



Master in Artificial Intelligence
Master of Science Thesis

HUMAN MULTI-ROBOT INTERACTION BASED ON GESTURE RECOGNITION

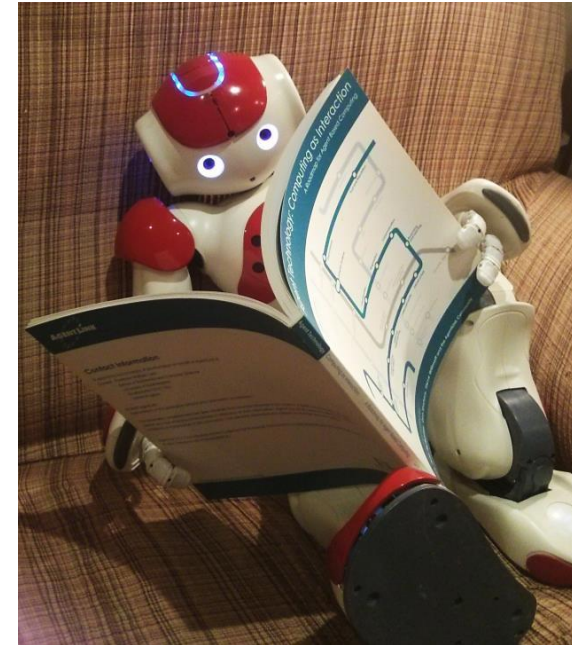
Gerard Canal Camprodon

Supervisors: Dr. Cecilio Angulo and Dr. Sergio Escalera



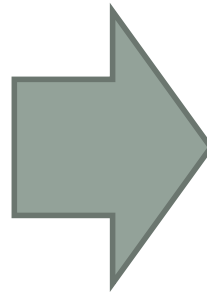
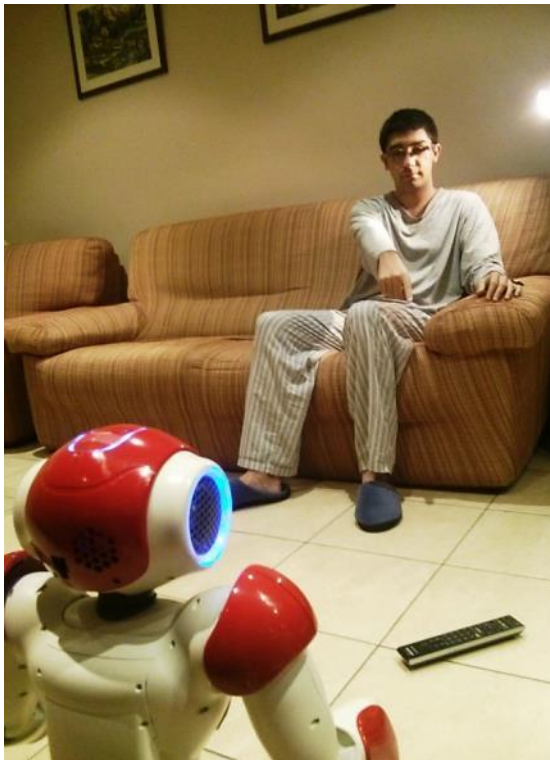
Introduction

- Nowadays robots are able to perform many useful tasks.
- Most of the human communication is non-verbal.
- HRI research on a gesture-based interaction system.



Motivation

- Elderly or handicapped person case.



Outline

- Goals
- Resources
- System overview
- Gesture Recognition
- Robot navigation
- HRI methods
- Results: Gesture recognition performance
- Results: User evaluation
- Conclusions
- Future work

Motivation

- Vision sensor too large to be carried by the robot.
- DARPA Grand Challenge idea of a driving humanoid.



Goals

- Design of a system *easy* to use and *intuitive*.
- *Real time*, therefore, *fast* response.
 - *Static* and *dynamic* gestures recognition.
 - *Accuracy* in pointing at the location.
 - *Multi-robot* tasks
 - Solving *ambiguous* situations.

