

Benchmarks for evaluating socially assistive robotics

David Feil-Seifer, Kristine Skinner and Maja J. Matarić
University of Southern California

Socially assistive robotics (SAR) is a growing area of research. Evaluating SAR systems presents novel challenges. Using a robot for a socially assistive task can have various benefits and ethical implications. Many questions are important to understanding whether a robot is effective for a given application domain. This paper describes several benchmarks for evaluating SAR systems. There exist numerous methods for evaluating the many factors involved in a robot's design. Benchmarks from psychology, anthropology, medicine, and human-robot interaction are proposed as measures of success in evaluating a given SAR system and its impact on the user and broader population.

Keywords: assistive technology, human-robot interaction, psychological benchmarks, socially assistive robotics

Benchmarks for evaluating socially assistive robotics

Socially assistive robotics (SAR) is a growing area of research with potential benefits for elder care, education, rehabilitation, and people with social and cognitive disorders, among others. SAR is the intersection of assistive robotics and socially interactive robotics (Fong, Nourbakhsh, & Dautenhahn, 2003); it arose from the large and growing body of applications suitable for robot assistance that involves social rather than physical interaction (Lord & McGee, 2001; Taub et al., 2006; Wolf et al., 2005). While identified rather recently as a specific field, SAR is already an active community (Feil-Seifer & Matarić, 2005) providing companionship, monitoring care, encouragement in rehabilitation exercises, and social mediation. However, as such systems are increasingly being developed, the overall goals and the long-term effects of socially assistive robotics need to be considered. Although it is difficult to compare robotic systems designed for different tasks, it is important to do so to establish benchmarks for effective and ethical SAR design.

Table 1. A list of proposed benchmarks

Robotics Technology	Social Interaction	Assistive Technology
Safety	Autonomy	Impact on User's Care
Scalability	Imitation	Impact on Caregivers
	Privacy	Impact on User's Life
	Understanding of Domain	
	Social Success	

A similar argument for relative comparison has been made by proposing benchmarks for human–robot interaction (HRI) for better evaluating and understanding interactive robots (Kahn, et al., 2007). In that work, a series of questions were asked that clarified the role of the robot and how well that robot performed its function. In this paper, the authors refine existing benchmarks and propose additional human–robot interaction benchmarks pertaining to socially assistive robotics. These benchmarks take into account not only aspects of HRI, such as imitation and autonomy, but also the applications of a robot system and its consequences on the intended user population.

SAR addresses a compelling subclass of HRI, because the robot is goal-oriented, which creates additional challenges for system evaluation. For example, when evaluating a robot designed entirely for entertainment, one may focus on whether the user is having fun. However, when evaluating a goal-oriented system, one needs to consider whether the robot's behavior results in good care.

SAR systems need to be evaluated from a technical (robotics) perspective, a social perspective, and an assistive perspective. In addition, each benchmark should be applied with respect to the user(s) and the caregiver(s) involved. Each attribute examined by a benchmark also may have ramifications on who is or is not helped through socially assistive robotics, because SAR may extend care to people who otherwise would not have access. In this paper, the authors outline benchmarks (see Table 1) for the three aspects of socially assistive robotics: robotics technology, social interaction, and assistive technology. These benchmarks are described in a user-centered fashion. An evaluation of the impact of a robotic system on its users is critical to the success of SAR.

Background and related work

When discussing benchmarks for SAR, the authors first define the range of roles that a robot can play in a human user's life. Also important are existing benchmarks for robots and how they are currently applied. This section describes socially assistive robotics and cites pre-existing benchmarks for robots in society.

Examples of SAR

SAR has a wide range of application domains. Robotics researchers have studied the advantages and challenges of using social robots as assistants to the elderly, in the home or in hospitals. They work to automate some physical tasks that an elderly person may not be able to do, including applying makeup, brushing teeth, getting in and out of bed, getting into and out of a wheelchair, and adjusting a bed for maximum comfort. These robots can be used as part of a ubiquitous computing system (Jung et al., 2005), which combines cameras and other sensors and computer controlled appliances (such as light switches, doors, and televisions), or as a single unit (Baltus et al., 2000).

SAR systems have also been used as companion robots in the public areas of nursing homes, aimed at increasing resident socialization. These robots are not designed to provide a specific therapy but to be a focus of attention. Paro, a synthetic actuated stuffed seal toy, behaves in response to touch and sound (Wada, Shibata, Saito, Sakamoto, & Tanie, 2005). Its goal is to provide the benefits of pet-assisted therapy, which can affect resident quality of life in nursing homes that cannot support pets (Edwards & Beck, 2002). Initial studies have shown lowered stress levels in residents interacting with this robot and an overall increase in the amount of socialization among residents.

In rehabilitation robotics, non-contact assistive robots are now being developed. Such robots provide their assistance through social means with little or no physical contact. The authors' lab is currently studying applications for post-operative cardiac recovery (Kang, Freedman, Matarić, Cunningham, & Lopez, 2005) and post-stroke rehabilitation (Matarić, Eriksson, Feil-Seifer, & Winstein, 2007). These robots act as rehabilitation coaches, observing a user's progress and providing encouragement and guidance.

