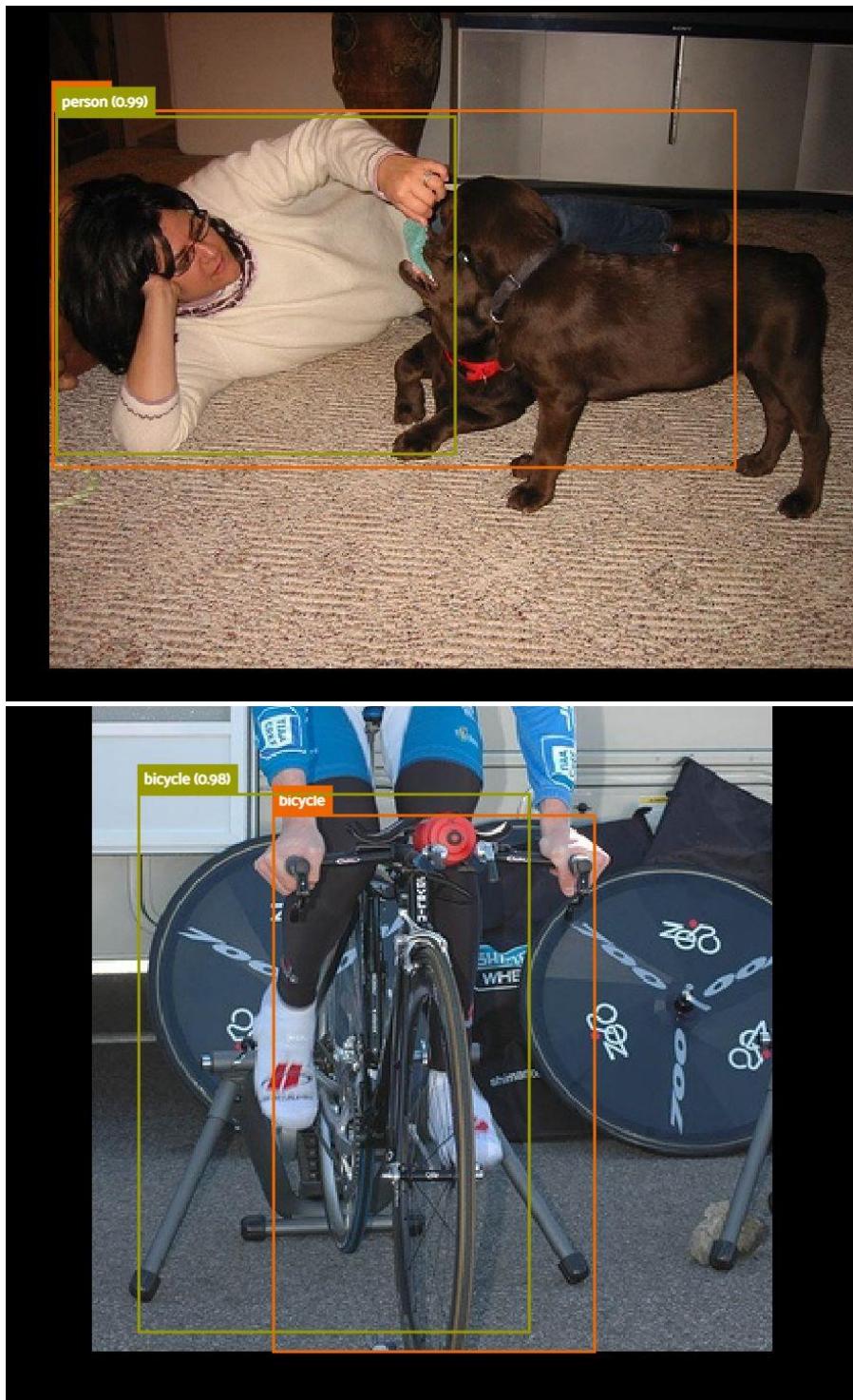
For the whole validation set the **mAP = 0.344** and **IOU = 0.407**

## Precision vs Recall Graph



## COCO IOU



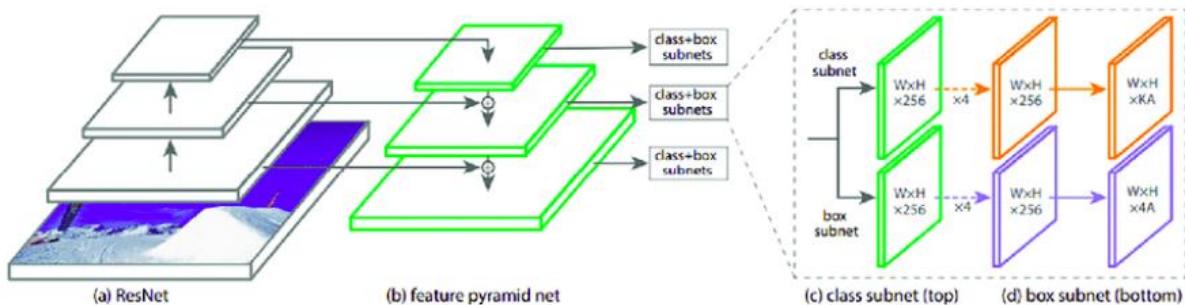


## mAP and IOU

For the whole validation set the **mAP = 0.348** and **IOU = 0.252**

## 2- RetinaNet

Architecture diagram



### Discussion

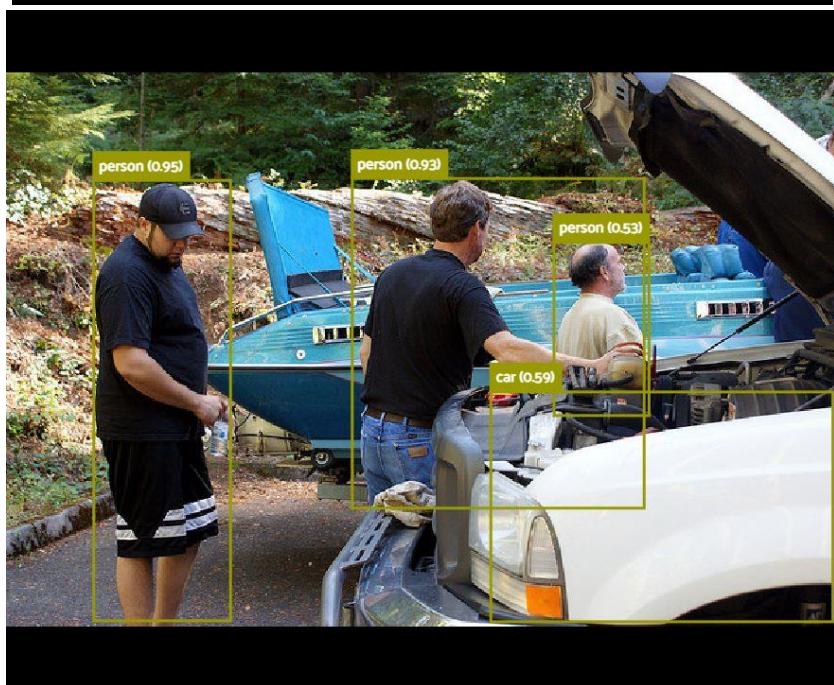
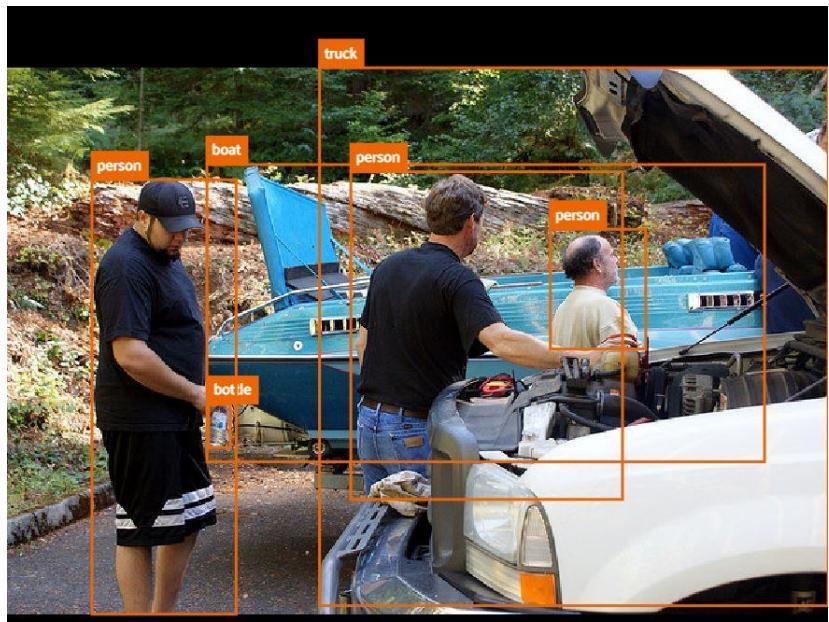
RetinaNet is one of the best one-stage object detection models that has proven to work well with **dense and small scale objects**. For this reason, it has become a popular object detection model to be used with aerial and satellite imagery.

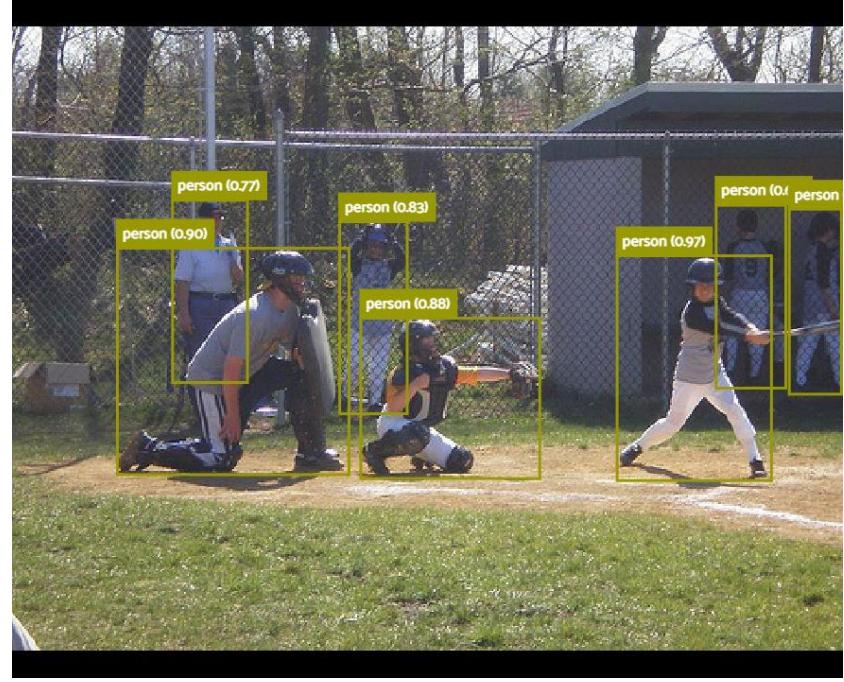
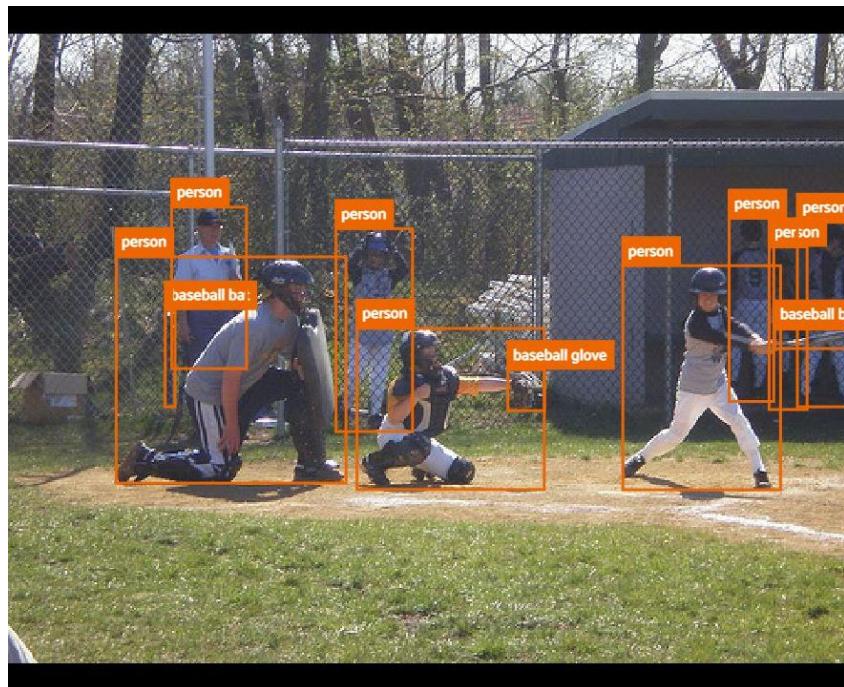
There are four major components of a RetinaNet model architecture :

