

Probabilistic Models of Proxemics for Spatially Situated Communication in HRI

Ross Mead and Maja J Matarić
Interaction Lab
Computer Science Department
University of Southern California
Los Angeles, California 90089–0781
{rossmead, mataric}@usc.edu

Abstract—To enable socially situated human-robot interaction, a robot must both understand and control *proxemics*—the social use of space—in order to employ communication mechanisms analogous to those used by humans. In this work, we focus on social speech and gesture production and recognition during both human-human and human-robot interactions. We conducted a data collection to model these factors as a function of the relative distance and orientation between sociable agents. The resulting models were used to implement a spatially situated autonomous proxemic robot controller. The controller utilizes a sampling-based approach, wherein each sample represents interagent distance and orientation, as well as agent speech and gesture production and recognition estimates; a particle filter (with resampling) uses these estimates to maximize the performance of both the robot and the human during the interaction. This functional approach yields pose, speech, and gesture estimates consistent with related work in human-human and human-robot interactions. This work contributes to the understanding of the underlying pre-cultural processes that govern proxemic behavior, and has implications for the development of robust proxemic controllers for robots situated in complex interactions (e.g., with more than two agents, or with individuals with hearing or visual impairments) and environments (e.g., with loud noises, low light, or visual occlusions).

I. INTRODUCTION AND BACKGROUND

Why do people position and orient themselves the way they do during social encounters? How does use of space impact their perception and production of social stimuli, such as speech and gesture? How can these considerations influence the design of autonomous control systems for sociable robots?

These questions are investigated in the study of *proxemics*, which is concerned with the interpretation, manipulation, and dynamics of spatial behavior in face-to-face social encounters. Proxemic behavior is governed by sociocultural norms, which determine the perception and production of social stimuli (e.g., speech and gesture) for each interacting participant [1].

The effect of sensory experience on social proxemics is evident in environments that interfere with the perception of these social stimuli, such as in rooms that are noisy, [2], poorly lit, or containing visual occlusions [3]. This effect is further highlighted in the analysis of the personal space requirements of people with visual and hearing impairments [4]–[6]. Such factors must be considered in the design of socially intelligent [7]–[9] and socially assistive robotic technologies [10], [11].

To facilitate socially situated human-robot interaction (HRI), a robot must often employ multimodal communication mechanisms analogous to those used by humans: speech production (via speakers), speech recognition (via microphones), gesture production (via physical embodiment), and gesture recognition (via cameras or motion trackers). This research focuses on answering the question: How do speech and gesture influence spatially situated communication between humans and robots in social encounters? Within that question there are several others, including: How does the relative pose (position and orientation) between a robot and a person affect speech and gesture recognition for each of them? How can a robot dynamically adjust its communication mechanisms to maximize human perceptions of its social signals?

The anthropologist Edward T. Hall [12] coined the term “proxemics”, and proposed that *cultural* norms define zones of intimate, personal, social, and public space [13]. These zones are characterized by the *pre-cultural* and *psychophysical* visual, auditory (voice loudness), olfactory, thermal, touch, and kinesthetic experiences of each interacting participant [14].

Contemporary probabilistic modeling techniques have been applied to socially appropriate person-aware robot navigation in dynamic crowded environments [15], [16], to calculate a robot approach trajectory to initiate interaction with a walking person [17], and to position the robot for user comfort [18].

Our previous work utilized advancements in markerless motion capture (specifically, the Microsoft Kinect) to automatically extract proxemic features based on metrics from the social sciences [19], [20]. These features were then used to recognize spatiotemporal interaction behaviors, such as the initiation, acceptance, aversion, and termination of an interaction [20]–[22]. These investigations offered insights into the development of spatially situated controllers for autonomous sociable robots, and suggested an alternative approach to the representation of proxemic behavior that goes beyond simple distance and orientation [23]. In [24], [25], we proposed a probabilistic framework for autonomous proxemic control that considered the sensory experience of each agent (human or robot) in a co-present social encounter. We now formally investigate this framework and model the relationship between spatial, speech, and gesture parameters to develop a controller for spatially situated communication in HRI.