Other rehabilitation projects have explored using a robot as a means of motivating rehabilitation through mutual storytelling (Lathan et al., 2001). In these experiments, a robot and a user construct a story, which, when acted out, requires the user to perform physical therapy exercises. Rehabilitation applications help the user spend more time on moderated therapy.

Finally, SAR is being studied as a tool for diagnosis (Scassellati, 2005) and socialization of children with autism spectrum disorder (ASD) (Dautenhahn & Werry, 2002; Kozima, Nakagawa, & Yasuda, 2005). When used for diagnosis, robots can observe children in ways that people cannot. In particular, eye-tracking studies have shown remarkable promise when evaluating children with ASD. Robots have been found to be more comfortable social partners for children with ASD than people. These robots encourage social behavior, such as dancing, singing, or playing with the robot and with other children or parents in the hope of making such behavior more natural.

Each of the above examples shows aspects that should be investigated to gauge the effectiveness of a system. Usually users benefit most from SAR technology, although sometimes caregivers do. Some of these robots bring a type of care to a population that would not have it otherwise or they make existing care cheaper or more effective. These are just some of the ways SAR systems can be evaluated.

Examples of existing benchmarks and theory

Isaac Asimov (1950, pp. 44–45) proposed some of the original benchmarks for robotics when he wrote the now-famous three laws. These “laws” largely evaluate a robot from a safety perspective. Currently robotics researchers assume safety is a “low-level” capability and are considering evaluation from an ethical standpoint. In such an evaluation, the robot must behave in an inherently ethical manner.

While no specific ethical guidelines have yet been established, active discussions have begun in the field. Turkle (2006) addressed the attachment that occurs between humans and robots when residents of a nursing home are asked to “care for” a baby-like robot. The users in the experiment ascribed human qualities to the robot user in the experiment. This creates interesting social side-effects with ethical ramifications (Turkle, 2007). What happens when the robots break down, or get taken away? Previous benchmarks addressed the disparity between machines that exist only to serve a human “master” and those that exist in cooperation with their users and act with autonomy (Kahn et al., 2007). Is it harmful for people to treat a social being like a slave?

Wallach and Allen (2005) distinguish top-down and bottom-up approaches to robot morality. A top-down approach takes an ethical theory and guides the design and implementation of algorithms and subsystems to implement that theory. A bottom-up approach treats values as implicit to robot design. Yanco (2002) described the evaluation of an assistive robot, stating that such evaluation can be performed through user tests and comparison to a human in the same assistive role. Long-term studies were recommended to evaluate effectiveness in real-world settings. Others advocated a human-centered approach to design, suggesting ecological studies of the use of the robots in the intended environment rather than long-term user studies (Forlizzi, DiSalvo, & Gemperle, 2004).

Robot evaluation

When vulnerable populations are to use the robot, it is especially important to evaluate it as a physical and technological platform. This section highlights two benchmarks of particular concern to SAR: safety and scalability.

Safety

A robot's safety in its given domain is the primary concern when evaluating a SAR system. If a robot is not safe, it could harm the very users it is designed to help. A key advantage of SAR over physically assistive robots is the minimization of the risk associated with physical contact. The safety of a mobile platform involves the ability to maneuver about a scene without unwanted contact or collisions. Safety also refers to the protection of a robot's user and of the robot itself. This concept, as a benchmark, refers to safety in a bottom-up fashion, rather than Asimov's laws which refer to the concept in a more top-down fashion.

Safety for assistive robots has previously been addressed (Baker & Yanco, 2005; Rentschler, Cooper, Blasch & Boninger, 2003). Safety in these situations primarily refers to obstacle avoidance for guide-canes and wheelchairs. Robots have been designed to help a user navigate through a nursing home (Glover et al., 2003; Morris et al., 2003). Safety is an important benchmark for SAR: *How safe is the robot itself, and how safe can the robot make life for its user?*

Scalability

Currently, the majority of human-robot interaction work occurs in research laboratories where systems are engineered for one environment and a pre-determined user population. Even within a group that needs assistance, there is a great difference between a prototypical user or environment and the range of potential users and environments. As SAR systems become more widespread in homes, schools, and hospitals, the question of scalability and adaptability arises: *How well would such robots perform outside of the lab? How well does a robot perform with real users as opposed to prototypical users?*

Scalability of user populations. When considering a robot for human use, there is a critical question to ask: *How many people can the robot help?* Consider, for example, a robot that uses speech recognition for understanding a user's intentions. Can the robot interact with someone who cannot speak or has recently had a stroke? If the robot is meant to be a companion, can it adapt to different users? How difficult is it to modify the robot for different needs?

