# Person Detection
# from a
# Top-View Perspective

Mameli Marco
Paolanti Marina
Pazzaglia Giulia
Frontoni Emanuele

Baldascino Giovanni - 1097405
Squarcella Loisi - 1096539

UNIVERSITÀ
POLITECNICA
DELLE MARCHE

VRAi

# INTRODUCTION

**Person Detection** from a Top-View Perspective by using most recent object detection frameworks.
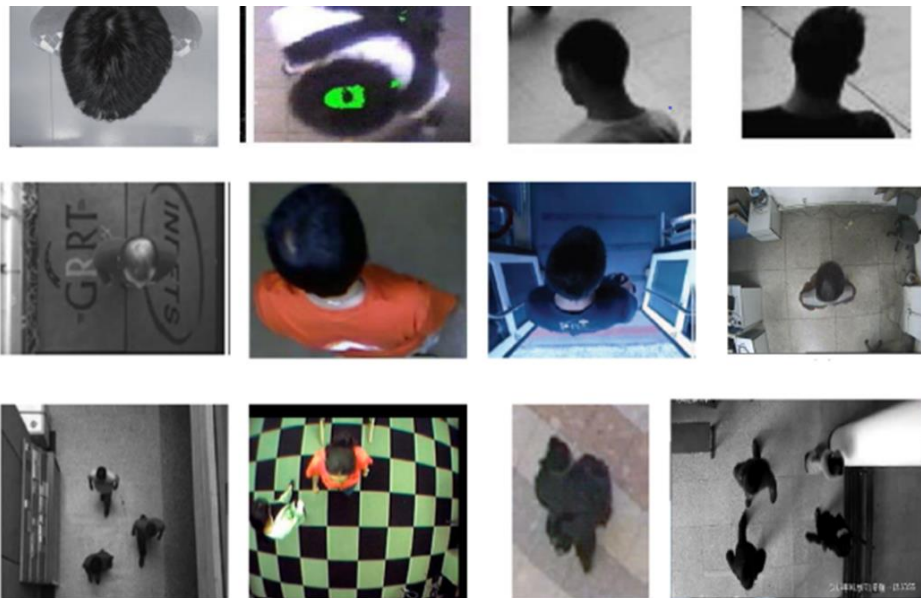