II. FRAMEWORK FOR HUMAN-ROBOT PROXEMICS

Our probabilistic proxemic framework considers how all represented social agents—in this case, the person and the robot—experience a co-present interaction. To do this, we must model the perception (input) and production (output) of speech and gesture conditioned on interagent pose (position and orientation). In this section, we describe the framework parameters, as well as the dependencies between these parameters, for their calculation and evaluation. In the next section, we discuss a data collection to model the calculation and evaluation of these framework parameters for their implementation in a socially situated autonomous proxemic control system.

A. Definition of Framework Parameters

Some number of sociable agents are co-located in a room and would like to interact. At any point in time and from any location in the room, each agent must be capable of estimating:

- 1) **an interagent pose**, POS_{AB} – where agent B is located with respect to agent A ;
- 2) **a speech output level**, SOL_{BA} – how loudly agent B will speak for agent A to hear;
- 3) **a speech input level**, SIL_{AB} – how loudly agent A will likely perceive speech produced by agent B ;
- 4) **a gesture output level**, GOL_{BA} – in what space agent B 's gestures will occupy for agent A to see;
- 5) **a gesture input level**, GIL_{AB} – in what space will agent A likely perceive gestures produced by agent B .

Note that these speech and gesture input/output parameters are not concerned with the *meaning* of the behavior, but, rather, the *manner* in which the behavior is produced.

The interagent pose (POS_{AB}) is expressed as a distance (d), and two orientations—one from agent A to agent B (α), and one from agent B to agent A (β). In robot navigation and trajectory planning, one typically works in world coordinates; however, for proxemics, it is beneficial to consider agent poses relative to one another.

Speech input and output levels (SIL_{AB} , SOL_{BA}) are represented as sound pressure levels, logarithmic measures of sound pressure relative to a reference value. A sound pressure level can be treated as a signal-to-noise ratio. This relationship is important when considering the impact of environmental auditory interference (detected as the same type of pressure signal) on SIL_{AB} and SOL_{BA} (saved for future work).

Gesture input and output levels (GIL_{AB} , GOL_{BA}) are each represented as a 3D region of space, called a *locus* [26]. For gesture production, we model the GOL_{BA} locus as the locations of agent B 's body parts in physical space. The GIL_{AB} can then be modeled the region of agent A 's visual field taken up by the body parts associated with the gesture output of agent B (i.e., GOL_{BA}) [24].

B. Modeling Framework Parameters

We model the relationships between these pose, speech, and gesture parameters as a Bayesian network to represent: (1) how people place themselves with respect to each other and a robot; (2) how interagent spacing influences human speech and gesture production (output); and (3) how interagent spacing influences speech and gesture perception (input) for both humans and robots. Formally, the models are:

- 1) $p(POS_{AB})$
- 2) $p(SOL_{BA}, GOL_{BA} | POS_{AB})$
- 3) $p(SIL_{AB}, GIL_{AB} | SOL_{BA}, GOL_{BA}, POS_{AB})$

III. DATA COLLECTION

We conducted a data collection to inform these models of pose, speech, and gesture parameters for humans and robots in our probabilistic framework. A comprehensive data collection for the large number of independent variables involved necessitates a large number of interaction scenarios. This work utilized a 2 (human vs. robot; between-participants) x 5 (choice distance vs. set distance {0.5m, 2.0m, 3.5m, 5.0m}); within-participants) design, and motivates future work in adaptive models for situated HRI.

A. Materials

Mead *et al.* [19] describes the software and hardware infrastructure developed for data collection and visualization. Microsoft Kinects were used to monitor the relative pose between agents (POS); six of these sensors were placed strategically in the environment to provide nearly complete coverage of the 6-meter-by-6-meter space. Mead *et al.* [19] demonstrated the ability to extract 6-DOF pose information from Kinect joint positions, which will serve as our representation of GOL . The Kinect field-of-view also allows for the extraction of the perceived GIL based on the body features (e.g., head, shoulders, arms, torso, hips, stance, etc.) observed by a mobile autonomous robot [24]. Speech output levels (SOL) were detected by microphone headsets (for people) and were produced by speakers (for the robot); both were calibrated to a sound pressure level meter for comparison across agents. Human speech input levels (SIL) were simulated using speaker directivity (based on α) [27] and a head-related transfer function (based on β) [28]. Robot SIL were recorded using the Kinect microphone array.

