

# CATIE Robotics @Home 2020 Team Description Paper

Rémi FABRE, Boris ALBAR, Clément DUSSIEUX, Christine JAUREGUIBERRY, Ludwig JOFFROY, Jean-Baptiste HOREL, Zhe LI, Alexandre PERROT, Clément PINET, Stéphane POUYET, Florian LARRUE, and Sébastien LOTY

Centre Aquitain des Technologies de l'Information et Electroniques (CATIE)  
1 Avenue du Dr Albert Schweitzer, 33400 Talence, France  
[r.fabre@catie.fr](mailto:r.fabre@catie.fr)  
<https://robotics.catie.fr/>

**Abstract.** This paper provides an overview of the CATIE Robotics Team's activities for participating at the RoboCup@Home 2020 in Bordeaux. A TIAGo<sup>1</sup> called Epock is used as the main robotic platform. Safe and autonomous navigation is achieved by integrating ROS compatible components. Snips, an offline, free tool is used for both speech recognition and natural language understanding. State-of-the-art neural networks are used for face recognition, person tracking and object detection. Problem specific developments are successfully tested for object localization. A Human-centered approach is used to enhance Epock's interactions. A simple but functional grasping pipeline is implemented. In 2019, good results have been achieved in three international robotics competitions validating the strong technical foundation the team has built. This document describes selected components, shares some feedback and discusses possible improvements.

## 1 Introduction

The CATIE Robotics Team has been formed at the beginning of 2018 and is part of CATIE, a digital technology transfer center at the crossroads between research and industry. CATIE is a non-profit organization supported by the Nouvelle-Aquitaine French region whose mission is to assist companies willing to adopt and integrate digital technologies in their technological and economic development. By creating a RoboCup@Home team, our ambition is to be part of an experts community, nurture our robotic knowledge and share it to foster progress towards a tangible goal. The competition offers an objective benchmark to measure progress and is a stimulating means to unite efforts locally.

In 2018, we focused on integrating proven technologies and achieving simple but robust behaviours in key fields such as safe navigation, Simultaneous Localization And Mapping (SLAM), force based grasping, person and object

---

<sup>1</sup> <https://tiago.pal-robotics.com/>

recognition. TIAGo, an off-the-shelf service robot, was chosen as our main development platform and Robot Operating System (ROS) as our middleware framework. This approach worked well in our first international competitions: we ranked second in GermanOpen’s 2019 RoboCup@Home<sup>2</sup>, third in Sydney’s 2019 RoboCup@Home and second in SciRoc’s Episode 7 challenge. As we now enter year two, our goal is to reach integrated, robust and safe high level behaviors for our robot. We currently work to enhance our vision pipeline and make it shareable with the community. We also explore flexible grasping and redesigning Epock’s head to enhance user interactions, particularly in order to make the robot’s intentions clearer to laymen.

In this paper, developments related to the navigation, perception, communication and object manipulation capabilities of Epock, as well as hardware modifications are presented. For each section, an overview of the approach, the results and the future work is given.

## 2 Navigation and SLAM capabilities

A robust, reliable navigation has always been our priority. We have established a full and working pipeline from map creation to autonomous navigation and localization. No collision happened during the total, long distance covered by the robot in different environments throughout the past year. We use AMCL for the localization and ROS move\_base for the navigation. A static map is created with Gmapping. We work on an approach using AMCL and Gmapping/Cartographer simultaneously. The SLAM algorithm task is to perform a low frequency mapping to update the map during navigation. At the same time, the localisation is acquired using AMCL.

## 3 Perception capabilities

Perception problems are mainly tackled by several different state-of-the-art neural networks detailed below. For each neural network, we have tried to find the best compromise between accuracy and performance. We mostly rely on siamese networks. This type of neural network can learn a similarity measure between objects [1]. Every object is mapped by the network to a vector. Object comparison can then be achieved using the distance between vectors. This enables these networks to naturally generalize to unknown instances. We also sometimes use some more problem specific approaches, described in dedicated subsections.