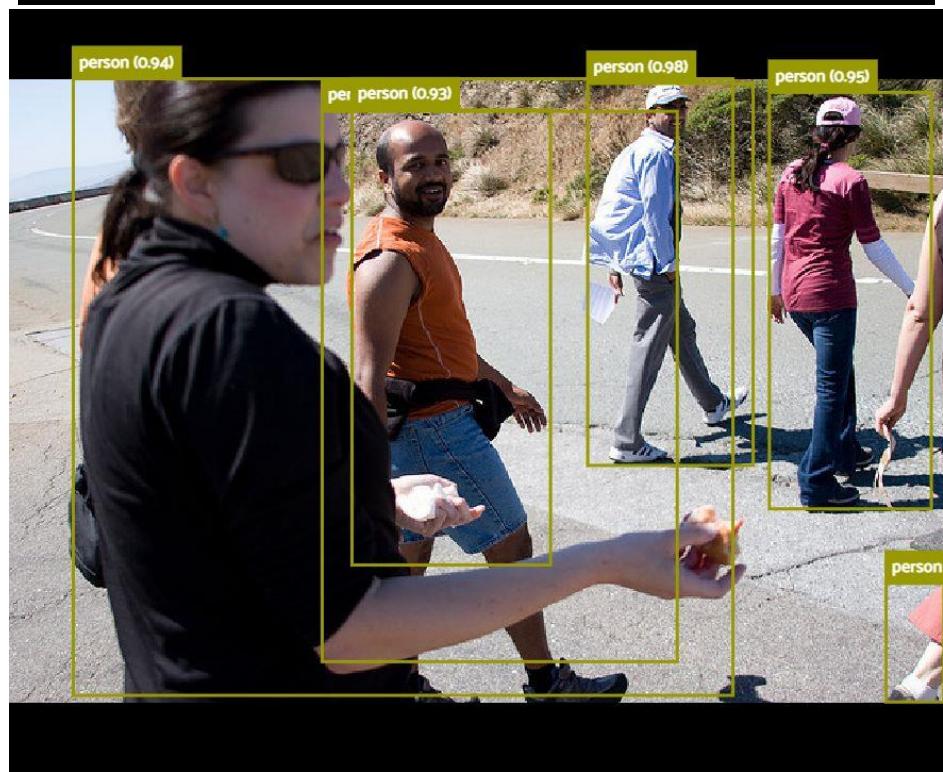
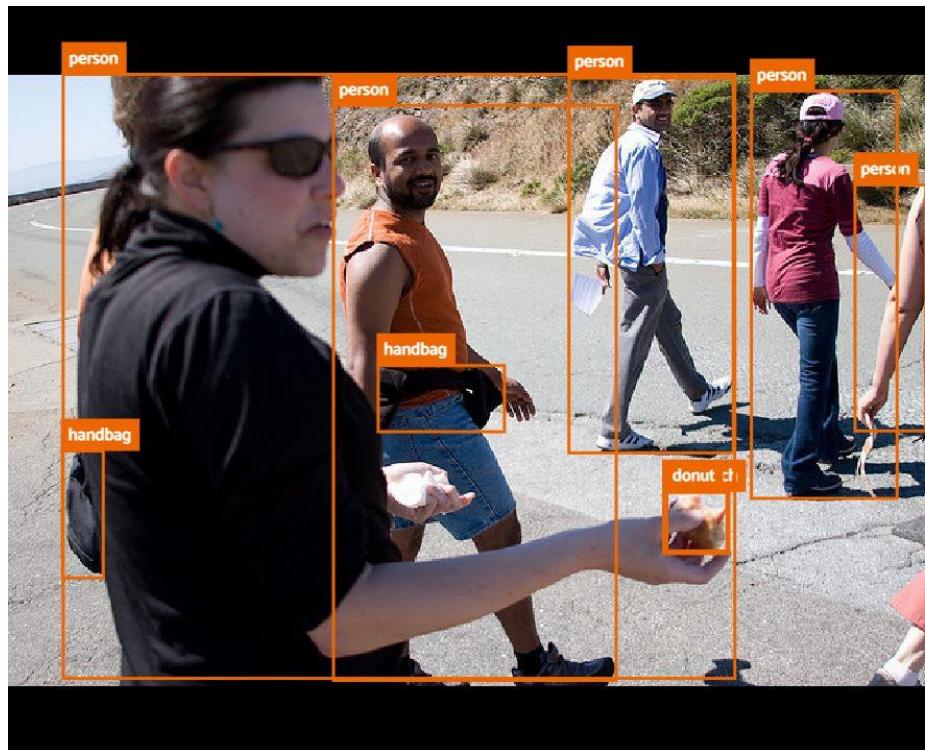
- Bottom-up Pathway - **The backbone network (e.g. ResNet)** which calculates the feature maps at different scales, irrespective of the input image size or the backbone.
- Top-down pathway and Lateral connections - The top down pathway upsamples the spatially coarser feature maps from higher pyramid levels, and the lateral connections merge the top-down layers and the bottom-up layers with the same spatial size.
- Classification subnetwork - It predicts the probability of an object being present at each spatial location for each anchor box and object class.
- Regression subnetwork - It regresses the offset for the bounding boxes from the anchor boxes for each ground-truth object.

## Results on COCO dataset

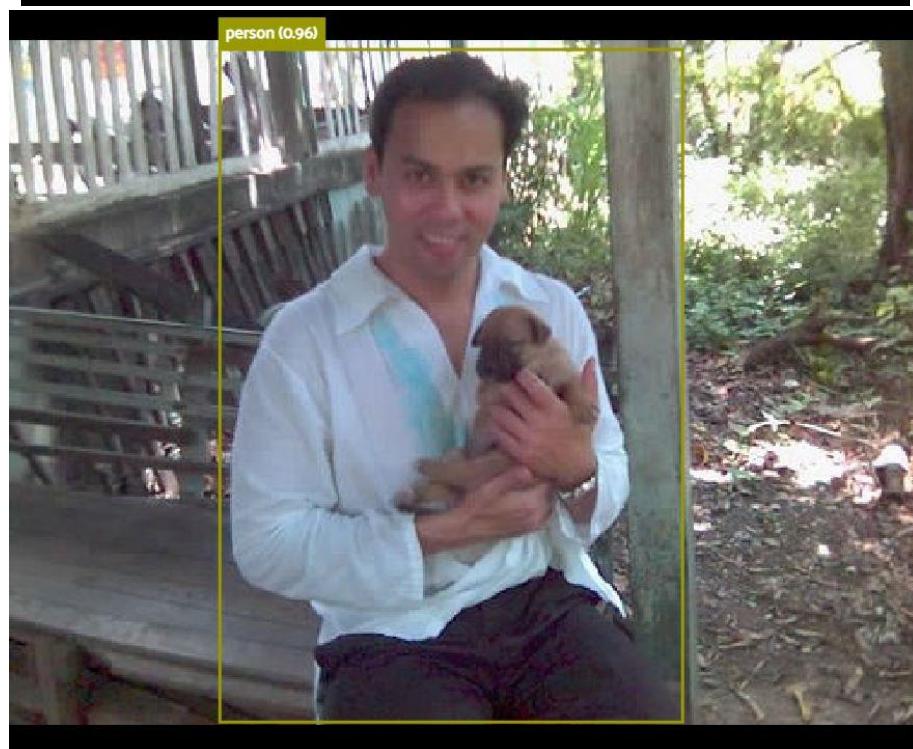
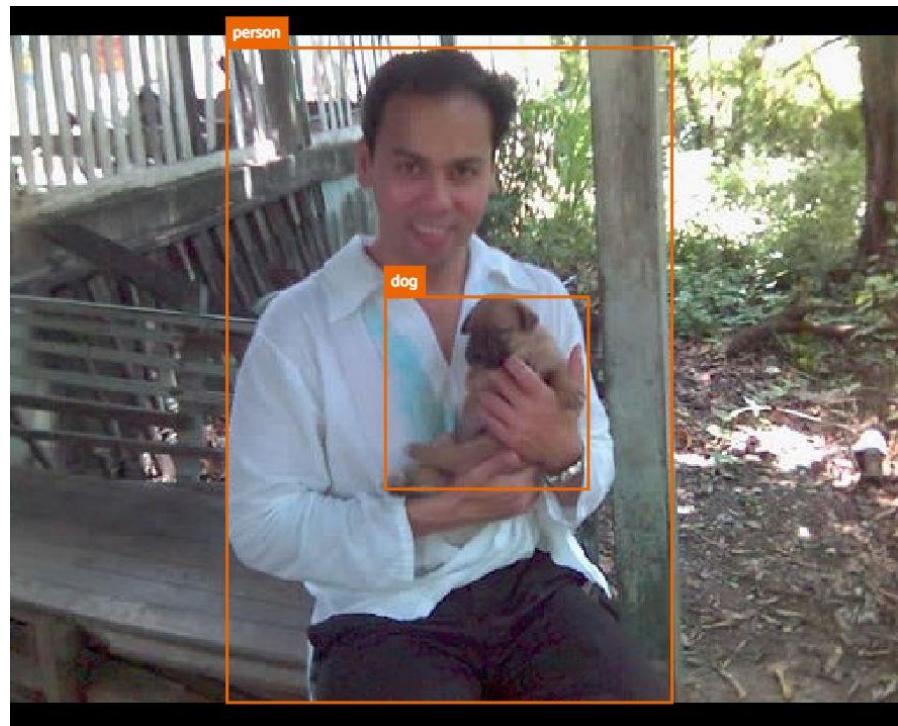


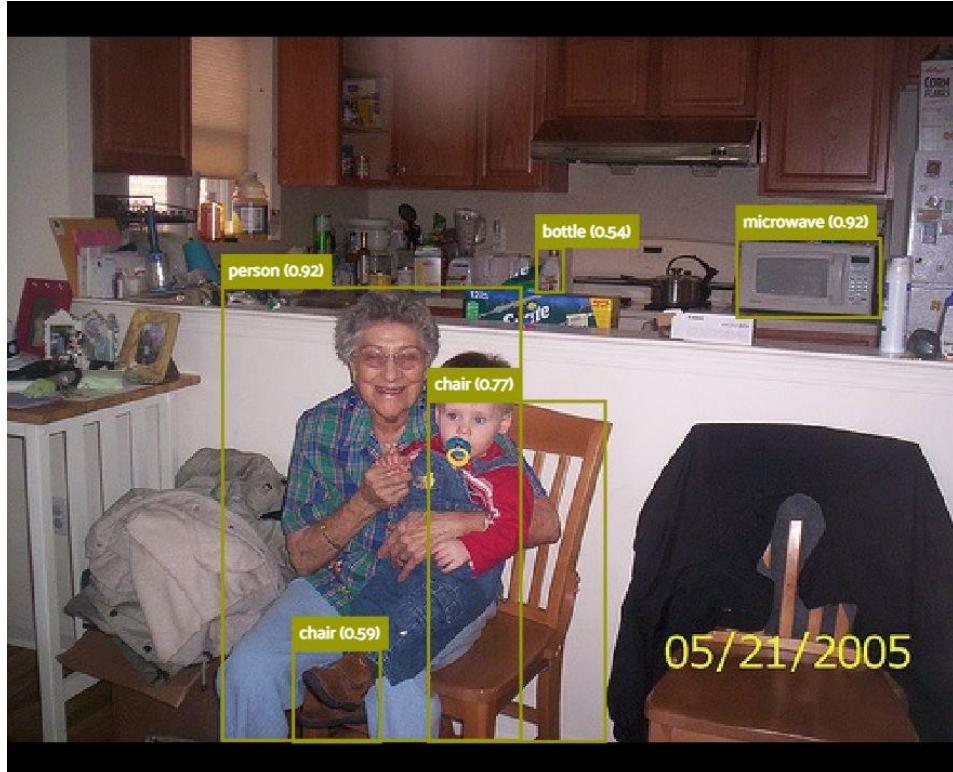
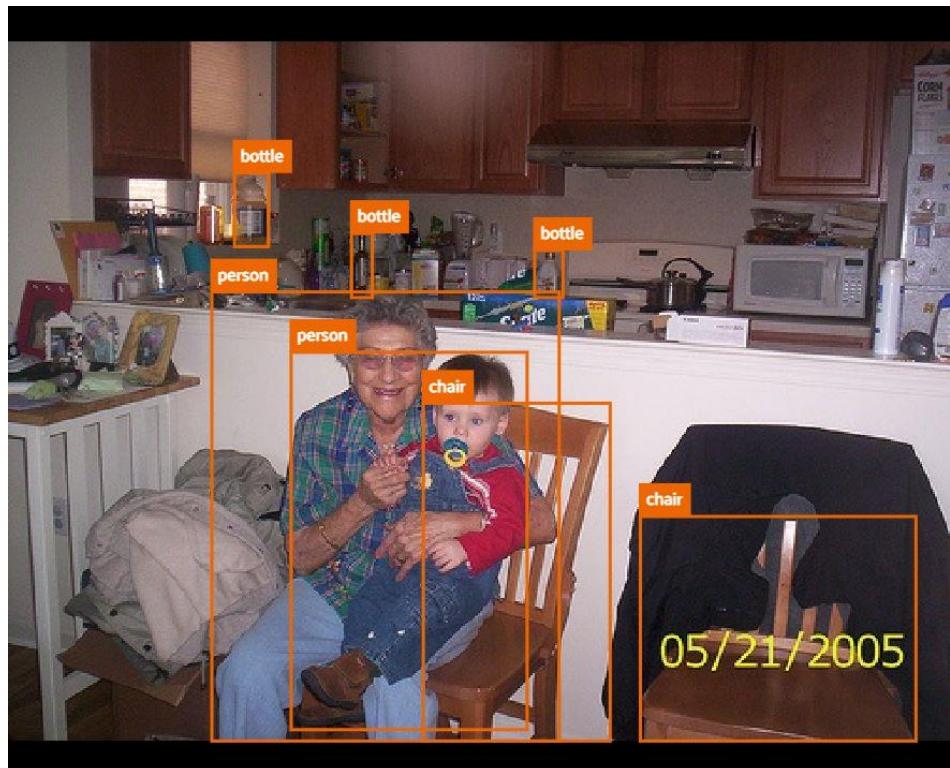




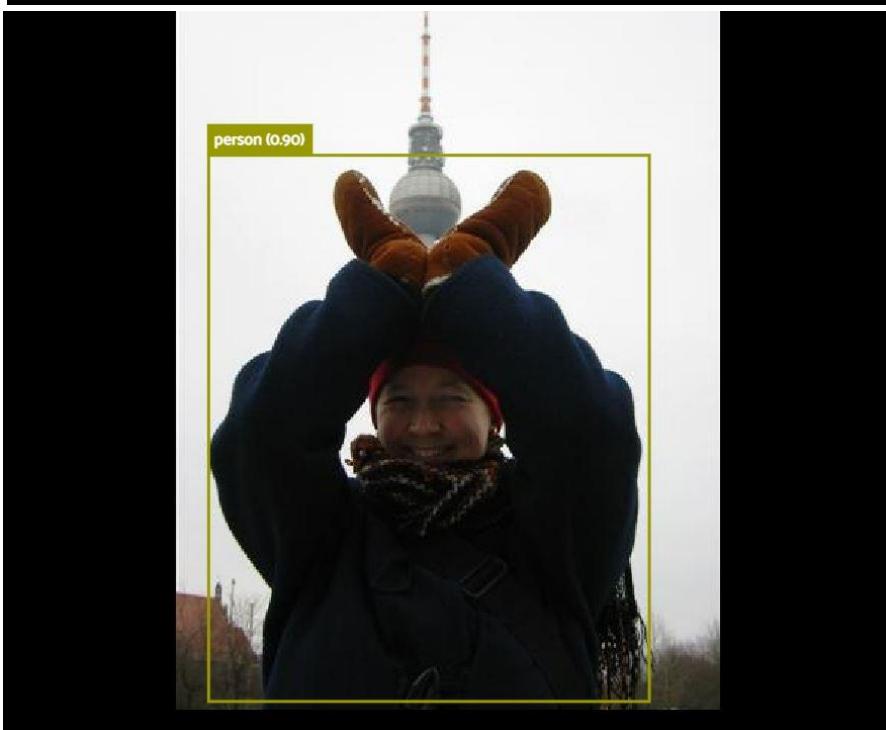
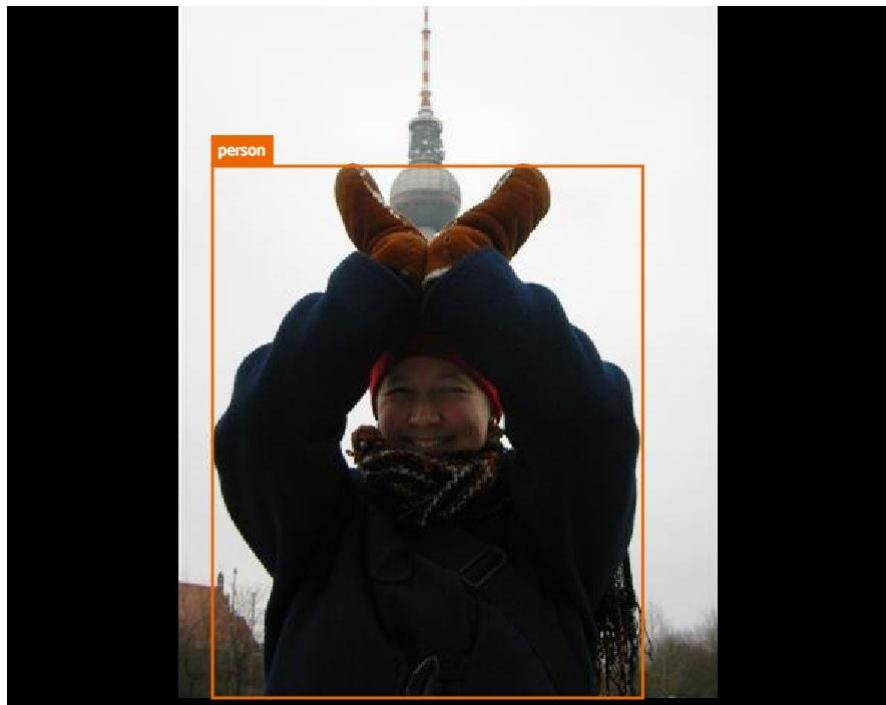