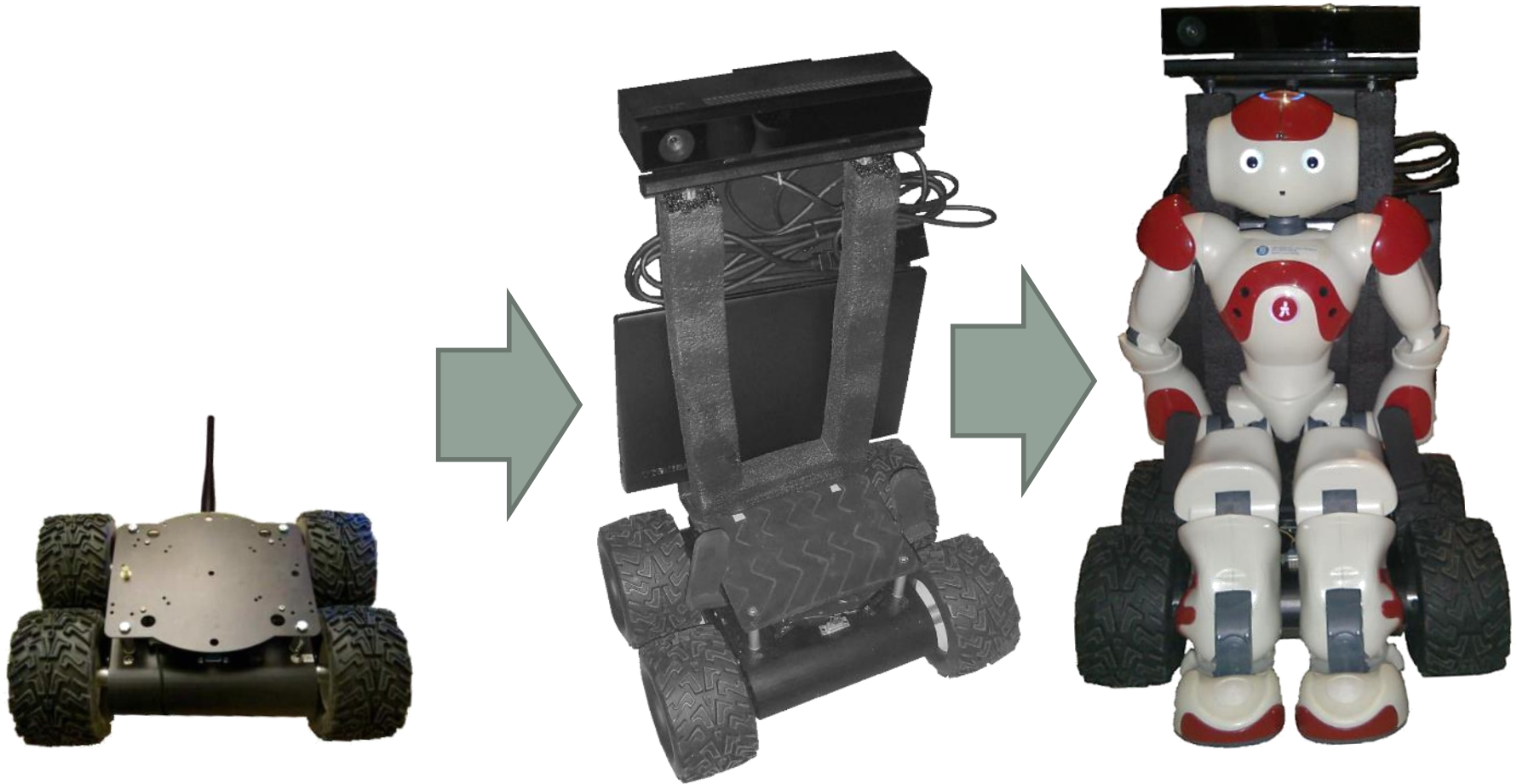
Hardware resources

- Microsoft Kinect version 2.
 - Windows 8.1 driver and USB 3.0.
- NAO.
 - CPU Geode.
 - NoaQi OS.
- Wifibot.
 - Intel Atom.
 - Ubuntu 12.04.