The PR2 mobile robot platform was used for human-robot interactions. The robot features two 8-DOF arms, a pan/tilt head, and a telescoping spine. The PR2 also features a semi-holonomic mobile base, which affords it the ability to strafe, a valuable motion for proxemic control. A Kinect was mounted atop the PR2's head to directly measure interagent pose and speech/gesture input levels. For the purposes of this data collection, the arms were placed at an "idle" position at the robot's sides, the head was pointed directly forward, and the spine was elevated so that the robot stood at 1.5 meters tall.

B. Participants

Participants were recruited via mailing list and word-of-mouth on the campus of the University of Southern California; thus, these results might not apply in other regions. A total of 20 participants (10 men, 10 women) were recruited for the data collection. All participants were university students between the ages of 18 and 27, and most had technical backgrounds; 10 of these participants had prior experiences with robots.

C. Procedure

Each participant was greeted at the door entering the private experimental space, and was informed of the nature of the data collection and their rights as a participant, which included a statement that the data collection could be halted at any time. Once participants agreed to proceed, they were fitted with a microphone headset; each participant was guided to a calibration station and asked to vocalize a relatively constant tone for a few seconds while their microphone configuration was calibrated to a sound pressure level meter. After this process was completed, the experiment could begin.

Each participant interacted with a co-present agent—either the PR2 robot or another participant. For human-human interactions, participants were acquainted prior to the study to eliminate many of the social behaviors associated with introductions (e.g., handshakes, which require special consideration during proxemic behavior analysis [14]). For human-robot interactions, participants were informed that the robot was learning how to communicate, and could see and hear them via the Kinect mounted atop its head.

Each participant was guided to a separate station to watch a short (1-2 minute) cartoon. After the cartoon had finished, each participant was instructed to enter the experimental area and discuss the cartoon (or any other topic of their choosing) with the other co-present agent. Cartoons were chosen because they are commonly used in speech and gesture analysis studies, as people appear to gesture frequently when describing them, most likely due to their unreal nature [29]. Participants would interact for 4-6 minutes, and then return to the station to watch another cartoon. Participants watched a total of six cartoons during the session; for human-human interactions, participants would watch and discuss three different cartoons and three of the same cartoons to elicit different kinds of interactions.

Participants were given no instructions as to where to stand during the first two interactions with the co-present agent. This “choice distance” condition sought to elicit the “desired” pose, speech, and gesture behaviors for the interaction. For the remaining four interactions, participants were instructed to stand at particular distances with respect to the other co-present agent. These distances were 0.5, 2.0, 3.5, and 5.0 meters, the order of which was randomly selected for each interaction. This “set distance” condition sought to expose how people would modulate speech and gesture to compensate for changes in interagent pose.

After all six interactions, participants completed a questionnaire, the details of which are beyond the scope of this paper. The data collection was then discussed with each participant.

D. Dataset

Each data collection session lasted one hour, which included 20-45 minutes of interaction time. We collected pose, speech, and gesture data for 5 human-human interactions (10 participants) and 10 human-robot interactions (10 participants).

Recorded audio and visual data were each independently annotated by two coders for occurrences of speech and gesture production. Interrater reliability was high in both conditions, with $r_{speech} = 0.92$ and $r_{gesture} = 0.81$. This yielded 2457 and 2232 spoken phrases, and 1613 and 638 gesture frames for human-human and human-robot interactions, respectively. Interagent pose data was available throughout each interaction.

E. Analysis and Discussion

The collected data were modeled to inform pose, speech, and gesture parameters for both humans and the robot in our proxemics framework. We now analyze and discuss these models as they were formally proposed in Section IV.