Scalability of environments. Most systems to date have been tested in research labs or hospitals. In the future, however, SAR will be used in homes and other unpredictable environments. In such domains, there exists an important benchmark for success: *Can a robot go wherever its intended user can?*

Social interaction evaluation

A critical benchmark of SAR is the evaluation of the robot as a social platform. Social interaction and engagement are both the primary means of interaction and the driving force behind the design. When assessing a robot's social performance, it is also necessary to consider the larger goal of the robot in the assistive context.

Existing HRI benchmarks

Previously proposed HRI benchmarks (Kahn et al., 2007) are directly relevant to SAR. In many respects, the same comparisons and evaluations that hold for humanoid robotics also hold for SAR. However, the goal of SAR is not to make as humanlike a robot as possible, but to make a robot that can best assist a user population. It is important, therefore, to evaluate SAR not only from the perspective of modeling human characteristics but also from a user-oriented perspective. The following sections describe how some of the previously identified HRI and humanoid benchmarks relate to SAR in particular.

Autonomy. Autonomy is a complex trait in the SAR context. It is favorable, when constructing a system that is designed to stand in for a human in a given situation to have a degree of autonomy which allows it to perform well in its desired tasks. SAR systems, such as robots for proactive social interaction with children with ASD (Dautenhahn & Werry, 2002) and motivational robot tools, require some autonomy (Wainer, Feil-Seifer, Shell, & Matarić, 2006; Matarić et al., 2007; Tapus & Matarić, 2006).

However, in situations such as dispensing medication and monitoring therapy (Fortescue et al., 2003), only as much autonomy is desired as is appropriate. The robot should dispense medication, but a doctor should monitor it. A human supervisor should monitor therapy to verify it is going according to plan. In general, SAR contexts require a robot to engage in believable social interaction, but the user or human caregiver must retain authority. For example, rehabilitation should terminate if the user is in pain. Social interaction should only occur when it is tolerable. Partial or adjustable autonomy on the part of the SAR system allows for an appropriate degree of authority and autonomy.

The ultimate benchmark for autonomy from a SAR perspective is related to the potential behavior of the robot: *Is the robot able to participate in activities necessary for proper assistance? Can a user and caregiver put the necessary trust in a robot system for that robot to be able to perform effectively?*

Imitation. Alan Turing proposed a test of artificial intelligence (AI), whereby a system is evaluated by whether it could fool a human user communicating with it through a teletype (Turing, 1950). This test was later elaborated as the Total

Turing Test (Harnad, 1989), in which a system communicating in humanlike ways (speech, facial expressions) tries to fool a human user into believing it is human. Since 1950, one of the benchmarks for success in AI and HRI has been how well the system can imitate human behavior. However, when dealing with goal-oriented systems for assistance and treatment, it is not necessarily desirable to imitate human behavior. Thus, imitation has varying degrees of importance and appropriateness for SAR.

A robot's personality can affect a user's compliance with that robot (Kiesler & Goetz, 2002). A robot, when exhibiting a serious personality, could provoke a greater degree of compliance than a more playful robot personality. When the robot's degree of extroversion or introversion matches the user's, task performance can be improved (Tapus & Matarić, 2006). Thus, imitation and social relatedness could be significant factors in task performance. The imitation benchmark proposed by Kahn et al. (2007) could be revised for SAR: *How does imitation (and reciprocity) affect task performance?*

Realistic robotics introduces new complications for social robot design (Duffy, 2003). Further, it has been implied that anthropomorphism has a negative influence on social interaction when the robot's behavior does not meet a user's expectations (Shneiderman, 1989). The "uncanny valley" hypothesis suggests that as a robot becomes very similar in appearance to a human, that robot appears less, rather than more, familiar (Mori, 1970; MacDorman & Ishiguro, 2006). Physical similarity or an attempt at very faithful humanlike imitation of behavior could cause such discord. It could therefore be important to ensure that any imitation of human likeness or behavior creates a believable portrayal of human behavior.

A benchmark based on imitation is troubling because of this relationship between imitation and user impressions. When a robot imitates a human, or vice versa, interaction between them can seem more natural because of a sense of familiarity. However, that same similarity can create a negative impression if the robot's capabilities do not match its humanlike appearance or behavior (Goetz, Kiesler, & Powers, 2003). Therefore, there exists a possible benchmark for imitation: *Does the interaction between the human and the robot reflect an accurate and effective impression of the robot's capabilities?* Also, since there is a large range of social awareness and proficiency in user populations, there is a related question: *Does the interaction between the human and the robot reflect an accurate and effective impression of the human's capabilities?* Robots that misrepresent their abilities can create an ethical dilemma (Turkle, 2007).

