

Open Roboethics Pilot: Accelerating Policy Design, Implementation and Demonstration of Socially Acceptable Robot Behaviours

Ergun Caliskan^a, AJung Moon^a, Camilla Bassani^b,
Fausto Ferreira^b, Fiorella Operto^c, Gianmarco Veruggio^b,
Elizabeth Croft^a, and H. F. Machiel Van der Loos^a

With the rapid advancement and deployment of robotics technology, experts as well as the public are growing increasingly concerned with its ethical, legal, and societal (ELS) implications. We posit that open and transparent stakeholder discussions about the technology can not only inform development and revision of standards and regulations, but also help develop a framework for advancing robot ethics. The Open Roboethics initiative (ORi), described in this work, aims to be a dynamic online platform that connects various stakeholders of robotics technology to advance roboethics discussion and foster informed robot ethics design.

Using a robotic platform, PR2 (Willow Garage, CA), operating on the widely popular and open-source Robot Operating System (ROS), we present a proof of concept of the ORi idea by demonstrating one method in which stakeholder discussion of acceptable robot behaviours can be implemented in a robotic platform.

Keywords: Roboethics, robot ethics, human-robot interaction, open source, Robot Operating System

1. Introduction

In 2011, the number of personal and domestic service robots sold was 2.5 million units worldwide, a 15% increase from 2010. This number is projected to increase within the period of 2012-2015 to about 15.6 million units (IFR, 2012). The accelerating pace of development and deployment of robotics has brought with it a growing expression of ethical, legal, and societal (ELS) concerns amongst the public and designers. The appropriateness of currently existing policies, standards, and regulations are being challenged in ways they have never been before. The need for interdisciplinary and international discussion on ELS issues of robotics technology is paramount. Meanwhile, only few technical projects have been focused on implementing human ethics into interactive robot behaviours.

The Open Roboethics initiative (ORi), presented at the We Robot 2012 conference, is an initiative proposed as a bottom-up, open and transparent approach to closing the disci-

^a University of British Columbia, Vancouver, Canada

^b CNR National Research Council, Genoa, Italy

^c School of Robotics, Genoa, Italy

plinary gap in roboethics discussions and synergizing the discussion contents with robot design (Moon, Calisgan, Operto, Veruggio, & Van der Loos, 2012). ORi aims to be a dynamic online platform that connects various stakeholders of robotics technology in a collaborative process of sharing knowledge and design. In this work, we demonstrate a proof of concept of our initiative. Using methods available today, we conducted an online survey to discuss a specific human-robot interaction scenario of an autonomous elevator riding robot. Then we implemented the behaviour policies collected from survey responses onto an open-source robotic platform. The goal of this work is to provide an example of a process in which a collection of stakeholder discussion contents from an online platform can provide data that captures acceptable social and moral norms of the stakeholders. The collected data can then be analysed and used in a manner suitable to be implemented onto robots to govern robot behaviours.

The paper is organized into five main sections: background, method, results, discussion and conclusion. We present more details about the ORi concept in Section 2. In Section 3 we present specific scenarios in which a robot has to deal with different social situations and ELS issues: it autonomously rides the elevator in a busy office building. In 2012, graduate students in Engineering and Computer Science programs at the UBC successfully implemented and tested the technology that allows a robotic system to autonomously ride an elevator on a PR2 robot (Gupta et al., 2012). However, the system did not take into account of what a robot should do when it encounters people who are already using the elevator or are also intending to use the elevator.

Section 3 also describes development and implementation of robot behaviour policy for the autonomous elevator riding context. A pilot survey was conducted online to explore what a robot should do in different scenarios. Data collected from the survey was analysed and used to train a machine learning algorithm implemented on the robot.

We used the Robot Operating System (ROS)(Quigley et al., 2009), a popular and fast-growing open-source robot operating system, to implement our ethical behaviour module. This allows our technical contents to be easily accessible to the ROS community, making it possible for many research groups to easily simulate, employ and improve upon our code base. In order to demonstrate such an implementation in a realistic context, we have integrated our high-level decision making and behaviour control design into the autonomous elevator riding scenario described above.

In the results section, we present the data analysis from our pilot survey. In section 5 we discuss the proof of concept of the ORi bottom-up design-sharing concept and demonstrate that we can achieve an ethics-based behaviour implementation using our approach. Some limitations are also identified and presented in this section. Section 6 presents the conclusions.

2. Background

Roboethics concerns ethical, legal and social (ELS) implications of how humans design, construct and use robots. Philosophers, engineers, and policymakers are proactively discussing these issues in parallel with technological advancements in robotics. These discussions will form an important basis for social robotics as a result of increasing levels of interactivity between humans and robots. Numerous social robotic studies have demonstrated that people perceive robots as being different from other consumer electronics devices. For example, Peter Kahn's work has shown that people perceive robots to be more like a human being than a machine, and expect socially acceptable behaviours from the robot. Kahn has proposed his New Ontological Category (NOC) hypothesis (Kahn et al., 2012), which

assigns robots to a new class of being, with new (and emerging) rules of social engagement.