- Two laptops:
 - Intel i5
 - Intel Core 2 duo

Hardware resources modifications

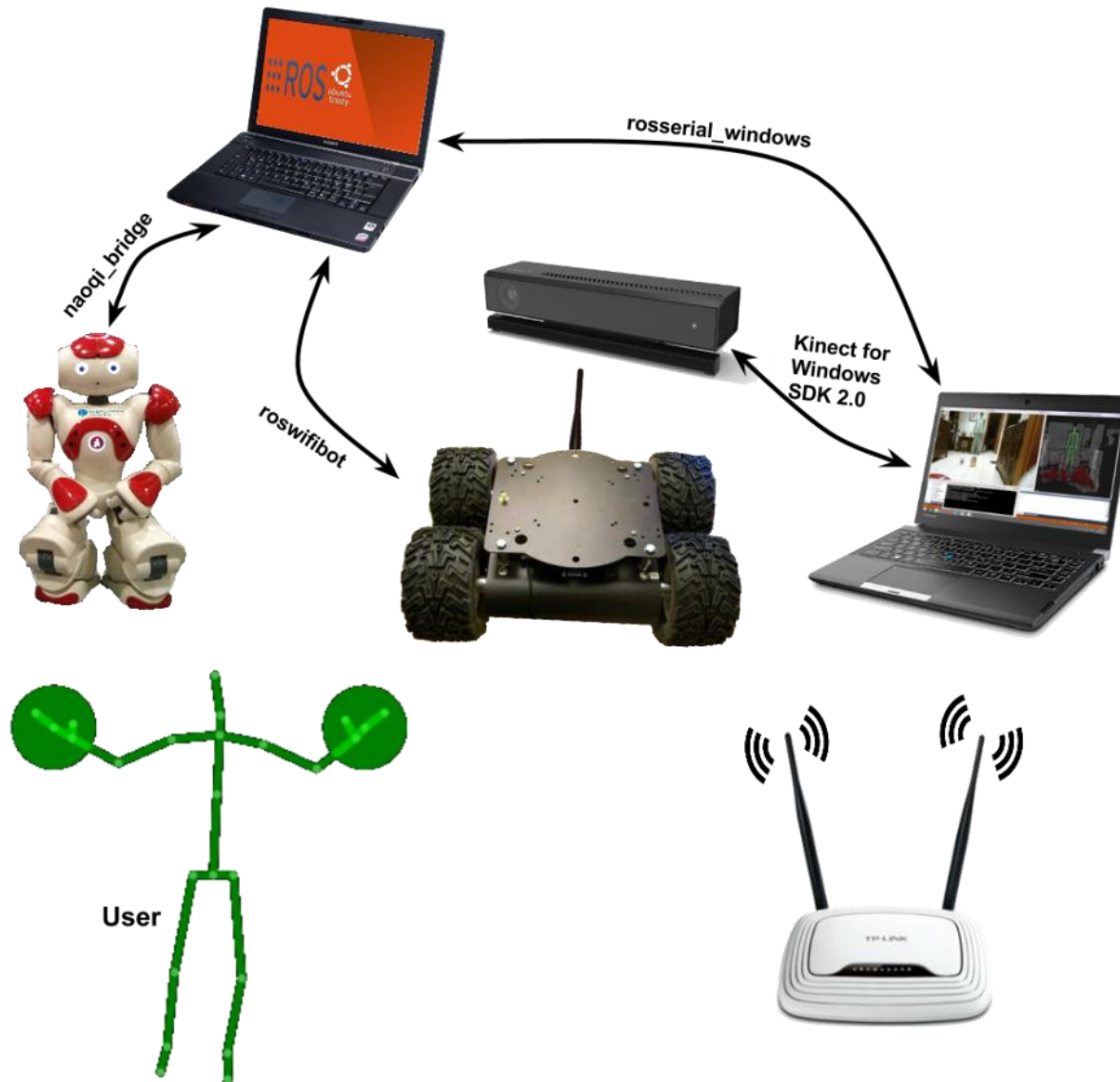


Software resources

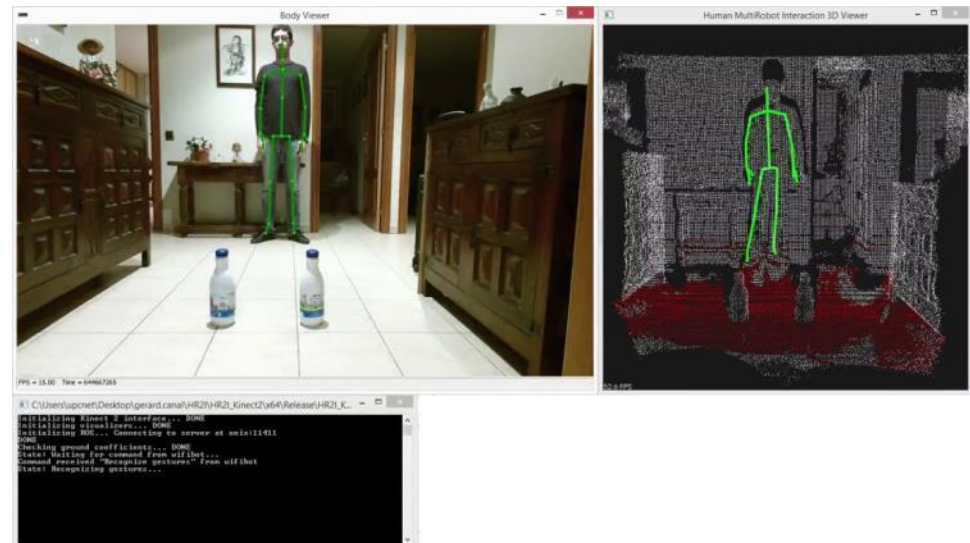
- ROS: Robot Operating System.
 - To program the robots.
 - SMACH to implement the Finite State Machines in Python.
 - Indigo Igloo version in Ubuntu 14.04.
- Kinect for Windows SDK 2.0.
 - C++ mode.
- PCL: Point Cloud Library.
 - Implemented in C++.



System overview



The diagram illustrates the robot's task sequence and gesture recognition. At the top, a blue box labeled 'Approach to pointing location' leads to a green box labeled 'Goal reached'. From 'Goal reached', a blue arrow points to a box labeled 'Segments objects', and an orange arrow points to a box labeled 'Approaches to object'. The 'Approaches to object' box leads to an orange box labeled 'Points object'. To the right, a stick figure shows two gestures: 'Wave gesture' (indicated by a green circle with a checkmark) and 'Point At gesture' (indicated by a green circle with a cross). A small inset image shows a robot in a kitchen environment.



Computer Vision: Gesture Recognition

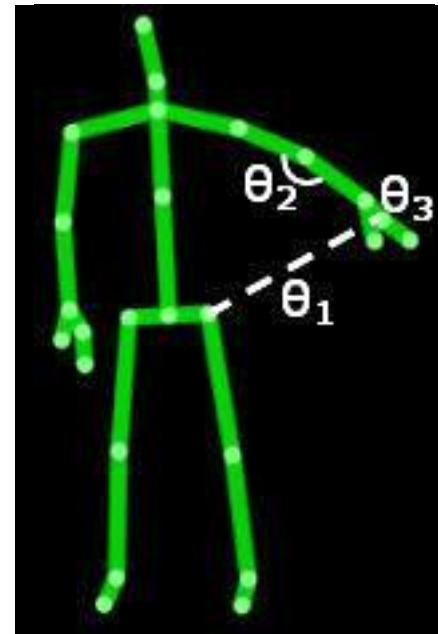
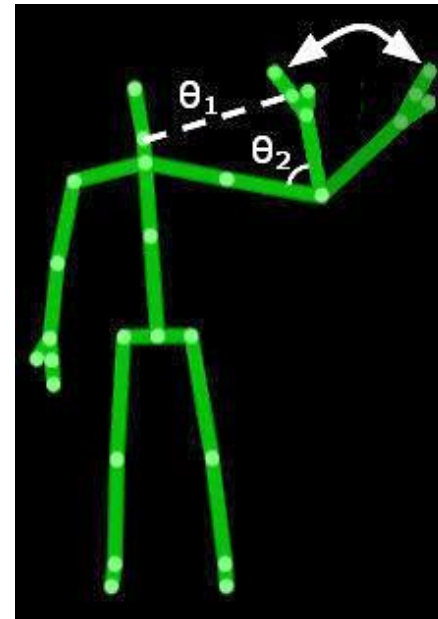
- Two types of gestures:
 - Static
 - Dynamic
- One gesture of each type:
 - Wave
 - Point at
- Described by means of skeletal features [1].



[1] J. Shotton, A. Fitzgibbon, M. Cook, T. Sharp, M. Finocchio, R. Moore, A. Kipman, and A. Blake. Real-time human pose recognition in parts from single depth images. In *Proceedings of the 2011 IEEE Conference on Computer Vision and Pattern Recognition, CVPR '11*, pages 1297– 1304, Washington, DC, USA, 2011. IEEE Computer Society.

Gesture recognition: Skeletal features

- Wave gesture:
 - θ_1 : Neck – Hand distance
 - θ_2 : Elbow angle
- Point at gesture:
 - θ_1 : Hand – Hip distance
 - θ_2 : Elbow angle
 - θ_3 : Hand 3D position



Gesture recognition: Dynamic Time Warping

- Used for sequence alignment.
- Applied for dynamic gesture recognition.
- An (infinite) input sequence is aligned with a gesture model.

