

# Markovito Team Description Paper for RoboCup@Home 2023

L. Enrique Sucar<sup>1</sup>, Eduardo F. Morales<sup>1</sup>, José Martínez-Carranza<sup>1</sup>, Alejandro Gutiérrez-Giles<sup>1</sup>, Sergio A. Serrano<sup>1</sup>, Arquímedes Méndez-Molina<sup>1</sup>, Jorge Gutiérrez<sup>1</sup>, and David Carrillo-López<sup>1</sup>

Instituto Nacional de Astrofísica, Óptica y Electrónica (INAOE), Luis Enrique Erro  
# 1, C.P. 72840, Tonantzintla, Puebla, México  
<http://robotic.inaoep.mx/markovito/>

**Abstract.** In this paper we present the current status of Markovito Jr., a service robot developed by the Markovito team at INAOE. Our robot is based on a Robotnik RB-1 robot platform and incorporates a set of general purpose modules that are integrated in a layered behavior-based architecture implemented on the Robot Operating System (ROS). Our current research focuses on the development of novel algorithms for the detection, recognition and manipulation of objects, people tracking, and human-robot interaction, as well as improvements on our modular software architecture. Important results were obtained by our group in robust indoor navigation, grasp confirmation and reinforcement learning published at the Mexican Inter. Conf. on Artificial Intelligence 2019, and the RoboCup Symposium 2021. In the context of robotics competitions, our team has achieved good results in the national contests and also has previous participation at the international RoboCup@Home competition.

## 1 Introduction

Robotics in its evolution has been leaving aside the repetitive tasks of controlled industrial environments and, with a new generation of skilled robots in the execution of domestic tasks known as service robots, is getting involved in dynamic environments with human interaction. When introducing service robots into society, they are expected to have skillful mechanisms to perform most of the tasks a human can perform.

Service robotics is a challenging, and exciting field that for several years has motivated our research group Markovito<sup>1</sup> at INAOE. With our robotic platform, we actively develop research in many subfields of service robotics. The Mexican Tournament of Robotics (TMR) has been the main competitive platform to test our results annually. Besides the good historical results obtained in the robotics competitions, the philosophy of our group is mostly focused on research in the various areas related to service robotics. As a result, our group has several

---

<sup>1</sup> Markovito Home Page: <http://robotic.inaoep.mx/~markovito/>

publications in international journals and conferences, and has served for several students from INAOE to successfully develop their master’s and doctoral theses.

In this Team Description Paper we present the recent research carried out by our group, which is directly related to achieving the fundamental skills required by a service robot in the context of the RoboCup@Home competition. The rest of the paper is structured as follows: in Section 2 we describe our team in more detail; Section 3 presents the general software architecture we use; an overview of the latest research developments in our laboratory in the different robot skills is presented in Section 4; finally the conclusions and some lines of future work are presented.

## 2 Group’s description and research focus

Markovito’s team is composed of researchers and students from the Robotics Laboratory of INAOE. The laboratory was created in 2007, however, the first research in the field of robotics dates back to 2001. Before our current robotic platform Markovito Jr. (2017 - ), we used Sabina (2012 - 2017) and the original Markovito (2007 - 2012).

The team integrates a variety of researchers and students. This year the team is constituted by four senior researchers, one post-doctoral researcher, three PhD students, four B.Sc. students on research internships, and one technician.

The Markovito team has participated in the RoboCup@Home category at previous RoboCup competitions in 2009, 2011, 2012 and 2016; in Turkey 2011 our team qualified for the second stage of the competition. We have participated in almost all Mexican RoboCup@Home competitions; in 2015 and 2016, the Markovito team won the 1st place.

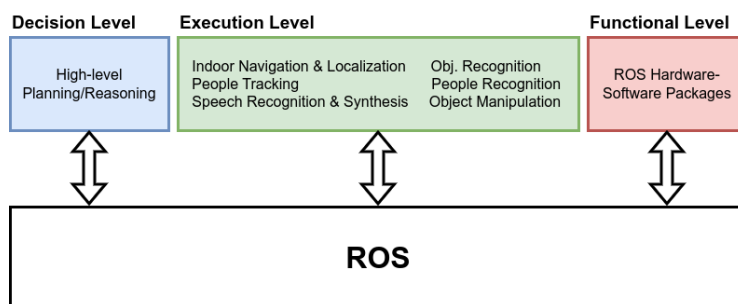
Currently, our group’s main research topics are computer vision, robust object manipulation, speech recognition, indoor navigation within dynamic environments, task planning, reinforcement learning, and human-robot interaction. A more detailed description of our current work is presented in Section 4.