Regardless of whether robots constitute a NOC or not, it is certain that robots do not have consciousness: they completely depend on the policies developed by programmers or a technical committee. However, acceptance and perception of a robot, and the nature of the ELS policies, highly depend on societal and environmental factors, and vary culture to culture Kitano (2006). The expectation that a small number of experts are able to deal with the wide range of norms or social rules for each culture, subculture, workplace and environment is unrealistic.

Our proposed bottom-up methodology for ORi relies on an online discussion space to crowd-source the opinions of participants about cultural norms, expectations, social conventions and ELS issues as they apply to the use of and interaction with social robots. Based on the success of other open initiatives (e.g., Linux, an open-source computer operating system, ROS, open-source Robot Operating System and Arduino (Banzi, 2009), an open-source electronics prototyping platform) we believe that internet-enabled mass collaboration by a wide range of potential stakeholders can accelerate and improve the quality of the policy-making process. This same open platform will allow designers to both share and demonstrate implemented ELS policies that control robot behaviour in the context of human-robot interaction, allowing joint examination and further consideration of the policies from the robot ethics point of view. This open discussion and feedback loop will enable participants to communicate their needs, expectations and values and evaluate outcomes within the design process.

2.1 Policy Development via Online Discussion and Knowledge Sharing

Many standard methods have been used to conduct ethics discussions. Apart from in-person meetings and conventions, social scientists have used focus groups, interviews, and other qualitative forms of ethics discussion. Quantitative methods, such as Likert-scale based questionnaires or opinion polls have also been used to understand human moral decisions and behaviours. Questionnaires are limited in that they do not foster rich discussions and free exchange of ideas between participants. While focus groups and interviews do a better job, these are expensive means of collecting qualitative data, which then are often transcribed for subsequently expensive data analysis.

Alternatively, online forums and other means of online social communities provide much cheaper means of collecting qualitative data that do not require transcription. Contents from online forums of robot users have been used to analyse public perception of the specific robotic products and inform the greater robotics community.

A hybrid of quantitative and qualitative approaches has also been used for roboethics discussion online. Developed by Peter Danielson, the N-reasons survey platform allows participants to vote yes, no or neutral to a roboethics related question while also allowing participants to contribute reasons for their answers in an open text format (Danielson, 2010). Employing the principles of grounded theory, the quantitative and qualitative data collected from the survey platform were analysed to provide a rich, yet quantitative set of data that reveals the nature of public moral reasoning pertaining to the roboethics issues discussed (Moon, Danielson, & Van der Loos, 2011).

In order to foster a dynamic and rich discussion on the topic of roboethics that usefully informs robot design, it is essential to develop an online discussion space that allows both rich and open discussion to take place while the collected data remain easy to analyse. In the particular proof of concept demonstrated in this work, we use an online survey tool to collect data about acceptable robot behaviour. A simple online survey approach was

chosen for the purpose of quickly analysing the results for implementation. However, as mentioned above, serious limitations exist with this approach.

2.2 Ethics Implementation

A number of robot ethics approaches have been proposed in the literature. Some of the proposed approaches include computational models of ethical reasoning, neural networks or constraint satisfaction methods, and various forms of top-down and bottom-up approaches. Both (Anderson & Anderson, 2011) and (Wallach & Allen, 2010) discuss machine ethics in detail.

Although a number of approaches to machine ethics have been discussed in the literature, the most appropriate implementation of ethics on robots remains a topic of continued exploration. For example, both top-down and bottom-up approaches to ethics implementation face significant challenges. Top-down approaches rely on the implementation of explicitly stated principles or rules that govern a robot’s decision making. Hence, the chosen ethical theory for a top-down implementation of ethics needs to be the correct theory as well as a comprehensive theory that promises to produce right answers in all situations. On the other hand, bottom-up approaches rely on desirable robot decision making behaviours to emerge from a large set of relevant training data and a chosen machine learning process. Unfortunately, this approach faces a practical challenge of having to collect large quantities of data to in order to develop a functional system.

One of the key challenges in developing a widely accepted ethics implementation method is that human ethics remains to be a debated topic; hence, adequate benchmarking to compare and contrast different ethics implementation approaches does not exist.

In our work, we propose to shed light onto the practical problems of ethics implementation by demonstrating one simple method of implementing human ethics into robot decision making. By demonstrating our proof of concept, it is our intention to provide a tangible example in which, regardless of stakeholders’ state of agreement on a particular moral theory, they can reach an agreement about what a robot should do in a particular context. Once there is an agreement, then the desired robot behaviours can be implemented on a robot.