### 3.1 Object Recognition

Object recognition is mainly done using YOLO. We trained our own PyTorch implementation of this neural network from scratch on the COCO dataset. Using

<sup>2</sup> <https://www.youtube.com/watch?v=7Y4RjxWRqxE&t=5s>

$320 \times 320px^2$  images as input, we were able to obtain a  $mAP_{50}$  of 50 without overfitting, close to the 51.5 value indicated in the original YOLO paper [2].

A new dataset is created for each new object we have to deal with and only the last YOLO layer is retrained. The dataset creation process is pretty heavy, so we designed a mechanical system to take photos of the objects, consisting of a rotating podium and a chroma-key background (green or blue) to partially automatize the process. The objects are then included in various fake backgrounds and augmented in several ways (luminosity, chromaticity, etc.). A COCO-like dataset is then generated and used for training. But the result of object recognition during the inference in the real environment could be improved. We therefore considered the solutions like Generative Adversarial Networks or Domain Randomization to overcome the gap between the synthetic and real data. These approaches may increase the quality of the generated dataset and increase performances in object recognition.

We recently added support for the MaskRCNN architecture, which enables the object mask extraction and not only a bounding box. Training is done using the same dataset described above. Recent work tries to extend this architecture in a siamese-like way. If it is successful, this approach could be used for one-shot object recognition and could enable to skip retraining the neural network each time an object is added.

Pragmatic approaches sometimes work better than neural networks. In the *SciRoc Pick and Pack* challenge, we used point cloud depth filtering to detect objects, based on the fact that they are on the shelves. We also took advantage of objects not being too close to one another to identify distinct blob of points for each object, with basic color filtering. This code was written with OpenCV.

### 3.2 Face detection and recognition

Face detection and recognition is done by a combination of two neural networks. The first one is a network that is used out of the box and that performs the face detection using a state-of-the-art Multi-Task Cascaded Convolutional Neural Network (MTCNN). Once the face bounding boxes are extracted as well as face features (eyes, nose, ...), the image is cropped and normalized. Faces are then processed through a siamese network. This network has been trained on more than 8000 different identities using the publicly available VGG 2 dataset. The distance between the faces and the previously detected people are computed and the closest person according to the L2 distance is returned if the value is below a predefined threshold. Examples of uses include checking if the right person is following Epopk or during the *Find my mate* task in Robocup@Home for identifying people faces.

### 3.3 Person attributes detection

We trained a neural network to identify a person's attributes from a picture. The network was trained on the CelebA public dataset using a resnet architecture.

The training of this network relies heavily on the described above face detection process for normalization. Attributes include glasses, hat, hair color, face shape, etc. This is used to generate descriptions for the *Find my mate* task in Robocup@Home. We wish to add attributes like hair size, skin color and some other useful features that are missing from the dataset using transfer learning from the current network.

### 3.4 Person tracking and *follow me* behavior

We worked on a solution to track people using a siamese network combined with a Kalman filter. First, we use YOLO to extract the person to be followed. We then get a vector representing the target vector using a siamese neural network trained on the Market1501 public dataset. This dataset contains a large amount of identities with images from six different cameras. During the tracking phase, several different people can be present. Each vector corresponding to a detected person is compared to the target vector in order to keep following the right person. We currently use a moving average of the previously detected vectors as a reference for distance calculation. This way, the target person’s orientation or size can change during long tracking sessions without disturbing the system. We combined a Kalman filter<sup>3</sup> to use the available detections and the previous predictions to obtain the best guess of the person’s current position. Assuming constant velocity motion, we can easily predict the target person’s next position. With the measurements and the guess of the position given by the Kalman filter, we can filter the wrong people with more confidence. This approach could be used to track any type of object by generalizing the network to learn features independently of the object type.

During the *follow me* phase, we control the robot’s head to always be looking in the direction of the target’s last known position. Using the target’s position as a goal position, we use Dijkstra’s algorithm to create the path to the target taking into account the obstacles registered on the costmap by the LIDAR and the RGB-D camera. The velocity commands are computed by the local\_planner.

### 3.5 Pose Recognition

The pose recognition is done with the OpenPose library<sup>4</sup> which computes face, body and hands 2D keypoint detection in real-time if equipped with a recent graphics card.