The **aim** is to build an efficient model to detect the people in the scene. The provided scenario is a retail environment.
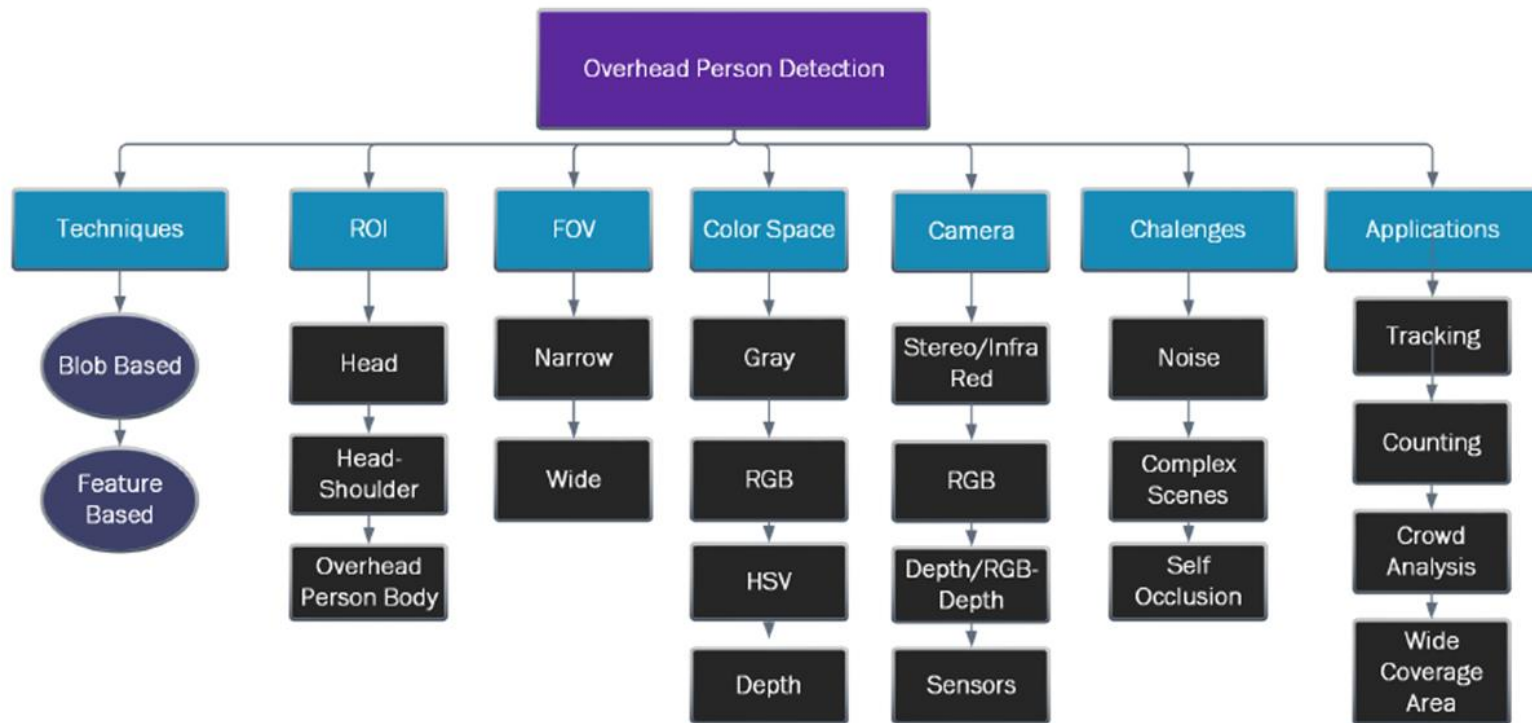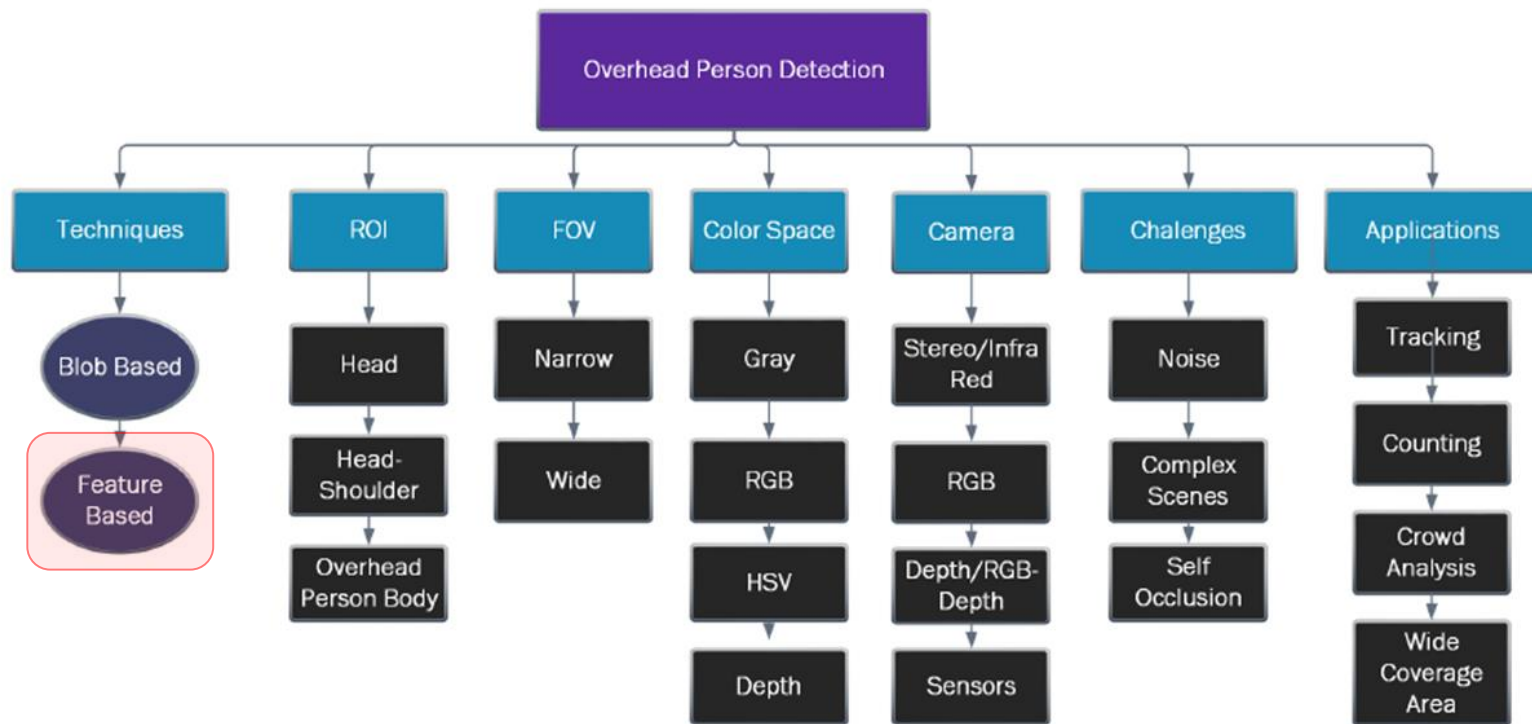
# INTRODUCTION