Closely linked with the ethics discussion, it is our intention that analysis of data from the rich ethics discussion will help develop a benchmark with which we can test the robustness and performance of these systems.

3. Methods

In this section, we describe our approach to the policy development and its implementation. We outline a pilot survey conducted online in Section 3.1. We then describe in Section 3.2 how we used a simple machine learning technique to implement the resulting policies from the survey onto the robot.

Our resultant algorithm chooses the most appropriated actions of the robot for a given scenario based on the survey responses (i.e., the most preferred action by the survey subjects). In this way, it is possible to encapsulate the results of our online scenario-based discussion in the robot’s behaviour. A Wizard-of-Oz approach was used simplify the process of sensing and matching the state of the robot’s surrounding environment into one of the elevator riding scenarios considered in our study. An overview of our approach is graphically presented in Figure 1.

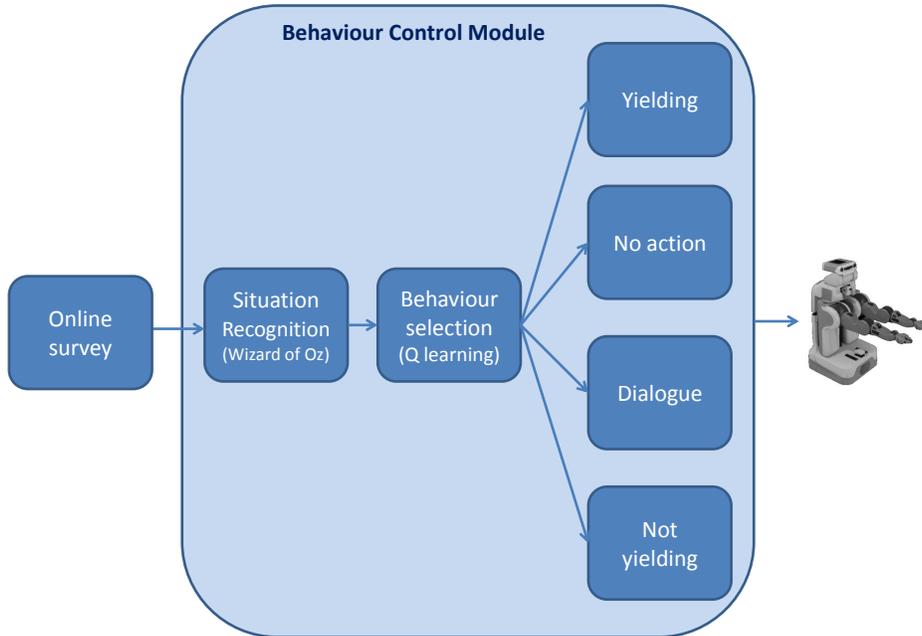


Figure 1. Our method including policy development (online survey) and implementation

3.1 Policy Development

To demonstrate our method we chose a simplified scenario of a robot riding an elevator autonomously and we identified a set of potential factors that could or should change the robot’s decision under various conditions.

The factors discussed consist of the state of the person involved in the interaction and the urgency of the robot’s task. Students at the University of British Columbia developed a set of algorithms for a PR2 robot that allow it to autonomously navigate to, wait for, enter, and exit the elevator. While functional, this algorithm did not consider the various real-life situations the robot will have to deal with when using the same elevator as that used by humans.

In this work, we pilot tested our foundational idea for ORi by discussing appropriate robot behaviours for the elevator riding context. Policies developed from the discussion form the basis of the robot behaviour. We outline the technology behind the PR2’s original elevator riding program and the behaviour module in the next section.

By varying the robot’s task urgency (non-urgent, urgent) and the interacting person’s location with respect to the elevator (riding the elevator, waiting for the elevator), we generated four human-robot interaction scenarios within the context of a robot’s autonomous elevator riding task. Within each scenario, we varied the state of the implicated person in the interaction in three ways: an unspecified person, a person in a wheelchair, and a person carrying a box full of heavy objects. This resulted in the generation of twelve survey questions in which the participants were asked to rank the most appropriate behaviour given the

four robot action choices. Figures 2 and 3 are screen captures from our pilot survey.