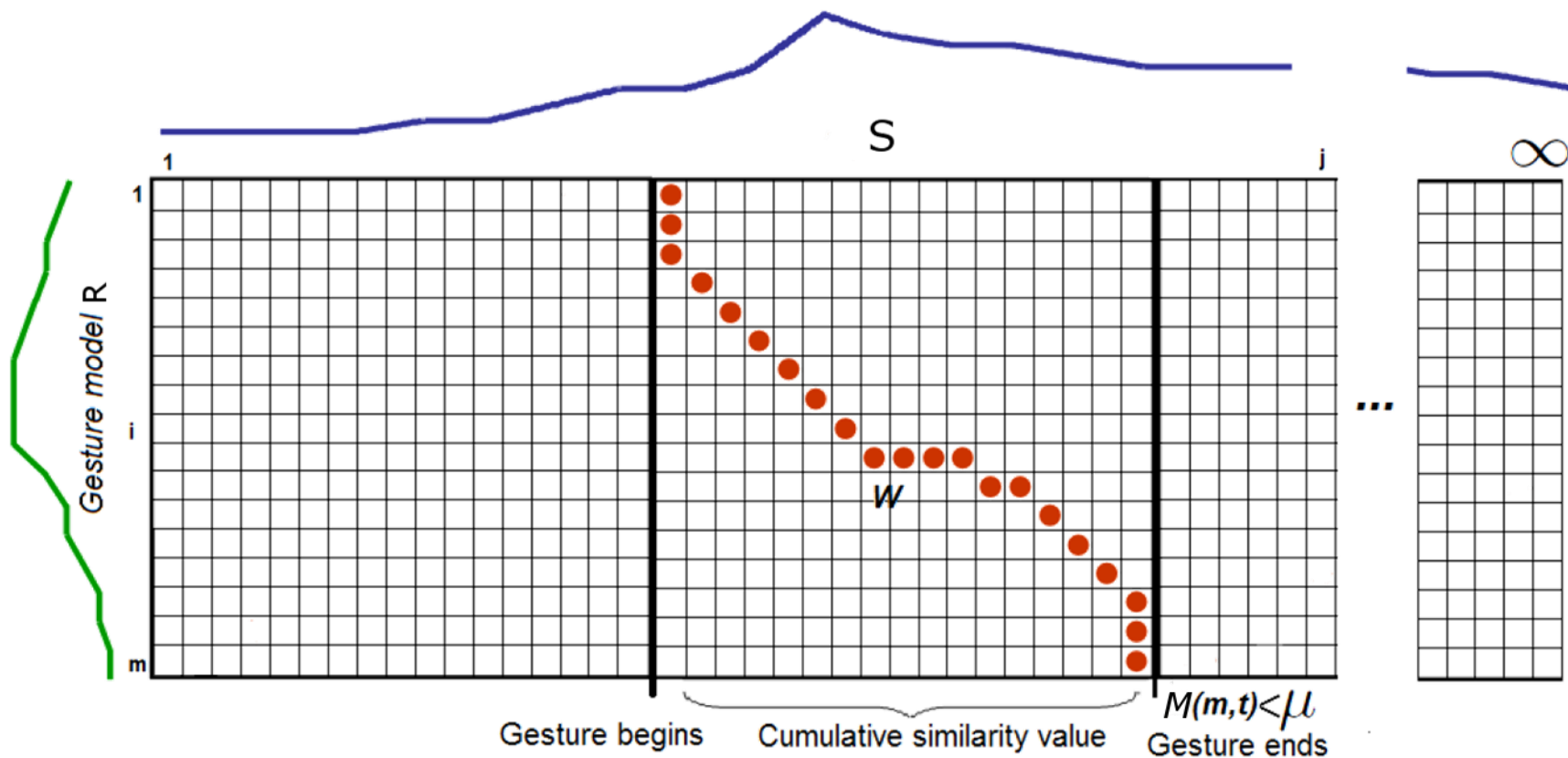
- Using a weighted L1 distance measure:

$$d_1(r, s) = \sum_{i=1}^k \alpha_i |r_i - s_i|$$

- A gesture is recognized when the input sequence is close enough to the model: $M_{m,k} < \mu, k \in [1, \dots, \infty]$.

Gesture recognition: Dynamic Time Warping

- $M_{i,j} = d_1(r_i, s_j) + \min\{M_{i-1,j}, M_{i-1,j-1}, M_{i,j-1}\}$



Static gesture recognition

- No dynamic time warping used.
- Checking features are within some thresholds.
- Checking involved limb is not moving.
- All during a certain number of frames.

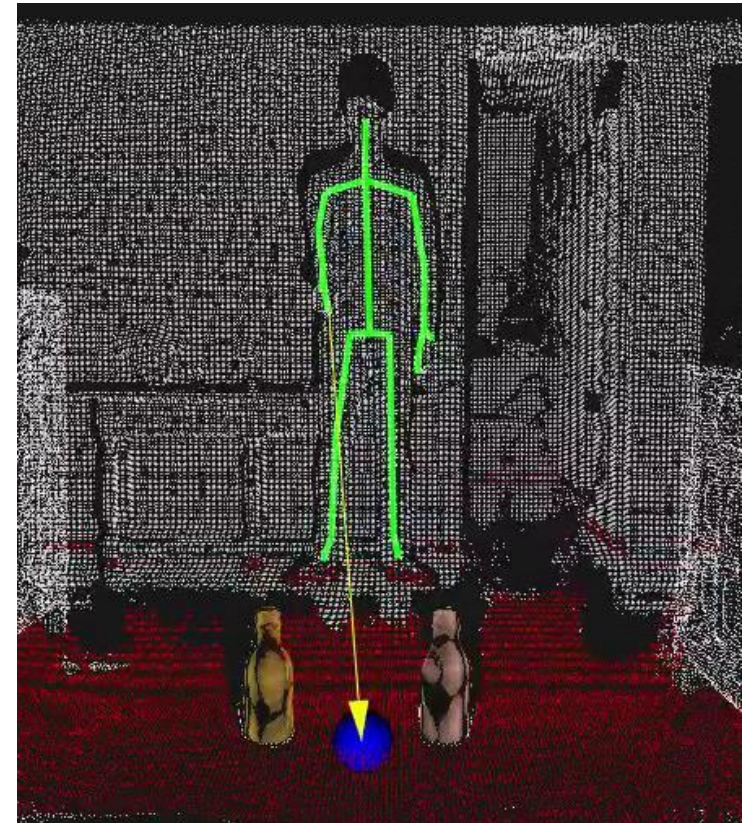
Static and Dynamic Gesture Recognition Algorithm (SDGRA)

- Dynamic and Static recognition in a joint algorithm.
- Multi-threaded to ensure real time.
- Possible multiple recognition in the same frame solved by keeping the one with less cost.

Gesture recognition:

Pointing gesture related methods

- Ground plane detection by RANSAC model fitting [2].
- Pointed point extraction using skeletal joints information.
- Object segmentation by Euclidean Cluster Extraction [3].

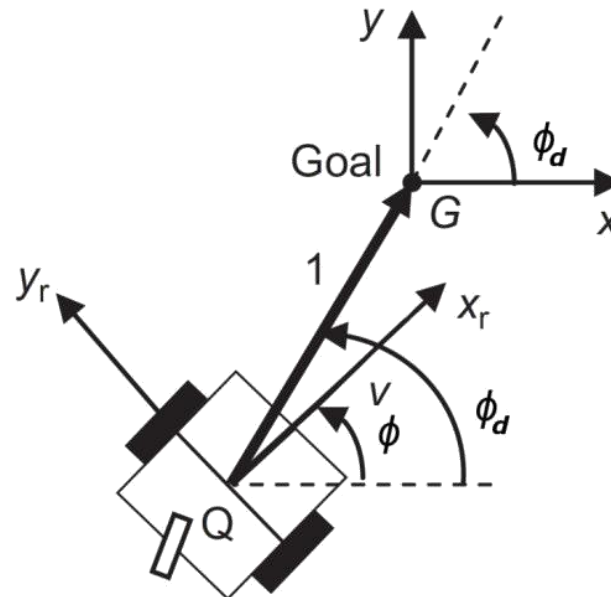


[2] M. A. Fischler and R. C. Bolles. Random sample consensus: A paradigm for model fitting with applications to image analysis and automated cartography. *Communications of the ACM*, 24(6):381–395, June 1981.