1) $p(POS_{AB})$: Interagent pose estimates from the Kinect were used to generate the mean (μ) and standard deviation (σ) distance (in meters) in both human-human interactions ($\mu = 1.44$, $\sigma = 0.34$) and human-robot interactions ($\mu = 1.94$, $\sigma = 0.61$). With such a small sample of the population in each condition, we can make no significant claims about these two distributions. However, we can refer to the abundance of social science literature on human-human spacing, noting that participants in our data collection positioned themselves near the boundary between personal and social distance [1], which is consistent with their interpersonal relationship (recall that participants were acquainted before the study began) [14]. Similarly, we can refer to the literature on proxemics in HRI. Walters *et al.* [30] consolidates and normalizes the results of many human-robot spacing studies, reporting mean distances of 0.49–0.71 meters. Takayama & Pantofaru [31] investigated spacing between humans and the PR2 (the same robot used in our study), reporting mean distances of 0.25–0.52 meters. We suspect that the apparent differences between our results and that of related work can be attributed to the experimental procedure—in many of these studies, participants are told to respond to a distance or comfort cue; however, in our study, participants were more focused on the interaction itself, so the spacing might have been less conscious. We will use the collected data in our own proxemic framework, and will continue to gather more data to better inform the model.

2) $p(SOL_{BA}, GOL_{BA} | POS_{AB})$: Results of the collected speech output levels are shown in (Fig. 2). In the “choice” condition, human-human and human-robot speech output levels were ($\mu = 65.80$, $\sigma = 6.68$) and ($\mu = 67.23$, $\sigma = 5.90$) dB SPL, respectively. We performed an ANOVA F -test for our 2x5 between-by-within study, and found significance ($p < 0.05$) in speech output levels at different distances across human-human vs. human-robot conditions. Furthermore, we found a trend towards significance ($p = 0.073$) for the statistical interaction between distance and human-human vs. human-robot conditions; this suggests that the way in which a person uses speech to compensate for distance might be

different between a robot and another person. These results are consistent with findings that people speak more loudly with a robot when it is trying to recognize their speech [32], and merit further investigation.

Gesture output levels were modeled as uniform distributions over body part locations with an agent’s workspace (as measured by the Kinect or based on robot kinematics). We are developing non-uniform distributions of these locations.

3) $p(SIL_{AB}, GIL_{AB} | SOL_{BA}, GOL_{BA}, POS_{AB})$: Estimates of human speech input levels (SIL) were simulated using models of human *speaker directivity* (based on α) [27] and a *head-related transfer function* (HRTF) of the human listener (based on β) [28]; each of these transforms is validated and commonly used in the audio processing community. Robot speech input levels were recorded using the Kinect microphone array, and were consistent with an HRTF we developed using the technique described in [28].

Human gesture input levels (GIL) were estimated based on a 60° horizontal field-of-view with respect to the human head frame as tracked by Kinects placed in the environment. Robot GIL was modeled based on human body features (e.g., head, shoulders, arms, torso, hips, stance, etc.) tracked by the Kinect on the PR2 [24] (Fig. 1).

We utilized this model SIL and GIL to estimate the performance rates of automated speech and gesture recognition systems (Figs. 3–5). For speech recognition, we trained a model on a corpus of annotated spoken phrases (2 male; 2 female; 1500 phrases each); audio files in this corpus were automatically attenuated based on the SOL and POS . For gesture recognition, we counted the number of times a particular body part appeared in the Kinect visual field for different values of POS ; this provides us with a more situated model of GIL that can inform a gesture *classification* system as to whether or not relevant gesture features will occur under the conditions detected by its sensors.

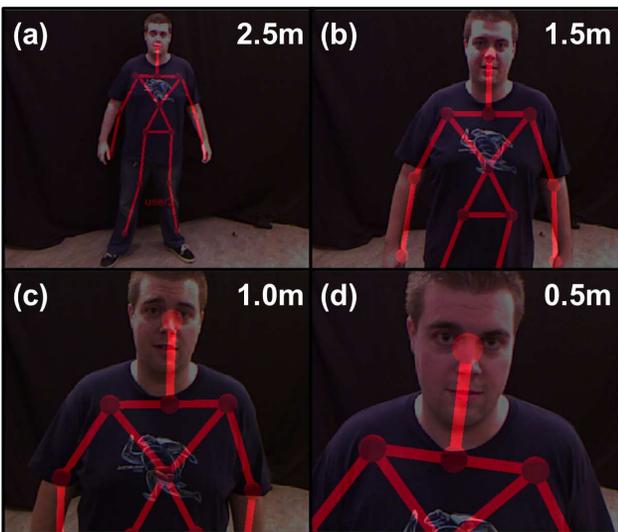


Fig. 1. Body features that fall into the Kinect field-of-view, depicted at four distances: (a) 2.5 meters, (b) 1.5m, (c) 1.0m, and (d) 0.5m. [24]