## Results on different dataset

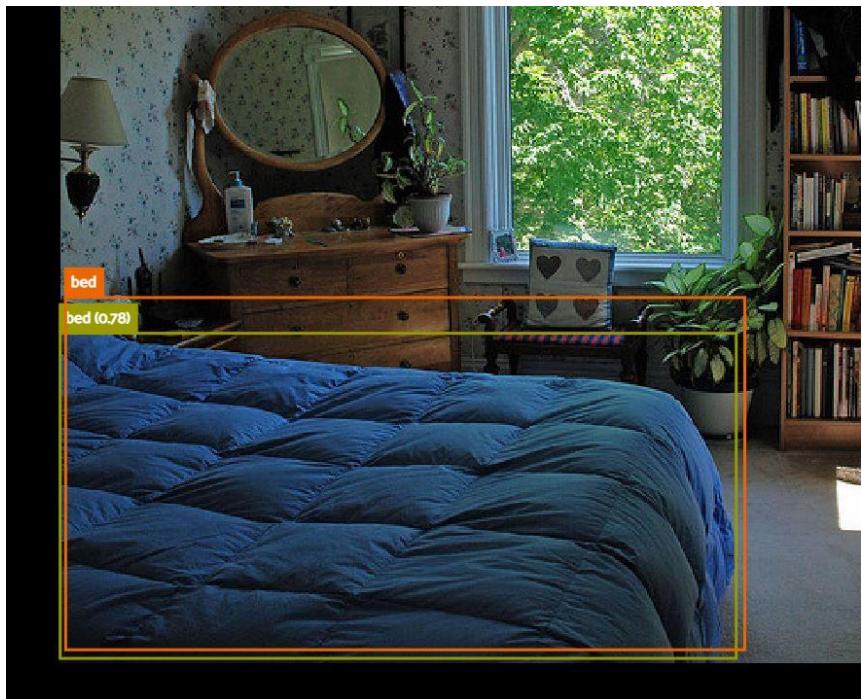








## COCO IOU

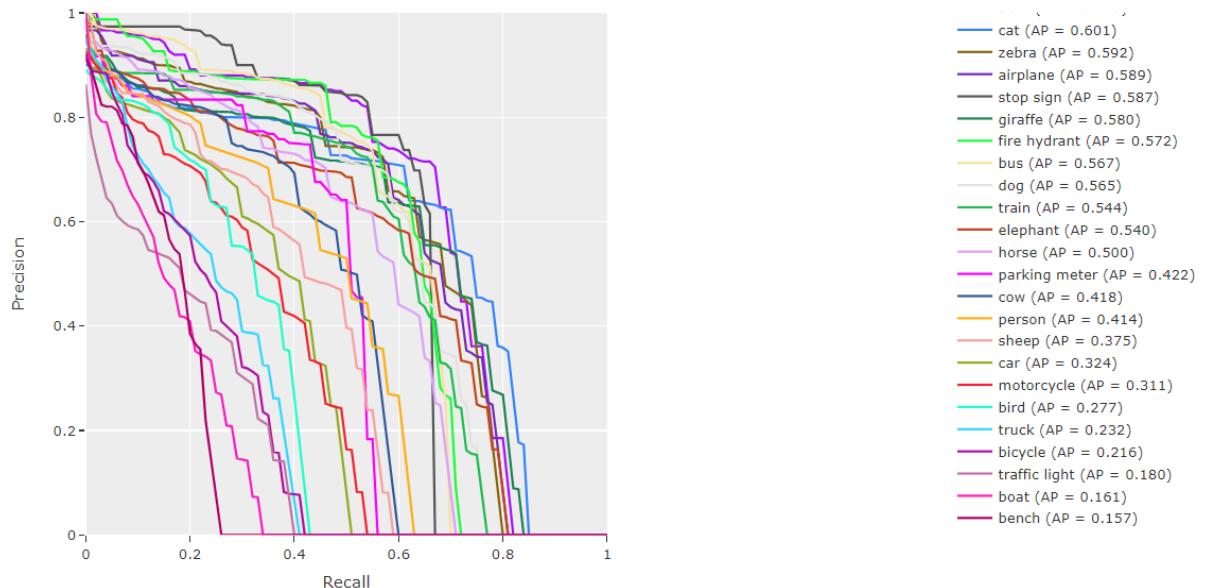




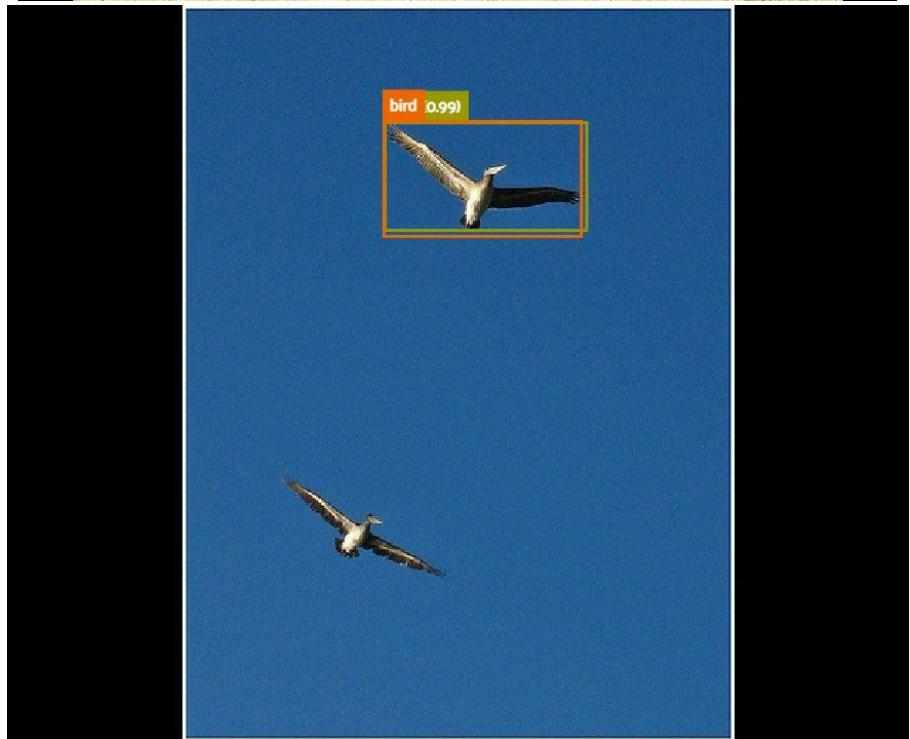
## mAP and IOU

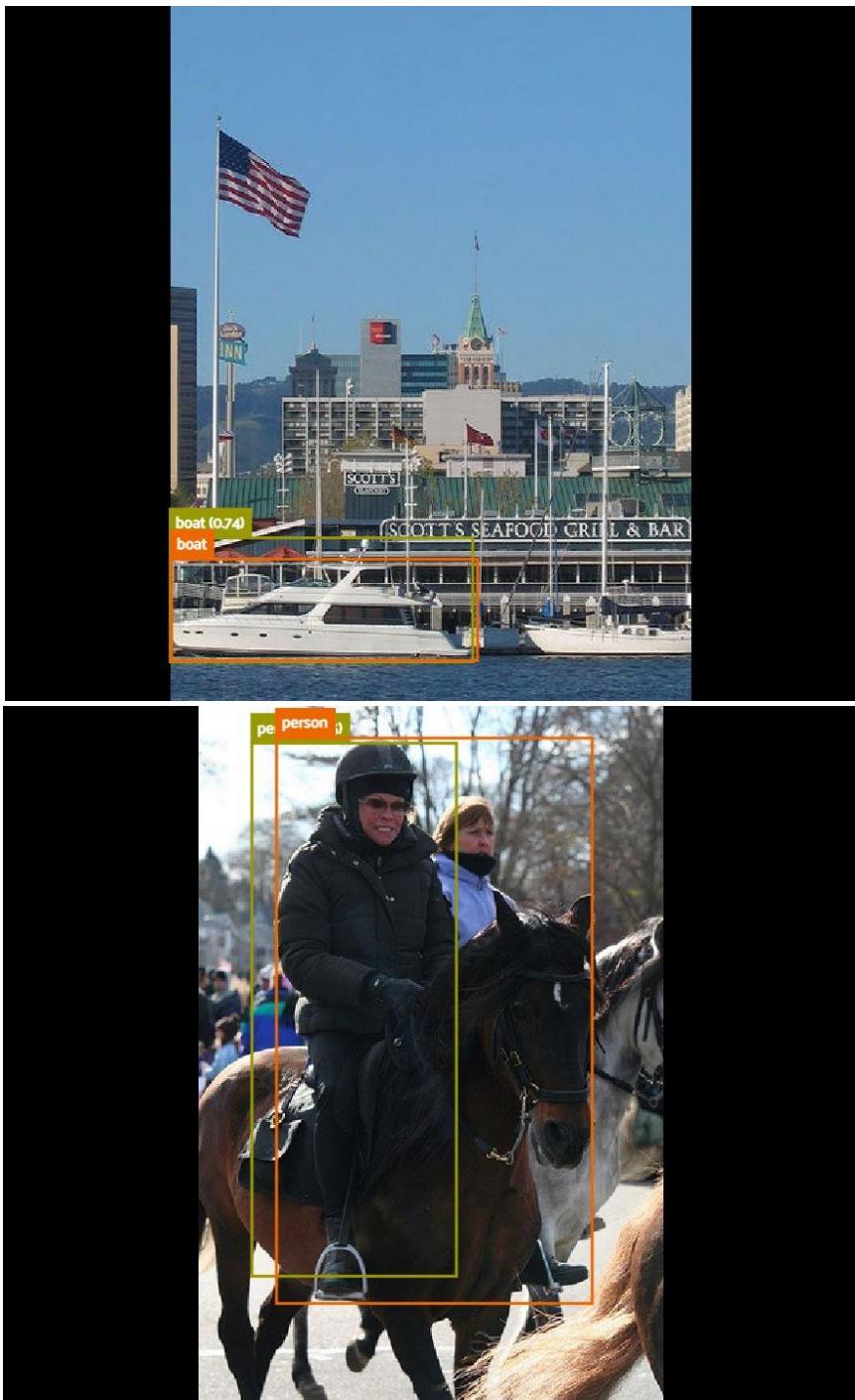
For the whole validation set the **mAP = 0.297 and IOU = 0.368**