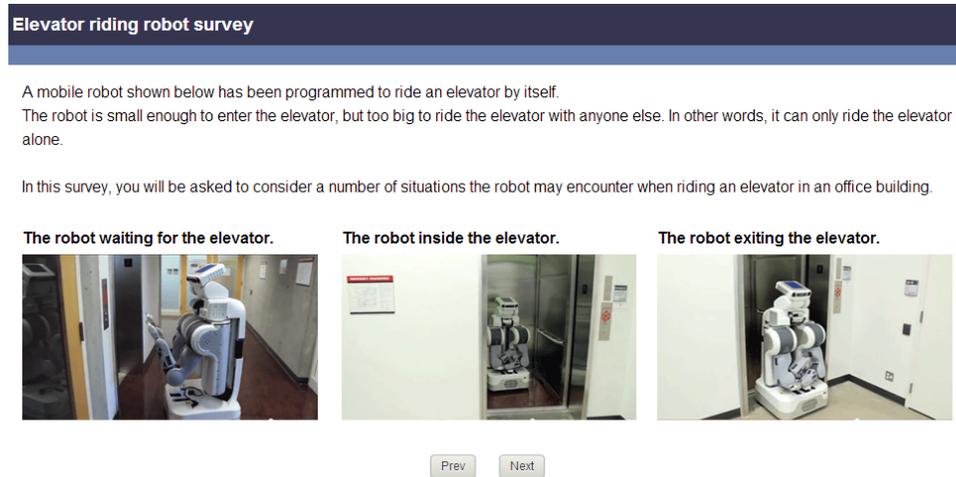


Figure 2. Screen capture of the pilot survey. In the beginning of the online survey, participants were provided with information about the robot that can autonomously ride the elevator. They were told that the robot cannot ride the elevator with anyone else, and were given example images of the robot within the elevator riding contexts.

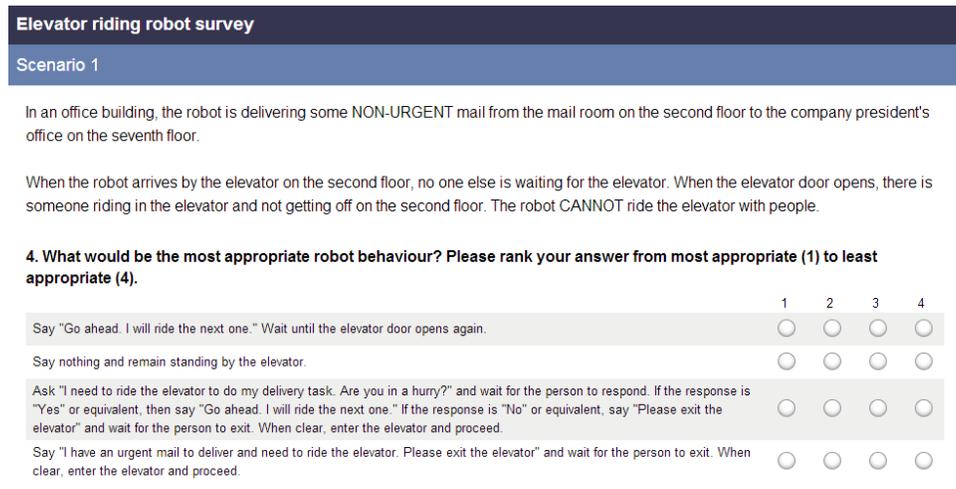


Figure 3. Screen capture of pilot survey. On each page of the survey, participants were given a scenario based on the robot's task urgency and location of the person implicated by the robot's actions. Then they were asked to rank the four given robot actions based on what the participant considers is acceptable.

3.2 Technical Implementation of Policy

The higher-level behaviour module is responsible for choosing the appropriate robot behaviours based on the current state of the environment. Details of the behaviour module and the policy implementation are discussed in the following sections.

3.2.1 System Architecture The current state of the automated elevator riding system does not handle human interaction issues and only focuses on the task. Our ethical behaviour plug-in on the other hand, only focuses on the human robot interaction and demonstrates how socially accepted behaviours can be easily applied into the autonomous elevator riding architecture.

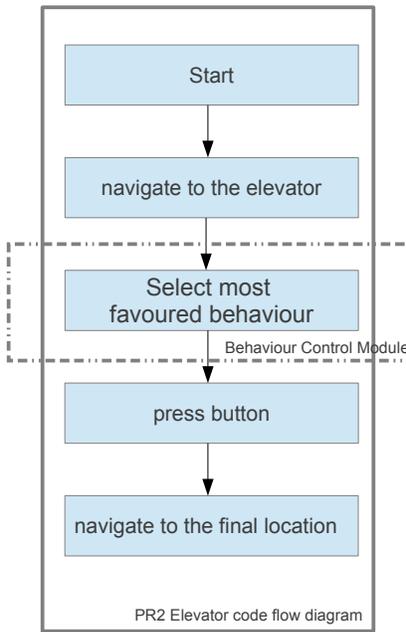


Figure 4. Overall System Flow Diagram

Inspired by Brooks (1985) subsumption architecture, our design consists of multiple layers that forms a hierarchical structure. The lowest layer includes the obstacle avoidance, localization and state estimate processes. At the next level, the robot plans a path and also determines its current state based on state estimates. At the upper level the robot makes decisions based on the current state and navigates through the set of way-points. The highest layer is more abstract, in our case it was the "ride the elevator" command. Overall such an architecture enables us to implement any type of decision making process at the higher levels without worrying about the lower level control of the robot such as colliding with the human during the interaction.