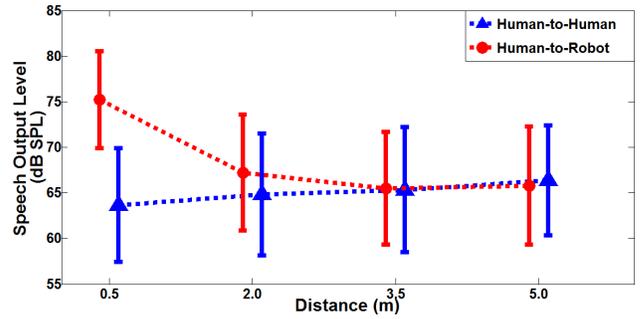


Fig. 2. Human speech output levels changing with distance with respect to both a robot and another human.

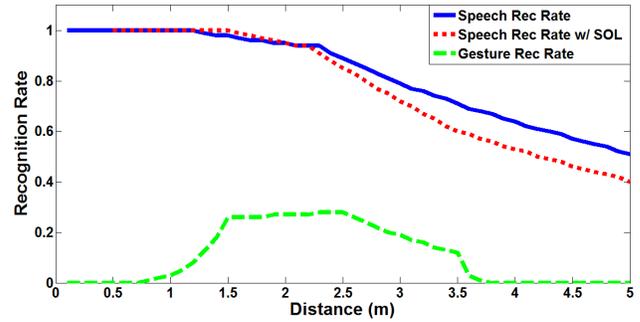


Fig. 3. Speech and gesture recognition rates as a function of distance; note the performance decrease in automated speech recognition when the model of human-robot SOL is applied, highlighting its utility for inference.

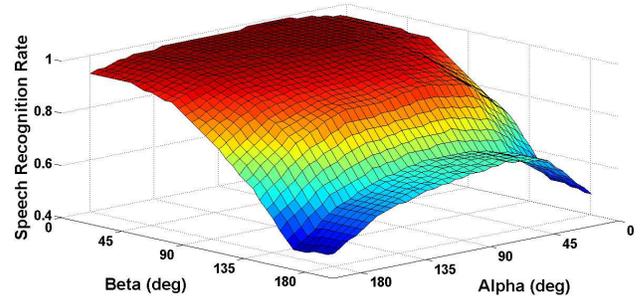


Fig. 4. Speech recognition rates as a function of speaker/listener orientations.

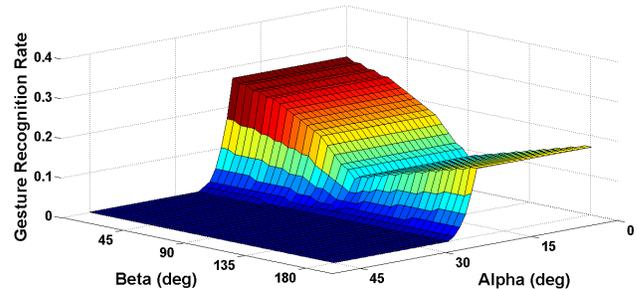


Fig. 5. Gesture recognition rates as a function of speaker/listener orientations.

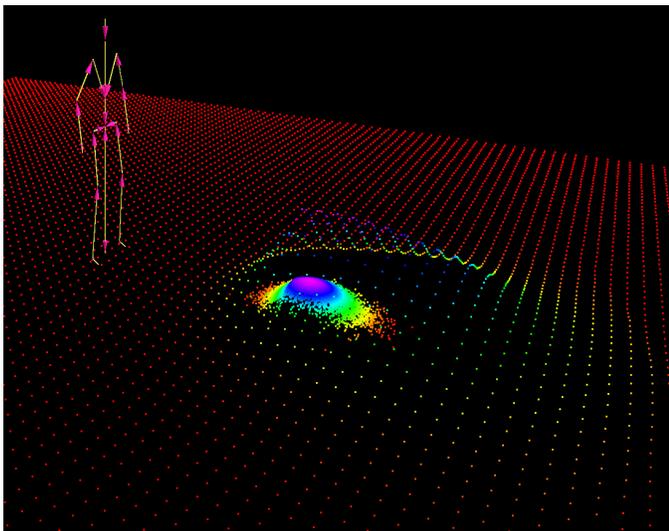


Fig. 6. Our autonomous proxemic controller uses a sample-based approach to estimate how interacting agents would experience a social interaction based on models in our socially situated proxemic framework; a particle filter (with resampling) then maximizes the social performance of all interacting agents.

