

# Object Detection-Based Behaviour using YOLO on Thymio

Francesco Saverio Zuppichini<sup>1</sup> and Alessia Ruggeri<sup>2</sup>

**Abstract**—The ability to classify objects in both space and domain is encoded in an enormous variety of species on the planet. Recently, deep learning approaches have been developed to make machines able to mimic the same task with interesting results. With this paper, we propose a minimal architecture to archive the same goal using a two wheels robots, the Thymio, with a mounted frontal camera in order to make it act based on the objects detected in the surrounding.

## I. INTRODUCTION

Finding objects in space by identifying their identity and their coordinates is a basic task that almost each animal is able to accomplish. The ability to solve this problem is extremely important in robotics, where a machine needs to act accordingly based on its surrounding. With our architecture, only the raw data from the camera are used as input to an Object Detection model and the decision are based directly on the prediction, without the need of any other sensor.

Specifically, our robot will wait for an user-specified class, e.g. *dog*, and, when found, it moves towards the object by keeping it in the middle of its view. We also implemented a PID based algorithm to avoid obstacles and properly stop in front of founded item.

The paper is structured as follows: we will first describe the architecture by showing each individual part of the topology and how they interact; then, we will discuss the challenges during the development process and finally we will show the results.

## II. ARCHITECTURE OVERVIEW

We used ROS[ROS] as operating system to develop and deploy our code on the Thymio. Our architecture is composed by:

- 1) A robot with a mounted camera and an on-board computer
- 2) A web server running the Object Detection model
- 3) A ROS node running on a local machine that moves the robot

As a robot, we used the *MightyThymio* [2], which is a robot developed on top of the Thymio, a small mass-produced mobile robot usually used in primary and secondary education. The MightyThymio has a powerful Odroid C1 on board, a 1GB quad-core ARM Linux machine, with a HD camera mounted on top. The authors also implemented a ROS node to communicate with it and open-sourced the code.

### *MightyThymio*

In our project, the MightyThymio has the same configuration proposed in [2]. Briefly, the on-board computer expose the Thymios sensors through ROS. Thus, another ROS node

can publish and subscribe in order to communicate with its sensors. For example, we can subscribe to the *camera sensors* to get information from the proximity sensors to avoid hitting obstacles. Also, we can simply subscribe to the camera's topic to get both compressed and uncompressed images.

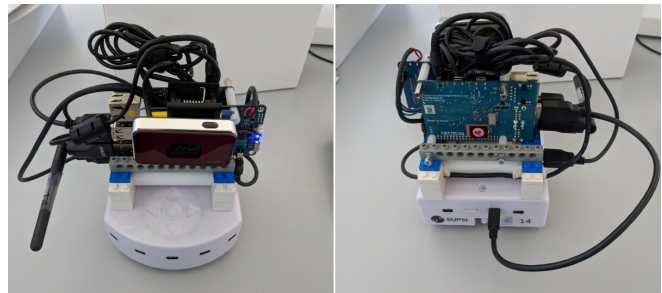


Fig. 1. The MightyThymio

### *Local ROS node*

We created a new ROS node to communicate with the running node on the MightyThymio in order to make it act accordingly with a specified behaviour based on the object detected by the model. The user specifies a set of target classes, for example *dog* and *cat*. Then our node, after having subscribed to the camera topic of the MightyThymio, sends the fetched image to the model trough an HTTP request and gets the bounding boxes for the detected classes. If there is no target prediction, the robot enters in *exploration* mode and starts spinning around. When one or more of the target classes are found, then our robot will start going forward while keeping the bounding box of the target in the centre of the camera in order to follow it. If we detect more classes, then the closest, which is represented by the biggest bounding box area, will be followed. In addition, we have an obstacle-avoiding algorithm in order to properly stop the Thymio before it hits the target.

### *Object Detection REST-API Web Server*

We used pre-trained YOLO[4], the state-of-the-art real-time object detection model, to predict the bounding boxes for each class. The model was trained on the COCO(Common Object in Context)[1] Dataset, that contains 80 classes.

On top of that, we built a REST-API Web Server using Python3 and Flask, a micro-framework to create Web Servers, to make the model accessible from the Web. We decided to decouple the logic behind the robot behaviour and the model in order to make our architecture more scalable.

For instance, if needed, we can deploy our model in the cloud while keeping the running ROS code intact.

