

PQMEC@HOME Team Description Paper

Alexandre C. F. Filho, Ana L. C. M. Frayne, Eduardo D. Peixoto, Franco B. G. Junior,
Gabriel P. O. Ruotolo, José G. R. Teles, José R. R. Teles, Letícia L. Mendes,
Lucas R. S. Gris, Luis F. F. de O. Freitas, Márcio G. B. Junior, Nabila de P. e Silva,
Rodrigo M. de Carvalho, Victor G. Pimenta, Werisson E. S. Pereira

August 4, 2023

Abstract

This paper describes the service robot Miss Piggy of team Pequí Mecânico that will participate on the Robocup@Home competition which takes place annually in Brazil. This competition has influenced the development of research in natural language processing, computer vision, robotic manipulation, simultaneous localization and mapping.

performing several other tasks that a person might request. Our robot's name is Miss Piggy (see Figure 1).

1 Introduction

Pequí Mecânico robotic team exists since 2011 and took part in Latin American and Brazilian Robotics Competition in various categories: IEEE Standard Educational Kit (SEK), IEEE Open, RoboCup Small Size Soccer (F180), IEEE Humanoid Robot Racing (HRR), IEEE Very Small Size Soccer (VSSS) and RoboCup Soccer Simulation 2D. In 2019 we decided to compete for the first time in the Robocup@Home league.

Service robots are hardware and software systems that assist humans to perform daily tasks in complex environments. In order to achieve this, they have to be able to understand spoken or gesture commands from humans; to avoid static and dynamic obstacles while navigating in known and unknown environments; to recognize and manipulate objects and



Figure 1: Miss Piggy in the LARC 2022 in São Bernardo do Campo/SP - Brazil

2 Miss Piggy’s Robotic Architecture

2.1 Manipulator

The manipulator model that accompanies Miss Piggy is a ViperX 6DoF (Trossen Robotics), consisting of 8 Dynamixel servos (both XM540 and XM430) coupled in a way that 6 degrees of freedom were obtained, with a maximum payload of 750g. The automation approach varies between classic control applications and reinforcement learning algorithms, depending on the desired activity.

2.2 Human Interaction

The main way to interact with the robot is through voice commands, we use custom enhanced versions for embedded systems of Neural Networks, deployed into our Jetson Xavier via the *NVIDIA Riva SDK* [8] for the speech related capabilities (Automatic Speech Recognition and Speech Synthesis) that is able to provide state-of-the-art performance. To handle the NLU (Natural Language Understanding) core, the approach consists of a Large Language Model (a custom Llama 2) to handle its conversational operations, which consists of being able to help the operator based on the actual environment and also understand actions/commands to be executed on the context of a personal domestic robot.

Aside from the voice command is possible to input commands through it’s Graphical User Interface (GUI) on his main screen in a simple text format. The GUI is in a web-based format that allow user’s to access it remotely on any mobile device, it provides control and feedback over the actions of the robot through text and audio, the interface also show a friendly animated face that reacts to the user input.

3 Current Research

3.1 Human Robot Interaction

There are several types of information that are important for the robot to perform an action, visual,

sound, language. The classic way method to handle all these different types is combining individual processes of each input with some heuristic. Our current work in this area aims to replicate multi-modal perception similar to *Kosmos-2* [11]. To improve the precision of this network is necessary to acquired a larger dataset of the actual environment instead of using related ones, thus we are storing all the data received by the robot during our tests and this will improve his capability of understanding its surroundings over time.

3.2 Computer Vision

3.3 Human detection

At the moment our main vision sensor is the Intel realsense D435i camera. Although there is the OpenNI API for gestures recognition and human body movements for a sensor, we decide to use the OpenPose [2] algorithm for these tasks.

The OpenPose it’s a Deep Learning based human pose estimation, providing real-time multi-person detection. The model takes as input a color image and produces, as output, the 2D locations of keypoints for each person in the image. The authors also provide two pre-trained models, one trained on the Multi-Person Dataset (MPII) [1] and the other trained on the COCO dataset [4], producing 18 and 15 points, respectively.

In order to get a better frame rate in the Jetson Xavier GPU, we have adopted the Tensorflow implementation provided by [3]. The solution provides several variants that have some changes to the network structure for real-time processing on the CPU or lower-power embedded devices. We achieve good accuracy and real-time performance, approximately 8 FPS.