IV. SPATIALLY SITUATED AUTONOMOUS CONTROLLER

The data-driven probabilistic framework of proxemics has been implemented as an autonomous proxemic controller for a sociable robot. The controller utilizes a sampling-base approach (similar to [18]) to infer information about latent variables (e.g., the perceived speech and gesture input levels of a person) based on sensed data (e.g., interagent distance), and uses these inferences to select parameters for robot social behavior (e.g., speech and gesture output levels) from some interagent pose. The robot considers (samples) possible interagent poses, estimating how all interacting agents would perceive and produce social stimuli (speech and gesture) if the interaction occurred with the robot at that pose—these estimates are maintained in a “particle”. At initialization, particles are sampled from the space uniformly. For each sample i , interagent pose (POS_i) and the corresponding speech input/output levels (SIL_i, SOL_i) and gesture input/output levels (GIL_i, GOL_i) are inferred based on the model as if the interaction were to occur at that pose; the sample is then assigned a weight (w_i) based on the likelihood (“desirability”) of an interaction actually occurring with these parameters (discussed in Section III). The interagent pose and input/output levels with the maximum weight are selected as the desirable parameters for robot pose, speech, and gesture behaviors; a particle filter resamples based on weight, conditioning these parameters on the previous state, ensuring that they will not change drastically between timesteps (Fig. 6). For interactions between two agents, this level of inference might be excessive; however, for interactions between three or more agents, such inference is necessary to determine appropriate parameters.

This controller has been implemented on the PR2 mobile robot (Fig. 7). The robot selects pose, speech, and gesture parameters to maximize its performance in the interaction.

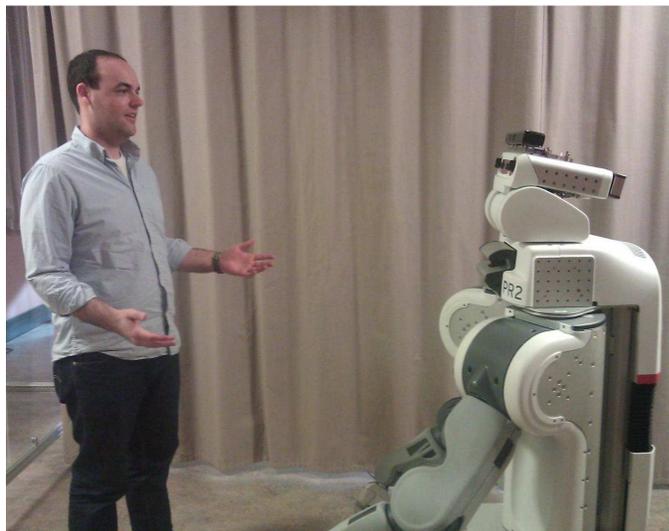


Fig. 7. The PR2 uses our controller to position itself in a conversation.

V. SUMMARY AND CONTRIBUTIONS

In this work, we considered how proxemic behavior is influenced by speech and gesture production and recognition during both human-human and human-robot interactions. We conducted a data collection to model these factors conditioned on interagent distance and orientation. The resulting models were used to implement a spatially situated autonomous proxemic robot controller. The controller utilizes a sampling-based approach to maximize the performance of both the robot and the human during the interaction. This functional approach yields pose, speech, and gesture estimates consistent with related work in human-human and human-robot interactions. This research contributes to the understanding of the underlying pre-cultural processes that govern proxemic behavior, and has implications for the development of robust proxemic controllers for robots situated in complex interactions and environments. This work further contributes to HRI community by providing data-driven models and software as part of the Social Behavior Library (SBL) in the USC ROS Packages Repository (<https://code.google.com/p/usc-interaction-software>).

VI. FUTURE WORK

Proxemic behavior is *not* cross-cultural [1], and the desired sensory experience of each interacting human participant varies all over the world [13]. We did not constrain our data collection to focus on a particular culture, as the robot cannot currently make any assumptions or inferences about with whom it will be interacting. However, in the future, it would be beneficial to develop separate models conditioned on different cultures, which could be used in conjunction with contextual knowledge and about the user and/or region.