[3] R. B. Rusu. Clustering and segmentation. In *Semantic 3D Object Maps for Everyday Robot Manipulation*, volume 85 of *Springer Tracts in Advanced Robotics*, chapter 6, pages 75–85. Springer Berlin Heidelberg, 2013.

Mobile Robotics: Wifibot's navigation

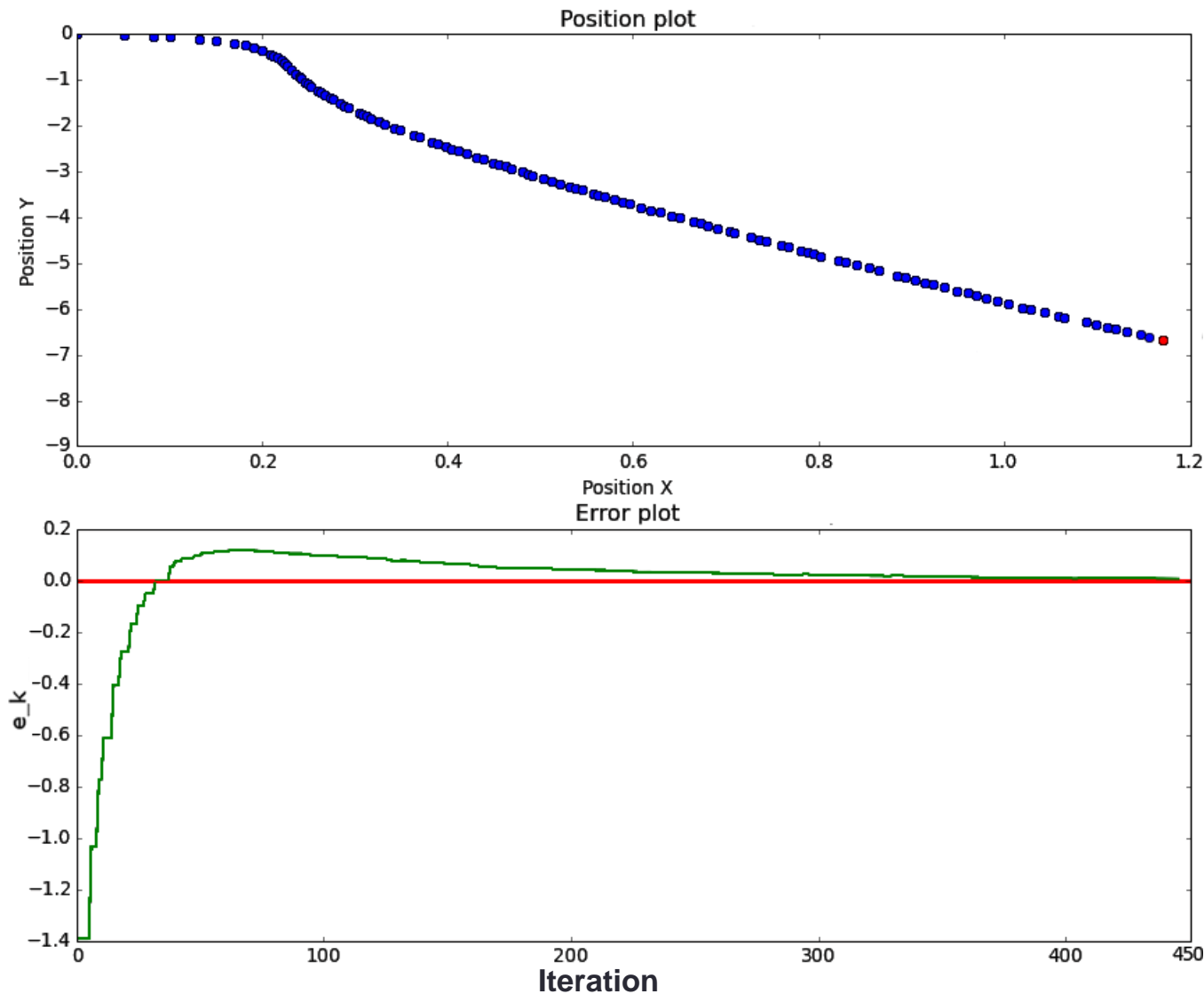
- Simple PID controller to move the robot towards a goal.
- No obstacles taken into account (free space assumption).
- Heading angle of the robot is the controlled variable.



Robot navigation: PID controller

- Differential drive model $\begin{cases} \dot{x} = v \cos \phi \\ \dot{y} = v \sin \phi \\ \dot{\phi} = \omega \end{cases}$
- Desired heading is $\phi_g = \arctan \frac{y_g - y}{x_g - x}$
- Heading error $e = \phi_g - \phi$.
- So the error is minimized $\omega = PID(e)$.

Robot navigation: PID example



HRI methods:

NAO going down the wifibot



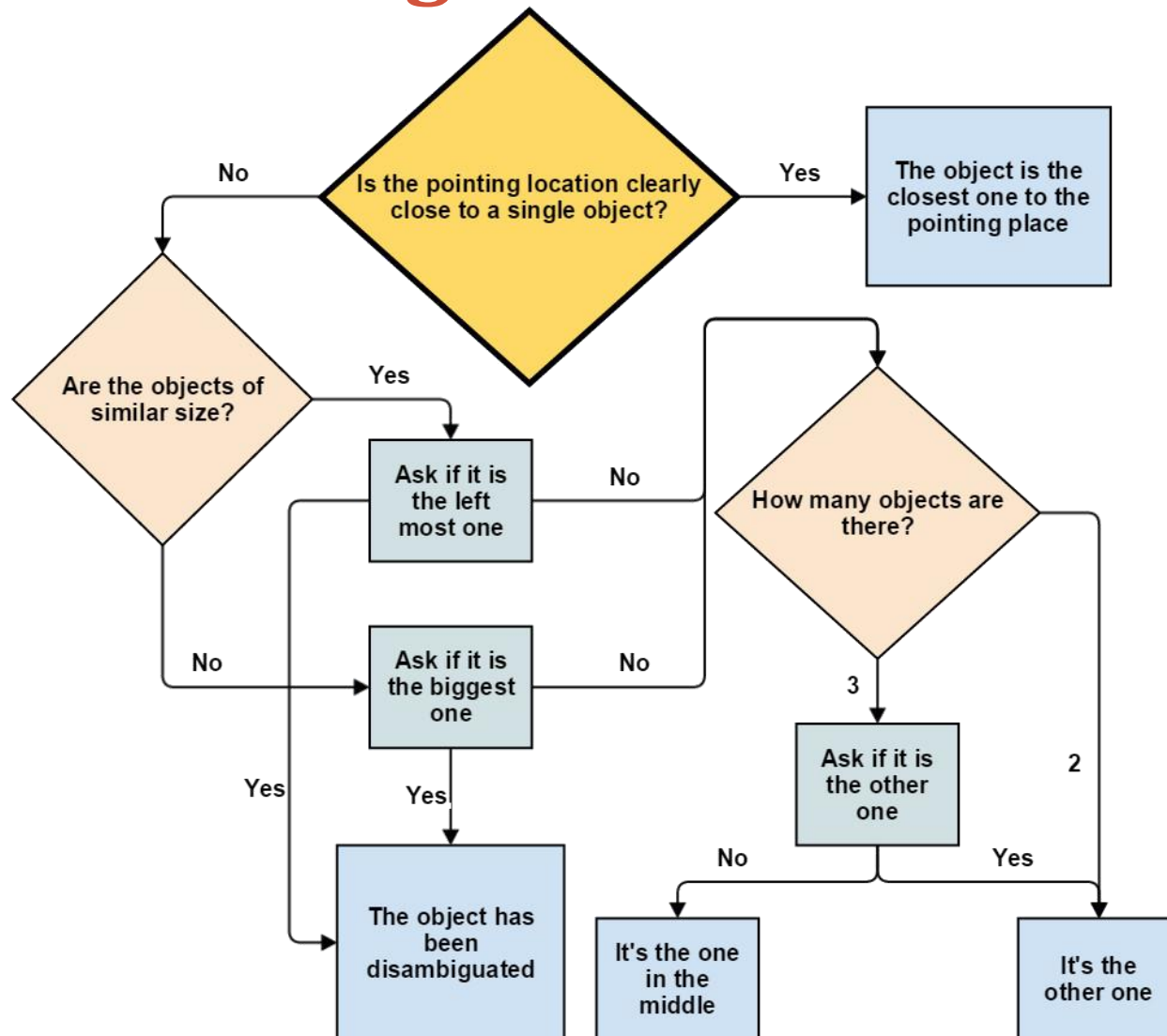
HRI methods:

Object disambiguation

- Extra information may be needed in case of doubt.
- Solve it by means of a small spoken dialogue.
- Use of simple questions about object's features like size and position.

HRI methods:

Object disambiguation



HRI methods: Interaction techniques

- The robot performs human-like gestures.
- Non-repetitive verbalization of its actions to enhance understanding.
- Eye color information to inform the user about its speech recognition state.

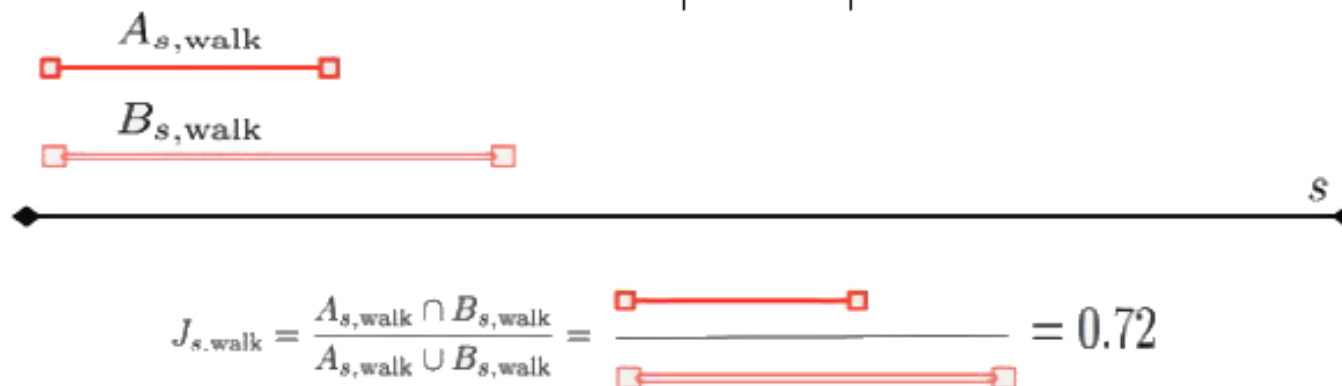


Results: Recognition performance.

Jaccard index

- Performance measured on a labeled set:
 - 7 sequences made by 3 different users
 - 61 gestures, 27 static and 34 dynamic
 - 2082 gesture frames
- Overlap / Jaccard index as performance metric:

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|}$$



Results: Recognition performance.

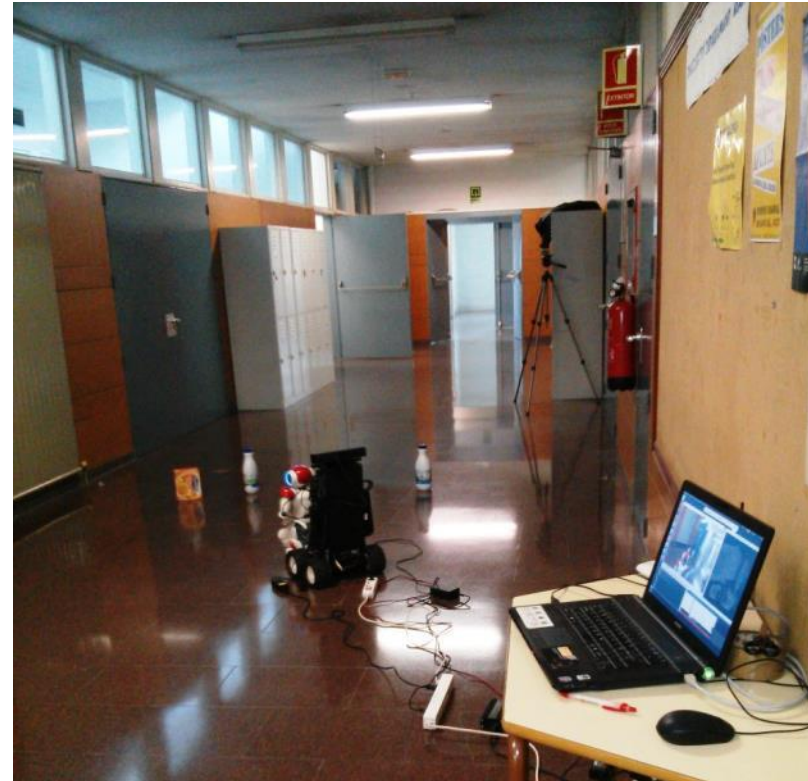
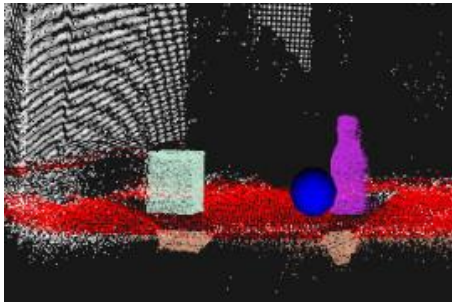
Jaccard index

- LOOCV test mean:
 - Static gestures: 0.463
 - Dynamic gestures: 0.492
 - Mean: 0.489



Results: User experience evaluation

- Testing environment.
- Implied some issues.