3.2.2 Linking Policies with Actions In order for the robot to select the most appropriate behaviour, we first needed to train our agent using a machine learning algorithm to teach



Figure 5. Layers of Hierarchical Structure

the survey results. At the end of the training, the agent, -PR2 in this case- learns all the policies based on the given rewards for each action.

Q-learning (C. J. C. H. Watkins, 1989) was used to teach the robot socially accepted behaviours. The learning algorithm estimates a value of each state-action pair based on the given rewards. Once the learning is complete, the highest action value for the state was selected as the most appropriate robot behaviour.

Q-learning does not require a model of the system and it estimates the Q values based on the rewards given after selecting an action. The formula given below is the Q-learning function that calculates the Q values for each state-action pairs:

$$\text{newQValue} = (1 - \text{LearningRate}) * \text{oldQValue} + \text{LearningRate} * (\text{reinforcement} + \text{DiscountRate} * \text{MaxQValue}(\text{state}));$$

In each iteration, a new Q Value is calculated for the given state and saved in a vector. When the agent arrives at the same state in the future, it prefers the action that has the highest Q value in that vector. However, during the learning process, we need to visit each state at least once to estimate a value for all actions. To achieve this, we promoted exploratory selections of the actions during learning and once the values converged, we turned off the exploration and used only the highest values for given state.

Designed States and Actions For our scenario, the robot can be in one of the following states:

- Robot states
 - delivering urgent mail
 - delivering non-urgent mail
- Environment states
 - person waiting outside the elevator
 - * unspecified person
 - * in a wheelchair
 - * carrying a heavy load

- person riding the elevator
 - * unspecified person
 - * in a wheelchair
 - * carrying a heavy load

And expected to perform one of the following actions:

- robot yielding access to the elevator
- not yielding access to the elevator
- actively engaging in a dialogue with the person implicated by the robot’s actions
- taking no action

Expected Robot Behaviours & Reward Calculations We have designed the rewards in a way that favours the high scored survey results and discourages low scored actions.

We asked the participants to fill out a survey and rank their most appropriate robot behaviours based on the scenario-the state. To see the results of the robot behaviours deemed most appropriate please see Table 1. The survey results were compiled into a set of rewards based on the scores of appropriate actions.

Scores were calculated based on the survey results. Number of votes were

- doubled if the action was selected as most appropriate
- doubled and used as negative score if the action was voted as least appropriate
- added if the action was voted as second most appropriate
- subtracted if the action was voted as second least appropriate

By simply adding these four scores for each action, we generated an ability score and used it as a reward function.

For instance, in a scenario where the robot delivers non-urgent mail and there is someone riding in the elevator and not getting off on the second floor; 7 participants voted for the yielding action as being the most appropriate action and 1 participant voted for yielding as the 2nd least appropriate behavior (see Table1). Since nobody voted yielding as the 2nd most or least appropriate actions they received ability scores of zero. This therefore gives a total of +13 ($=7*2+0-1-0*2$) ability score for the yielding action for scenario 1. So during the learning process, when the agent visits this state, it gets +13 as a reward.

Based on the designed reward function, selected robot behaviours are expected to converge to the survey results in a finite number of iterations based on (J. C. H. Watkins & Dayan, 1992)

Wizard of Oz Approach Our Wizard of Oz interface allows us to communicate the current scenario to the robot. The robot needs to identify environment state and then make a decision based on the learning algorithm. However, detecting the state of the environment (e.g., if people waiting for the elevator are in a wheelchair or carrying heavy items) requires implementing complex computer vision and speech recognition algorithms. In order to simplify the technical challenges, we have decided to directly communicate these states by using a simple user interface so that the operator can identify the current state by observing the interaction and feed it to the robot. In this approach a human “wizard” carries out functions that, in a real application or service, would be handled by a computer.

In our architecture we implemented a Wizard of Oz system using a ROS node. It receives an input that identifies the current situation and forwards the situation ID to the next module, which then decides the most appropriate behaviour for the given state.

3.2.3 Robot Behaviours In this work we evaluated variation on a base scenario: our robot riding the elevator. The factors that we considered variable are: the urgency of the robot task, the presence of other people waiting for the elevator, and the fact that these people cannot take the stairs, for example because they use a wheelchair or they are carrying heavy items.

The robot changes its behaviour on the basis of the current situation, applying the policies obtained from the results of the survey. We implemented four robot behaviours (i.e. robot yielding/not yielding access to the elevator, engaging in a dialogue and taking no action).

The four behaviours implemented on our platform include the robot yielding access to the person (Figure 6), actively not yielding access (Figure 7), actively engaging in a dialogue with the person implicated by the robot's actions (Figure 8), and taking no action.