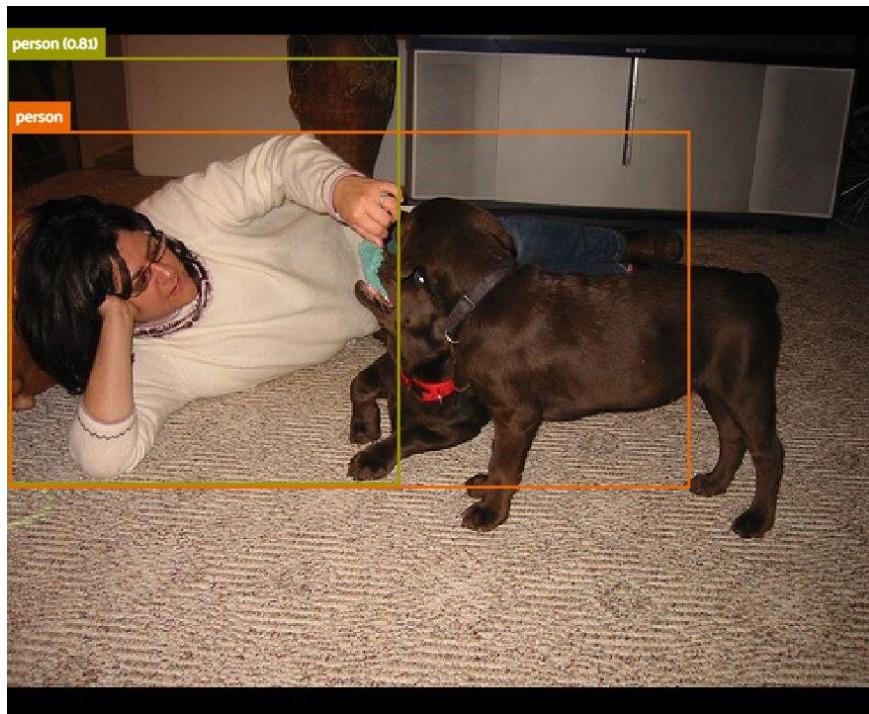
## Precision vs Recall Graph



## VOC IOU





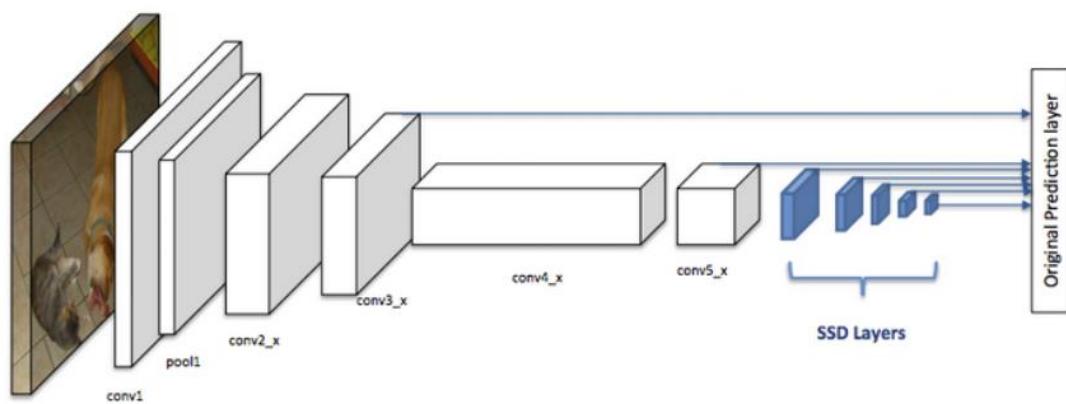


## mAP and IOU

For the whole validation set the **mAP = 0.35 and IOU = 0.36**

## 3- Single Shot Detector

Architecture diagram

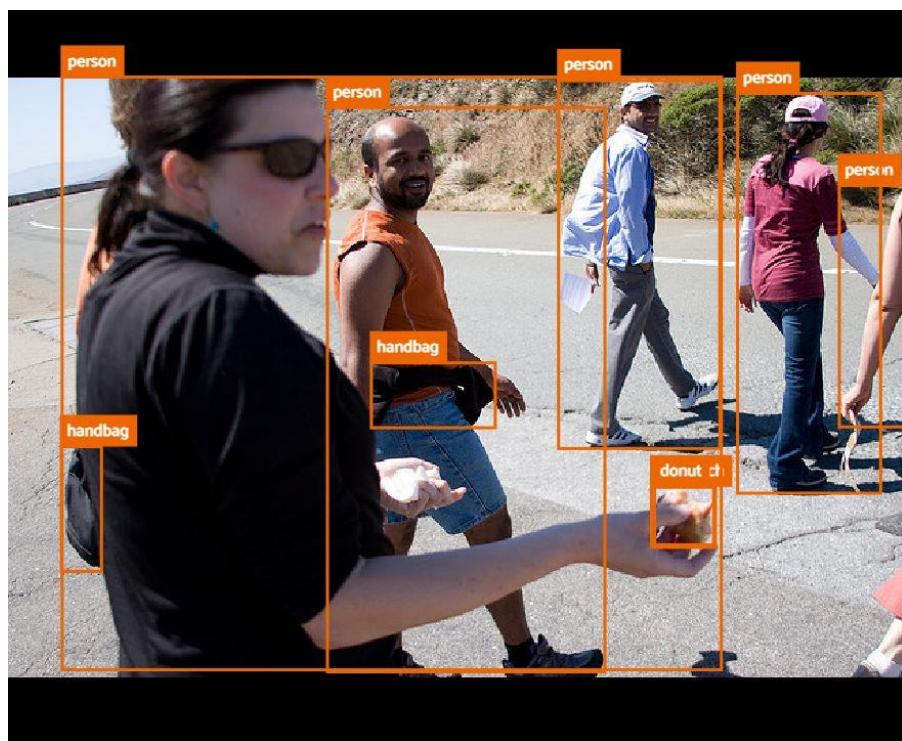


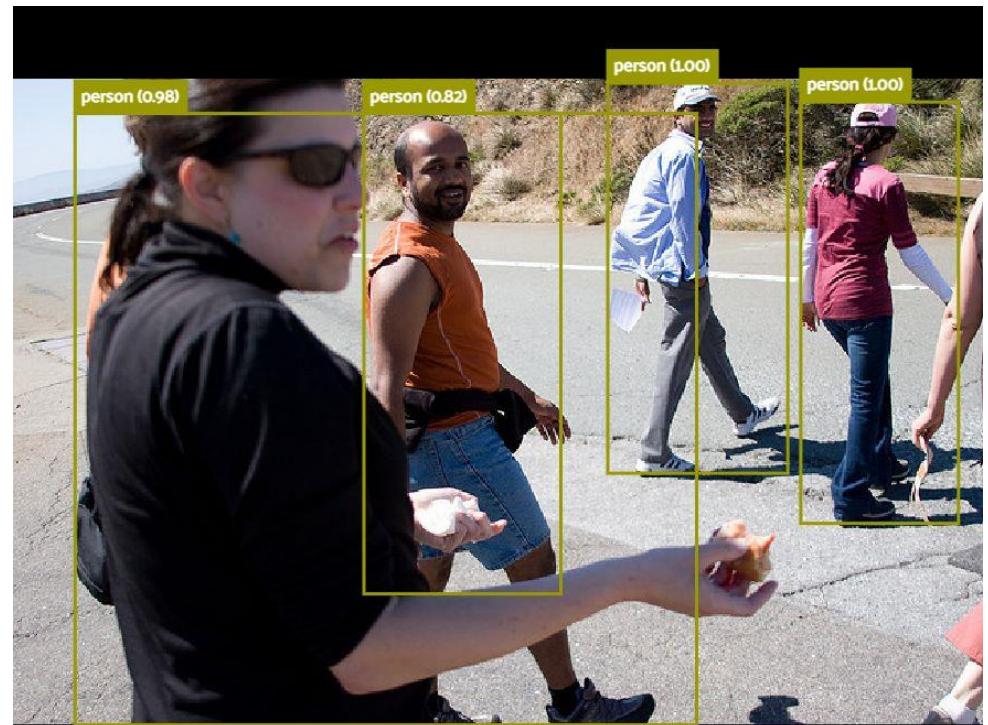
## Discussion

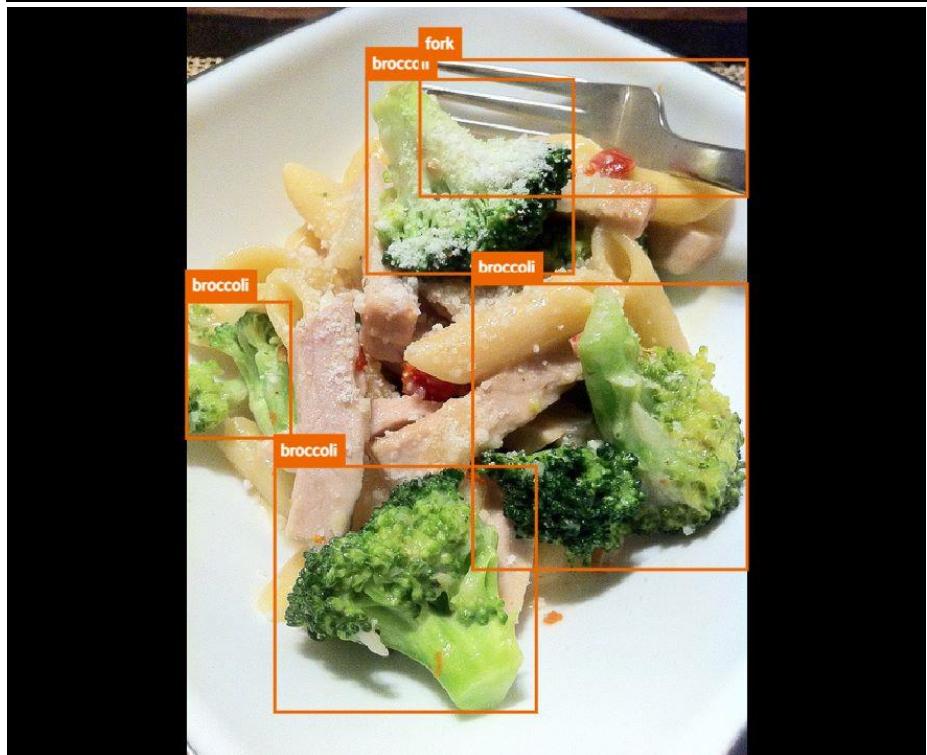
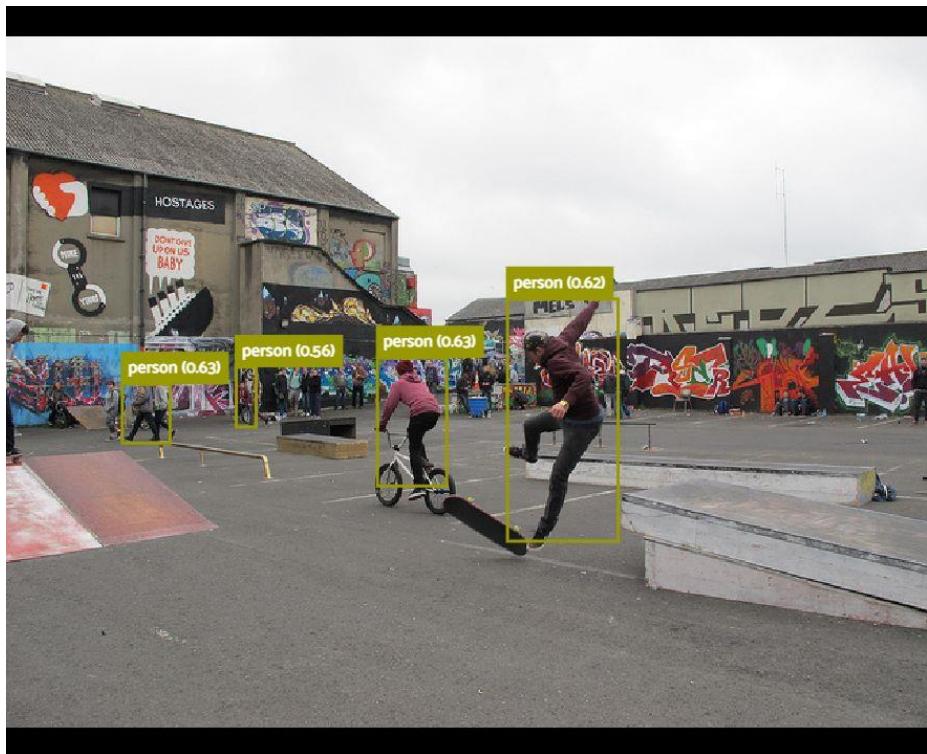
SSD has two components: a backbone model and SSD head.

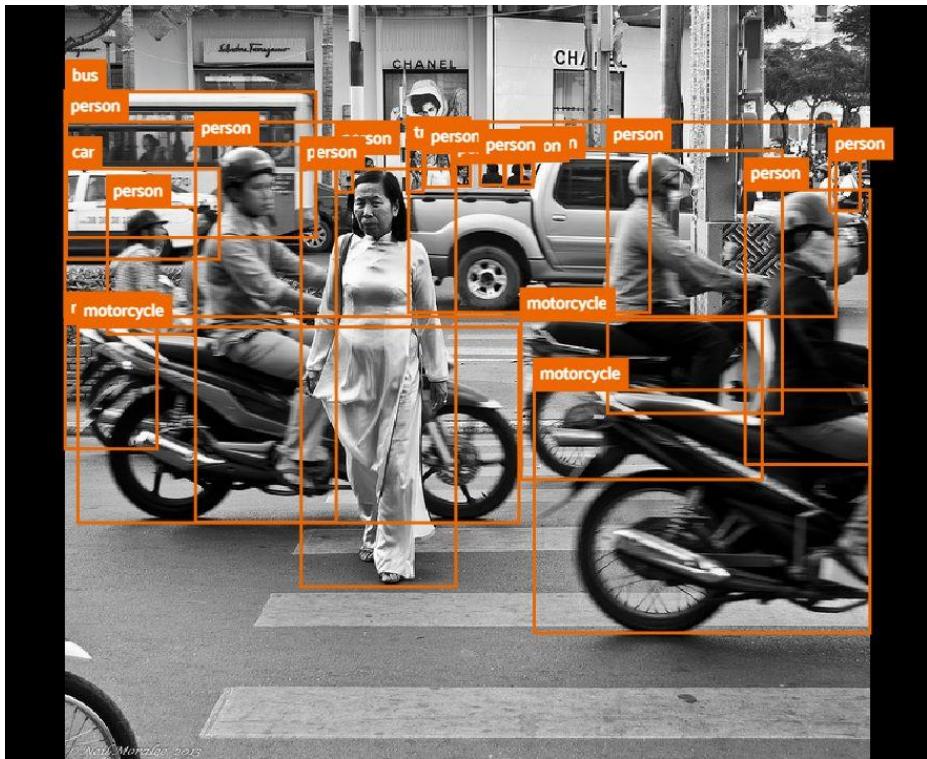
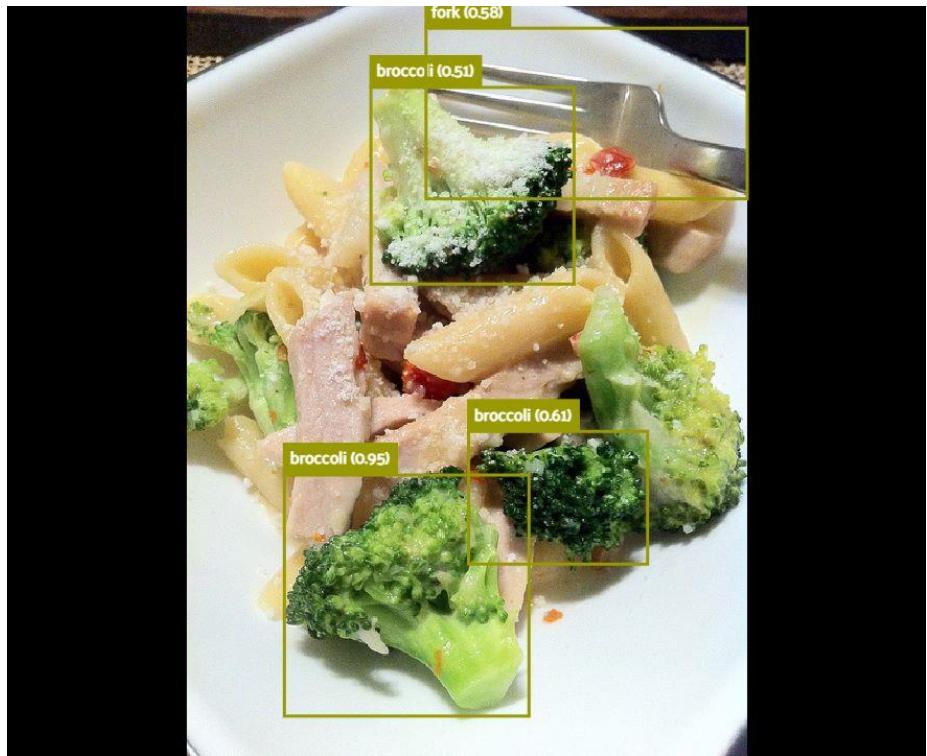
*Backbone* model usually is a pre-trained image classification network as a feature extractor. This is typically a network like **ResNet trained on ImageNet** from which the final fully connected classification layer has been removed. We are thus left with a deep neural network that is able to extract semantic meaning from the input image. The *SSD head* is just one or more convolutional layers added to this backbone and the outputs are interpreted as the **bounding boxes and classes of objects in the spatial location of the final layers activations**.