## 3 Markovito’s Architecture

Markovito’s software architecture has been designed and developed as a layered behavior-based [1] architecture that uses ROS [2] for communication (see Fig. 1). The architecture has three different levels:

1. **Decision level:** This is the highest level in the architecture, at which planning and reasoning is performed by manipulating logical symbols that represent high-level concepts, such as relative-position relations (*e.g.* in front of) and causal effects of actions (*e.g.* picking-up and moving objects causes a change in their location). The symbolic representation of the environment and the planning is performed with a sorted answer set programming approach [3].

2. **Execution level:** Modules in this level interact with the functional level through ROS architecture. This level includes the modules to perform basic tasks such as navigation, object manipulation, speech recognition, etc.
3. **Functional level:** At this level the modules interact with the robot’s sensors and actuators, relaying commands to the motor or retrieving information from the sensors.



**Fig. 1.** Software architecture of Markovito: a layered behavior-based architecture that uses ROS for communication.

At the execution level, the robot’s skills are encapsulated into general purpose modules. By doing so, programming the robot to solve different tasks becomes less time-consuming and endows the architecture with scalability, which becomes more important as new modules are required.

## 4 Markovito’s skills and current work

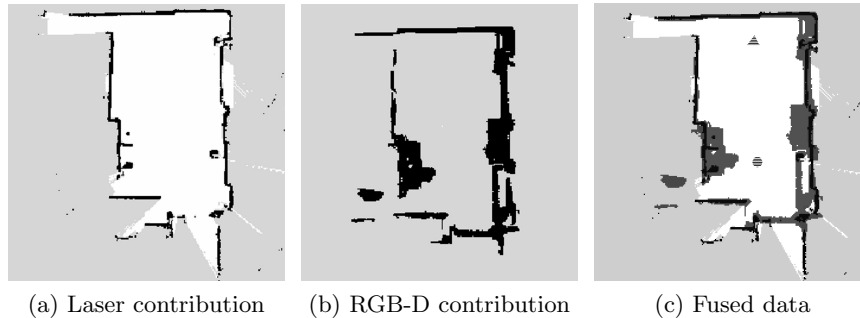
In this section we present the most recent work on Markovito’s skill set, which encompasses indoor navigation in dynamic environments, people tracking, speech recognition, object recognition, and robust object grasping and manipulation.

### 4.1 Indoor Navigation

The navigation algorithm that we are currently using (with some modifications) was designed by the Active Vision Group (AGAS) [4]. They present and explain the use of their software which they call Homer GUI, Homer mapping and Homer navigation corresponding to the graphical user interface, mapping algorithm and navigation algorithm, respectively. This software uses a particle filter to solve the SLAM problem and the A\* algorithm for path planning.

The navigation system developed by AGAS was modified by integrating a depth image sensor, which gathered point clouds of the robot’s surroundings

and projected them onto the two-dimensional plane (that corresponds to the floor) [5], as shown in Fig. 2. The projected points were combined with laser data so that the robot could identify obstacles at different heights above the floor, unlike the original system which could only detect obstacles at the level of the laser sensor (usually installed in the robot’s mobile base). Additionally, our navigation system is able to re-plan a path to its target location when static (*e.g., furniture*) and dynamic obstacles (*e.g., people, pets*) are hindering, which is important for the safety of the robot and people in a real scenario.



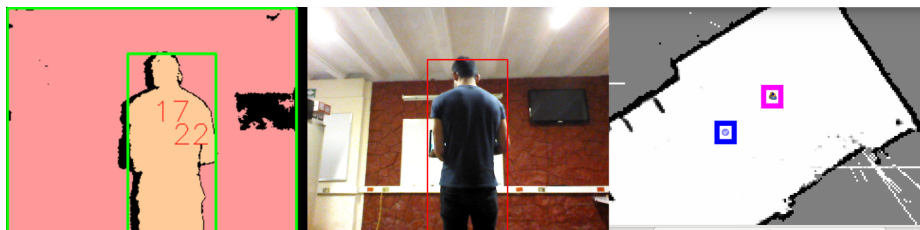
**Fig. 2.** Static occupancy grid map used for Small Scattered Obstacles scenario integrating laser and RGB-D (figure taken from [5]).