**MAIN STEPS:**

- Definition of the **Region of Interest** (ROI);
- People localization.

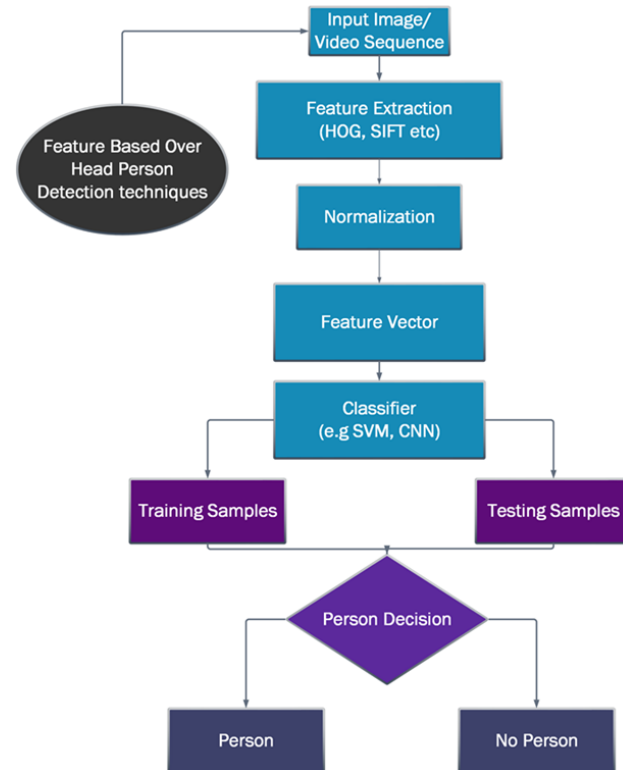# STATE OF ART – Person Detection

# STATE OF ART – Person Detection
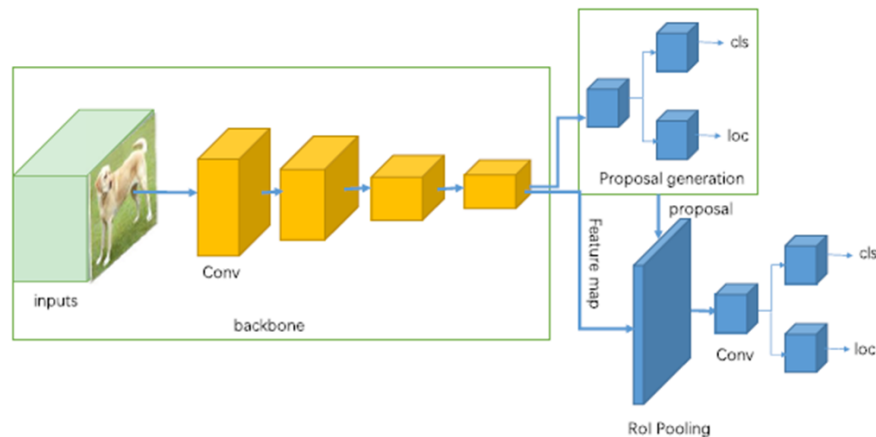
# STATE OF ART – Person Detection

These **Feature based Techniques** operate on *features extracted* from overhead view videos and images.

The extracted features contain shape, color, texture, etc…. The images are often divided into samples for training and testing.
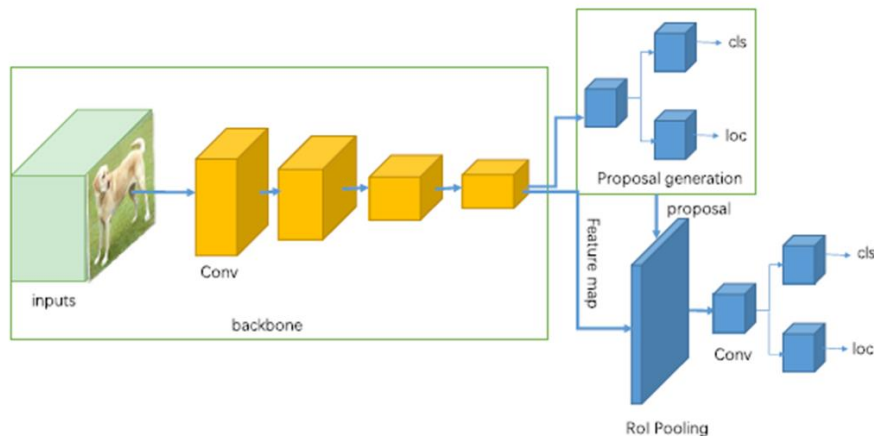
# STATE OF ART – Feature based Techniques