We following endpoints are available:

```
GET /model
POST /prediction --image --size --compress
```

Each of them returns a JSON serialised response. The `/model` endpoint is used to get informations about the model on the server. These informations are the classes and the colour associated to each one of them. It follows a snapshot of the response:

```
{
  "classes": [ "person", "bicycle", "car", "motorbike", "
    aeroplane", ... ]
  "colors": [[255, 38, 0], [0, 63, 255], ... ]
}
```

The client can get a prediction by sending an image to `/prediction`. Also, it can specify whatever the image is compressed or no and if the server should resize it before feed it to the model. In our case, we send a compressed image directly taken from the MightyThymio after subscribed to the `/camera/compressed` topic and we define the resize size to (190,190). The small size is due low computational power since we had no access to a GPU. The response looks like:

```
{'res': [{'boxes': [96.21580505371094,
  219.6134033203125,
  437.66387939453125,
  594.16748046875],
  'class': 'bicycle',
  'class_idx': 1,
  'score': 0.48406994342803955}, ...
]
'time': 0.174407958984375}
```

For each detected class we return the bounding boxes, the `boxes` key of the json, the class name, its index and the score, the confidence of the model. Each box is composed by the top, left, bottom and right coordinate of the bounding box with respect to the left of the image and the top. In the above example, the bicycle's bounding starts 96 pixels from the left and 219 from the top.

### III.

#### A. YOLO: Real-Time Object Detection

YOLO [4] is a trained model for real-time object detection. While prior work on object detection repurposes classifiers to perform detection, YOLO handle the object detection as a regression problem to spatially separated bounding boxes and associated class probabilities. YOLO avoids resorting to complex pipelines that are slow and hard to optimize; instead, it uses a single convolutional neural network that simultaneously predicts multiple bounding boxes and class probabilities for those boxes. Using YOLO, You Only Look Once at an image to predict what objects are present and where they are.

YOLO neural network uses features from the entire image to predict each bounding box simultaneously across

all classes. The convolutional neural network contains 24 convolutional layers, that extract features from the image, followed by 2 fully connected layers, whose aim is to predict the output probabilities and coordinates. The system divides the input image into an  $S \times S$  grid; if the centre of an object falls into a grid cell, that grid cell is responsible for detecting that object. Each grid cell predicts  $B$  bounding boxes and confidence scores for those boxes. These confidence scores reflect how confident the model is that the box contains an object and also how accurate it thinks the box is that it predicts. There is also a fast version of YOLO, designed with a neural network with fewer convolutional layers and fewer filters in those layers.

YOLO is extremely fast, it reason globally about the image when making predictions, it learns generalizable representations of objects and it is open source. YOLO still lags behind other detection systems in accuracy; while it can quickly identify object in images, it struggles to precisely localize some objects, especially small ones. However, it is still the state-of-the-art model for object detection in real-time.

For our project, we used the third version of YOLO, or YOLOv3 [3], which has been improved by making it a little bigger than before, but more accurate.

### IV. IMPLEMENTATION

We design our ROS node to find object and follow them until it is close enough. In addition, we also implemented a PID based algorithm to avoid obstacle and to stop before crash into the target object. From an high level view, the robot can be represented as a three-state machine.

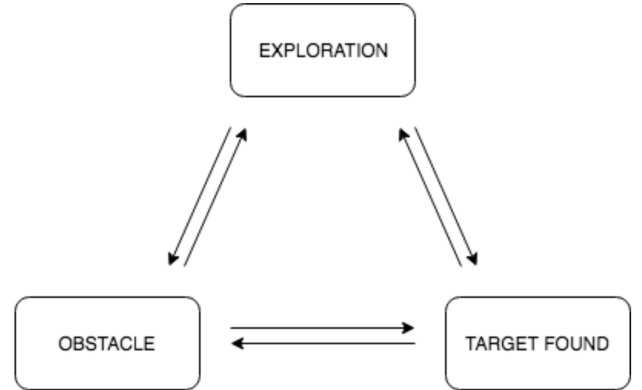


Fig. 2. State Machine

Thus, at any time the machine can be **only** in one of the three states: *exploration*, *obstacle* and *target found*.

a) *Exploration*: The robot starts to spin on itself with a fixed angular velocity to scan the surrounding until it finds the target or it detect an obstacle.

b) *Target Found*: When we find the object we are looking for, we start moving towards it by keeping the its bounding box in the camera's center. We used a Proportional Integral Derivative, PID, controller to select the correct angular velocity. Given the camera's image and the bounding

box we define the error at time  $t$  as:

$$e_t = (l_b + \frac{w_b}{2}) - \frac{w_c}{2} \quad (1)$$

Where,  $w_c$  is the camera's image width, in our case 640,  $l_b$  and  $w_b$  are the left component and the width of target's bounding box respectively. To find  $w_b$  we need to just subtract the the left and right component of the box. We are not taking the absolute value of the error since the sign is used to decide in which direction to turn. Before feeding the error to the PID controller we normalise it between zero and one by dividing by  $w_c$ , otherwise we should have defined a very small proportional parameter for the PID.

After we calculate the error for the current box, we feed it to a PID with a proportional parameter of 0.5, and integral of 0 and a derivative of 3. However, even if we did not do any hyperparameter search, we notice that even with smaller values of the derivative component of the PID the robot works properly.

The following figure shows the error for each frame given different position of the target class *bottle*.