## 4.2 People tracking

Our people tracking system is based on visual information obtained from the robot’s RGB and depth image sensors. The system is able to keep track of multiple people within its range of view by performing at every instant of time a procedure that follows four main steps:

1. Segmentation on the depth image is performed to obtain a vector of blobs and their position with respect to the robot’s framework.
2. In order to track the blobs, those from the current instant are matched to those from the previous one, based on the Euclidean distance.
3. The Single Shot Detection (SSD) neural network [6] is employed to recognize people in the RGB image of the current instant, if a person is found then its bounding box is extracted.
4. The bounding boxes of people found in the previous step are matched against the depth image in order to filter out those blobs that are not people. In Fig. 3 is shown an example of the robot locating a person within its map.

In order to follow somebody, after the system projects the position of the person of interest in the map the robot uses to navigate, the robot knows where



**Fig. 3.** Example of the tracking system locating a person within its map. The left image shows the depth image segmented into two blobs: a person (orange) and the background (pink). Next to the depth image, in the RGB image the person’s bounding box is extracted using a pre-trained CNN. On the right of the image, the location of the robot (blue square) and the tracked person (fuchsia square) are shown within the robot’s map.

the person is within the map and can safely get closer without crashing into potential obstacles, for instance, tables and chairs. Additionally, because of the simplicity of our approach, it does not require of a GPU to work in real time.

### 4.3 Speech Recognition and Synthesis

We use the Amazon’s Alexa Service as our speech-to-text engine. The vocabulary of Markovito is described and implemented as an Alexa Skill [7], in which we specify what questions and commands Markovito can address. The vocabulary of Markovito is task-independent, which means that regardless of the task that is being solved, Markovito can address voice-queries at any time, *e.g.* answering a question.

Moreover, in case Markovito has no access to internet, it employs a speech-to-text model to process voice commands locally [8]. After the model transcripts the audio data into a sentence, parts of speech (POS) tagging is performed to retrieve the command/question and verify whether it is within Markovito’s vocabulary or not before it attempts to address it.

### 4.4 Object and people recognition

Among the sensors Markovito is equipped with, imagery (both RGB and depth) provide the most descriptive information about the robot’s immediate surroundings. However, in order to exploit such rich data and extract useful information from it, we employ multiple deep neural networks, such as convolutional neural networks (CNN). Currently, our robot identifies two main categories of entities in images:

1. **People:** For tasks that require Markovito to identify people’s location within the room, it uses the single shot detection (SSD) network [6], which returns

a list of boxes that enclose people in the RGB image. Moreover, our robot is also able to perform face recognition, based on the tiny face detector network [9] which is capable of identifying faces even with low resolution. Similar to SSD, the tiny face detector returns boxes that enclose the faces found in an image. In addition to locating where people and their face are with respect to Markovito’s RGBD sensor, it uses Google’s MediaPipe model<sup>2</sup> for pose estimation [10]. This enables the robot to understand non-verbal commands such as pointing at things and raising a hand to call the robot’s attention.

2. **Objects:** Unlike for people recognition, in which we use a CNN to crop areas of an image that contains a person or its face, for object recognition our system performs segmentation on depth images to locate where objects are with respect to the robot’s framework. Currently, an important constraint is that the objects must be over a flat surface. Then, it extracts the bounding boxes in the RGB image that correspond to the location of the objects in the depth image. Finally, the RGB snippets are classified with a pre-trained CNN (VGG [11]) that has been retrained to classify a set of known objects. Currently, one research line of the group is focused on the creation of an incremental learning scheme of objects in the scene (not necessarily on flat surfaces) so that a human instructor can progressively teach new objects to the robot. Once the robot incorporates the new objects learned, the idea is that it will be able to detect any appearance of these in the scene simultaneously. We are currently reviewing YOLO[12] and SSD[6] approaches to include some modifications to add new objects to the database to make the retraining as inexpensive as possible. We are working on an application with a graphical interface to speed up the process of capturing images of objects and also to re-train and manage existing models in a simple way.