Similarly, these models of proxemic behavior could be personalized to an individual. The use of a spatially situated communication has particular implications in socially assistive robotics domains [10], [11], such as when interacting with a person with a visual or hearing impairment. The models

of speech and gesture input would be conditioned on the sensory experience of such the individual, and the robot would dynamically adjust its communication modalities to maximize the perception of its social stimuli.

The current implementation of the system has the robot producing speech and gesture based on models of how humans produce such behaviors. This is based on the assumption that humans actually want the robot to behave this way. Future work will investigate the production of these behaviors from the ground up by conducting a human psychometric study of robot speech and gesture output.

Complex environments introduce extrinsic sensory interference that requires socially situated proxemic behavior to be dynamic. For example, if one is speaking in a quiet room, listeners do not need to be nearby to hear; however, if one is speaking in a noisy room, listeners must be much closer to hear at the same volume and, thus, perceive the vocal cues contained in the utterance. Similarly, if one is speaking in a small, well-lit, and uncrowded or uncluttered room, observers may view the speaker from a number of different locations; however, if the room is large, poorly lit, or contains visual occlusions, observers must select locations strategically to properly perceive the speech and body language of the speaker. In future work, we will integrate extrinsic interference into our probabilistic framework of human-robot proxemics.

ACKNOWLEDGMENTS

This work is supported in part by an NSF Graduate Research Fellowship, as well as NSF NRI-1208500, IIS-1117279, and CNS-0709296 and ONR MURI N00014-09-1-1031 grants, and the Willow Garage PR2 Beta Program. We thank Amin Atrash for his insights into probabilistic modeling techniques; Tian Zhu and Dhruv Monga for their development of the audio processing pipeline; Edward Kaszubski, Akanxa Padhi, and Divya Sharma for their assistance in the capture and analysis of visual data from the Microsoft Kinect; Elaine Short for her input on the statistical analysis; and Cliff Nass and Leila Takayama for their design considerations for the data collection.

This work is also in submission for the 2014 *International Symposium on Experimental Robotics (ISER 2014)*.