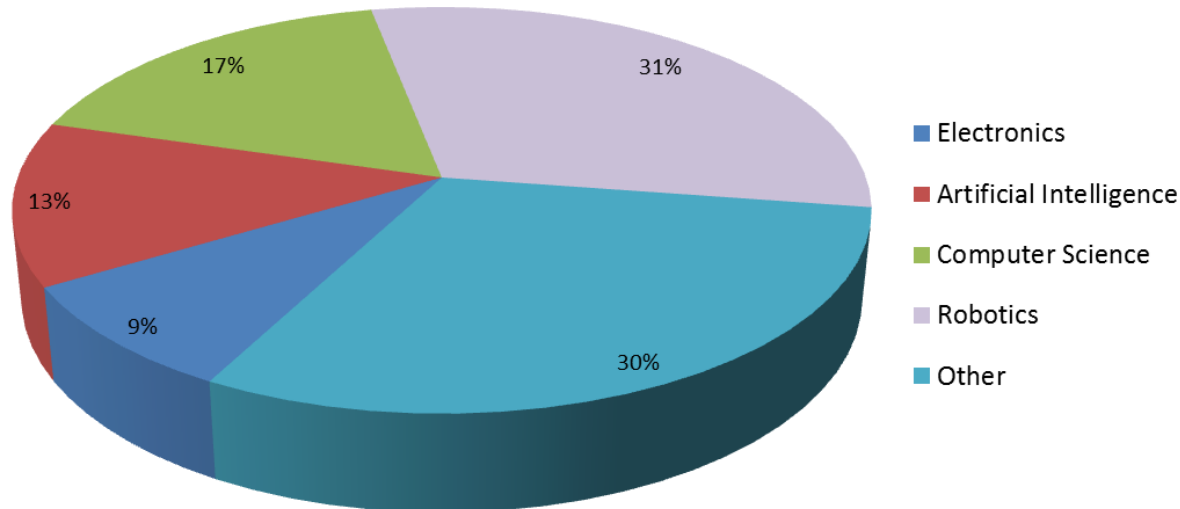


Results: User experience evaluation.