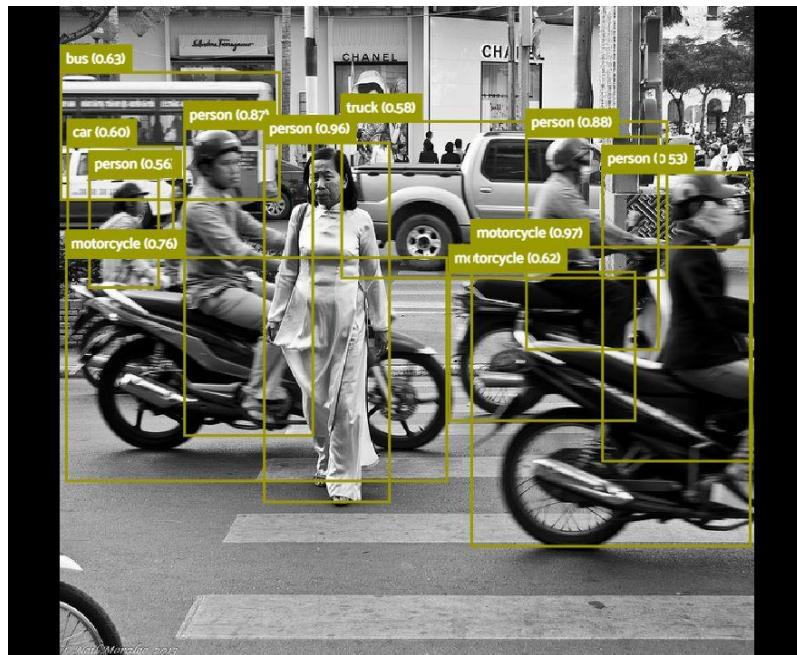
## Results on COCO dataset







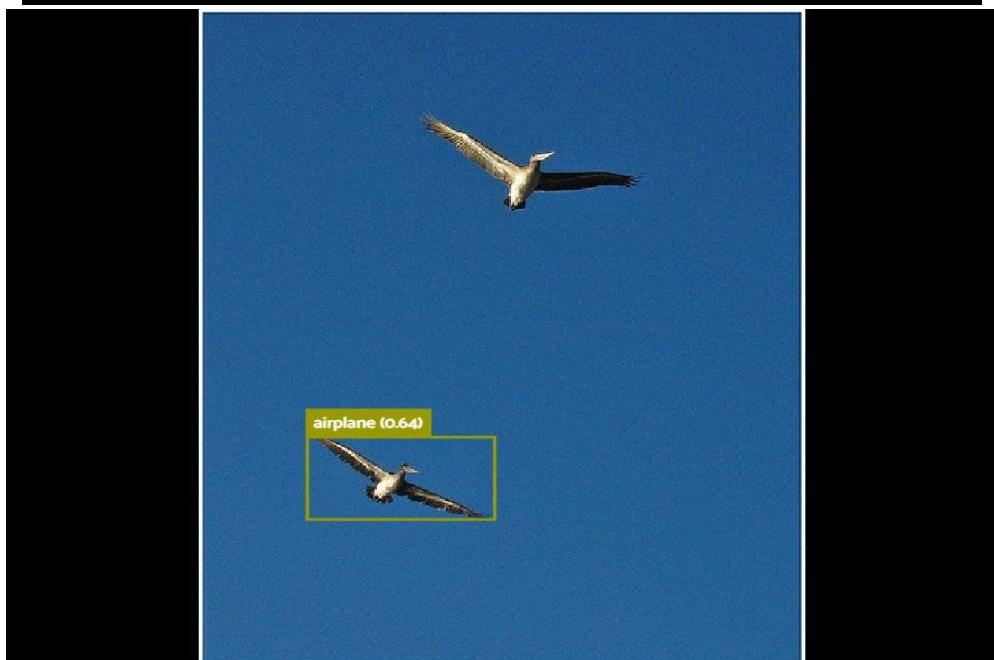


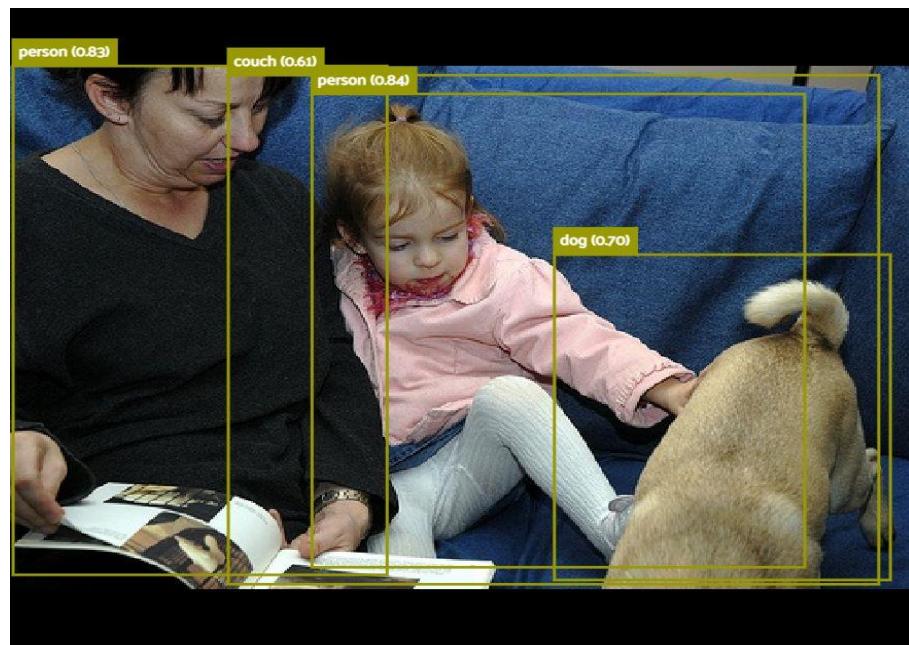
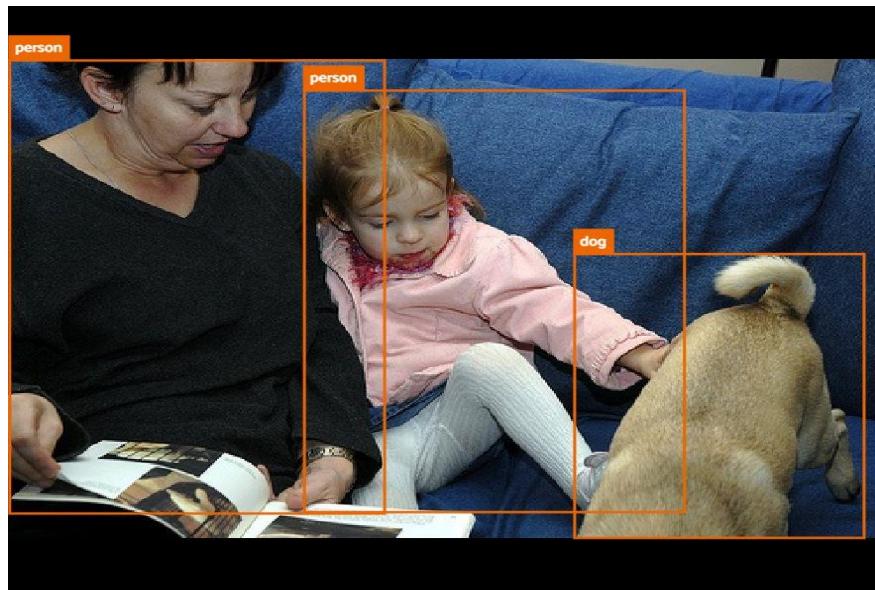


Results on different dataset



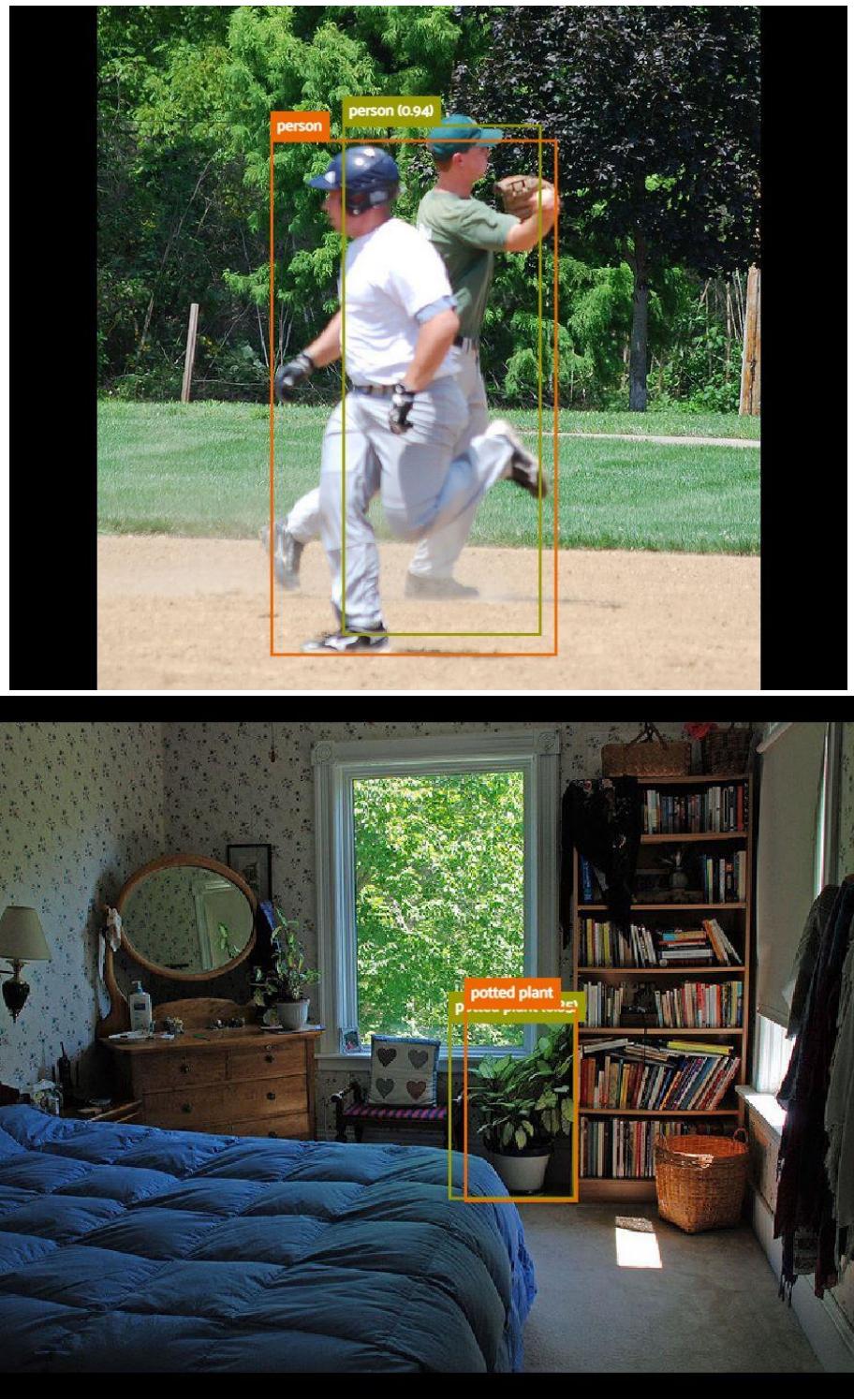






## COCO IOU

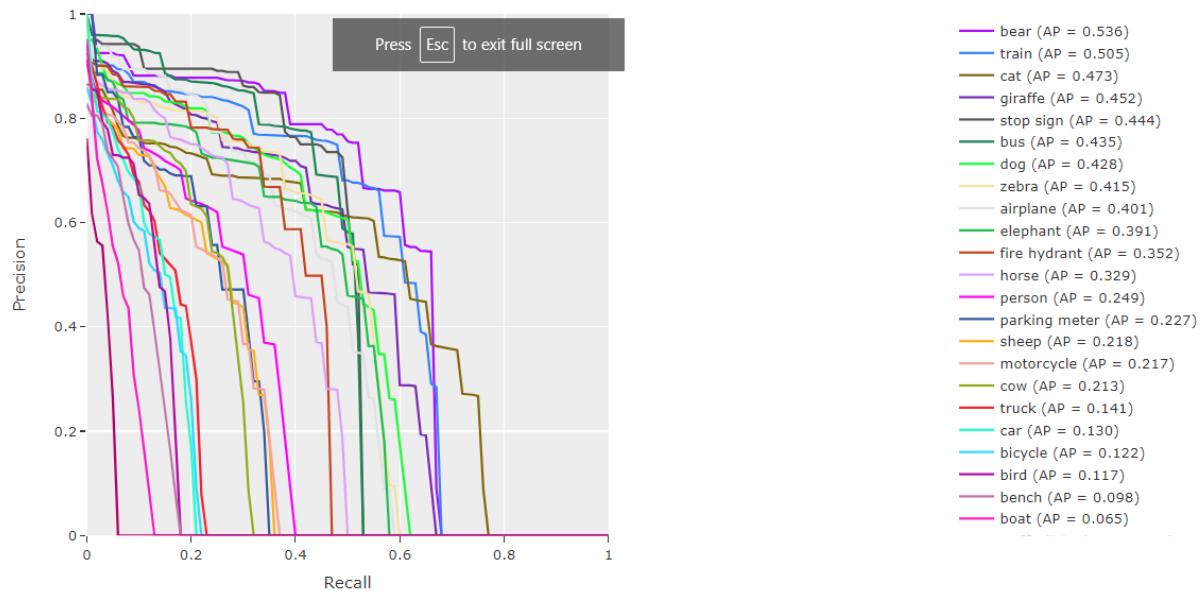




## mAP and IOU

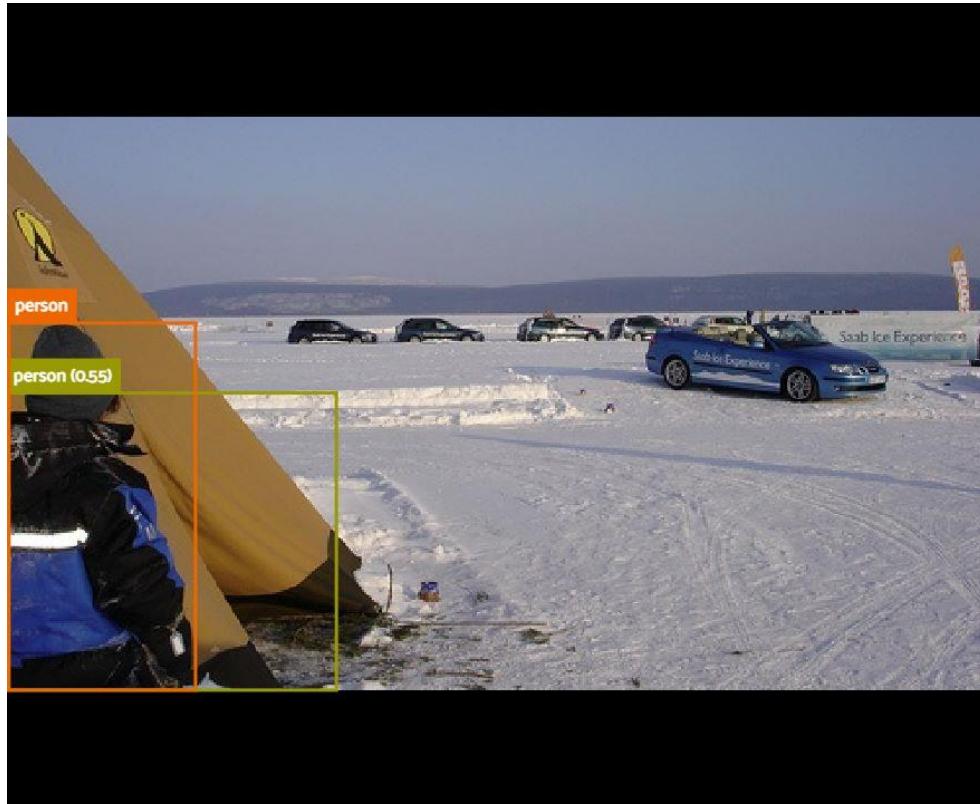
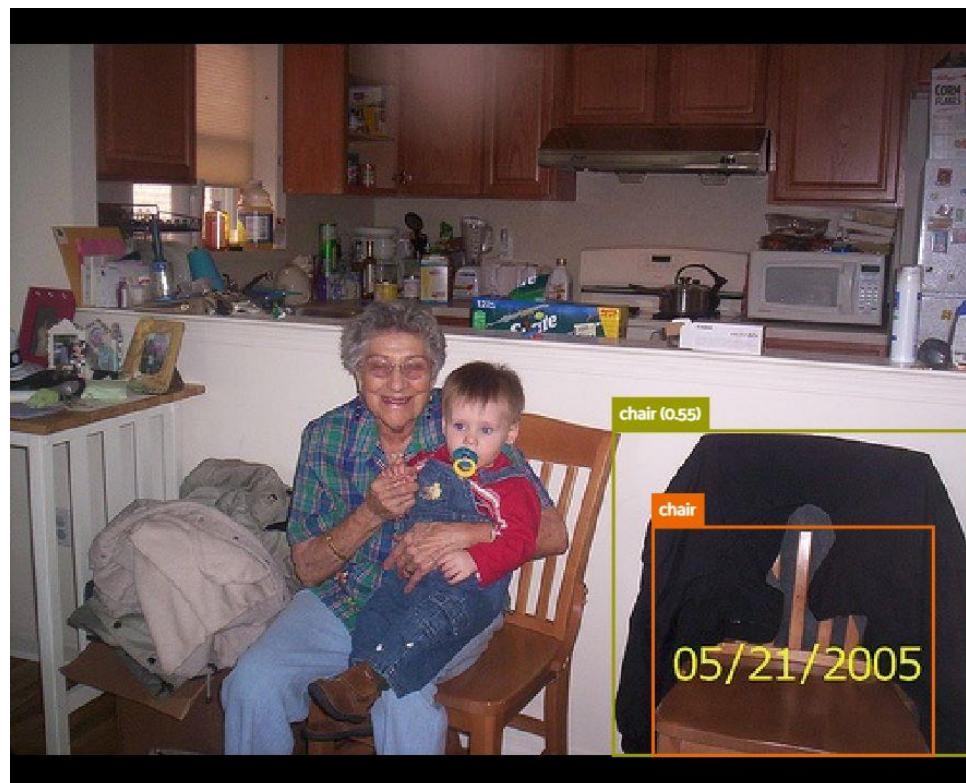
For the whole validation set the **mAP = 0.182 and IOU = 0.221**