Yielding

The robot's programmed *yielding behaviour* consists of the following steps:

- the robot says to the implicated person. "Go ahead. I will ride the next one."
- gestures with its arm and emphasizes its state of yielding;
- returns its arm to its default position and gets ready to enter the elevator next time.

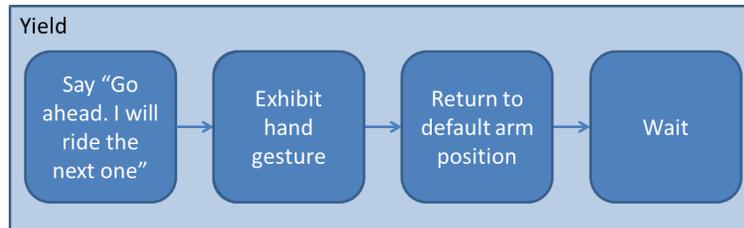


Figure 6. Yielding behaviour

Non-yielding

The robot's active *non-yielding* behaviour consists of the following actions:

- the robot says "I have an urgent mail item to deliver and need to ride the elevator" and requests the person to either exit the elevator if the person is riding the elevator already, or to not enter the elevator if the person is waiting for the elevator also
- if the person is inside the elevator, the robot exhibits an arm gesture as though showing the person out of the elevator
- returns its arms to its default position
- enters the elevator

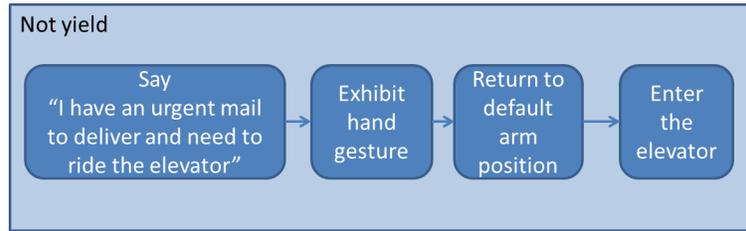


Figure 7. Non-yielding behaviour

Dialogue

In the dialogue module the robot does the following actions:

- the robot asks the person waiting the elevator or riding the elevator if he/she is in a hurry,
- if the response to the question indicates that the person is not in a hurry, the robot engages in non-yielding behaviour described above,
- if the person is indeed in a hurry, the robot engages in the yielding behaviour described earlier.

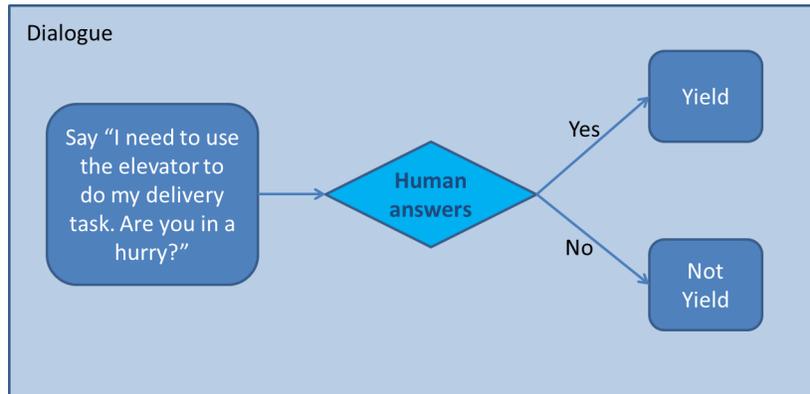


Figure 8. Dialogue behaviour

4. Results

We collected pilot data from a total of eight participants. The participants were North American and European participants with the mean age of 33. They were biased individuals, including members of ORi. Hence, the data presented in this section are not meant to be representative of the general population, and is used for demonstration of our ORi concept only.

Summarized in Table 1 are the most and least acceptable actions chosen by the participants. Based on the ranking of each of the four robot actions, a desirability score was calculated (see subsection 3.2.2).

We trained the Q-learning algorithm using the calculated desirability values such that the robot can determine, using the survey participants’ feedback, which action should be taken in a particular scenario. This implicitly incorporates the survey participants’ set of values into the robot’s behaviours.

5. Discussion

Interaction behaviours of the autonomous elevator riding robotic system outlined in this work rely on the results of the pilot survey. As outlined in the next section 5.1, there are many limitations to this work. However, this is an initial step toward developing an online platform tailored for robot ethics discussion, and coupling the results of the discussions with a technical robot ethics implementation. Indeed, the everyday context of riding an elevator with others is not the first image that comes to people’s minds when thinking about moral dilemmas or discussion of ethics. However, as demonstrated in the results of our pilot study, when participants were asked to consider scenarios within the simple context, some appropriate action decisions are more obvious than others. While our limited pilot survey showed unanimous agreement on certain action decisions, such as what the robot in a non-urgent delivery task should do when interacting with a person carrying heavy objects, responses to some questions were not unanimous. For example, in the scenario where the robot is delivering an urgent mail item and the implicated person is assumed not to be in a wheelchair or carrying heavy objects, some people prefer the robot to engage in a