**Two-stage Detectors** (R-CNN, Faster R-CNN, etc..) use a *Region Proposal Network* to generate regions of interests.
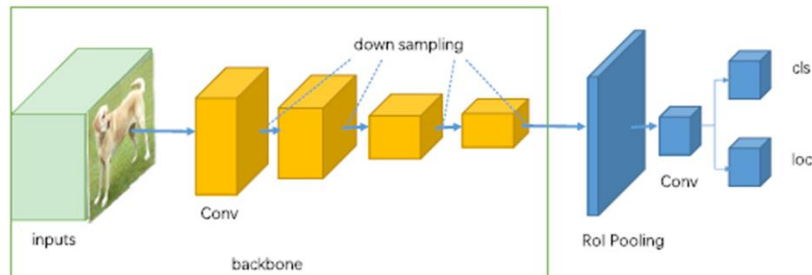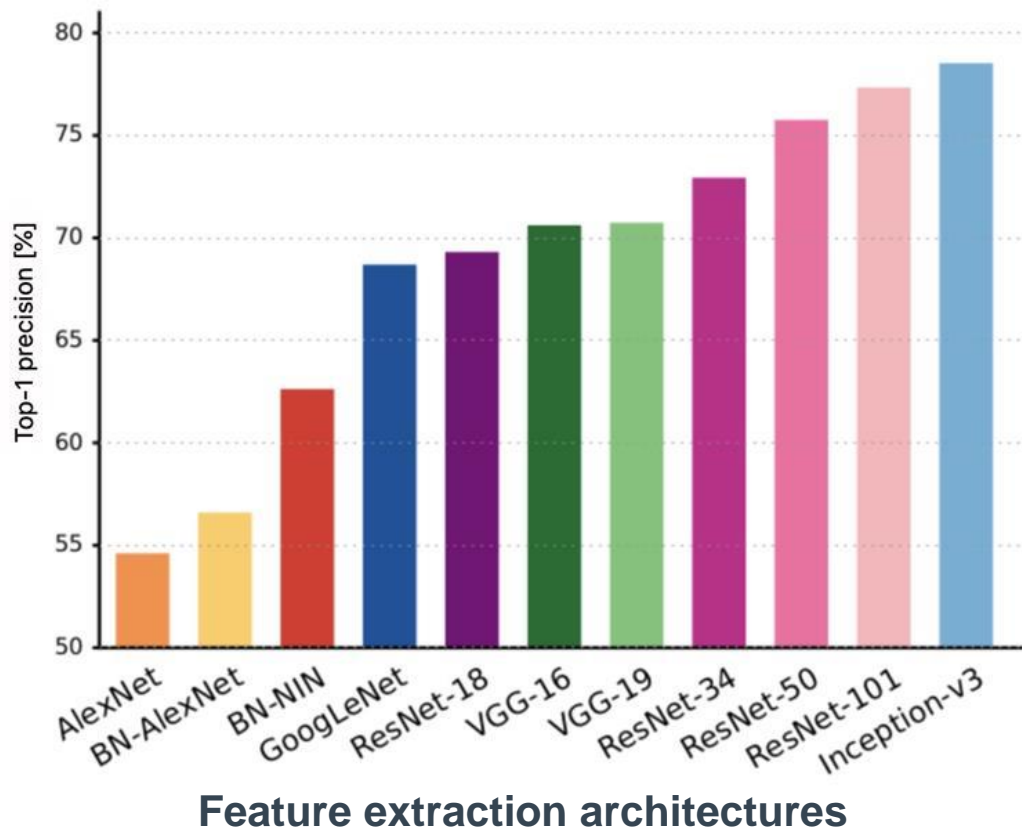
# STATE OF ART – Feature based Techniques

**Two-stage Detectors** (R-CNN, Faster R-CNN, etc..) use a *Region Proposal Network* to generate regions of interests.



**One-stage Detectors** (YOLO, ***SSD***, etc..) treat object detection as a *simple regression problem*.

# STATE OF ART – Feature based Techniques



**Feature extraction architectures**

# STATE OF ART – Feature based Techniques

# STATE OF ART – One-stage Detector

**SSD** (Single-shot Detector) discretizes the output space of *bounding boxes* into a set of default boxes over different aspect ratios and scales per feature map location.

In SSD the *prediction layer* is acting on fused features of different levels. Head module consists of a series of *convolutional layers* followed by several classification layers and localization layers.

# STATE OF ART – Single-shot Detector



loc : $\Delta(cx, cy, w, h)$
conf : $(c_1, c_2, \cdots, c_p)$

Each prediction is composed of:
- Bounding box with shape offset ($\Delta cx$, $\Delta cy$, $w$ and $h$);
- *Confidences* for all object categories or all the classes.

UNIVERSITÀ
POLITECNICA
DELLE MARCHE

VRAi

# MATERIALS AND METHODS – Labelling

*Dataset* annotation using **LabelImg**. LabelImg is a *graphical image annotation tool*.