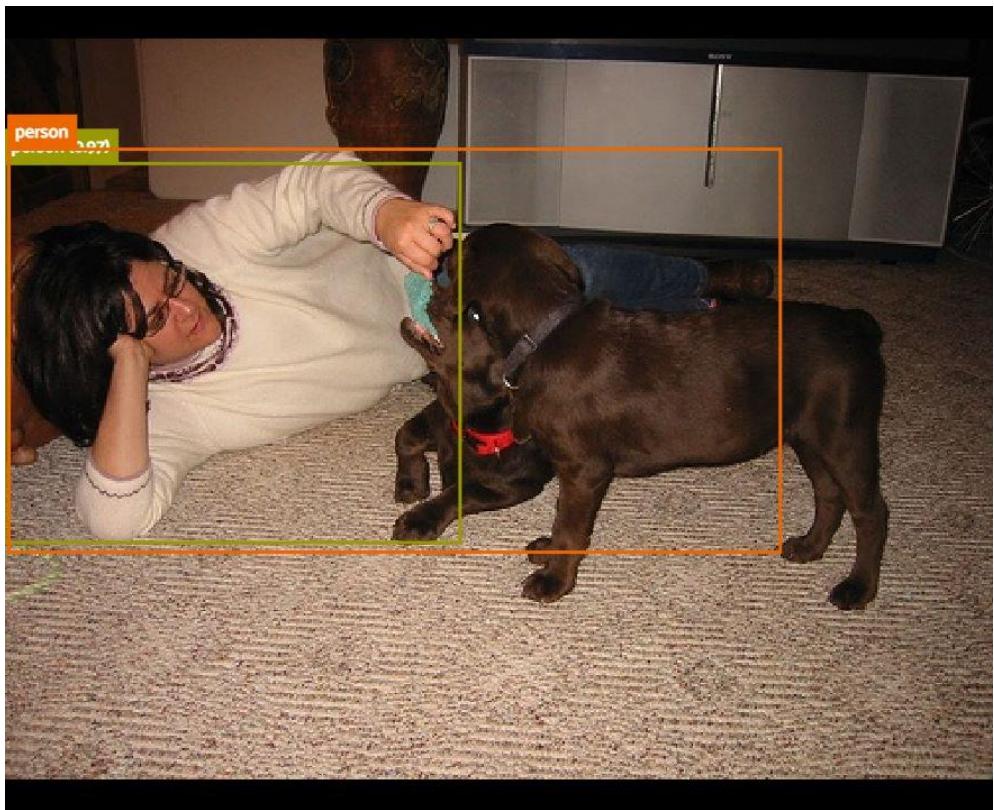
## Precision vs Recall Graph



## Voc IOU



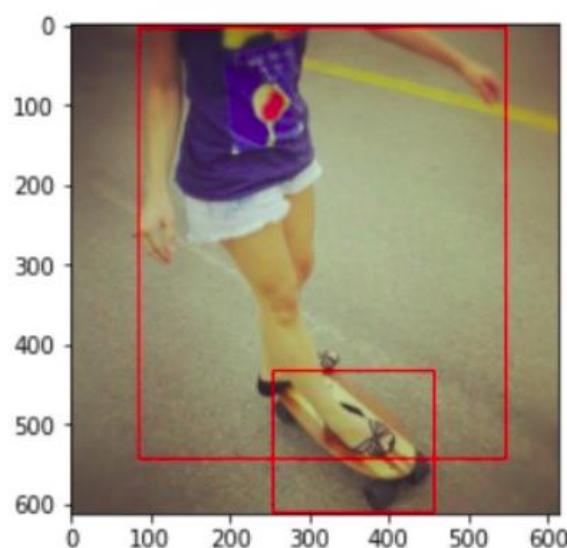


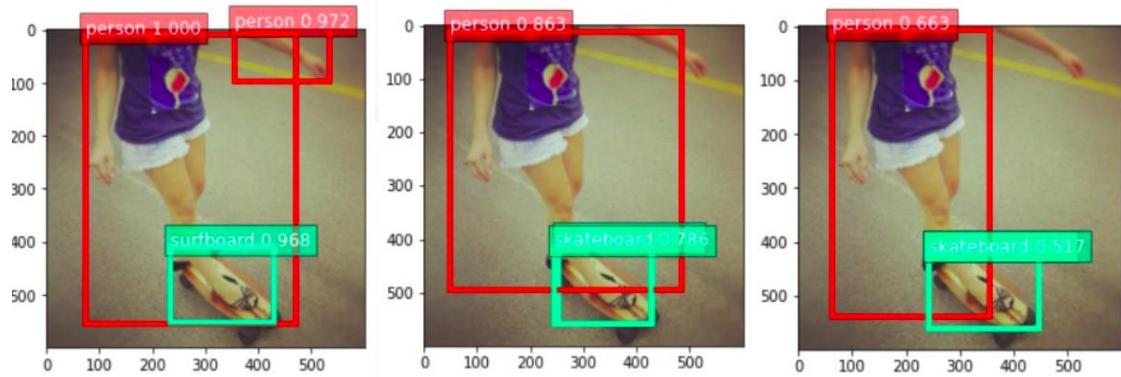


## mAP and IOU

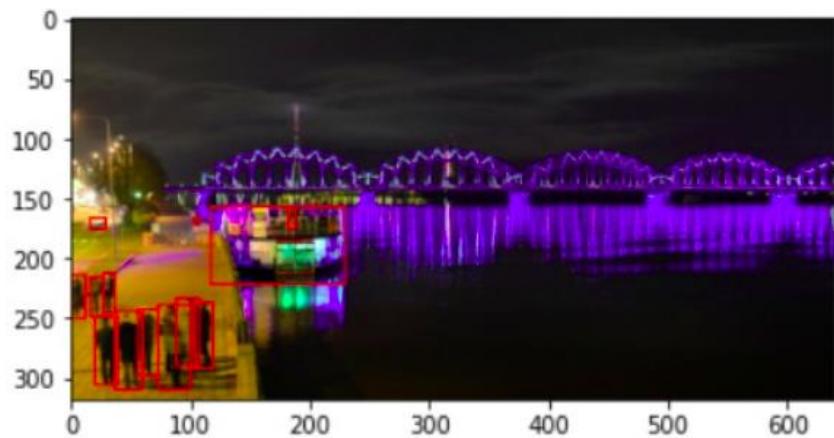
For the whole validation set the **mAP = 0.27** and **IOU = 0.323**

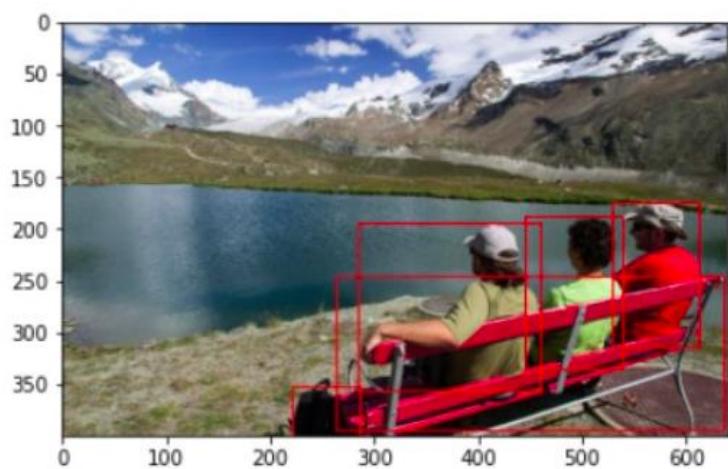
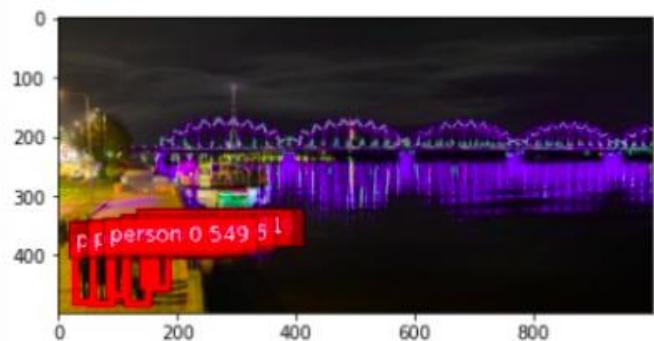
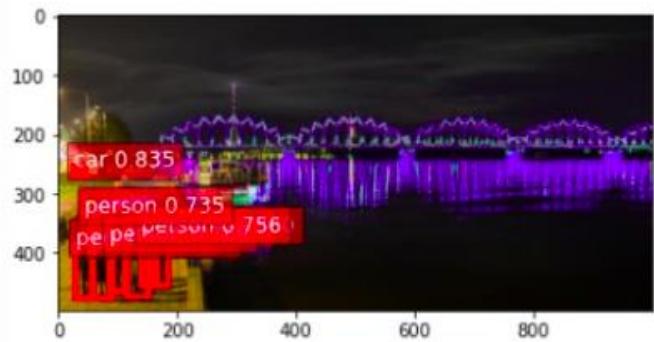
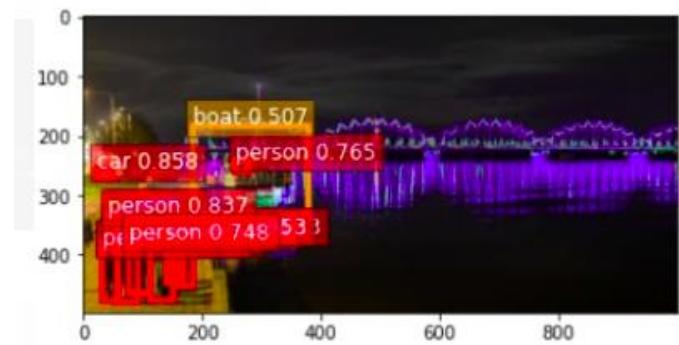
## Test Models on Random Images

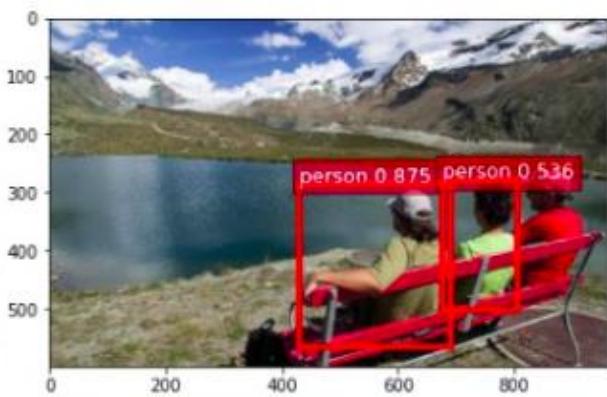
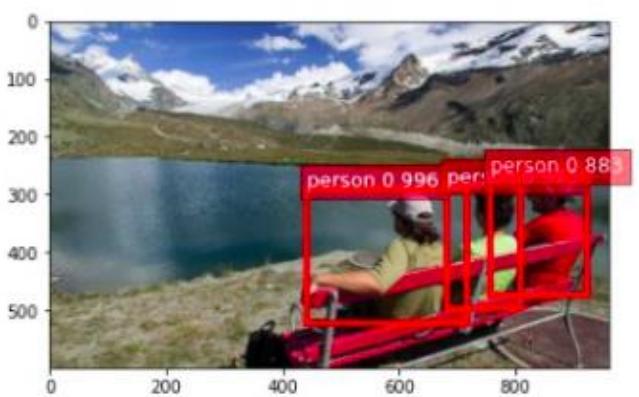
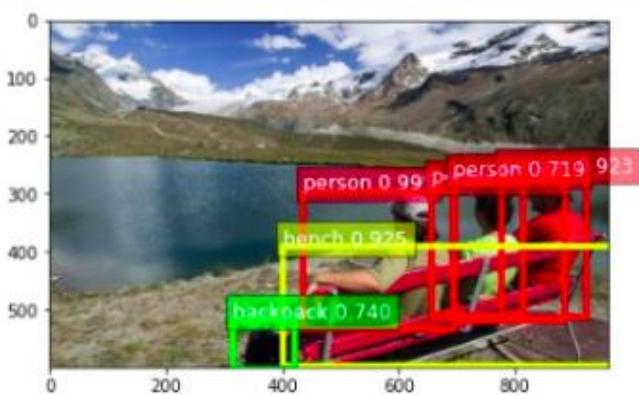


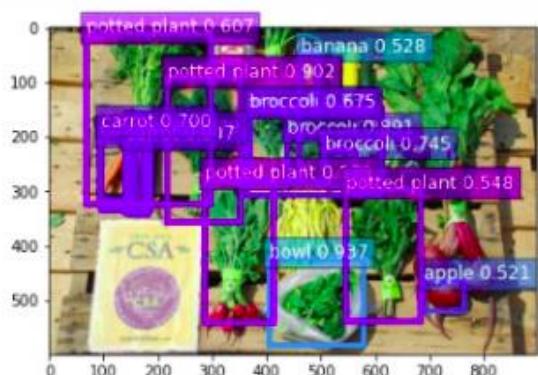


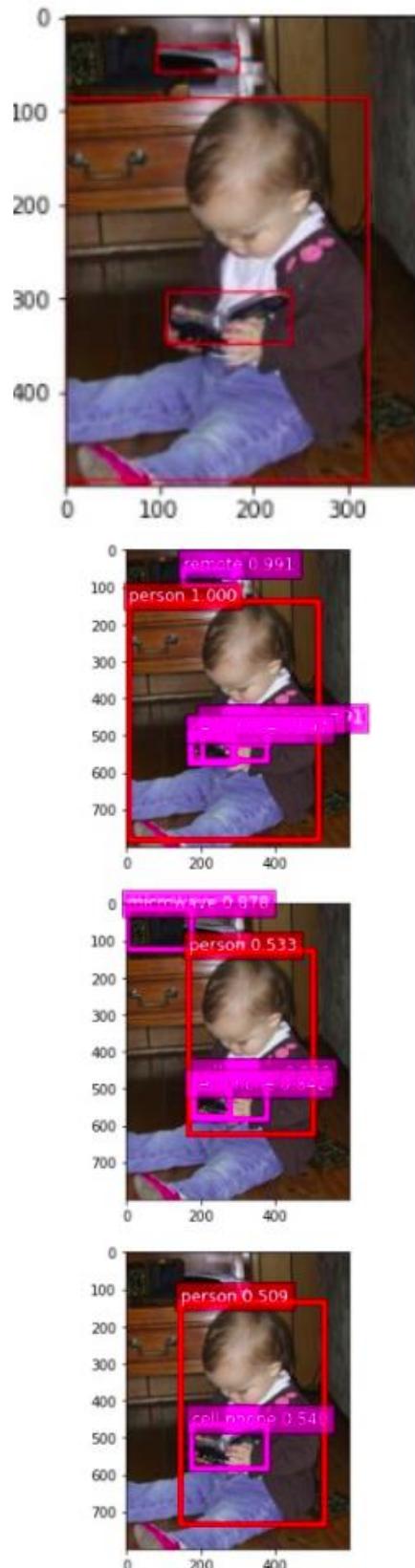
Note: Faster R-CNN gets better confidence and good localization, the SSD has some overlapping but detected all the objects with accepted confidence, Retina has good localization but low confidence.