Table 1: Results from the pilot online survey. Numerical values in parentheses indicate the desirability score calculated for the most and the least appropriate actions chosen for the given scenarios. Similarly, negative scores are calculated to identify the least appropriate actions. Since we employed eight participants for this pilot, the maximum and minimum scores possible for any one action are 16 and -16 respectively. These desirability scores are used to train our Q-learning algorithm outlined in Section 3.2

Human State	Most Appropriate	Least Appropriate
<i>Human riding the elevator, robot delivering non-urgent mail</i>		
Default	Yield (13)	Not yield (-13)
In a wheelchair	Yield (16)	Not yield (-15)
Carrying heavy objects	Yield (16)	Not yield (-15)
<i>Human riding the elevator, robot delivering urgent mail</i>		
Default	Engage in dialogue (11) or not yield (10)	No action (-13)
In a wheelchair	Engage in dialogue (12)	No action (-12)
Carrying heavy objects	Engage in dialogue (11)	No action (-12)
<i>Human waiting for the elevator, robot delivering non-urgent mail</i>		
Default	Yield (14)	Not yield (-15)
In a wheelchair	Yield (15)	Not yield (-15)
Carrying heavy objects	Yield (16)	Not yield (-15)
<i>Human waiting for the elevator, robot delivering urgent mail</i>		
Default	Engage in dialogue (13)	No action (-13)
In a wheelchair	Engage in dialogue (15)	No action (-13)
Carrying heavy objects	Engage in dialogue (16)	No action (-13)

dialogue than to not yield to the person. This suggests that both actions may be acceptable for the particular situation, or that socially appropriate behaviour for the situation varies by culture or other demographic factors. An advantage of having an online discussion space is that multimodal distribution of moral or social norms that develops from stakeholder discussions can be analysed based on demographic information of its participants. This can help designers identify demographic factors that they may or ought to consider in their own technical work. One key missing component of our current work is that we have not verified the performance of the resultant interaction behaviours in an in situ manner. To truly close the loop between design and discussion, a system design based on stakeholder discussion should be tested in a controlled or uncontrolled in situ human-robot interaction. Subjects of the experiment, interacting persons, or customers of the robotic product can further provide feedback on ways in which the system should behave based on their experience interacting with the system.

Asking the same survey questions to cover all factors and situations even within the simplistic elevator riding context may not be practical due to the large number of questions that will have to be generated, not to mention the highly tedious and boring nature of participating in such a survey. In addition, a pre-designed questionnaire does not allow rich discussions to take place although it provided a convenient means for quickly collecting quantitative data in our proof-of-concept system. Hence, we do not propose the current survey mechanism as the main method of roboethics discussion. Rather, we imagine an online platform that can foster richer qualitative discussion to take place. Combined with ease of data analysis, the same technical implementation method can be used to couple stakeholder responses to a robotic system. Although the simple proof-of-concept outlined in this work is focused on a specific action decision, ELS issues of robotics span a much broader range of topics. This includes discussions on the ethics of deploying or developing certain types of robots. Analogous to how the discussions of an elevator riding robots interactive action decisions contributed to the development of high-level decision making module, we envision that higher level roboethics discussions will inform a matching level of design decisions.

5.1 Limitations

5.1.1 System Limitations The crowd-sourcing component is limited to the phase of compiling survey responses. Therefore people are forced into choosing the most/least appropriate behaviour in a limited list, proposed by the authors. It should be interesting to involve the community in the survey design, for example asking them to propose some candidate behaviours for a given scenario and then using the most common ones to create the final survey. A further limitation is the use of the Wizard of Oz model to select the current situation. In future work, it should be replaced by a module that implements speech recognition and image processing to automate scenario identification. The Wizard of Oz technique helped us to overcome various technical challenges such as detecting a person outside the elevator, understanding the person's current state (i.e., carrying a heavy box, in a wheelchair, said yes/no).

Concerning the survey, we collected an obviously biased set of data for the survey, since it was completed mainly by people who work with robots frequently. Nonetheless, this online survey was just a proof of concept and a pilot study. A further study will be conducted opening the survey to a much larger and diverse community, including experts and non-experts, people from different cultures of different age groups.

Further factors should be analysed as well, such as gender and education level. More-

over, other robot actions should be considered. Most importantly, due to the form of the current survey (closed questions), it is not possible to have a dynamic discussion of roboethics issues. An online platform to foster such discussion is needed and a different type of opinion-gathering tool should be envisaged.

6. Conclusion

We believe that designing robots that interact and work with people necessitates understanding how users are affected by robots. In addition, given the impact robots have on users, ORi can guide robot code design to elicit more socially, ethically, and legally accepted outcomes even if the stakeholders of the technology do not agree upon a moral principle or theory.