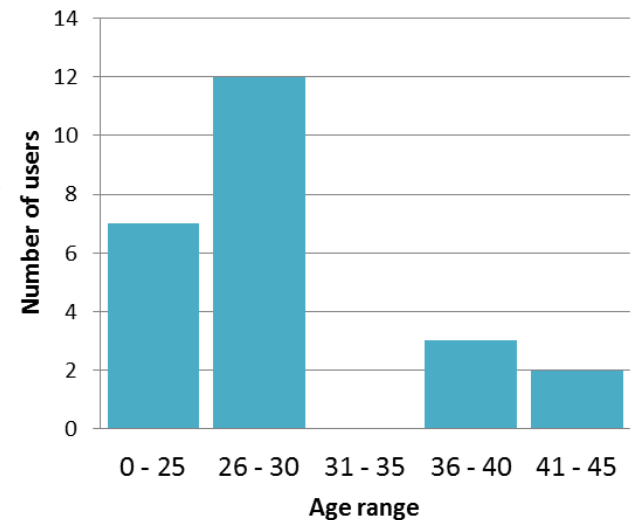
Users survey

- 24 users tested the system

User's background



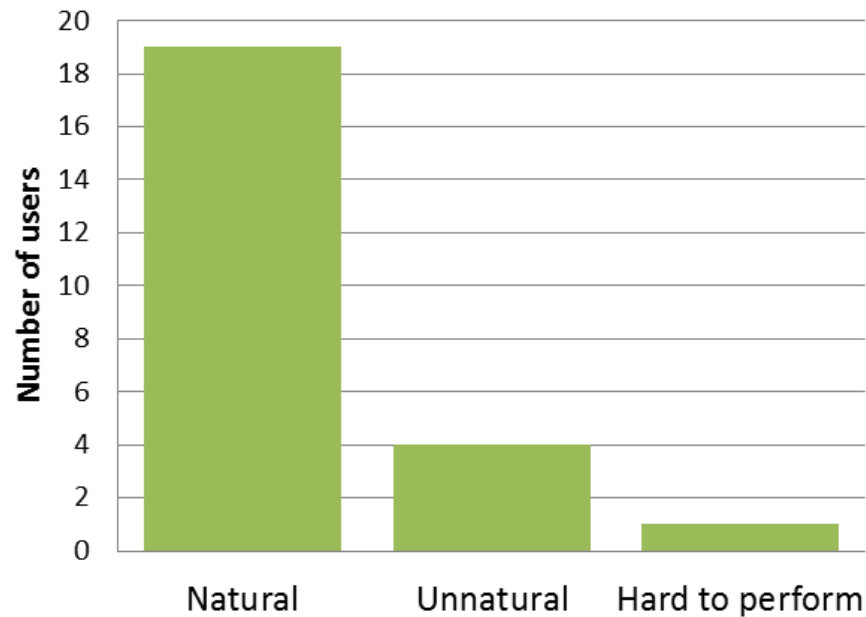
Age distribution



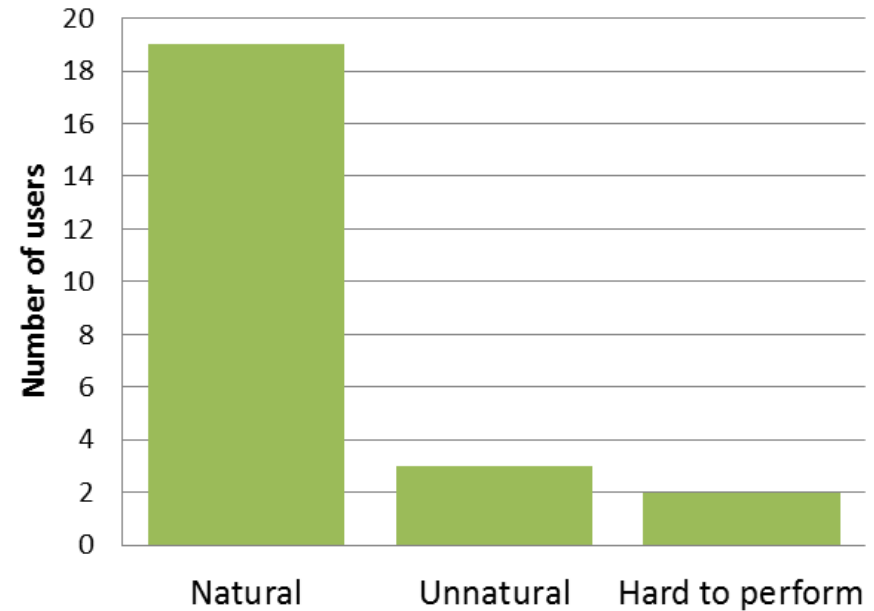
Results: User experience evaluation.

Users survey

Wave gesture naturalness



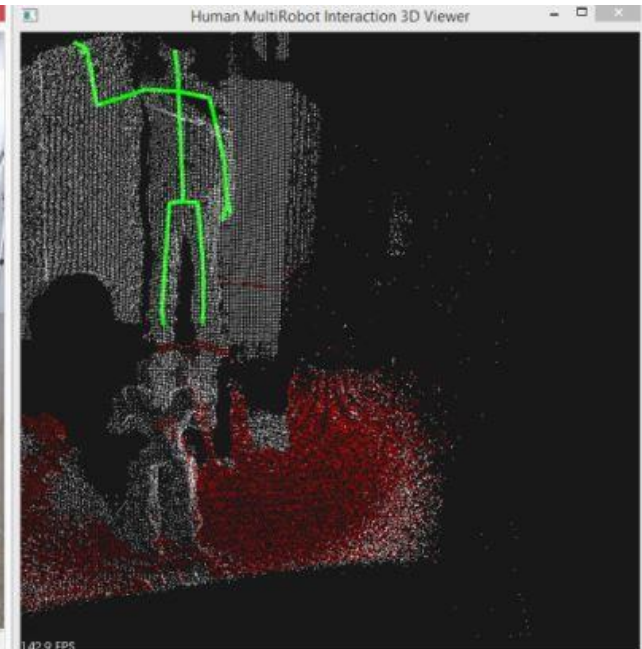
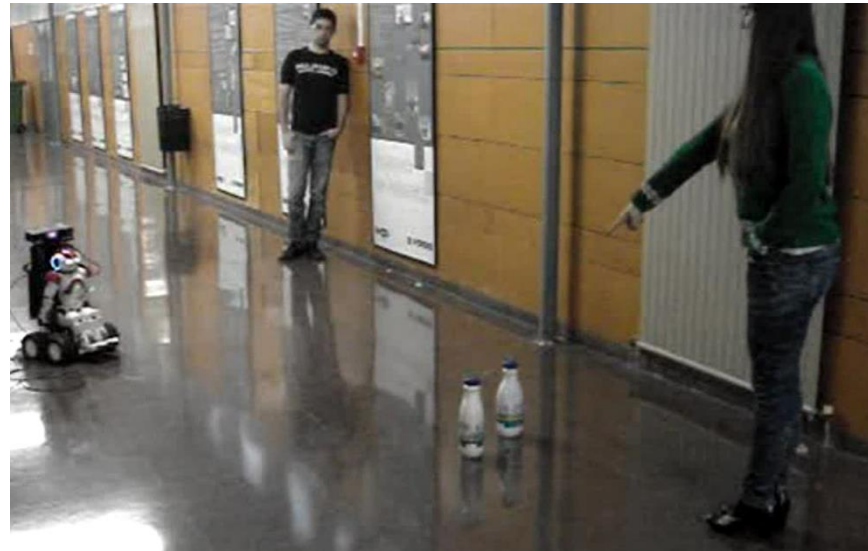
Point At gesture naturalness



Demonstration



User tests examples

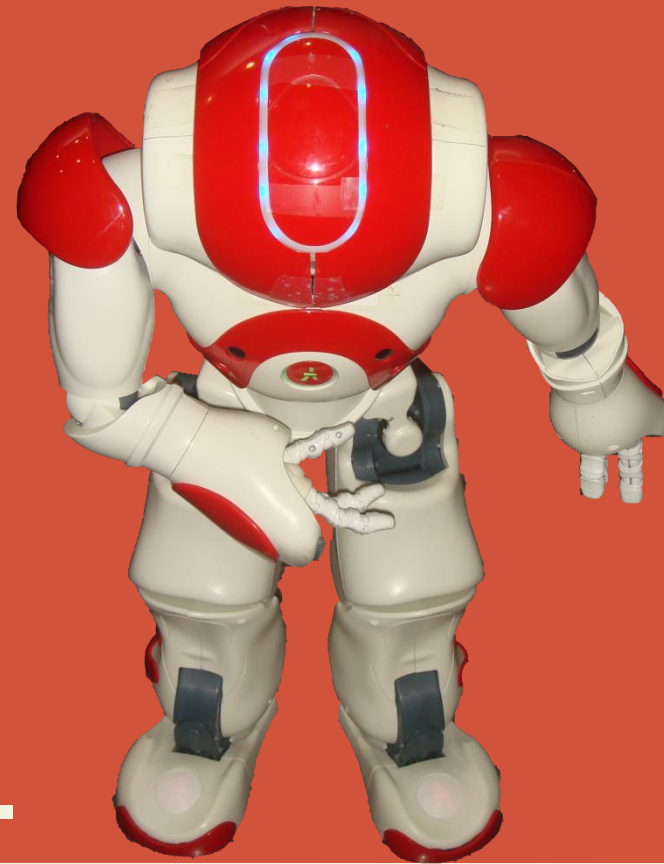


Conclusions

- Potential utility in household environments.
- Natural gestures as said by the test users.
- Easy to interact with the system and able to fulfill a task successfully in most of the cases.
- Working in real time, with correct response times.
- Generic and scalable framework.

Future improvements

- Addition of face gestures such as nodding and refusing.
- Face checking to improve tracking and feedback in case of failures.
- Bring the objects to the user.
- Object aware navigation.
- NAO tracking for walking correction.
- Cognitive interaction.
- More affective interaction.



THANK YOU.

***No robot was harmed in the making of this Master Thesis.*