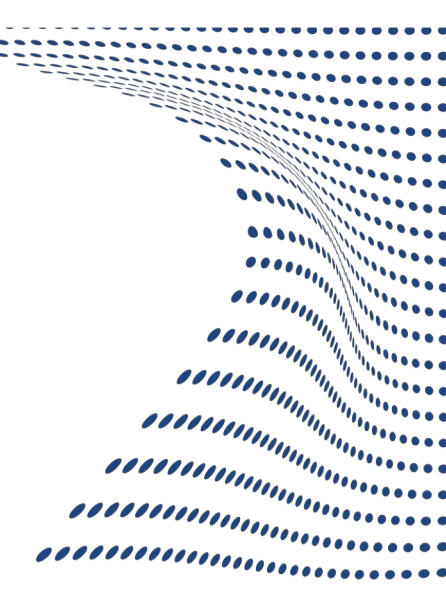# Incremental Learning for Object Detection on Embedded Systems using Machine Generated Buonding Boxes

**Vittorio Mazzia**

**PhD Supervisors:**

Marcello Chiaberge          Mario Roberto Casu

## Motivation and background

Training object class detectors typically requires large amount of data in which images are **manually annotated** with bounding boxes (bbox) for every instance of each class. This is particularly true for lightweight object class detectors that progressively improve their mean average precision ($mAP$) increasing the number of examples available. The presented research suggests a methodology to exploit **generated data** from the field and a **collaboration** with multiple independent deep neural networks to obtain an increasingly more performing **embedded model** for the designated tasks.

## Materials and methods

- **Dataset:** (OIDv4_ToolKit)

| | | Apple | Grape | Lemon | Orange | Pear |
|---|---|---|---|---|---|---|
| **O I D V 4** | Train | 624 | 755 | 367 | 583 | 204 |
| | Validation | 24 | 44 | 41 | 25 | 4 |
| | Test | 57 | 124 | 79 | 95 | 27 |
| | **Videos** | 5' 4'' | 12' 25'' | 34' 3'' | 7' 9'' | 9' 43'' |

- **Hardware**:
  - Tesla K80 (4992 Cuda Cores)
- **Networks**:
  - Faster R-CNN (with ROI-align)
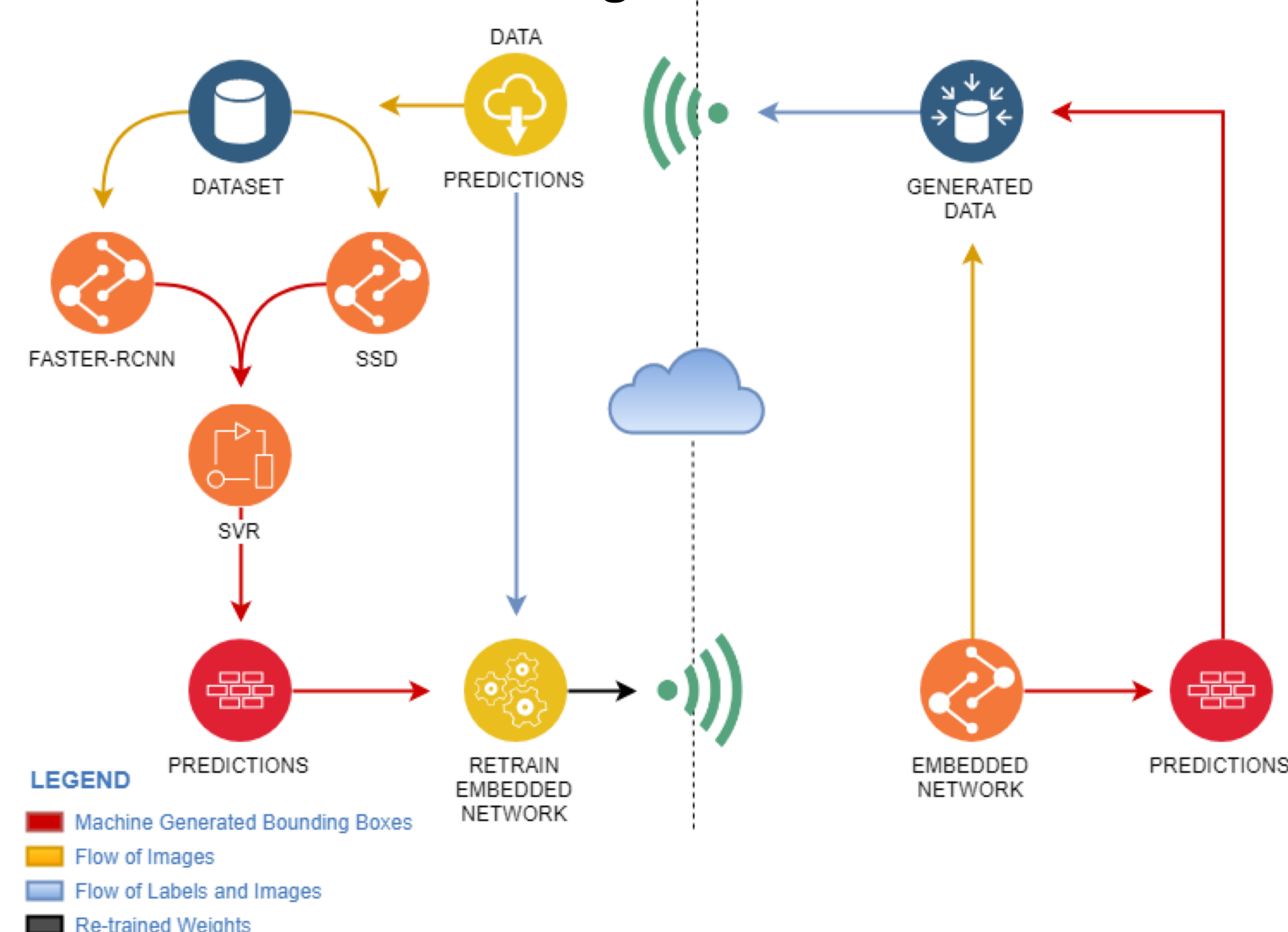  - SSD (with Focal-Loss (1.1))