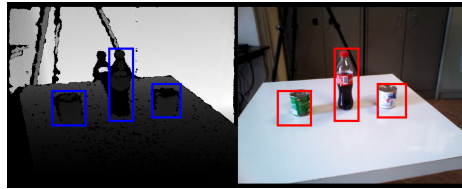


Fig. 4. Objects are located by performing segmentation on the depth image.

#### 4.5 Object manipulation and grasping

For tasks that require the robot to pick up and move objects, we employ *GraspIt!* [13] in tandem with *MoveIt* [14], which are two libraries developed for object

<sup>2</sup> <https://ai.googleblog.com/2020/12/mediapipe-holistic-simultaneous-face.html>

manipulation tasks. The former is for identifying the best grasp position for a particular object and the robot’s end effector, while the latter performs motion planning.

Although *GraspIt!* and *MoveIt* have shown to be quite effective and accurate, a robot remains prone to fail at grasping objects as a consequence of exogenous factors, *e.g.*, somebody moving the target object after the manipulator started approaching, or the object slips from the robot’s gripper. Moreover, since in many cases other sub-tasks will depend on the robot’s success in grasping an object (as part of a larger task), failing at picking up an object might have a high cost, therefore, feedback is required.

In order to endow the robot with a source of feedback, we developed a grasping confirmation mechanism that is low time-consuming and quite robust [15]. We specified a position for the robot’s gripper in which it could be observed by the robot without occlusion. Fine-tuning was performed to a pre-trained neural network by training the fully connected layers to identify two classes of images: i) empty gripper and ii) gripper holding something. Thus, after our robot attempts to grasp an object, it takes the gripper to the predefined position and verifies if it succeeded, if not, it can decide whether to try it again or move onto the next sub-task.

In addition to the confirmation method, the robot also uses visual feedback of its end-effector position when it approaches an object that it is trying to pick up. In a visual-servo style, Markovito stays aware of its surroundings to increase its odds of picking an object up, as unexpected events might occur as a consequence of its end effector getting closer to the object. For instance, push other objects in a way that end up hindering reaching the desired object, or someone changing the object’s position while the robot attempts to grasp it.

On the other hand, taking into account the human-robot interaction environments when the human gives some object to the robot to take it or vice versa, we have developed a program to automatically close and open the fingers of the manipulator arm. Our system uses the current sensors of the robotic arm, which allow it to measure the force exerted in each direction and decide to open or close taking into account threshold values.

## 5 Conclusions and future work

We have developed a set of general purpose modules, integrated in a layered behavior-based architecture that communicates its modules with ROS. These features enable Markovito to perform on tasks that are made of sub-tasks that can be solved with Markovito’s skills, such as the RoboCup@Home tasks. Based on this framework and different platforms, we have participated in the Mexican Robotic Tournament since 2007 achieving top positions each year. For instance, in 2015 and 2016 we won a 1st place in the Mexican RoboCup@home competition. We have also participated in the international Robocup@Home competition in 2009, 2011, 2012 and 2016. In order to suffice the needs of a real world domestic scenario, we believe that by integrating reactive and deliberative behaviors

becomes a must, so that a robot can act towards solving a task (planning) while it adapts to unexpected situations (reactive behavior).