Privacy. The presence of robots affects a user's sense of privacy (Kahn et al., 2007). This is especially important to SAR. Privacy is an area of concern for designers of assistive systems (Baillie, Pucher, & Kpesi, 2004). In contrast to ubiquitous systems (Lee & Keating, 1994), users know when the robot can see them. A robot

can be told to leave when privacy is desired. In some ways, a robot could be seen as providing more privacy for potentially embarrassing environments than a human caregiver in a similar situation. Therefore, there exists an appropriate benchmark from a SAR perspective: *Does the user's perceived sense of privacy relate to better robot performance as an assistive presence? Does user privacy impact user satisfaction?*

Task-oriented benchmarks

In addition to existing HRI benchmarks, the task-oriented nature of SAR suggests some additional task-oriented benchmarks for SAR, such as rehabilitation therapy, socialization, and tutoring. These benchmarks pertain to how the social aspects of the robot can affect the overall task performance of the robot. As with the previous HRI benchmarks, these could apply to all social robots, but when put into an assistive context, the task-related effects make these features stand out. Task performance is described as the robot's ability to assist a user in a given task.

Social Success. Considering the importance of the social aspect of a robot, a question should be asked: *Does the robot successfully achieve the desired social identity?* This is the most amorphous of the benchmarks, but its evaluation is simple. When the robot is intended to be playful, do users find the robot playful? If the robot is intended to be a social peer, do users act as if it were a social peer? The social success of the robot is a fundamental component of SAR applications. As discussed in the Imitation and Cost/Benefit sections, the social identity of the robot (both the personality and the role of the robot) affects the user's task performance.

Understanding of Domain. The understanding of social dynamics is a critical component of SAR. Robotics researchers employ user and activity modeling to achieve it. Modeling efforts include emotion recognition (Cassell, Sullivan, Prevost, & Churchill, 2000; Cowie et al., 2001), and integration of vocalizations, speech, language, motor actions, and gestures (Busso et al., 2004; C. Lee & Narayanan, 2005). Sensing for social understanding does not only have to relate to traditional social channels. For example, robotics researchers used radio frequency identification (RFID) tags and position tracking to observe children in school hallways. This enabled the robot to determine more appropriate social behavior by detecting who was in social range and with whom the robot was interacting (Kanda, Hirano, Eaton, & Ishiguro, 2003). Thus, social understanding can come from both human-oriented social perception (such as speech recognition or face recognition), and purely non-human-oriented social perception (such as galvanic skin response evaluation of emotional state). Therefore, there exists an understanding-related benchmark for social interaction: *Does a robot's social understanding of human behavior help task performance?*

This benchmark is especially hard to employ because our own understanding of social interaction may not be sufficient to judge an artificial being's social proficiency. Theory of mind and simulation theory are currently popular in the social science literature; however, such models are only theoretical and may not be correct (Gallagher, 2004). To evaluate robot performance, it is necessary to have an accurate model of social interaction.

Assistive evaluation

There are many ways to view an assistive tool. For the domains of SAR, the authors have selected impact on the user's care, impact on caregivers, impact on the user's life, and the role of the robot as benchmarks of an assistive platform. An important way to view how an assistive robot performs when caring for people is by observing how people care for other people in similar situations.

The role of a robot in assistive care has a great impact on how benchmarks are to be applied. The role of the robot differs from one application to another. In some cases, the robot is meant to stand in for a human caregiver while in other cases the robot is meant to assist a human caregiver.

Impact on the user's care

The most prominent performance evaluation is how an assistive robot changes the over-all care of the user. For evaluation, the metrics for success include the comparative success of a robot system relative to a human caregiver and the cost to benefit ratio of that care.

Success Relative to a Human Caregiver. A good place to start when evaluating the effect a robot has on a user's care is to compare the results of care with a robot caregiver to those of care with a human caregiver: *How does the robot perform relative to a human performing the same task?* In such situations, it is possible to evaluate changes in care by using existing metrics. Most of these effects deal directly with the subject of the robot's interactions. For rehabilitation tasks, overall rehabilitation can be a metric. For learning tasks, overall learning measures such as grades, tests, or evaluations can be used.

In other cases, when the actions of the robot are not mimicking human job performance, new metrics can be created. In the cardiac-spirometry task (see Background and Related Work section), compliance with the robot compared to compliance with a human is a suitable metric. For companion robots, user satisfaction could be used. A key role of SAR is to provide care where human care

is not available, so comparison to human care may not be the most appropriate evaluation measure.

Cost/Benefit Analysis. The robot can perform in several different capacities for any given task. For example, in a rehabilitation setting a robot could serve as a therapist, giving advice on specific movements, a motivational coach, giving encouragement and monitoring progress, a cognitive orthotic, reminding the users of important items, a companion robot, a learning aid, or as a demonstration, showing a user how to do specific exercises. The role of the robot for a given task can inform the sophistication of the robot and its social and assistive capacities.