OpenPose gives the Body-Foot Estimation (default skeleton), but in the @Home competition, we need a more high-level information like "the person is seated" or "the person is pointing to the left". This is determined by simple rules based on the skeleton, such as: "If the knees are approximately at the same height than the hips and the neck is significantly higher, then the person is sitting". We wrote this code on top of OpenPose. It can categorize the person’s

<sup>3</sup> [https://en.wikipedia.org/wiki/Kalman\\_filter](https://en.wikipedia.org/wiki/Kalman_filter)

<sup>4</sup> <https://github.com/CMU-Perceptual-Computing-Lab/openpose>

pose (sitting, standing, raising arm, etc.). A way to improve the pose recognition would be to use 3D keypoints.

### 3.6 Precise object localization

When Epock has to interact with an object at a known position in the environment, it cannot only rely on its own localization.

Sometimes the actual object position can be slightly different from the theoretical one and Epock's localization alone is not accurate enough for precision tasks such as grasping.

For example, in the *take out the garbage* challenge from Robocup@Home, robots need to detect garbage bins to take trash bags out of them. In *SciRoc Pick and Pack* challenge, robots need to drop the items in a crate. These tasks require a precise localization of both the bins and the crate.

A LIDAR based pattern recognition is used to detect objects with a specific geometry. Indeed, a trash bin seen by the LIDAR is a semicircle. To add more robustness to this detection, the object theoretical position is also taken into account. An object matching the geometry, but located far from its expected position is probably not the seeked one.

## 4 Communication capabilities

### 4.1 Text-To-Speech, Natural Language Understanding and speech synthesis

For the speech-to-text and Natural Language Understanding (NLU) capabilities, we are using Snips <sup>5</sup>. Snips is an off-the-shelf solution, free for non-commercial use software. Although its accuracy may be lower than other well-known cloud-based solutions, both the speech recognition and the NLU engine run offline with low latency, which is vital for the RoboCup competition. It is very lightweight, it has been designed to run on a Raspberry Pi. However, one drawback is that the speech-to-text only works with examples that have been learnt, it does not work with a new sentence pattern or new words. Every step of the Snips process is linked to a specific MQTT topic, this is a very convenient feature to keep control over the stack<sup>6</sup>. Custom code will then take over and call the appropriate functions.

To fix some problems encountered during previous competitions and match its needs, we tweaked Snips to get some kind of continuous speech-to-text. Listening is programmatically triggered every few seconds and the robot tries to match each phrase heard with what it has learnt. The confidence score must be high enough to avoid triggering false positives. The continuous text-to-speech is quite experimental.

<sup>5</sup> <https://snips.ai>

<sup>6</sup> <https://snips.gitbook.io/documentation/ressources/hermes-protocol>

One observed issue in real conditions was some variability in speech recognition quality. In a robustness driven approach, we scripted automatic record of every audio interaction to double check every fail. This enabled us to detect some noises by hardware at different levels of the sound acquisition pipeline.

We are also starting to explore ways to do text-to-speech without Snips, to gain more control and freedom in this process.

For text-to-speech, we are using TIAGO’s default text-to-speech module: Acapela<sup>7</sup>.

## 4.2 Improvement of Epock’s Collaboration with Humans

The different challenges encountered during our competitions were an opportunity to see how Epock interacts with people. We noticed that some people had difficulties understanding when to speak to Epock or what Epock was doing. To address this issue, we worked on interactions with a Human-centered approach by adding a screen where the current state of Epock is displayed (help needed, listening, moving, in action). A user study was conducted to see if these pictograms were well understood and to get some feedback. Impressions and comments of people not involved in the project were collected. Now, with the additional pictograms, users can adapt their behavior according to the robot state shown on the screen as they would with another human. Our method was based on Bastien and Scapin’s Criteria<sup>8</sup>.

This first study should help initiate an in-depth reflection on the creation of a global system thought to improve human/robot collaboration in a non-experimental environment.