Note: Overall Faster R-CNN is more accurate while SSD and RetinaNet are faster which gives more convenient in real-time detection.

But for localization and precision Faster R-CNN gives better results.

# Comparison

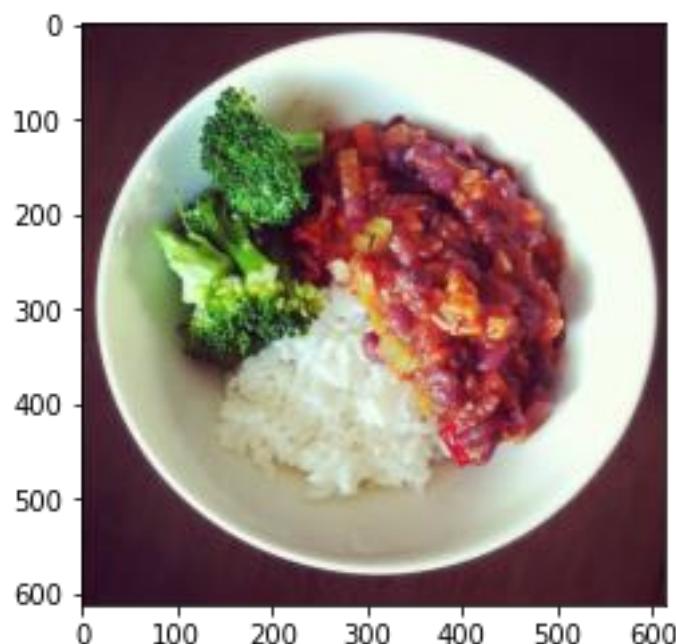
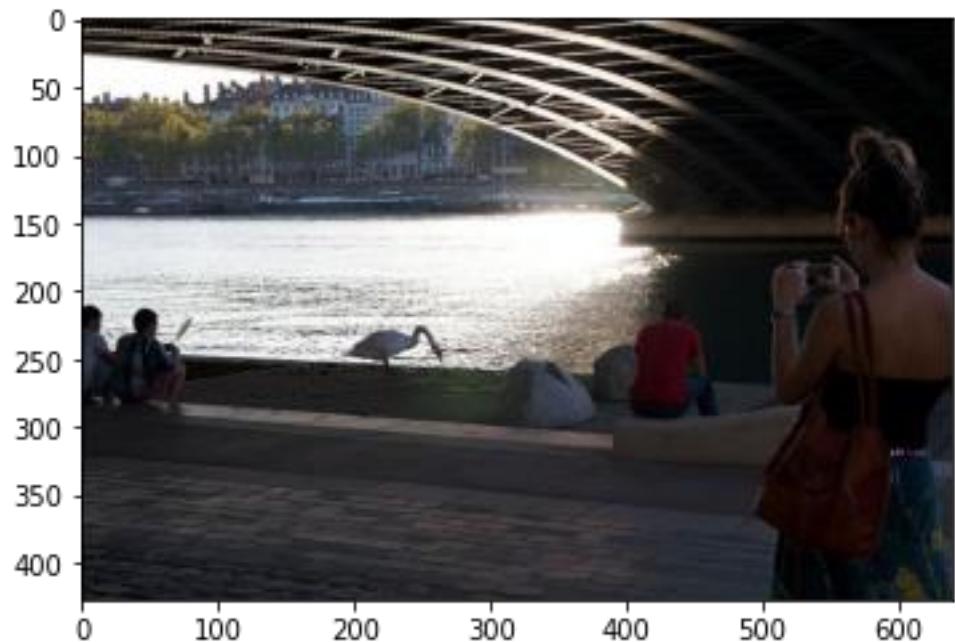
	Coco-RCNN	COCO-SSD	COCO-Retina	VOC-RCNN	VOC-SSD	VOC-Retina
mAP	0.344	0.182	0.297	0.348	0.273	0.351
IOU	0.406	0.221	0.368	0.252	0.323	0.361

# Notebooks

1. Coco dataset:
  - 1.1. <https://colab.research.google.com/drive/1BQKgLPMClvozxpPFRUFEJgwW6EJf8NbK?usp=sharing>
  - 1.2. [https://colab.research.google.com/drive/1N5uDqRMsunGsCu6a\\_OeGhpRYDISunfnK?usp=sharing](https://colab.research.google.com/drive/1N5uDqRMsunGsCu6a_OeGhpRYDISunfnK?usp=sharing)
  - 1.3. [https://colab.research.google.com/drive/1ovCPIEO4hiVo4XrmjdDYdLe\\_ivHLP6Hk?usp=sharing#scrollTo=zS8OVB8p3lqz](https://colab.research.google.com/drive/1ovCPIEO4hiVo4XrmjdDYdLe_ivHLP6Hk?usp=sharing#scrollTo=zS8OVB8p3lqz)
2. Voc dataset:
  - 2.1. [https://colab.research.google.com/drive/1YaWDKbePhhfOkMSyz9kB\\_K4KIAQwgNWP?usp=sharing](https://colab.research.google.com/drive/1YaWDKbePhhfOkMSyz9kB_K4KIAQwgNWP?usp=sharing)
  - 2.2. <https://colab.research.google.com/drive/18XnoypdZQbbn98sQQkuPXUUgonaFmgu?usp=sharing>
  - 2.3. [https://colab.research.google.com/drive/1DCJSAOzIdC4O02EFEqwa0JSBkwOy\\_Ruv?usp=sharing](https://colab.research.google.com/drive/1DCJSAOzIdC4O02EFEqwa0JSBkwOy_Ruv?usp=sharing)

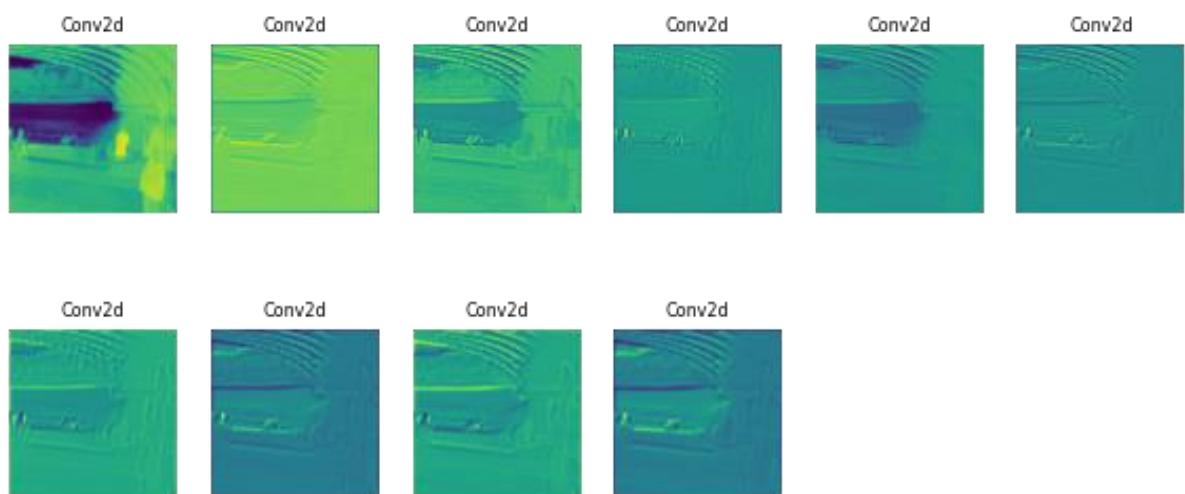
# Feature Maps

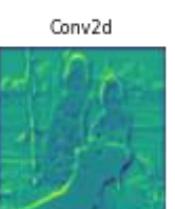
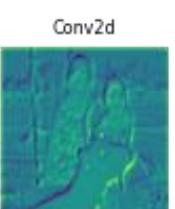
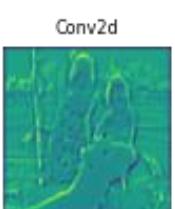
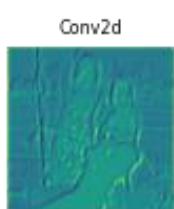
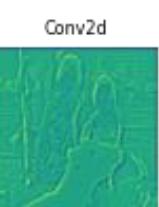
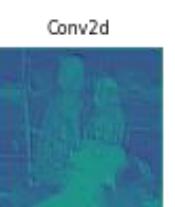
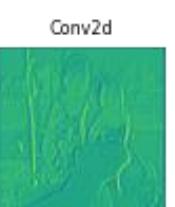
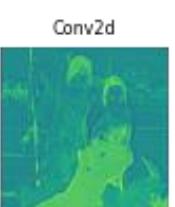
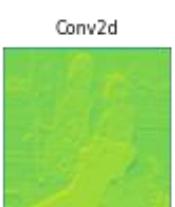
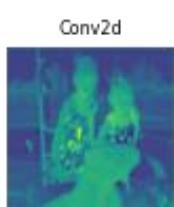
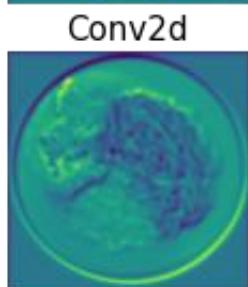
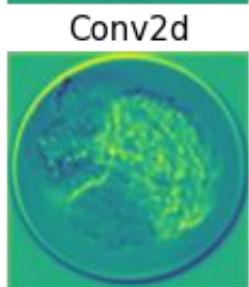
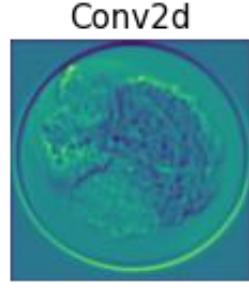
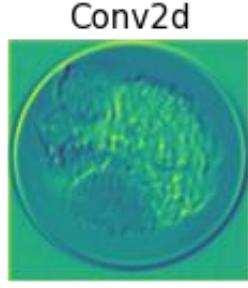
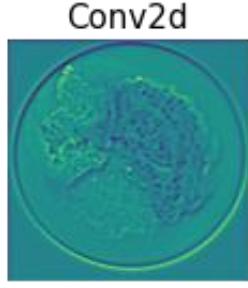
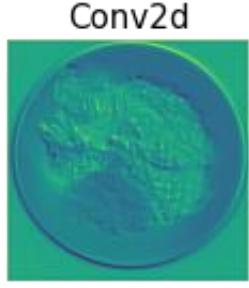
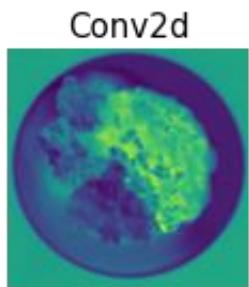
## Image samples





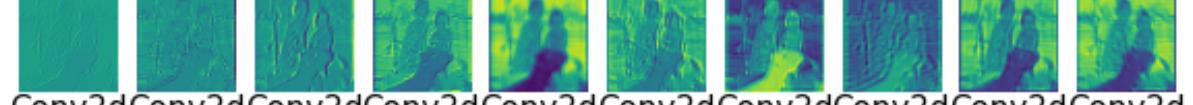
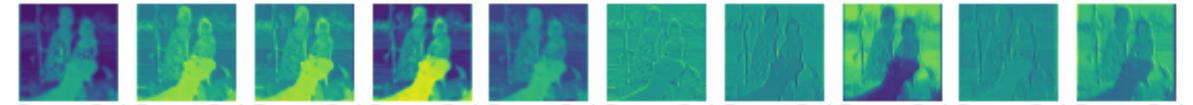
ssd\_model.backbone



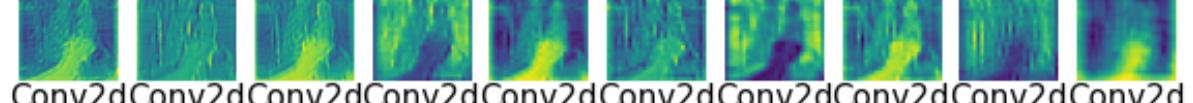


rcnn\_model.backbone

Conv2dConv2dConv2dConv2dConv2dConv2dConv2dConv2dConv2dConv2dConv2dConv2d



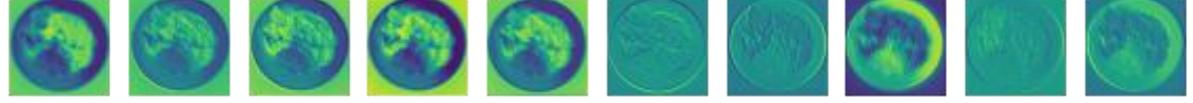
A horizontal row of 10 small square heatmaps, each representing a different layer or stage of a convolutional neural network. The heatmaps show increasing levels of complexity in the learned features, from simple edges to more complex shapes and finally to a detailed semantic segmentation map.



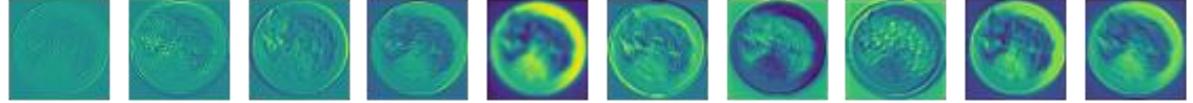
A horizontal row of nine square heatmaps, each representing a different feature or pattern. The colors range from dark blue to bright yellow, indicating varying values or intensities across the spatial domain.