While providing a comparable level of care, SAR may offer a solution for a lower cost than other options. In this case, a cost/benefit analysis can be a benchmark of success. Furthermore, the provision of care through SAR can also result in faster training and recovery and thus lower healthcare costs. This suggests a good measure for SAR: *Does the use of the robot (a) change the cost/benefit ratio of providing such care or (b) make such care available where it was not previously possible?*

Impact on caregivers

In some cases, the goal of a system is not to increase the standard of care, but to make the caregivers' jobs easier and more manageable. For example, in one study the goal of a system (Kang et al., 2005) was to reduce the overall workload for cardiac nurses, given the overall nurse shortage in the US and worldwide. The robot observed cardiac patients, post-surgery, who needed to perform breathing exercises (spirometry) ten times per hour. The robot approached the patient's bed, encouraged the patient to perform the exercise, monitored the number and depth of the breaths taken, and collected data on performance. By automating such repetitive and monotonous tasks, caregivers can attend to other tasks and provide more individualized services. However in that study, the robot did not provide any care not already given by a human caregiver; it just automated the simplest aspects of existing care.

Caregiver impact is a useful benchmark: *Does the job condition of the caregiver improve as a result of the robot?* Additionally, it is important to observe cooperation: *How well does the caregiver work with the robot?* This arises out of a concern that caregivers are not used to working with robots, and may need to adjust their work habits (Scholtz, 2002). A necessary caveat, however, is that since this benchmark does not directly refer to the user, success may not reflect improved care. This benchmark should be evaluated in concert with user-centered benchmarks.

Impact on the user's life

The most important area for benchmarking relates to how an assistive robot affects the user's life in general. This impact can be measured not only by user studies related to the robot's performance on a given task, but also by considering how the robot affects the user's overall quality of life, lifestyle, satisfaction with the robot, and role in the community.

Satisfaction with care. User satisfaction is an important aspect of assistive therapy success. Users' impression of a nurse robot's personality affected compliance with that robot (Kiesler & Goetz, 2002). Thus, satisfaction can be a useful benchmark for success. The authors currently are investigating ways to use questionnaires to measure satisfaction (Wainer et al., 2006), although the relation between satisfaction and task performance or user compliance requires further investigation. An important question when designing an assistive system is raised: *Does user satisfaction with a system affect assistive task performance or user compliance?*

Existing quality of life measurements. A way to evaluate the effects of a particular therapy regimen is to assess changes in the overall quality of life of the user. Quality of life, while a rather nebulous measurement, can be an effective tool for measuring overall effects of therapy (Wood-Dauphinee, 1999). Some recommend using repeated measures with the same survey to capture changes over time. The SF-36 survey is designed for a patient to rate health-related quality of life (Aaronson et al., 2004). It assesses the comprehensive quality of life from the patient's point of view. The 15-D survey produces quality of life numbers along several dimensions (Sintonen, 1994). Experiential measures, such as the dementia care mapping (DSM), are used in addition to other measures to assess a person's quality of life over time (Kitwood & Bredin, 1992). These kinds of measures are especially important, because they focus on the users of a particular service (Younger & Martin, 2000). Socially-sensitive care (e.g., involving eye-contact, favorable attention) is also important to the overall measure of care. Such measures for assessing quality of life constitute a benchmark for SAR: *Does the robot result in a general increase in the user's perceived quality of life?*

Impact on the User's Role in the Community. The introduction of new therapies can also affect the community. When fish tanks were introduced in a nursing home, residents gathered around the fish tanks in common areas and engaged in more conversation. Their nutrition also improved (Edwards & Beck, 2002). The introduction of objects of social interest can thus affect the community.

When robotics researchers introduced Paro in the common areas of a nursing home (Wada et al., 2005), residents were also in the common areas longer and socialized more. Stress proteins also decreased in their urine. The Robovie project used a robot to stimulate social interaction among a group of elementary school

students (Kanda et al., 2003). By telling “secrets” about itself, the same robot was able to elevate a student’s status in the group (Cowley & Kanda, 2005). When assessing a particular robot-assisted therapy, it is important to look beyond the robot’s immediate effects on the individual user: *Does the robot increase or decrease the amount of socialization in its user community? Is the robot’s overall impact on the community positive or negative?*

A potential critique of socially assistive robotics is that social robots could reduce the time their users spend in contact with other people (Turkle, 2007). Thus, the use of this benchmark is critical for ensuring such an effect does not occur.

Discussion

One of the main uses of a socially assistive robot is to provide the services that a human caregiver is unable to provide. The reason could be financial, in which case a socially assistive robot would be able to meet the needs of someone who otherwise could not afford the help required. The reason could be lack of personnel, in which case the robot could fill a void when there are not enough skilled workers to meet the demands of those needing their assistance. Alternatively, the reason could be the user’s discomfort with having human assistance. In these situations SAR can provide great benefits.

However, to determine whether socially assistive robotics is effective, clear benchmarks are needed to measure their success. While there are already benchmarks for HRI, these are insufficient for SAR, as researchers need to evaluate the performance of a robot’s assistance and social interaction. This paper is an attempt to jumpstart a broader dialog on the topic of benchmarks for socially assistive robotics, which should evaluate both the performance of a robot and its effects on the user population.

One of the most challenging aspects of establishing these benchmarks is that many aspects of SAR are difficult to measure. Establishing whether a robot can make eye contact with a person is easy, but evaluating how the person reacts to and is affected by the robot is more difficult. Does the person get bored or frustrated? Does the user consider the robot helpful and effective? Is the robot perceived as competent? Is it trusted to perform its intended tasks? A robotics researcher should ask these questions when evaluating how well a robot can interact with and assist people.