3.4 Human Tracking

With the OpenPose algorithm, we can detect multiple persons in the image, however, it does not provide means for people identification. In order to track the operator in the follow-me task, the time dived to use feature descriptor for person classification.

The OpenPose provides a set of body keypoints which can be used to create a local representation of the texture. This local representation is constructed using Local Binary Patterns (LBPs) [9]. With a neighborhood of size r surrounding the keypoints we compute the LBPs. Then we compute a histogram that tabulates the number of times each LBP pattern occurs. So, we treat this histogram as our feature vector. To handle multiresolution grayscale and rotation invariance, we use the extension to the original LBP proposed by [10].

In the follow-me task when the person stands in front of the robot we compute the LBPs histogram of the operator and register the vector in the memory. During the process we use a similarity metric to find the best match.

3.5 Navigation

In order to navigate through robocup@home environment, Miss Piggy uses *Navigation 2* [15]. Thus, the robot can move from wherever it is to a desired point avoiding obstacles in its map using three layers of `costmap2d ros2` package[12]:

- `obstacle_layer`: uses data from LIDAR as `laser_scan` and point cloud provided from `realsense`'s sensors as sources of observation;
- `static_layer`: map acquired from mapping stage (explained in next subsection);
- `inflation_layer`: propagating cost values out from occupied cells that decrease with distance as stated by environment.

To perform the best movement according to Miss Piggy's structure constraints (e.g. differential driver, acceleration limits, minimum velocity), we adopted a `smac` local planner. An implementation of this approach in Robot Operation System is available in [7]. We are using it to achieve optimization of global planner at during runtime and minimizing the trajectory execution time.

3.5.1 Mapping

Previous to the task's execution moment, the robot Miss Piggy will navigate the entire environment for the purpose of tracking objects and mapping the house. This way, using a Lidar as laser scanner and dead-reckoning sensors as odometry and IMU combined by Extended Kalman Filter, we can create a 2D map of occupancy performed by a *Simultaneous Localization and Mapping* approach.

By this day, we are searching for the best SLAM algorithm which attempts hardware limitations and sensors uncertainty. RTABMAP is accessible as a ROS package in [13]. It is a response that fit well into the faced problems and proved to be modularizable enough for testing and replicability in both virtual and real environments."

3.5.2 Localization

An implementation of Extended Kalman Filter (EKF) available in ROS is [5]. By using that, we can combine odometry of wheels, visual odometry and IMU data to provide more accurate relative position measurements. Collecting absolute position measurements like `laser_scan` data from Lidar, we merge information and pass it through a Localization based on map approach.

By now, ORB-SLAM is a good solution to acquire visual odometry data, and we are using a self-adapting ROS package inherited from RGB-D application [6]. Our solution changes some stacks of algorithm to obtain odometry messages and parameterize the code for `realsense` sensor.

There are a lot of solutions to the problem of localization with a map (e.g. EKF-Localization, Multi-Hypothesis Tracking, Grid Localization, Monte Carlo Localization explored by [16]). Nevertheless, each of them has great complexity measures in practice. Therefore, an Adaptive Monte Carlo Localization (AMCL) has been used as state-of-art localization problem. We are working with an AMCL available in [14] and parameterizing the code to hitch our sensors.

3.6 3D Modeling and Simulation

After the development of the first 3d modeling software, the sketchpad, in 1963, the creation of solid objects through the representation of a volumetric object became possible and increasingly used - in fields ranging from cinema to robotics.

The modeling in its most common aspect is performed by creating a mesh of segments that will give shape to the object, this is developed by several techniques, the most common being the polygon technique, the vertex technique and the edge technique.

In Robotics, the robot geometry description is based on the Unified Robot Description Format (URDF), which is a package with XML format to represent the robot model.

The new 3D modeling done for Miss Piggy was entirely based on the formulations of the Unified Robot Description Format (URDF), but completely implemented using the updates contained in the Foxy distribution of ROS2, aiming to enhance the transcription from the real robot to the virtual one and its performance in simulation.

4 Conclusion and future work

This is Pequi Mecânico's third attempt to compete in this category. Our robot was designed based on other teams' projects, information available at ROBOCUP @home wiki, and our own insights on the competition's challenges. The Deep Learning approaches we use will allow the robot to learn throughout the course while it receives commands and perform tasks. We already have several versions of Miss Piggy which perform different sets of tasks, and, until the competition, we will have many more.