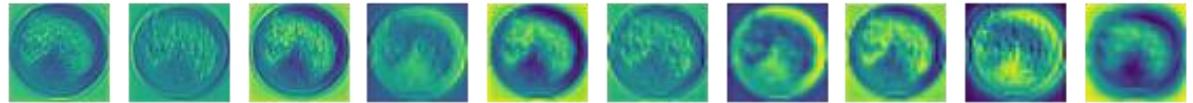
Conv2d Conv2d Conv2d Conv2d Conv2d Conv2d Conv2d Conv2d Conv2d Conv2d



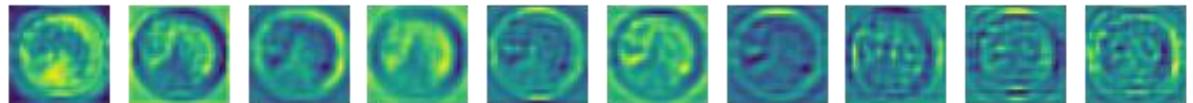
Conv2d Conv2d Conv2d Conv2d Conv2d Conv2d Conv2d Conv2d Conv2d Conv2d



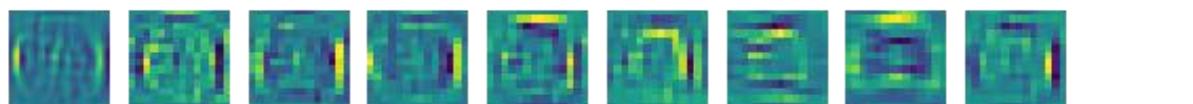
Conv2d Conv2d Conv2d Conv2d Conv2d Conv2d Conv2d Conv2d Conv2d Conv2d

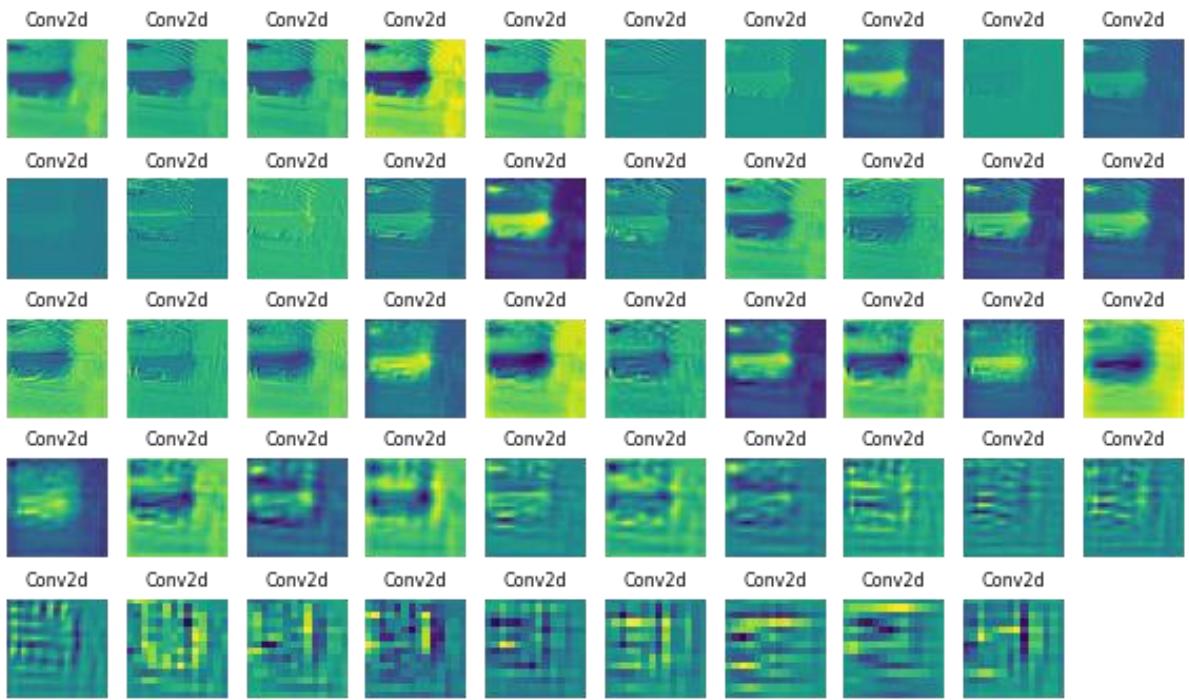


Conv2d Conv2d Conv2d Conv2d Conv2d Conv2d Conv2d Conv2d Conv2d Conv2d

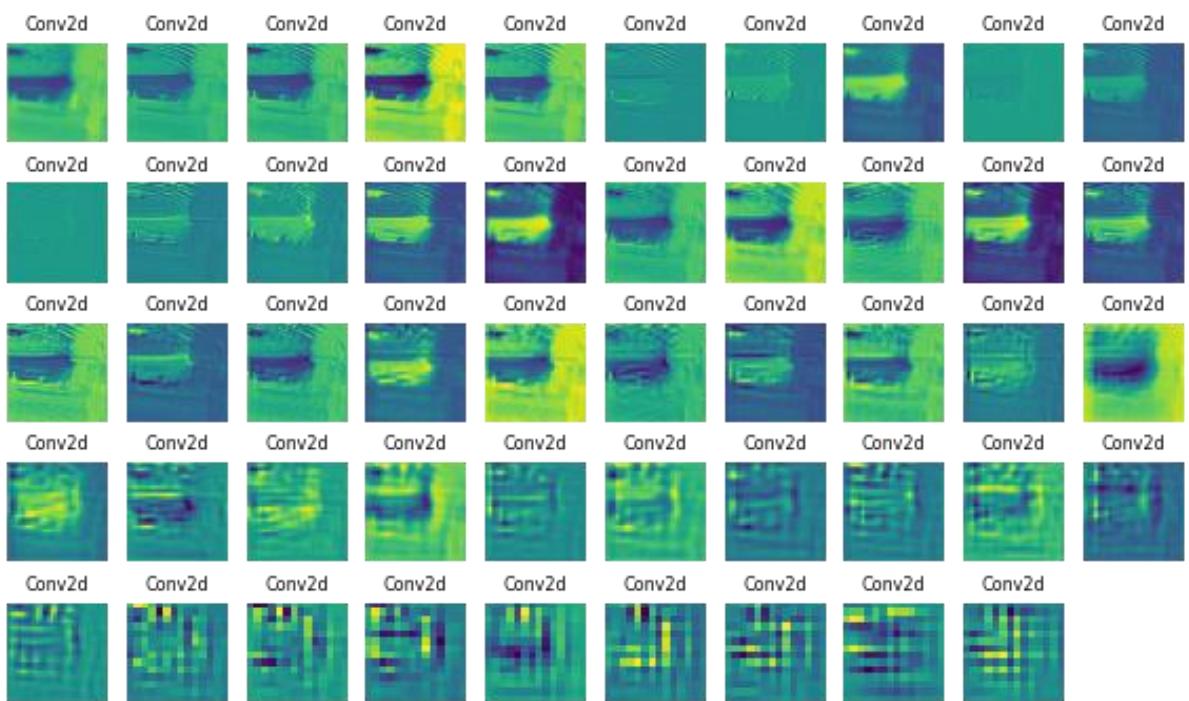


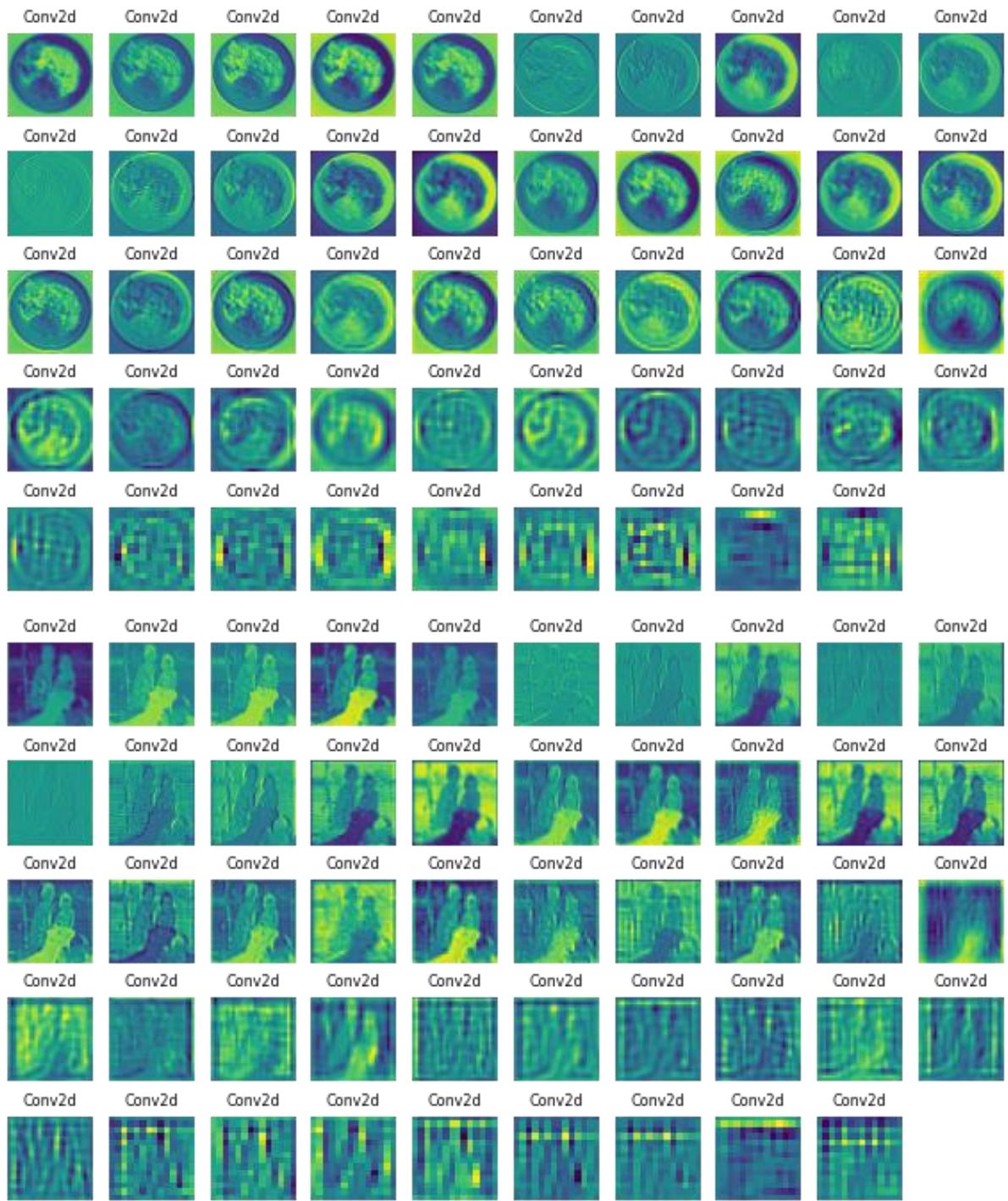
Conv2d Conv2d Conv2d Conv2d Conv2d Conv2d Conv2d Conv2d Conv2d





retina\_model.backbone





## Notebook:

<https://colab.research.google.com/drive/1i6pqgryGTpJXUmwYYh4mE3upPdQFdox?usp=sharing#scrollTo=LVv2vluXfaN3>

# Latex Report

```
\documentclass{article}

\begin{document}

\begin{titlepage}
\begin{centre}
\huge{CV Assignment 4} \\
[2cm]
\text{\LARGE Name 1: Abdelrahman Salem Mohamed} \\
[0.5cm]
\text{\LARGE ID 1: 6309} \\
[0.5cm]
\text{\LARGE Name 2: Reem Abdelhalim} \\
[0.5cm]
\text{\LARGE ID 2: 6114 } \\
[0.5cm]
\end{centre}
\end{titlepage}

\section{Object Detection}
\text{In this assignment, you will work on Cocco dataset which is a large-scale object detection,}
\text{segmentation, and captioning dataset. You are required to run 3 different object detectors on}
\text{this dataset. We will learn to use and differentiate between the architectures.}
\text{Either PyTorch or TensorFlow are allowed to run your models.}

\noindent
\section{Dataset}
\text{We test models using coco 2017 validation set which consists of 5000 different images with their meta data from segmentation points and object bounding boxes.}

\section{Models}
\text{We used three different models:}
\begin{enumerate}
\item Faster R-CNN (Two stage Model)
\item RetinaNet (One stage Model)
\item Single shot detector (One stage Model)
\end{enumerate}

\subsection{Faster R-CNN }
\paragraph{Discussion}
\text{Two stage based model where, first the model learn how to extract Region of proposals, then for every Region of proposal fe perform the normal Fast R-CNN network}
```

from ROI pooling then from fully connected layer compute the two losses one for object box (regression), and one for class type (classification).

\text Each network has its own loss functions and the whole system weights consider the both losses path, and during the test the network needs first to identify the ROI and apply Object detection technique, so Faster R-CNN consider slower than one stage models but have better mAP.

\text The backbone network is usually a dense convolutional network like ResNet or VGG16

#### \subsection{ RetinaNet }

##### \paragraph{Discussion}

\text RetinaNet is one of the best one-stage object detection models that has proven to work well with \textbf{dense} and small scale objects}. For this reason, it has become a popular object detection model to be used with aerial and satellite imagery.

\text There are four major components of a RetinaNet model architecture :

\text a) Bottom-up Pathway - \textbf{The backbone network (e.g. ResNet)} which calculates the feature maps at different scales, irrespective of the input image size or the backbone.

\text b) Top-down pathway and Lateral connections - The top down pathway upsamples the spatially coarser feature maps from higher pyramid levels, and the lateral connections merge the top-down layers and the bottom-up layers with the same spatial size.

\text c) Classification subnetwork - It predicts the probability of an object being present at each spatial location for each anchor box and object class.

\text d) Regression subnetwork - It regresses the offset for the bounding boxes from the anchor boxes for each ground-truth object.

#### \subsection{ Single Shot Detector}

##### \paragraph{Discussion }

\text SSD has two components: a backbone model and SSD head.

\text \textit{Backbone} model usually is a pre-trained image classification network as a feature extractor. This is typically a network like \textbf{ResNet} trained on ImageNet} from which the final fully connected classification layer has been removed. We are thus left with a deep neural network that is able to extract semantic meaning from the input image

\text The \textit{SSD head} is just one or more convolutional layers added to this backbone and the outputs are interpreted as the \textbf{bounding boxes and classes of objects in the spatial location of the final layers activations}.

### \section{Test Models on Random Images}

\text Note: Faster R-CNN gets better confidence and good localization, the SSD has some overlapping but detected all the objects with accepted confidence, Retina has good localization but low confidence.

\text Note: Overall Faster R-CNN is more accurate while SSD and RetinaNet are faster which gives more convenient in real-time detection.

\text But for localization and precision Faster R-CNN gives better results.

### \section{Comparison}

\begin{tabular}{|p{0.6in}|p{0.6in}|p{0.6in}|p{0.6in}|p{0.6in}|p{0.6in}|p{0.6in}|} \hline

& Coco-RCNN & COCO-SSD & COCO-Retina & VOC-RCNN & VOC-SSD & VOC-Retina \\ \\hline  
mAP & 0.344 & 0.182 & 0.297 & 0.348 & 0.273 & 0.351 \\ \\hline  
IOU & 0.406 & 0.221 & 0.368 & 0.252 & 0.323 & 0.361 \\ \\hline  
\\end{tabular}  
\\end{document}

Link for the zip file:

[https://drive.google.com/drive/folders/19Lwy4htajEA9VhwZxqQvD6hG5cpljrHK?usp=share\\_link](https://drive.google.com/drive/folders/19Lwy4htajEA9VhwZxqQvD6hG5cpljrHK?usp=share_link)