## References

1. Ronald C Arkin, Ronald C Arkin, et al. *Behavior-based robotics*. MIT press, 1998.
2. Morgan Quigley, Ken Conley, Brian Gerkey, Josh Faust, Tully Foote, Jeremy Leibs, Rob Wheeler, and Andrew Y Ng. Ros: an open-source robot operating system. In *ICRA workshop on open source software*, volume 3, page 5. Kobe, Japan, 2009.
3. Evgenii Balai, Michael Gelfond, and Yuanlin Zhang. Towards answer set programming with sorts. In *Logic Programming and Nonmonotonic Reasoning: 12th International Conference, LPNMR 2013, Corunna, Spain, September 15-19, 2013. Proceedings 12*, pages 135–147. Springer, 2013.
4. Viktor Seib, Raphael Memmesheimer, and Dietrich Paulus. A ros-based system for an autonomous service robot. In *Robot Operating System (ROS)*, pages 215–252. Springer, 2016.
5. Orlando Lara-Guzmán, Sergio A Serrano, David Carrillo-López, and L Enrique Sucar. Rgb-d camera and 2d laser integration for robot navigation in dynamic environments. In *Mexican International Conference on Artificial Intelligence*, pages 661–674. Springer, 2019.
6. Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and Alexander C Berg. Ssd: Single shot multibox detector. In *European conference on computer vision*, pages 21–37. Springer, 2016.
7. Henry Lee. *Voice user interface projects: build voice-enabled applications using dialogflow for google home and Alexa skills kit for Amazon Echo*. Packt Publishing Ltd, 2018.
8. Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. Robust speech recognition via large-scale weak supervision. *arXiv preprint arXiv:2212.04356*, 2022.
9. Peiyun Hu and Deva Ramanan. Finding tiny faces. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 951–959, 2017.
10. Rohit Josyula and Sarah Ostadabbas. A review on human pose estimation. *arXiv preprint arXiv:2110.06877*, 2021.
11. Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
12. Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You only look once: Unified, real-time object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 779–788, 2016.
13. Andrew T Miller and Peter K Allen. Graspit! a versatile simulator for robotic grasping. 2004.
14. Sachin Chitta, Ioan Sucan, and Steve Cousins. Moveit![ros topics]. *IEEE Robotics & Automation Magazine*, 19(1):18–19, 2012.
15. Sebastián Salazar-Colores, Arquímides Méndez-Molina, David Carrillo-López, Esaú Escobar-Juárez, Eduardo F Morales, and L Enrique Sucar. A fast and robust deep learning approach for hand object grasping confirmation. In *Mexican International Conference on Artificial Intelligence*, pages 601–612. Springer, 2019.



## Markovito Software and External Devices

The base platform is a Robotnik RB-1 model, modified with a series of after-market devices to improve the robot's sensing and computing capabilities.

### Robot's Software Description

For our robot's system, we are using the following software:

- Platform: Ubuntu 14.0 (onboard PC), Ubuntu 16.0 (upper torso PC), Ubuntu 18.0 (lower torso PC).
- Speech recognition: Amazon Alexa Service (online) and OpenAI Whisper model (locally).
- Face recognition: Tiny Face Detector.
- People recognition: Single Shot Detection.
- People pose/gesture recognition: Google's MediaPipe pose estimator.
- Object recognition: custom software based on CNN classification over depth images segmentation.
- Arm control: Kinova driver, MoveIt.
- Manipulation: GraspIt.
- Grasp Confirmation: custom software based on CNN classification over depth images segmentation.
- Obstacle avoidance: custom modification of the AGAS navigation package Homer.

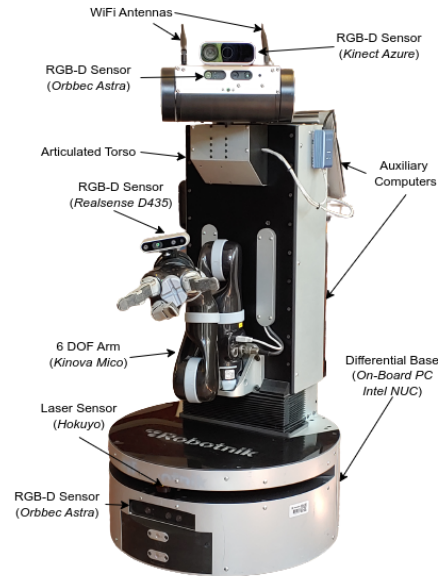


Fig. 5. Robot Markovito

### External Devices

In addition to the robot's original hardware, we have attached the following external devices to enhance its sensing and computing capabilities:

- Upper torso attached Intel Laptop.
- Lower torso attached MSI Laptop (with GPU capabilities).
- Kinect-Azure RGBD sensor<sup>3</sup> attached to the head.
- Intel realsense RGBD sensor<sup>4</sup> attached to the robot's wrist.

<sup>3</sup> <https://learn.microsoft.com/en-us/azure/Kinect-dk/hardware-specification>

<sup>4</sup> <https://www.intelrealsense.com/depth-camera-d435/>