About 17.000 images divided into:

- **Training set** → 70%

- **Testing set** → 30%

- **Validation set** → 30% of training set

UNIVERSITÀ
POLITECNICA
DELLE MARCHE

VRAi

# MATERIALS AND METHODS – Training

Three types of Single-shot Detectors were used: **SSD300**, **SSD512** and **SSD7**. The architecture of the first two is almost the same (9 and 10 layers respectively) while the SSD7 provides a simplified approach (7 layers).

| Method | mAP | FPS | batch size | Input resolution |
|--------|-----|-----|-----------|------------------|
| SSD300 | 74.3 | 46 | 1 | $300 \times 300$ |
| SSD512 | 76.8 | 19 | 1 | $512 \times 512$ |
| SSD300 | 74.3 | 59 | 8 | $300 \times 300$ |
| SSD512 | 76.8 | 22 | 8 | $512 \times 512$ |

The default *backbone* was based on **VGG16**, then replaced with **ResNet50**.

Furthermore, a subsequent approach was fine tuning starting from a pre-trained model.

# MATERIALS AND METHODS – Training

Three types of Single-shot Detectors were used: **SSD300**, **SSD512** and **SSD7**. The architecture of the first two is almost the same (9 and 10 layers respectively) while the SSD7 provides a simplified approach (7 layers).

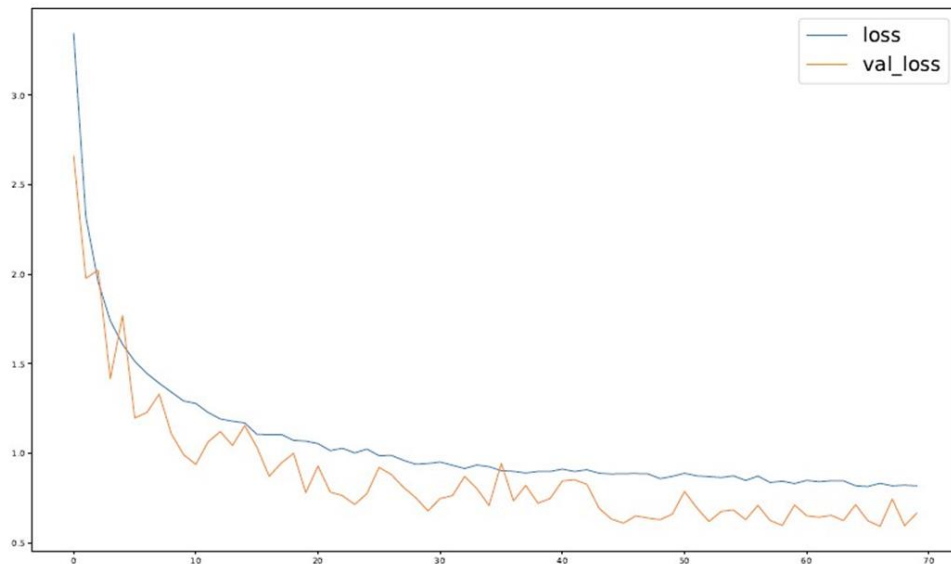|  |  | Batch size | Steps per epoch |
|---|---|---|---|
| VGG16 | SSD300 | 32 | 260 |
|  | SSD512 | 16 | 520 |
| ResNet50 | SSD300 | 32 | 260 |
|  | SSD512 | 12 | 690 |