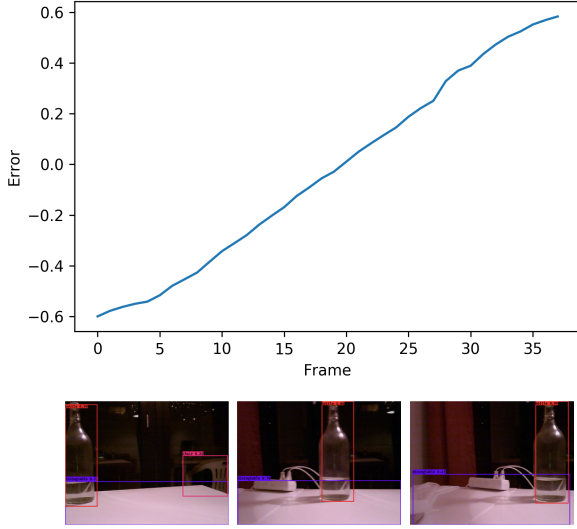


Fig. 3. Bounding Box Error

The error is bigger when the bottle is on the left on the right and it is zero when it is in the center of the image. We used it as angular velocity in order to keep the robot pointed to the target object. We move it forward with a fixed linear velocity of 0.1.

c) *Obstacle*: When the Thymio's proximity sensors of the Thymios, five in the front and two in the back, detect something it stops any activity and goes into the *Obstacle* state. We implemented a simple yet effective PID controllers based algorithm to move the Thymio away from any object in from of it. We defined two PIDs, one for the linear velocity and one for the angular. The errors are  $e_l$  for the linear PID and  $e_a$  for the angular PID, defined as follow. We denote the front sensor set as  $S_f$  and the rear sensor set as  $S_r$ , where

each element in the set is a double representing the value from the sensor. In our case,  $|S_f| = 5$  and  $|S_r| = 2$ .

$$e_l = \sum_{i=0}^{|S_f|} S_f(i) - \sum_{i=0}^{|R_f|} R_f(i) \quad (2)$$

$$e_a = (\sum_{i=0}^{\frac{|S_f|}{2}} S_f(i) - \sum_{i=\frac{|S_f|}{2}}^{|S_f|} S_f(i)) - (\sum_{i=0}^{\frac{|R_f|}{2}} R_f(i) - \sum_{i=\frac{|R_f|}{2}}^{|R_f|} R_f(i)) \quad (3)$$

Also, we turn on the five frontal Thymio's LED as user feedback. The number and position of the LEDs are based on the width of target's bounding box. For example, a big object on the left side will turn on two LEDs, while a smaller one only one. The follow picture shows an example where the robot is looking for a *cup*:

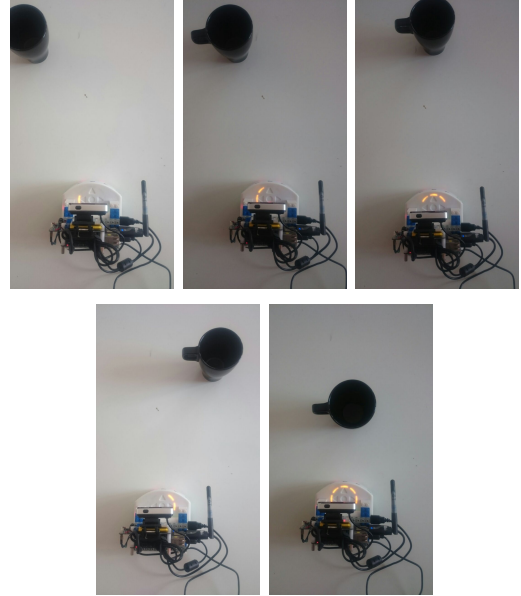


Fig. 4. On board LEDs

## V. SERVER

## VI. RESULTS

## VII. CONCLUSION

In this paper we presented a light and scalable architecture that uses a real-time Object Detection model to correctly find objects in the surrounding environment and to perform actions based on that. Since we only require a camera, our work can be used into similar robots running ROS with the correct hardware in order to extend their functionalities. Furthermore, since we used YOLO to perform the classification task, the model can be easily re-trained with other classes to expand the context domain. This symbiosis between Machine Learning and Robotic highlights the power of both areas by creating new interesting user cases and scenarios in which both disciplines can shine.

#### REFERENCES

- [1] Tsung-Yi Lin Michael Maire Serge Belongie Lubomir Bourdev Ross Girshick James Hays Pietro Perona Deva Ramanan C. Lawrence Zitnick Piotr Dollar. "Microsoft COCO: Common Objects in Context". In: 2015.
- [2] Jerome Guzzi et al. "Mighty Thymio for Higher-Level Robotics Education". In: *AAAI*. 2018.
- [3] Joseph Redmon and Ali Farhadi. "YOLOv3: An Incremental Improvement". In: 2018.
- [4] Joseph Redmon et al. "You Only Look Once: Unified, Real-Time Object Detection". In: 2016.