$$\text{CE}(p_t) = -\log(p_t) \qquad (1)$$
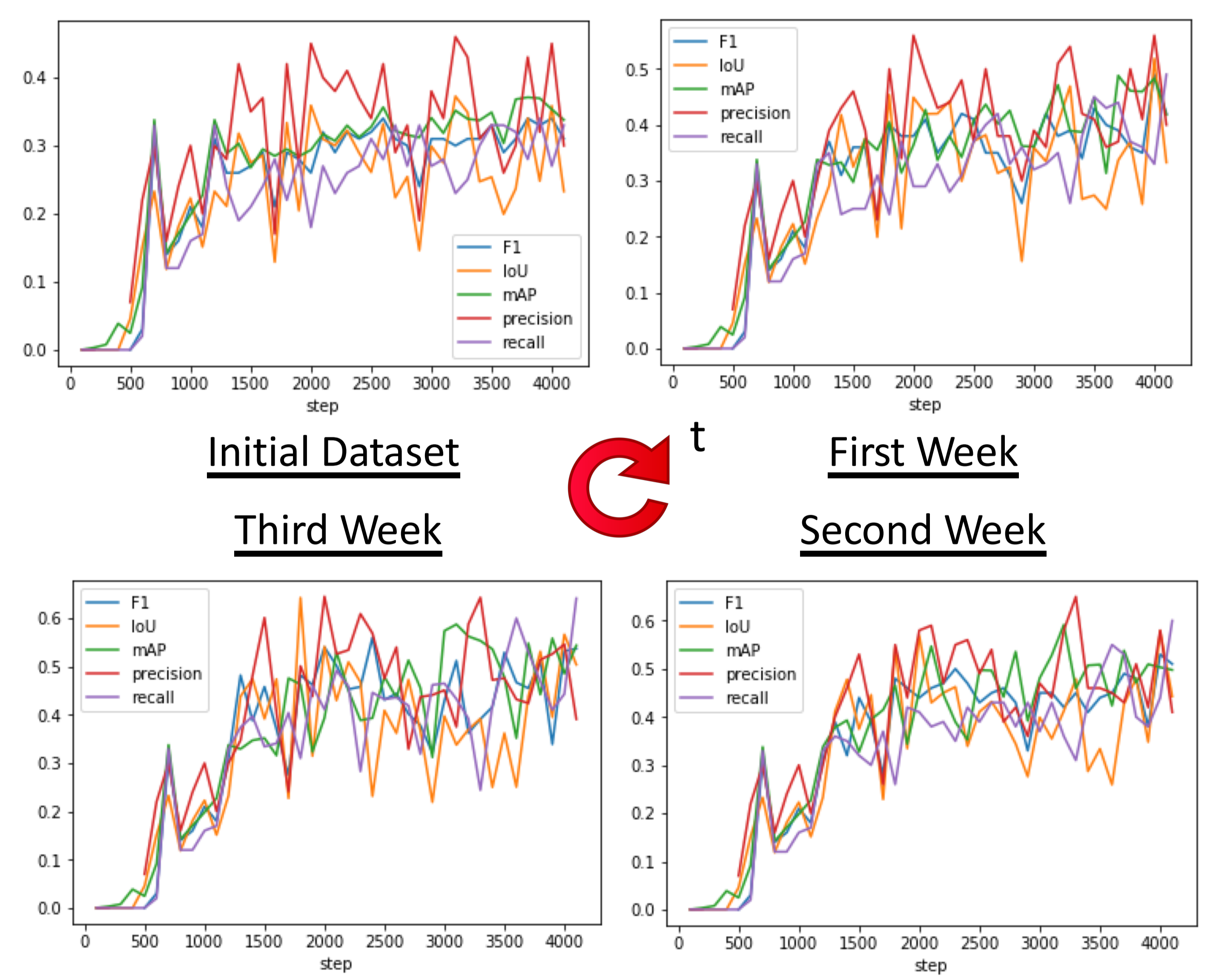$$\text{FL}(p_t) = -(1 - p_t)^\gamma \log(p_t) \qquad (1.1)$$

## Proposed Algorithm

A first architecture of the algorightm is shown in the graph at the bottom and it follows the following steps:

- An initial dataset is used to train a two-stage Faster-RCNN, a Single Shot Multibox Detector (SSD) and a **lightweight version** of it.
- Data generated by the embedded network (frames & predictions) is sent to the cloud.
- Received images are elaborated by the **ensemble network** that generates new bbox.
- New data are merged with the old one and, through a re-training, novel weights of the embedded SSD are generated



## Simulation Results



Initial Dataset          First Week

Third Week          Second Week

## Conclusions and future work

The methodology presented is the first of its kind and preliminary results have proven a remarkable effectiveness of the overall system. However, the proposed research requirese further studies to improve the algorithm and asesess its limitations and drawbacks.

- Substitute the SVR block with a FC layer that exploits backbone extracted features
- Look for saturation value of $mAP$