Likewise, while there are clear methods for measuring physical improvement, evaluating a robot’s social or emotional impact is a complex task. The authors want robots to be able to perform human caregiver tasks and to be able to interact positively with humans. However, robotic interaction should not replace human

interaction (Turkle, 2007) but rather improve it. A robot's emotional impact is equally important; a user should be able to trust a robot and feel safe using it.

Safety benchmarks are of the utmost importance, because a socially assistive robot is useless if it poses a threat to the user. Those interacting with the robot need to be assured of both their physical and emotional safety. Physical safety guidelines may vary depending on the user's condition, but are reasonably easy to establish with the aid of a healthcare professional. Emotional safety standards are more ambiguous but just as important. For example, many applications of SAR are in the medical field; if a robot records a patient's medical information, benchmarks are needed to ensure that that information is kept confidential.

This leads to ethical considerations and legal guidelines for developing these types of robots. Not only do robotics researchers need to act ethically, but so do their robots, even if they are not aware of morality. Some things to consider are the intended uses of the robot, the tasks it should be allowed to perform, and the situations it might encounter. For example, if the person being assisted suddenly needs emergency attention, what is the robot's responsibility? Additionally, guidelines should be established to determine who is in charge of the robot and when. Although most situations may warrant the user having control, there are important exceptions, such as dispensing medicine. In addition, robotics researchers need to be able to ensure that a doctor or other caregiver cannot use the robot to harm the person the robot was designed to help.

To address these issues, the authors have outlined an initial set of benchmarks that can be used to develop safe and effective robots that will have a positive impact on the user in both the short term and the long term. This list is by no means exhaustive and will evolve as SAR research grows more complex and develops new applications. Additionally, some questions will apply to some robots but not to others. This list is but a first step in the process of systematically measuring and evaluating the consequences of socially assistive robots and ensuring that they are safe, effective, and have a positive impact on the people they are designed to help.

Acknowledgments

The authors would like to thank Ronald Seifer and Leslie Feil for valuable insights into the issues surrounding ethics and medical economy. This work was supported by the USC Institute for Creative Technologies, the USC Provost's Center for Interdisciplinary Research, and the Okawa Foundation.

References

- Aaronson, N. K., Acquadro, C., Alonso, J., Apolone, G., Bucquet, D., Bullinger, M., et al. (2004). International quality of life assessment (IQOLA) project. *Quality of Life Research*, 1(5), 349–351.
- Asimov, I. (1950). *I, robot*. New York: Bantam Books.
- Baillie, L., Pucher, M., & Kpesi, M. (2004). A supportive multimodal mobile robot for the home. In C. Stary & C. Stephanidis (Eds.), *User-centered interaction paradigms for universal access in the information society* (Vol. 3196, pp. 375–383). Berlin: Springer.
- Baker, M., & Yanco, H. (2005). Automated street crossing for assistive robots. In *Proceedings of the International Conference on Rehabilitation Robotics* (pp. 187–192). Chicago, IL.
- Baltus, G., Fox, D., Gemperle, F., Goetz, J., Hirsh, T., Magaritis, D., et al. (2000). Towards personal service robots for the elderly. In *Proceedings of the Workshop on Interactive Robots and Entertainment (WIRE 2000)*, Pittsburgh, PA, USA.
- Busso, C., Deng, Z., Yildirim, S., Bulut, M., Lee, C., Kazemzadeh, A., et al. (2004). Analysis of emotion recognition using facial expressions, speech and multimodal information. In *Proceedings of the International Conference on Multimodal Interfaces* (pp. 205–211). State Park, PA.
- Cassell, J., Sullivan, J., Prevost, S., & Churchill, E. (2000). *Embodied conversational agents*. Cambridge, MA: MIT Press.
- Cowie, R., Douglas-Cowie, E., Tsapatsoulis, N., Votsis, G., Kollias, S., Fellenz, W., et al. (2001). Emotion recognition in human–computer interaction. *IEEE Signal Processing Magazine*, 18(1), 32–80.
- Cowley, S., & Kanda, H. (2005). Friendly machines: Interaction-oriented robots today and tomorrow. *Alternation*, 12(1a), 79–106.
- Dautenhahn, K., & Werry, I. (2002). A quantitative technique for analysing robot–human interactions. In *Proceedings of the IEEE/RSJ, International Conference on Intelligent Robots and Systems* (pp. 1132–1138). Lausanne, Switzerland.
- Duffy, B. (2003). Anthropomorphism and the social robot. *Robotics and Autonomous Systems*, 42(3), 177–190.
- Edwards, N., & Beck, A. (2002). Animal-assisted therapy and nutrition in Alzheimer’s disease. *Western Journal of Nursing Research*, 24(6), 697–712.
- Feil-Seifer, D., & Mataric, M. (2005). Defining socially assistive robotics. In *Proceedings of the International Conference on Rehabilitation Robotics* (pp. 465–468). Chicago, IL.
- Fong, T., Nourbakhsh, I., & Dautenhahn, K. (2003). A survey of socially interactive robots. *Robotics and Autonomous Systems*, 42(3–4), 143–166.
- Forlizzi, J., DiSalvo, C., & Gemperle, F. (2004). Assistive robotics and an ecology of elders living independently in their homes. *Human–Computer Interaction*, 19(1–2), 25–59.
- Fortescue, E., Kaushal, R., Landrigan, C., McKenna, K., Clapp, M., Federico, F., et al. (2003). Prioritizing strategies for preventing medication errors and adverse drug events in pediatric inpatients. *Pediatrics*, 111(4), 722–729.
- Gallagher, S. (2004). Understanding interpersonal problems in autism: Interaction theory as an alternative to theory of mind. *Interaction Theory*, 11(3), 199–217.
- Glover, J., Holstius, D., Manojlovich, M., Montgomery, K., Powers, A., Wu, J., et al. (2003). *A robotically-augmented walker for older adults* (Tech. Rep. No. CMU-CS-03–170). Pittsburgh, PA: Carnegie Mellon University, Computer Science Department.