5 Team Information

5.1 Robot Software Information

Operating System	Ubuntu 20.04
Middleware	ROS2 Foxy
Navigation	ROS2 <i>navigation stack</i>
Localization	<i>AMCL</i>
Mapping	<i>Gmapping</i>
Object Recognition	YOLOv8
Face Detection	YOLOv8
Human Detection	OpenPose
Gesture Recognition	OpenPose
Face Recognition	Eigenfaces
Speech Synthesis	FastSpeech & HiFi-GAN
Speech Recognition	Conformer-CTC
Natural Language Understanding	Llama 2

5.2 Robot Hardware Information

Base Motors	4x Hoverboard Wheel Motor 6,5"
Manipulators	Trossen Robotics Viper X 6 DOF Robotic Arm
Microphone	Rode Videomic Go
Display	Samsung Galaxy Tab S3
RGB-D Camera	Intel RealSense D435i
Speaker	JBL Charge 3
IMU	Intel RealSense D435i
LIDAR	RPLidar A2
Embedded Systems	NVIDIA Jetson Xavier and Intel NUC i7
Microcontroller	Teensy 4.0
Gimbal Setup	2x Dynamixel Servos MX28AT

References

- [1] Mykhaylo Andriluka et al. "2d human pose estimation: New benchmark and state of the art

- analysis". In: *Proceedings of the IEEE Conference on computer Vision and Pattern Recognition*. 2014, pp. 3686–3693.
- [2] Zhe Cao et al. "OpenPose: realtime multi-person 2D pose estimation using Part Affinity Fields". In: *arXiv preprint arXiv:1812.08008*. 2018.
- [3] Ildoo Kim. *Deep Pose Estimation implemented using Tensorflow with Custom Architectures for fast inference*. <https://github.com/ildoonet/tf-pose-estimation>.
- [4] Tsung-Yi Lin et al. "Microsoft coco: Common objects in context". In: *European conference on computer vision*. Springer. 2014, pp. 740–755.
- [5] Thomas Moore and Daniel Stouch. "A generalized extended kalman filter implementation for the robot operating system". In: *Intelligent autonomous systems 13*. Springer, 2016, pp. 335–348.
- [6] Raúl Mur-Artal and Juan D. Tardós. "ORB-SLAM2: an Open-Source SLAM System for Monocular, Stereo and RGB-D Cameras". In: *IEEE Transactions on Robotics* 33.5 (2017), pp. 1255–1262. DOI: 10.1109/TR0.2017.2705103.
- [7] *NAV2 smacplanner*. <https://navigation.ros.org/configuration/packages/configuring-smac-planner.html>. Accessed: 2023-06-30.
- [8] *NVIDIA Riva SDK*. <https://docs.nvidia.com/deeplearning/riva/user-guide/docs/index.html>. Accessed: 2023-06-30.
- [9] Timo Ojala, Matti Pietikainen, and David Harwood. "Performance evaluation of texture measures with classification based on Kullback discrimination of distributions". In: *Proceedings of 12th International Conference on Pattern Recognition*. Vol. 1. IEEE. 1994, pp. 582–585.
- [10] Timo Ojala, Matti Pietikainen, and Topi Mäenpää. "Multiresolution gray-scale and rotation invariant texture classification with local binary patterns". In: *IEEE Transactions on Pattern Analysis & Machine Intelligence* 7 (2002), pp. 971–987.
- [11] Zhiliang Peng et al. "Kosmos-2: Grounding Multimodal Large Language Models to the World". In: *arXiv preprint arXiv:2306.14824* (2023).
- [12] *ROS2 costmap_2d*. https://github.com/ros-planning/navigation2/tree/main/nav2_costmap_2d. Accessed: 2023-06-30.
- [13] *ROS2Wiki rtabmap*. https://github.com/introlab/rtabmap_ros. Accessed: 2023-06-30.
- [14] *ROSWiki AMCL*. <http://wiki.ros.org/amcl>. Accessed: 2019-06-17.
- [15] *The Marathon 2: A Navigation System*. <https://arxiv.org/abs/2003.00368/>. Accessed: 2023-06-30.
- [16] Sebastian Thrun, Wolfram Burgard, and Dieter Fox. *Probabilistic robotics*. MIT press, 2005.