## 5 Grasping capabilities

### 5.1 Arm control

Grasping as a whole is a new subject for the team this year but is of high interest and as such represents a substantial portion of the work done since 2019’s TDP. Our approach prioritizes security, then robustness and only then efficiency. Depending on the situation, the arm will either be controlled using predefined joint movements<sup>9</sup>, using MoveIt! planners<sup>10</sup> or using PAL’s closed implementation of the Whole Body Control (WBC)<sup>11</sup>.

To reach an acceptable security level, we use high level guidelines and low level behaviors. The guidelines include:

- Position the robot in order to maximize the surrounding space.

<sup>7</sup> <http://www.acapela-group.com/>

<sup>8</sup> [https://www.cocoaheads.fr/wp-content/uploads/files/Ergonomic\\_Criteria.pdf](https://www.cocoaheads.fr/wp-content/uploads/files/Ergonomic_Criteria.pdf)

<sup>9</sup> [http://wiki.ros.org/play\\_motion](http://wiki.ros.org/play_motion)

<sup>10</sup> <https://moveit.ros.org/>

<sup>11</sup> <http://blog.pal-robotics.com/whole-body-control-online-planner-tiago-robot/>

- Only move the arm if the robot is still.
- After grasping an item, try to retract the arm as best as possible before moving the base of the robot.
- Always announce vocally that the arm will move.

As for the low level, a "zero gravity" behavior was implemented<sup>12</sup> based on the force/torque sensor attached to the wrist. Its principle is simple: if an abnormal force or torque is felt, then the arm is moved in the direction that would reduce that force/torque. Once the disturbance disappears, the arm resumes its movement<sup>13</sup>. When grasping an item, its weight must be accounted for. If the weight is unknown (e.g. *Take Out The Garbage* test), this security is turned off when the item is grasped and turned on as soon as the item is weighed with the sensor. In its current state, the arm doesn't take into account the current measures of its motors and therefore won't adapt if an obstacle is touched with the elbow for example, this is a work in progress.

Epock's force sensing capabilities are also used to improve the overall robustness of the grasping pipeline. For example, if the weight measured by the wrist sensor goes below a predetermined threshold, the robot will assume that the item fell off its gripper and try to grasp it again if possible<sup>14</sup>. Also, some items are smaller than the precision we're capable of reaching at the end of the gripper along the Z axis (e.g. a fork on a table). In this case, we'll use a slow descending motion until a resistance is felt when the table is touched. The same approach was used to grasp the garbage whose size is unknown prior to the test.

We've tested the impact of using checkpoints to achieve a reproducible behavior and tried to evaluate the appropriate force to apply on objects depending on their resistance and stiffness. Formalizing these methods and making a generic module out of them is a work in progress. Finally, exploratory work is currently being pursued with soft grippers and will be integrated in Epock if the results are good enough.

## 5.2 Perception challenges

Grasping requires specific information from the object. This could include the position, the mask or the orientation of the object. We reviewed the state of the art on object recognition for grasping in 2018. The conclusion was that the 6D Object Pose Estimation methods seemed not accurate enough and too slow for real time, based on the Bop benchmark [3]. For now, we used either YOLO and the bounding box it provides or custom OpenCV code to identify specific features of the objects (color, height, ...) in the image. Using the latter, we achieved good results in the *Pick and Pack* final task of the SciRoc 2019 competition (6/6 items picked, 5/6 items packed).

<sup>12</sup> <https://www.youtube.com/watch?v=4auVU5Ifvpw&t=8s>

<sup>13</sup> <https://www.youtube.com/watch?v=7Y4RjxWRqx&t=6s>

<sup>14</sup> <https://www.youtube.com/watch?v=btS7S6dadN4>

Recent publications in the 6D pose estimation field (PoseCNN and PVNet) now show promising results, both in terms of accuracy and speed. We will try and test these new approaches in a near future.

## 6 RoboCup experience and community outreach

In addition to the results presented in the introduction, we were exhibitors at 2018 and 2019 Cap Sciences' Village des Sciences, that gathered more than 3000 people over the weekend around robotics and the RoboCup competition<sup>15</sup>. We took part in the following events in 2019: NAIA (Bordeaux) and Vivatech (Paris). We participated in RoboCup@Home Education Challenge @EURCJ 2018 and co-organized a similar workshop in early 2019 that gathered 30 students in Bordeaux.