REFERENCES

- [1] E. T. Hall, *Handbook for Proxemic Research*. Washington, D.C.: American Anthropology Association, 1974.
- [2] D. M. Lloyd, "Dont stand so close to me: The effect of auditory input on interpersonal space," in *Perception*, vol. 38, no. 4, London, England, 2009, pp. 617–620.
- [3] L. Adams and D. Zuckerman, "The effect of lighting conditions on personal space requirements," *The Journal of General Psychology*, vol. 118, no. 4, pp. 335–340, 1991.
- [4] J. D. Webb and M. J. Weber, "Influence of sensory abilities on the interpersonal distance of the elderly," in *Environment & Behavior*, vol. 35, no. 5, 2003, pp. 695–711.
- [5] L. Hayduk and S. Mainprize, "Personal space of the blind," *Social Psychology Quarterly*, vol. 43(2), pp. 216–223, 1980.
- [6] T. W. Mallenby, "The personal space of hard-of-hearing children after extended contact with "normals"," *British Journal of Social and Clinical Psychology*, vol. 14, no. 3, pp. 253–257, 1975.
- [7] C. Breazeal, "Socially intelligent robots," *Interactions*, vol. 12, no. 2, pp. 19–22, 2005.
- [8] K. Dautenhahn, "Socially intelligent robots: Dimensions of human-robot interaction," *Philosophical Transactions of the Royal Society B*, vol. 362, pp. 679–704, 2007.
- [9] C. Breazeal, "Social interactions in hri: The robot view," *IEEE Transactions on Man, Cybernetics and Systems*, vol. 34, no. 2, pp. 181–186, 2003.
- [10] A. Tapus, M. Matarić, and B. Scassellati, "The grand challenges in socially assistive robotics," *IEEE Robotics and Automation Magazine*, vol. 14, no. 1, pp. 35–42, 2007.
- [11] D. Feil-Seifer and M. Matarić, "Defining socially assistive robotics," in *International Conference on Rehabilitation Robotics*, ser. ICRR'05, Chicago, Illinois, 2005, pp. 465–468.
- [12] E. T. Hall, *The Silent Language*. New York, New York: Doubleday Company, 1959.
- [13] —, *the Hidden Dimension*. Chicago, Illinois: Doubleday Company, 1966.
- [14] E. Hall, "A system for notation of proxemic behavior," *American Anthropologist*, vol. 65, pp. 1003–1026, 1963.
- [15] D. Vasquez, P. Stein, J. Rios-Martinez, A. Escobedo, A. Spalanzani, and C. Laugier, "Human aware navigation for assistive robotics," in *Proceedings of the Thirteenth International Symposium on Experimental Robotics*, ser. ISER'12, Québec City, Canada, 2012.
- [16] P. Trautman and A. Krause, "Unfreezing the robot: Navigation in dense, interacting crowds," in *IROS*, Taipei, Republic of China, 2010.
- [17] S. Satake, T. Kanda, D. F. Glas, M. Imai, H. Ishiguro, and N. Hagita, "How to approach humans?: Strategies for social robots to initiate interaction," in *HRI*, 2009, pp. 109–116.
- [18] E. Torta, R. H. Cuijpers, J. F. Juola, and D. van der Pol, "Design of robust robotic proxemic behaviour," in *Proceedings of the Third International Conference on Social Robotics*, ser. ICSR'11, 2011, pp. 21–30.
- [19] R. Mead, A. Atrash, and M. J. Matarić, "Proxemic feature recognition for interactive robots: Automating metrics from the social sciences," in *International Conference on Social Robotics*, Amsterdam, Netherlands, 2011.
- [20] —, "Automated proxemic feature extraction and behavior recognition: Applications in human-robot interaction," *International Journal of Social Robotics*, vol. 5, no. 3, pp. 367–378, 2013.
- [21] —, "Recognition of spatial dynamics for predicting social interaction," in *HRI*, Lausanne, Switzerland, 2011.
- [22] D. Feil-Seifer and M. Matarić, "Automated detection and classification of positive vs. negative robot interactions with children with autism using distance-based features," in *HRI*, Lausanne, 2011, pp. 323–330.
- [23] R. Mead, A. Atrash, and M. J. Matarić, "Representations of proxemic behavior for human-machine interaction," in *NordiCHI 2012 Workshop on Proxemics in Human-Computer Interaction*, Copenhagen, Denmark, 2012.
- [24] R. Mead and M. J. Matarić, "A probabilistic framework for autonomous proxemic control in situated and mobile human-robot interaction," in *HRI*, Boston, Massachusetts, 2012, pp. 193–194.
- [25] R. Mead, "Space, speech, and gesture in human-robot interaction," in *Doctoral Consortium of ICMI 2012*, Santa Monica, California, 2012, pp. 333–336.
- [26] N. Rossini, "The analysis of gesture: establishing a set of parameters," in *Gesture-based Communication in Human-Computer Interaction*, vol. 2915. Springer-Verlag, 2004, pp. 463–464.
- [27] W. Chu and A. Warnock, *Detailed Directivity of Sound Fields Around Human Talkers*, ser. Research Report (Institute for Research in Construction (Canada)). Institute for Research in Construction, 2002.
- [28] B. Gardner and K. Martin, "Hrtf measurements of a kemar dummy-head microphone," MIT Media Lab Perceptual Computing, Boston, Massachusetts, Tech. Rep. 280, 1994.
- [29] D. McNeill, *Hand and Mind: What Gestures Reveal about Thought*. Chicago University Press, 1992.
- [30] M. Walters, K. Dautenhahn, R. Boekhorst, K. Koay, D. Syrdal, and C. Nehaniv, "An empirical framework for human-robot proxemics," in *New Frontiers in Human-Robot Interaction*, Edinburgh, 2009.
- [31] L. Takayama and C. Pantofaru, "Influences on proxemic behaviors in human-robot interaction," in *IROS*, 2009, pp. St. Louis, Missouri.
- [32] S. Kriz, G. Anderson, and J. G. Trafton, "Robot-directed speech: Using language to assess first-time users' conceptualizations of a robot," in *Proceedings of the 5th ACM/IEEE international conference on Human-robot interaction*, ser. HRI '10. IEEE Press, 2010, pp. 267–274.