- Goetz, J., Kiesler, S., & Powers, A. (2003). Matching robot appearance and behavior to tasks to improve human–robot cooperation. In *Proceedings of the 12th IEEE International Workshop on Robot and Human Interactive Communication (RO-MAN 2003)*, (pp. 55–60). Lisbon, Portugal.
- Harnad, S. (1989). Minds, machines and Searle. *Journal of Experimental and Theoretical Artificial Intelligence*, 1, 5–25.
- Jung, J., Do, J., Kim, Y., Suh, K., Kim, D., & Bien, Z. (2005). Advanced robotic residence for the elderly/the handicapped: Realization and user evaluation. In *Proceedings of the International Conference on Rehabilitation Robotics* (pp. 492–495). Chicago, IL.
- Kahn, P. H., Jr., Ishiguro, H., Friedman, B., Kanda, T., Freier, N. G., Severson, R. L., & Miller, J. (2007). What is a human? — Toward psychological benchmarks in the field of human–robot interaction. *Interaction Studies*, 8(3). (This issue)
- Kanda, T., Hirano, T., Eaton, D., & Ishiguro, H. (2003). Person identification and interaction of social robots by using wireless tags. In *International Conference on Intelligent Robots and Systems (IROS2003)* (pp. 1657–1664). Las Vegas, NV.
- Kang, K., Freedman, S., Matarić, M., Cunningham, M., & Lopez, B. (2005). Hands-off physical therapy assistance robot for cardiac patients. In *Proceedings of the International Conference on Rehabilitation Robotics* (pp. 337–340). Chicago, IL.
- Kiesler, S., & Goetz, J. (2002). Mental models and cooperation with robotic assistants. In *Proceedings of the Conference on Human Factors in Computing Systems* (pp. 576–577). Minneapolis, MI. New York: ACM Press.
- Kitwood, T., & Bredin, K. (1992). A new approach to the evaluation of dementia care. *Journal of Advances in Health and Nursing Care*, 1(5), 41–60.
- Kozima, H., Nakagawa, C., & Yasuda, Y. (2005). Interactive robots for communication-care: A case-study in autism therapy. In *Proceedings IEEE International Workshop on Robot and Human Interactive Communication (RO-MAN 2005)*, (pp. 341–346). Nashville, TN.
- Lathan, C., Vice, J., Tracey, M., Plaisant, C., Druin, A., Edward, K., et al. (2001). Therapeutic play with a storytelling robot. In *Conference on Human Factors in Computing Systems* (pp. 27–28). New York: ACM Press.
- Lee, C., & Narayanan, S. (2005). Towards detecting emotions in spoken dialogs. *IEEE Transactions on Speech and Audio Processing*, 13(2), 293–302.
- Lee, N., & Keating, D. (1994). Controllers for use by disabled people. *Computing & Control Engineering Journal*, 5(3), 121–124.
- Lord, C., & McGee, J. (Eds.). (2001). *Educating children with autism*. Washington, DC: National Academy Press.
- MacDorman, K. F. & Ishiguro, H. (2006). The uncanny advantage of using androids in social and cognitive science research. *Interaction Studies*, 7(3), 297–337.
- Matarić, M., Eriksson, J., Feil-Seifer, D., & Winstein, C. (2007). Socially assistive robotics for post-stroke rehabilitation. *Journal of NeuroEngineering and Rehabilitation*. 4(5), Feb 2007
- Mori, M. (1970). Bukimi no tani (The Uncanny Valley). *Energy*, 7, 33–35.
- Morris, A., Donamukkala, R., Kapuria, A., Steinfeld, A., Matthews, J., Dunbar-Jacob, J., et al. (2003). A robotic walker that provides guidance. In *Proceedings of the 2003 IEEE International Conference on Robotics and Automation, (ICRA)* (pp. 25–30). Taipei, Taiwan.
- Rentschler, A., Cooper, R., Blasch, B., & Boninger, M. (2003). Intelligent walkers for the elderly: Performance and safety testing of VA-PAMAID robotic walker. *Journal of Rehabilitation Research and Development*, 40(5), 423–431.