## 7 Conclusion

In this paper, we have given an overview of the approaches used by the CATIE Robotics team for the RoboCup@Home competition. We have detailed our approaches for navigation, detection, communication and grasping. In all these areas, we have made significant improvements, but we are still building a robust basis for the competition by catching up to the state of the art, which consumes most of our time. We have though shared our contribution to Vizbox and hope this is the first of many more.

## References

1. Gregory Koch, Richard Zemel, and Ruslan Salakhutdinov. Siamese neural networks for one-shot image recognition. In *ICML deep learning workshop*, volume 2, 2015.
2. Joseph Redmon and Ali Farhadi. Yolov3: An incremental improvement. *arXiv*, 2018.
3. Tomas Hodan, Frank Michel, Eric Brachmann, Wadim Kehl, Anders GlentBuch, Dirk Kraft, Bertram Drost, Joel Vidal, Stephan Ihrke, Xenophon Zabulis, et al. Bop: Benchmark for 6d object pose estimation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 19–34, 2018.

---

<sup>15</sup> <http://www.cap-sciences.net/au-programme/evenement/village-des-sciences-2018>



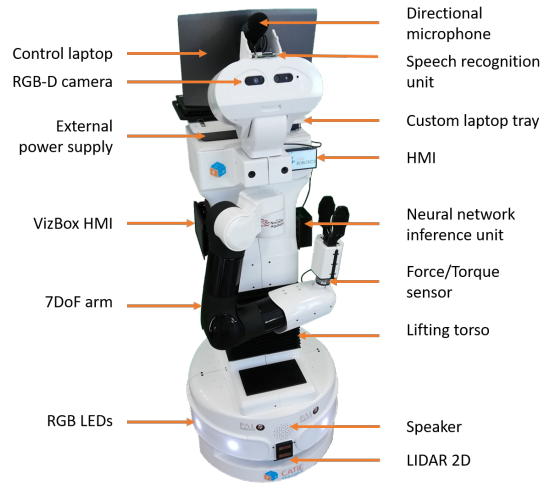
## Epock - Robot TIAGo Hardware Description

Robot TIAGo has been selected and is being customized for the @home competition purpose. Specifications are as follows:

- Base: differential drive base, 1m/s max speed.
- Torso: lifting torse (35cm lift stroke)
- One arm with a gripper (7 DoF). Maximum load: 2kg.
- Head: 2DoF (pan and tilt)
- Robot dimensions: height: 1.10m - 1.45m, base footprint: 54cm diameter
- Robot weight: 72kg.

*Our robot incorporates the following sensors:*

- RGB-D camera
- 2D LIDAR
- Stereo microphone
- Speaker
- Sonars
- IMU
- Motors current feedback
- Wrist force and torque sensor



**Fig. 1.** Robot TIAGo

## Robot's Software Description

*For our robot we are using the following software:*

- OS: Ubuntu 16.04
- Middleware: ROS Kinetic
- Simulation: Gazebo  
<http://gazebosim.org/>
- Visualisation: RViz  
<http://wiki.ros.org/rviz>
- Localization: AMCL  
<http://wiki.ros.org/amcl>
- SLAM: Cartographer and GMapping  
<https://github.com/googlecartographer/cartographer>  
<http://wiki.ros.org/gmapping>
- Navigation: move\_base  
[http://wiki.ros.org/move\\_base](http://wiki.ros.org/move_base)
- Arms control: moveIt! and play\_motion  
<http://moveit.ros.org/>  
[http://wiki.ros.org/play\\_motion](http://wiki.ros.org/play_motion)

- Face recognition: custom siamese neural network
- Object recognition: YOLO  
<https://pjreddie.com/darknet/yolo/>
- Pose detection: Open Pose  
<https://github.com/CMU-Perceptual-Computing-Lab/openpose>
- Speech recognition: Snips  
<https://snips.ai/>
- Speech generation: Acapela  
<http://www.acapela-group.com/>
- Task executor: SMACH  
<http://wiki.ros.org/smach>

## External Devices

*Our robot relies on the following external hardware:*

- Rode Videomic Pro external microphone
- External laptop
- 2 touch screens