$$Steps = \frac{Trainig\ set}{Batch\ size}$$

Epochs = 40

Learning rate = 0,001

# MATERIALS AND METHODS – Training

Three types of Single-shot Detectors were used: **SSD300**, **SSD512** and **SSD7**. The architecture of the first two is almost the same (9 and 10 layers respectively) while the SSD7 provides a simplified approach (7 layers).
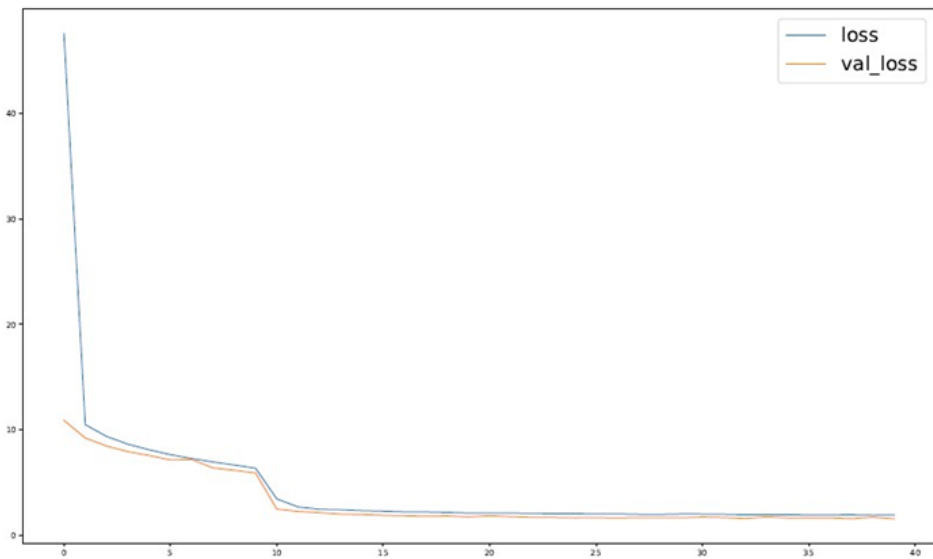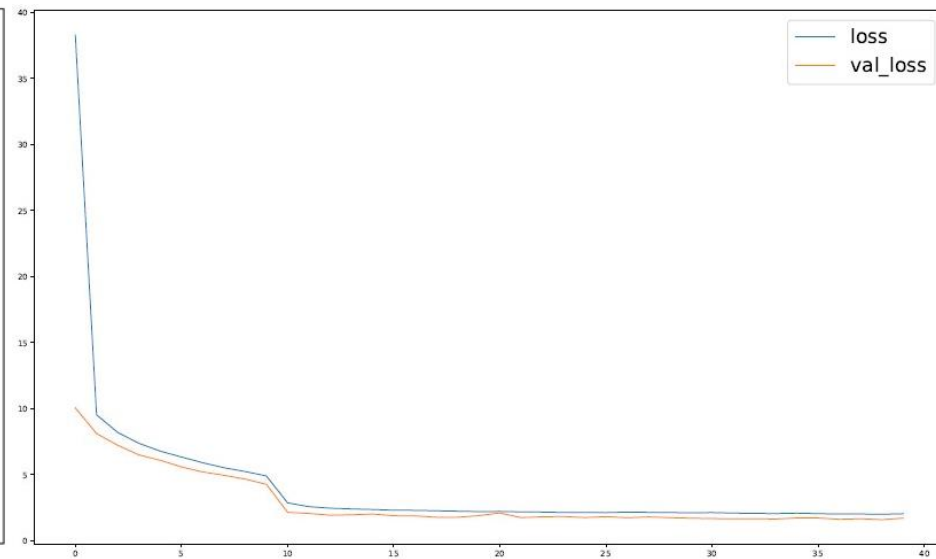


**SSD7**

Epochs = 70

Learning rate = 0,001

# MATERIALS AND METHODS – Training

*SSD300*

*SSD512*



*VGG16*

# MATERIALS AND METHODS – Training

The **Loss function** consists of two terms: $L_{conf}$ and $L_{loc}$ where N is the matched default boxes.