- Scassellati, B. (2005). Using social robots to study abnormal social development. In *Proceedings of the Fifth International Workshop on Epigenetic Robotics: Modeling Cognitive Development in Robotic Systems* (pp. 11–14). Nara, Japan.
- Scholtz, J. (2002). Evaluation methods for human-system performance of intelligent systems. In *Proceedings of the 2002 Performance Metrics for Intelligent Systems (PERMIS) workshop*. Gaithersburg, MD.
- Shneiderman, B. (1989). A nonanthropomorphic style guide: Overcoming the Humpty-Dumpty syndrome. *The Computing Teacher*, 16(7).
- Sintonen, H. (1994). *The 15-d measure of health related quality of life: Reliability, validity and sensitivity of its health state descriptive system* (Working Paper No. 41). West Heidelberg, Victoria, Australia: Center for Health Program Evaluation.
- Tapus, A. & Matarić, M. (2006). User personality matching with hands-off robot for post-stroke rehabilitation therapy. In *Proceedings, International Symposium on Experimental Robotics (ISER)*. Rio de Janeiro, Brazil.
- Taub, E., Uswatte, G., King, D., Morris, D., Crago, J., & Chatterjee, A. (2006). A placebo-controlled trial of constraint-induced movement therapy for upper extremity after stroke. *Stroke*, 37(4), 1045–1049.
- Turing, A. (1950). Computing machinery and intelligence. *Mind*, 49(236), 433–460.
- Turkle, S., Taggart, W., Kidd, C. D., Daste, O. (2006). Relational artifacts with children and elders: The complexities of cybercompanionship. *Connection Science*, 18(4), 347–361.
- Turkle, S. (2007). Authenticity in the age of digital companions. *Interaction Studies* (this issue).
- Wada, K., Shibata, T., Saito, T., Sakamoto, K., & Tanie, K. (2005). Psychological and social effects of one year robot assisted activity on elderly people at a health service facility for the aged. In *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)* (pp. 2785–2790).
- Wainer, J., Feil-Seifer, D., Shell, D., & Matarić, M. (2006). The role of physical embodiment in human-robot interaction. In *IEEE Proceedings of the International Workshop on Robot and Human Interactive Communication* (pp. 117–122). Hatfield, United Kingdom.
- Wallach, W., & Allen, C. (2005). Android ethics: Bottom-up and top-down approaches for modeling human moral faculties. In *Proceedings of the 2005 COGSCI workshop: Toward Social Mechanisms of Android Science* (pp. 149–159). Stresa, Italy.
- Wolf, S., Thompson, P., Morris, D., Rose, D., Winstein, C., Taub, E., et al. (2005). The EXCITE trial: Attributes of the wolf motor function test in patients with subacute stroke. *Neurorehabilitation and Neural Repair*, 19(3), 194–205.
- Wood-Dauphinee, S. (1999). Assessing quality of life in clinical research: From where have we come and where are we going? *Journal of Clinical Epidemiology*, 52(4), 355–363.
- Yanco, H. (2002). Evaluating the performance of assistive robotic systems. In *Proceedings of the Workshop on Performance Metrics for Intelligent Systems*. (page numbers) Gaithersburg, MD.
- Younger, D., & Martin, G. (2000). Dementia care mapping: An approach to quality audit of services for people with dementia in two health districts. *Journal of Advanced Nursing*, 32(5), 1206–1212.

Authors' address

David Feil-Seifer, Kristine Skinner and Maja J Matarić
Interaction Laboratory
Department of Computer Science
University of Southern California
Los Angeles, CA USA 90089-0781
{dfseifer, kmskinne, mataric}@usc.edu

About the authors

David Feil-Seifer is a doctoral student in the Department of Computer Science at the University of Southern California. He received his M.S. in Computer Science at the University of Southern California in 2007 and his B.S. in Computer Science at the University of Rochester in 2003. His research interests include assistive human-robot interaction and especially the design of robots to interact with people who have cognitive disabilities.

Kristine Skinner is an undergraduate research assistant, funded by Women in Science and Engineering, in the Human Interaction Lab at the University of Southern California, where she is a bachelor's student. She expects to receive her B.S. in Electrical Engineering in 2007. Her research interests include human-robot interaction.

Maja J. Matarić is Professor of Computer Science and Neuroscience, Director of the Center for Robotics and Embedded Systems (CRES), and the Viterbi School of Engineering Senior Associate Dean for Research at the University of Southern California. She received her Ph.D. in Computer Science and Artificial Intelligence at MIT in 1994. With the goal of getting robots to help people, her research interests include human-robot interaction and robot control and learning in complex environments.