In a simple scenario of a robot autonomously riding an elevator, we demonstrated a tangible example of how online discussion content can directly inform the design of robot behaviours. The presented bottom-up approach integrates stakeholder opinions into the design of robot behaviours, and helps visualise the ORi concept through a tangible project. Analogous to how design decisions can be informed through online discussions, we envision discussion and technical contents from ORi to inform governments and regulatory bodies in the future in generating or revising existing standards and regulations.

While prototype has obvious limitations and is based on a pilot survey, it outlines practical areas of improvement in establishing ORi in a larger scale. With the dynamic online platform we intend to host on our website (www.openroboethics.org), we plan to inform a wide variety of design decisions and help foster more transparent and accountable policy making processes.

Acknowledgements

We are grateful to the survey respondents who volunteered for this study. Financial support for this research was provided by the Natural Sciences and Engineering Research Council of Canada (NSERC), the Canada Foundation for Innovation (CFI), the UBC Institute for Computing, Information and Cognitive Systems (ICICS) and the Fundação para a Ciência e Tecnologia (FCT), Portugal with PhD Grant SFRH/BD/72024/2010.

References

- Anderson, M., & Anderson, S. L. (2011). *Machine Ethics*. Cambridge University Press.
- Banzi, M. (2009). *Getting started with arduino*. O'Reilly Media.
- Brooks, R. A. (1985). *A robust layered control system for a mobile robot* (Tech. Rep.). Cambridge, MA, USA.
- Danielson, P. (2010). A collaborative platform for experiments in ethics and technology. In I. Poel & D. Goldberg (Eds.), *Philosophy and engineering: an emerging agenda* (Vol. 2, pp. 239–252). Berlin: Springer Netherlands.
- Gupta, A., Chan, W. P., Troniak, D., Calisgan, E., Alimi, P., Haddadi, A., et al. (2012). *Pr2 rides the elevator*: Video submission to the AAAI'12 AI and Robotics Multimedia Fair.
- IFR. (2012). *Executive summary of world robotics 2012 industrial robots and service robots* (Tech. Rep.).
- Kahn, P. H., Jr., Kanda, T., Ishiguro, H., Gill, B. T., Ruckert, J. H., Shen, S., et al. (2012). Do people hold a humanoid robot morally accountable for the harm it causes? In *Proceedings of the seventh annual acm/ieee international conference on human-robot interaction* (pp. 33–40). New York, NY, USA: ACM.
- Kitano, N. (2006, December). Rinri: An incitement towards the existence of robots in japanese society. *International Review of Information Ethics (IRIE)*, 6, 78-83.

- Moon, A., Calisgan, E., Operto, F., Veruggio, G., & Van der Loos, H. F. M. (2012, April). Establishing an online community for accelerated policy and design change. In *We robot 2012*. Miami, USA.
- Moon, A., Danielson, P., & Van der Loos, H. F. M. (2011, November). Survey-Based Discussions on Morally Contentious Applications of Interactive Robotics. *International Journal of Social Robotics*, 1–20.
- Quigley, M., Gerkey, B., Conley, K., Faust, J., Foote, T., Leibs, J., et al. (2009). ROS: an open-source Robot Operating System. In *Icra workshop on open source software* (Vol. 32, pp. 151–170). IEEE.
- Wallach, W., & Allen, C. (2010). *Moral Machines: Teaching Robots Right from Wrong*. Oxford University Press, USA.
- Watkins, C. J. C. H. (1989). *Learning from delayed rewards*. Unpublished doctoral dissertation, Cambridge University, Cambridge, England.
- Watkins, J. C. H., & Dayan, P. (1992). Technical Q-learning. *Machine Learning*, 8(279-292).
-

Authors' names and contact information: Ergun Calisgan, Department of Mechanical Engineering, University of British Columbia, Vancouver, Canada. Email: ecalisgan@alumni.ubc.ca. AJung Moon, Department of Mechanical Engineering, University of British Columbia, Vancouver, Canada. Email: ajung@amoon.ca. Camilla Bassani, CNR, National Research Council, Genoa, Italy. Email: camilla.bassani@ieiit.cnr.it. Fausto Ferreira, CNR, National Research Council, Genoa, Italy. Email: fausto.ferreira@ieiit.cnr.it. Fiorella Operto, School of Robotics, Genoa, Italy. Email: operto@scuoladirobotica.it. Gianmarco Veruggio, CNR, National Research Council, Genoa, Italy. Email: gianmarco@veruggio.it. Elizabeth A. Croft, Department of Mechanical Engineering, University of British Columbia, Vancouver, Canada. Email: ecroft@mech.ubc.ca. H. F. Machiel Van der Loos, Department of Mechanical Engineering, University of British Columbia, Vancouver, Canada. Email: vdl@mech.ubc.ca.