$$L(x, c, l, g) = \frac{1}{N}(L_{conf}(x, c) + \alpha L_{loc}(x, l, g))$$

$L_{loc}$ is the **Localization Loss** which is the smooth L1 loss between the predicted box (l) and the ground-truth box (g) parameters. $L_{conf}$ is the **Confidence Loss** which is the softmax loss over confidences (c).
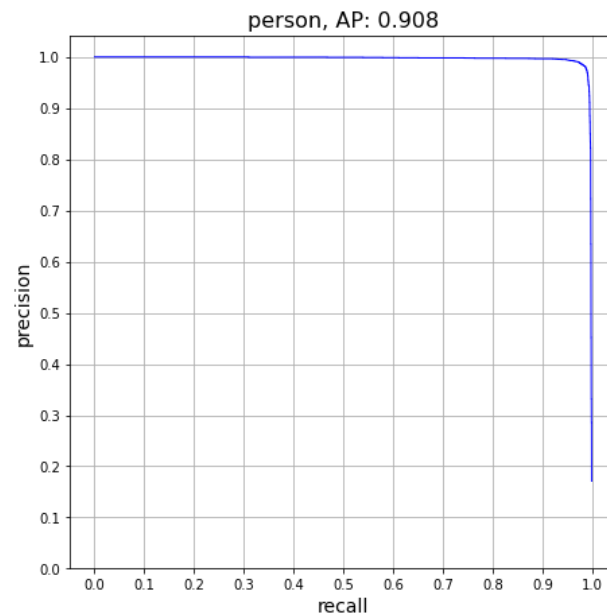
An attempt was made to modify the Loss function computation using *IoU* and F1 Score without having positive results.

# RESULTS AND DISCUSSIONS – Evaluate

The metrics used for the evaluation are **AP** (*Average Precision*), **Recall**, **F1 Score** and **IoU\***
*(Intersection over Union).*

| SSD300 | AP | Recall | F1 Score | IoU |
|---|---|---|---|---|
| VGG16 | 0,908 | 0,818 | 0,861 | 0,879 |
| ResNet50 | 0,909 | 0,618 | 0,736 | 0,842 |

| SSD512 | AP | Recall | F1 Score | IoU |
|---|---|---|---|---|
| VGG16 | 0,908 | 0,912 | 0,910 | 0,801 |
| ResNet50 | 0,909 | 0,872 | 0,890 | 0,846 |

*An attempt was made to use the IoU as localization loss in calculating the loss function during training



SSD512 (VGG16)

UNIVERSITÀ POLITECNICA DELLE MARCHE

VRAi

# RESULTS AND DISCUSSIONS – Fine Tuning

The SSD300 and SSD512 presented *weights* related to models already trained on the Pascal VOC 07 + 12 dataset; a **fine tuning** of the weights was made to go from 20 classes to the only <**person**> class.



person, AP: 0.744

SSD512

| FINE-TUNE | SSD300 | SSD512 |
|-----------|--------|--------|
| AP | 0,643 | 0,744 |
| Recall | 0,887 | 0,934 |
| F1 Score | 0,746 | 0,828 |
| IoU | 0,712 | 0,885 |



person, AP: 0.643

SSD300

# CONCLUSION AND FUTURE WORKS

As confirmed by the results, the proposed approach, i.e. replacing the feature extractor, proved to be excellent to solve the detection problem.

The best architecture remains the SSD512 with the VGG16.

- Changing the feature extractor with other architectures, for example DenseNet.

- Use of IoU metric as Localization Loss in the calculation of the Loss Function, using machines with more memory.

# REFERENCES

[1] Misbah Ahmad, Imran Ahmed, Kaleem Ullah, Iqbal khan, Ayesha Khattak, Awais Adnan, "Person Detection from Overhead View: A Survey".

[2] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, Alexander C. Berg, "SSD: Single Shot MultiBox Detector".

[3] Zhengxia Zou, Zhenwei Shi,Member, IEEE, Yuhong Guo, and Jieping Ye,Senior Member, "Object Detection in 20 Years: A Survey," IEEE.

[4] L. Jiao et al., "A Survey of Deep Learning-Based Object Detection," in IEEE Access, 2019, vol. 7, pp. 128837-128868.

[5] R. Girshick, J. Donahue, T. Darrell, and J. Malik, "Rich feature hierarchies for accurate object detection and semantic segmentation," in Proceedings of the IEEE conference on computer vision and pattern recognition, 2014, pp.580-587.

[6] R. Girshick, J. Donahue, T. Darrell, and J. Malik, "Region based convolution with deep convolutional neural networks," in Advances in neural information processing systems, 2012, pp. 1097-1105.

UNIVERSITÀ POLITECNICA DELLE MARCHE

VRAi